

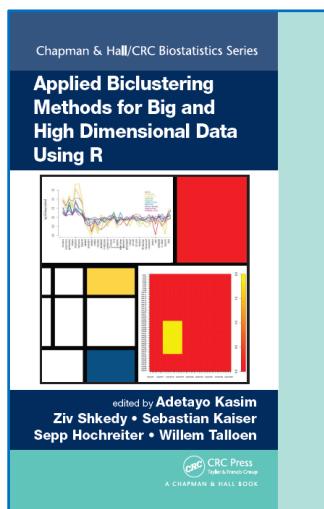
Computer intensive methods for bioinformatics

Applied Biclustering Methods for Big and High Dimensional Data Using R

Ziv Shkedy

2016-2017

Research Team



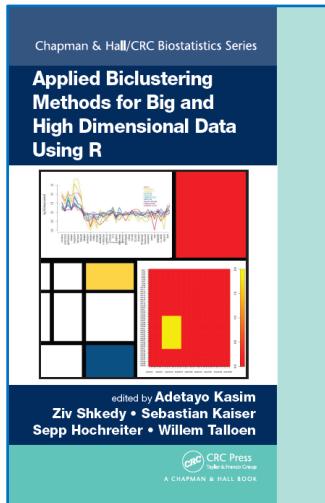
UHasselt:

Ewoud De Troyer
Rudradev Sengupta
Nolen Joy Perualila
Ziv Shkedy
Adetayo Kasim

and many others.....



Reference & R packages



R packages:

`biclust` (CRAN)
`biclustGUI` (CRAN)



Part 1

Introduction

Big data



- Everything is measurable.....
- We can collect a lot of data (and usually very quick).....
- How can we identify patterns in the data ?

Big data



- Today:
- Data analysis tool to discover local patterns in big data matrices.
- Case studies:
 - Sport.
 - Tourism
 - Drug discovery.

Part 2

Biclustering: local versus global patterns

Data structure

$$\mathbf{X} = \left(\begin{array}{cccc} X_{11} & X_{12} & \dots & X_{1m} \\ X_{21} & X_{22} & \dots & X_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \dots & X_{nm} \end{array} \right) . \quad \begin{array}{l} \text{Variables,} \\ \text{features...} \end{array}$$

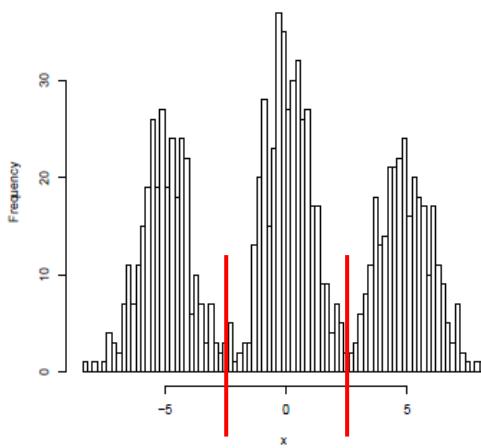

Observations, samples, conditions

Global patterns

- Find variables (observations) that can be grouped together due to a pattern in the data matrix.
- Examples:
 - All customers in a supermarket that have a tendency to buy the same products.
 - Genes with the same expression profiles in an expression matrix.

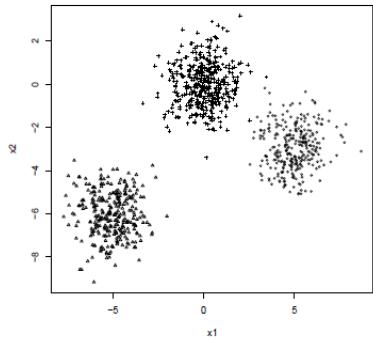
Clusters of observations

Example: three clusters of one variable

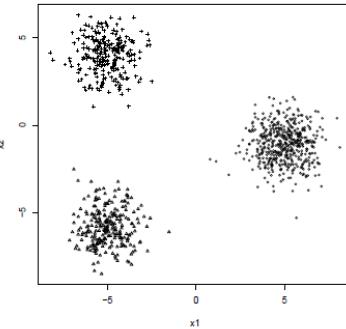


Clusters of observations

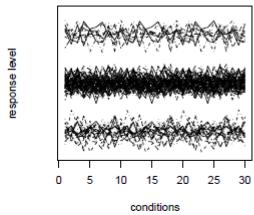
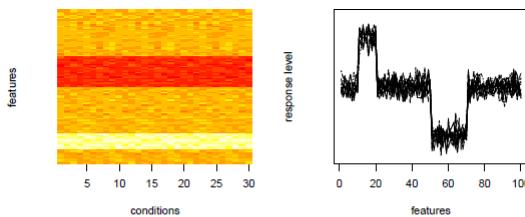
Example: three cluster in two variables



Example: how many clusters ?



Clustering and similarity measures



A data matrix with three clusters (of variables, rows).

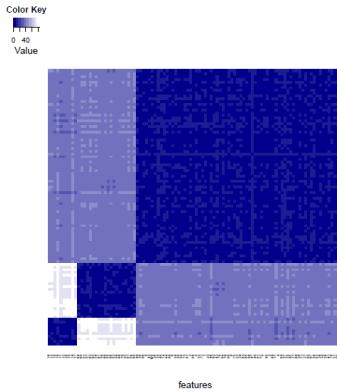
Clustering and similarity measures

The observe data matrix.
Clustering features.

$$\begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,n} \\ x_{2,1} & x_{2,2} & \dots & x_{2,n} \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \dots & x_{m,n} \end{bmatrix} \quad \text{features}$$

How similar are the features across
the samples (conditions).
Example: correlation across samples of
two features:

$$\rho(x_i, x_j)$$

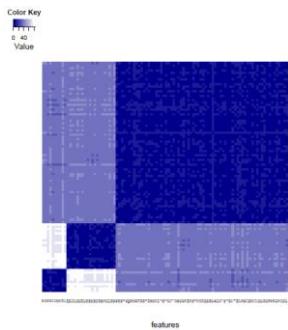


Example of three
clusters.

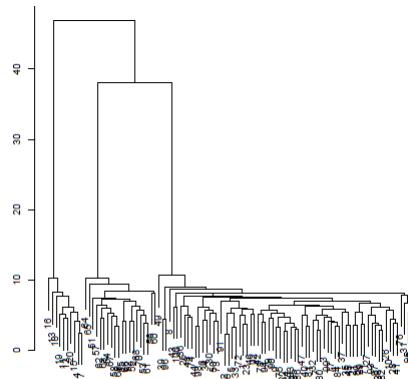
Hierarchical clustering

- Group variables according to their correlation (with each other).
- Variables are correlated across all observations.

Hierarchical clustering



Example of three clusters.

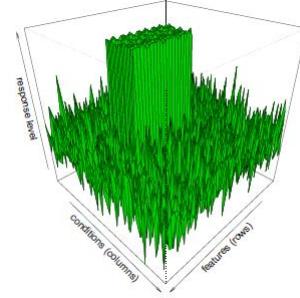
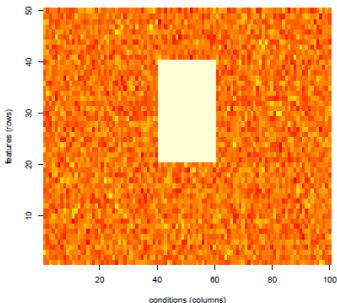


(b) Example 2.

Local patterns

- We are looking for:
 - A subset of features with the same characteristic across a subset of conditions.
 - Example:
 - A group of genes with the same expression patterns across a subset of samples.
 - A group of customers that buy the same products in a supermarket.
 - A group of students with the same results patterns across a group of subjects.
 -

Local patterns in a data matrix

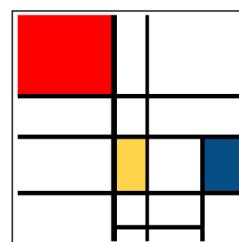


Example of a subset of features with high response level on a subset of conditions.

Local vs. global

A bicluster

$$\mathbf{X} = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1m} \\ X_{21} & X_{22} & \dots & X_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \dots & X_{nm} \end{pmatrix}.$$



A subset of features in a data matrix that have a similar response patterns across a subset of samples.

Example: a group of genes with a similar expression profiles across a group of samples

A bicluster: signal and noise

Within a bicluster: additive and multiplicative structures

$$Y = \text{signal} + \text{noise}$$

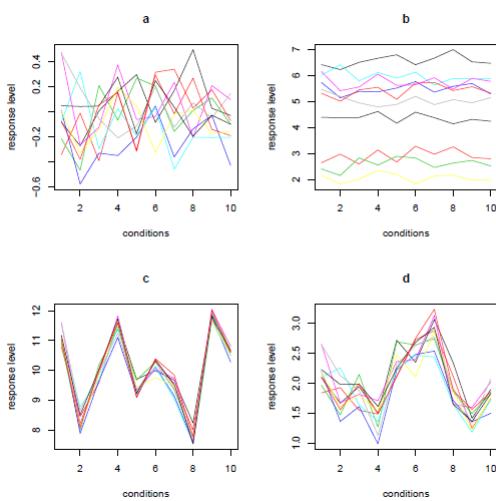
Signal structure: multiplicative or additive.

Outside a bicluster:

$$Y = \text{noise}$$

$$\begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,n} \\ x_{2,1} & x_{2,2} & \dots & x_{2,n} \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ x_{m,1} & x_{m,2} & \dots & x_{m,n} \end{bmatrix}$$

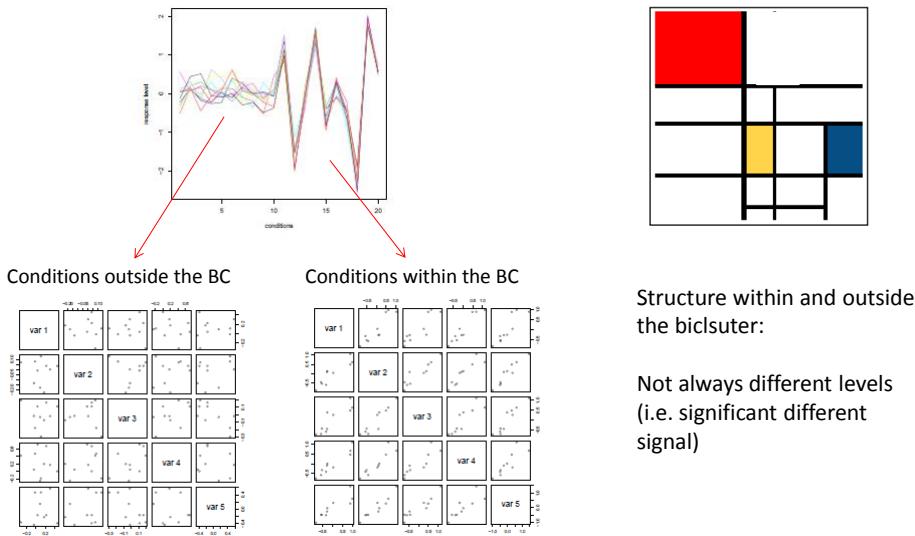
A bicluster: row and columns effects



Dominant effects:

- Rows ?
- Columns ?
- Rows and columns ?

A bicluster: correlation



Signal structure

$$Y = R + Error$$

$$Y = C + Error$$

Additive BCs

$$Y = R + C + Error$$

$$Y = R \times C \times Error$$

$$Y = R \times C^{Error}$$

$$Y = R \times C + Error$$

$$\log(Y) = \log(R) + \log(C) + Error$$

$$\log(Y) = Error \times (\log(R) + \log(C))$$

Multiplicative BCs

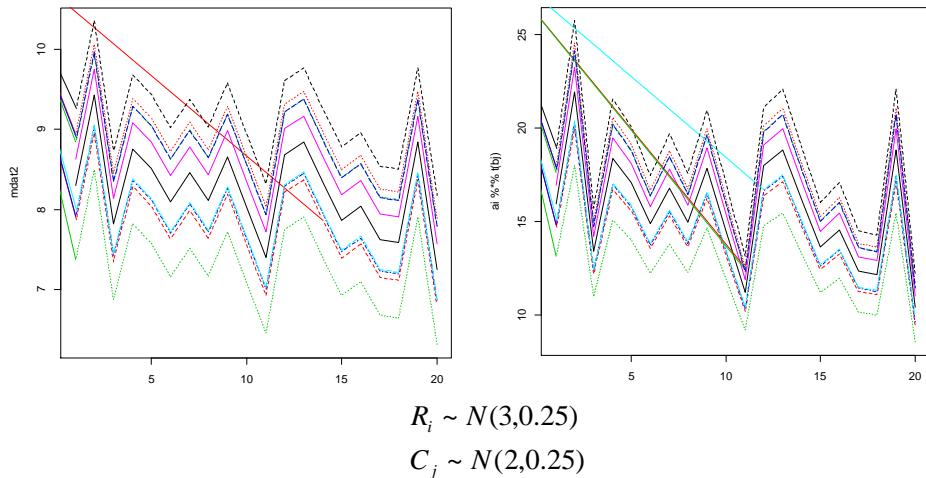
Signal structure

Additive BC:

$$Y = R + C$$

Multiplicative BC:

$$Y = R \times C$$



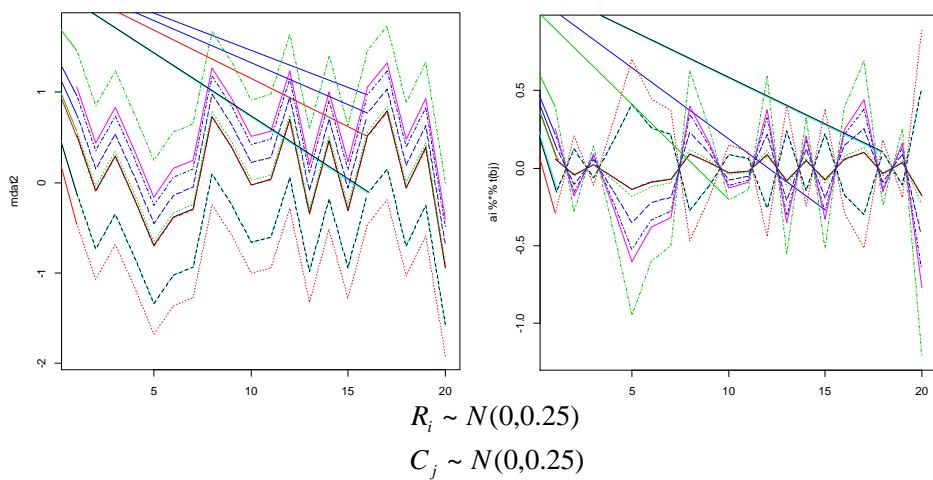
Signal structure

Additive BC:

$$Y = R + C$$

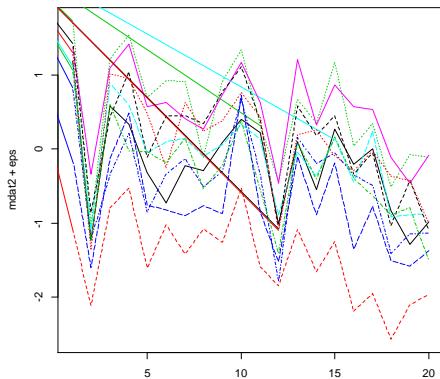
Multiplicative BC:

$$Y = R \times C$$

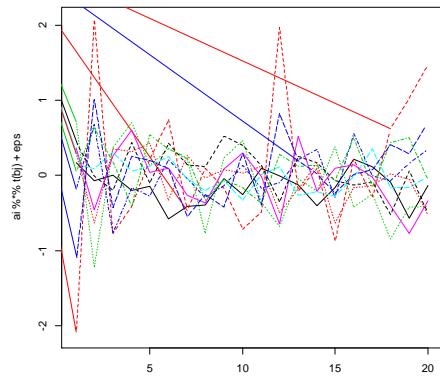


Signal + noise

Additive BC: $Y = R + C + E_{\text{error}}$



Multiplicative BC: $Y = R \times C + E_{\text{error}}$

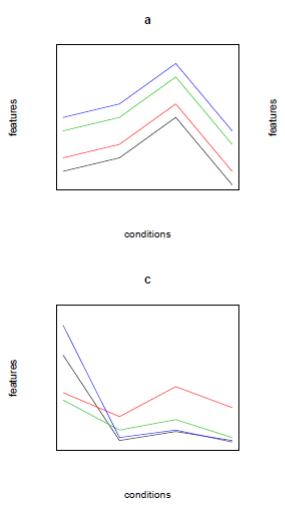


$$R_i \sim N(0, 0.25)$$

$$C_j \sim N(0, 0.25)$$

$$E_{ij} \sim N(0, 0.0625)$$

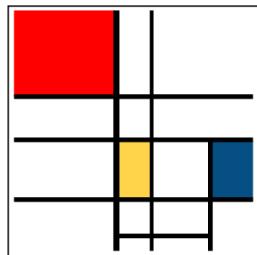
Types of biclusters



- Rows and columns effects.
- Coherent values .
- Coherent evolution.

Configurations of biclusters in the data matrix

Which structure we observed in the data matrix ?

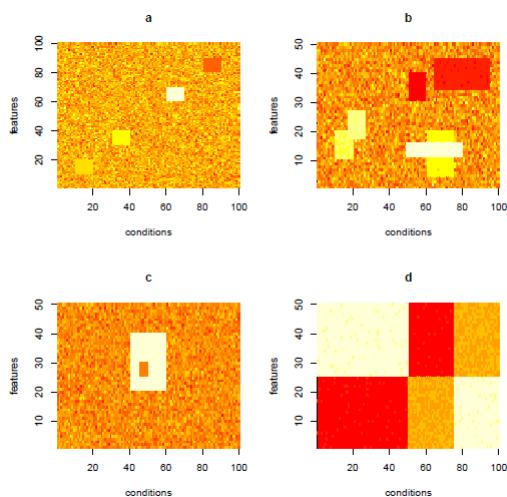


Piet Mondrian



Theo van Doesburg

Configurations of biclusters in the data matrix



Overlapping
Non overlapping

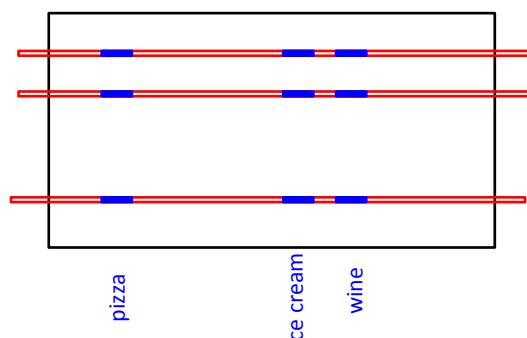
Local patterns

- Why local ?
- In a supermarket, if we know that a group of costumers have a tendency to buy : pizza, wine and ice cream we can help them to buy these products.
- In a holiday resort, if we know that costumers like to go to the sea and to have BBQ there....



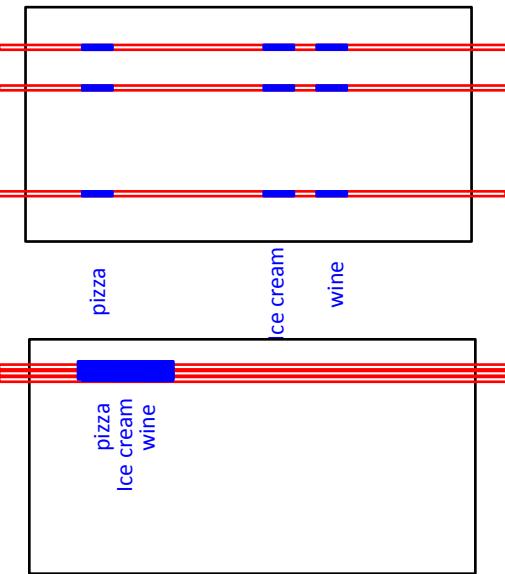
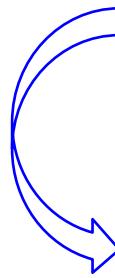
Local patterns

- In a supermarket, if we know that a group of costumers have a tendency to buy : pizza, wine and ice cream we can help them to buy these products.
- Why this is a bicluster ?

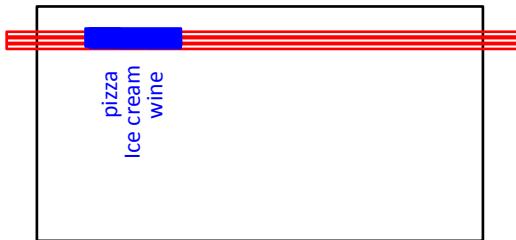


Local patterns

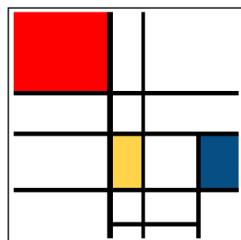
Re arrange
columns and rows



Many local patterns....



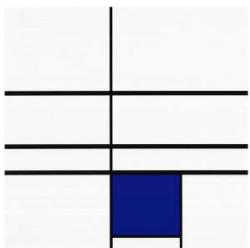
Pizza, ice cream &
wine



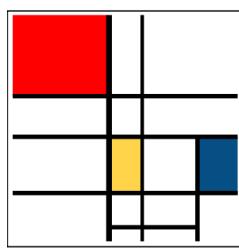
Organic
vegetables, free
ranged chicken
and cheese...

Exempels of BC

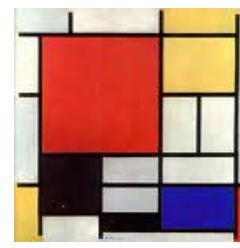
One bicluster



Three biclusters



Overlapping biclusters



Piet Mondriān

Overlapping only in one dimension (rows or columns).

Examples of overlapping biclusters

Many Overlapping biclusters



Theo van Doesburg

Overlapping only in one dimension (rows or columns).

Overlapping biclusters



Theo van Doesburg

Ben Nicholson



Other examples

Jean Arp



Paul Klee



Sonia Delaunay



Not everybody understood the concept of biclustering so good...

Part 3

Selection of Biclustering methods

- Computer science methods:
 - Bimax.
- Statistical methods:
 - The plaid model.
 - FABIA.

Part 3.1

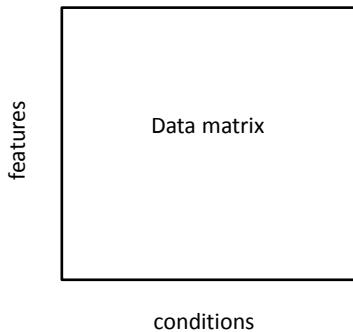
Bimax

Prelic, A., Bleuler, S., Zimmermann, P., Wil, A., Buhmann, P., Gruissem, W., Hennig, L., Thiele, L. and Zitzler, E. (2006) A systematic comparison and evaluation of biclustering methods for gene expression data. *Bioinformatics*, 22(9), 1122–1129.

Prelic et al. 2006

- Bimax (binary inclusion-maximal biclustering algorithm).
- The Bimax is a biclustering algorithm introduced by Prelic 2006 as a reference biclustering method for a comparison with different biclustering methods.

Data structure: binary data



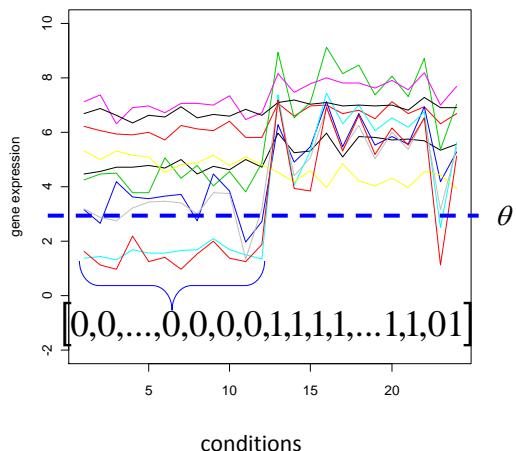
$$Z_{ij} = \begin{cases} 1 & \text{feature } i \text{ is active} \\ 0 & \text{on condition } j \\ & \text{otherwise} \end{cases}$$

Examples:

- In the supermarket: subject i buy product j
- In football: player i scores a goal in the last 10 minutes of the game.

Data structure: binary data

The original data is continuous.



Examples

- Gene i is expressed under condition j

$$Z_{ij} = \begin{cases} 1 & X_{ij} > \theta \\ 0 & X_{ij} \leq \theta \end{cases}$$

Data structure: binary data

$$\begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,n} \\ x_{2,1} & x_{2,2} & \dots & x_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \dots & x_{m,n} \end{bmatrix}$$

We are looking for subset of active features.

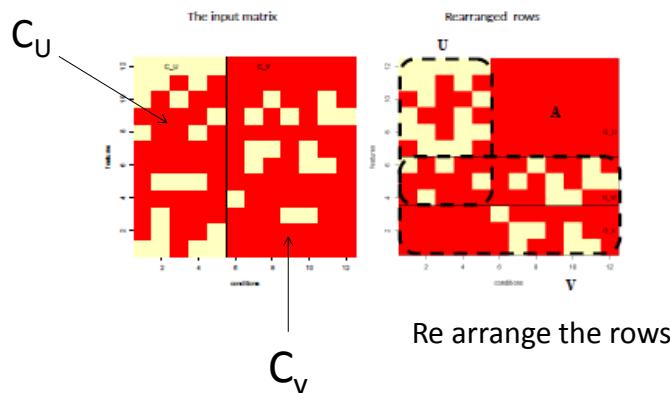
$$Z_{ij} = \begin{cases} 1 & \text{feature } i \text{ expressed (active) in condition } j \\ 0 & \text{otherwise,} \end{cases}$$

$$\begin{bmatrix} Z_{1,1} & 0 & \dots & 1 \\ Z_{2,1} & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ Z_{m,1} & 1 & \dots & 0 \end{bmatrix}$$

Can we find a sequence of 1s of features across the same conditions ?

The Bimax algorithm

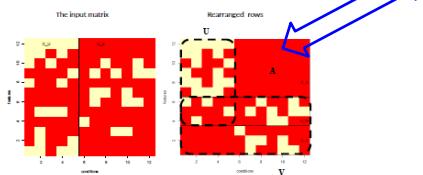
Divide the columns in two sets, C_U and C_V , based on the first row (in the first row, C_U contains only ones, while C_V contains only zeroes).



The Bimax algorithm

Step 1:

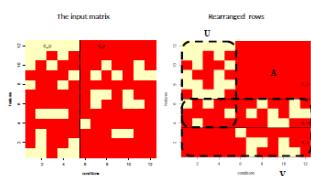
- Re arrange the data matrix.
- Exclude all rows/columns combinations with only zeros.



Step 2:

- Search for rows/columns combinations with ones.

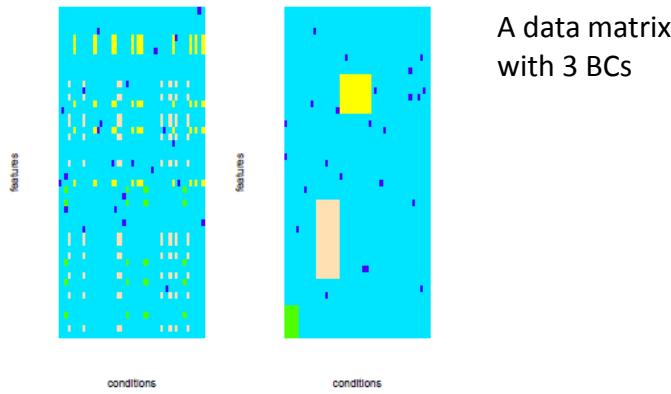
The Bimax algorithm: parameter setting



How many biclusters we are looking for ?

What is the minimum size of the bicluster (i.e. number of rows and columns) ?

The Bimax algorithm: an illustration

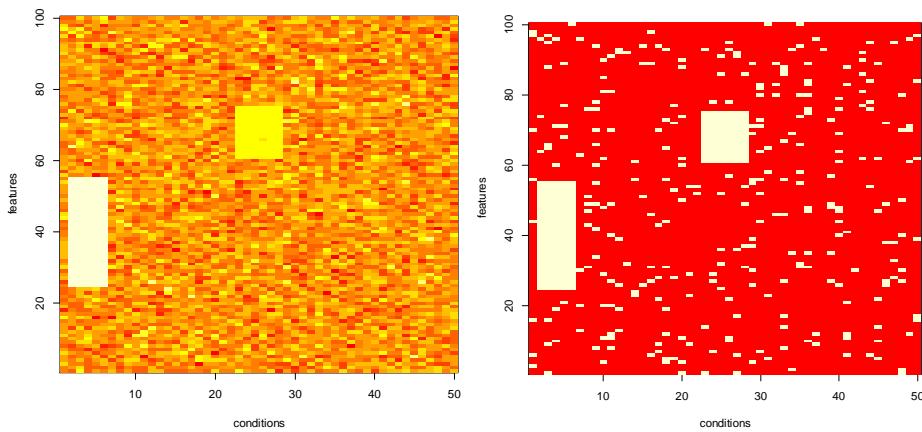


Example

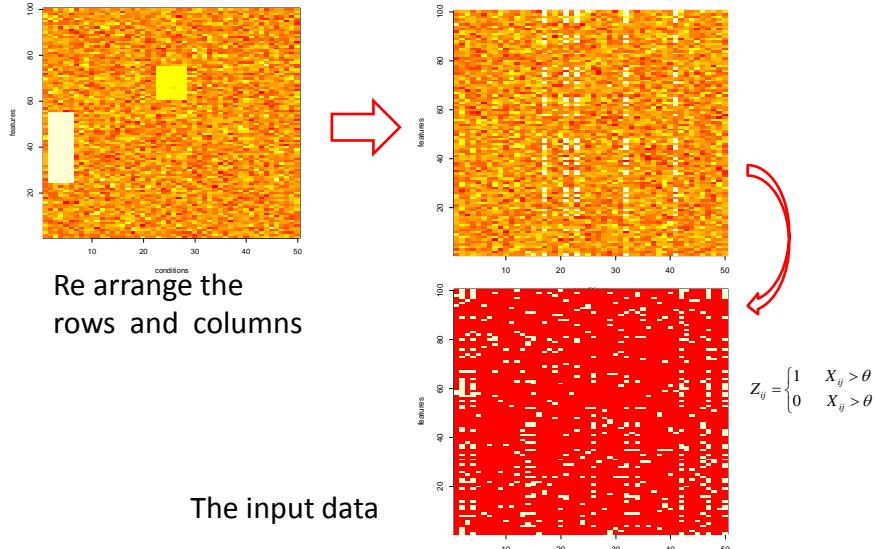
A 100 X 50 matrix with two BCs.

Before dichotomization (X).

After dichotomization (Z).



The same example: the input data



Results

Solutions

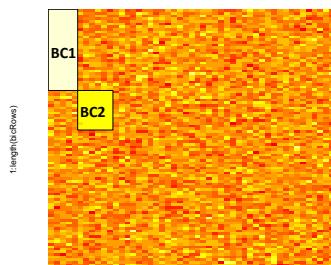
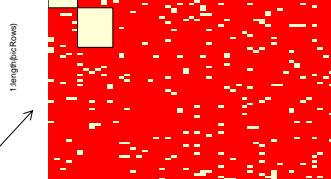
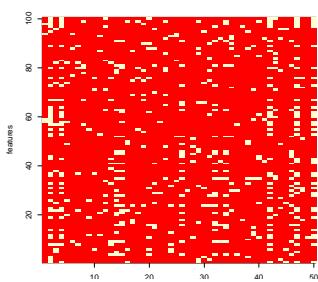
```
> test.b1<-binarize(test2,threshold=1.5)
> image(c(1:dim(test)[2]),c(1:dim(test)[1]),t(test.b1),ylab="features",xlab="conditions")
>
> bimaxbic<-biclust(test.b1,method=BCBimax(),minr=10,minc=5,number=2)
> summary(bimaxbic)
```

An object of class Biclust

```
call:
biclust(x = test.b1, method = BCBimax(), minr = 10, minc = 5,
number = 2)
```

Number of Clusters found: 2

```
Cluster sizes:
BC 1 BC 2
Number of Rows: 31 15
Number of Columns: 5 6
```

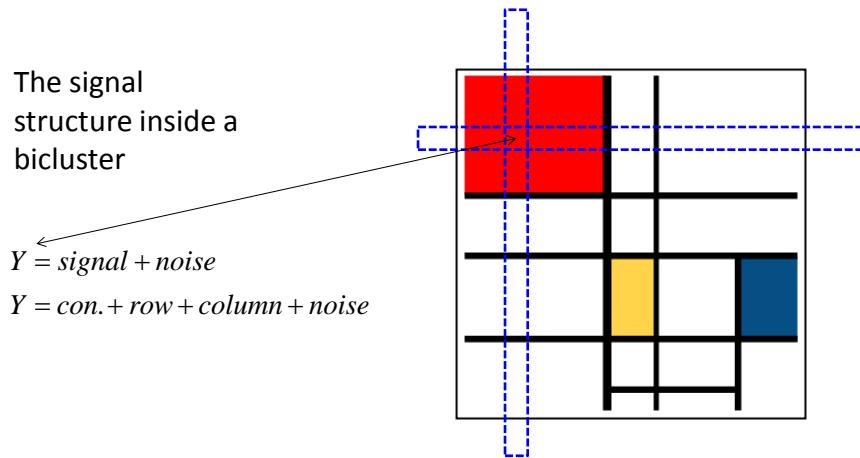


Part 3.2

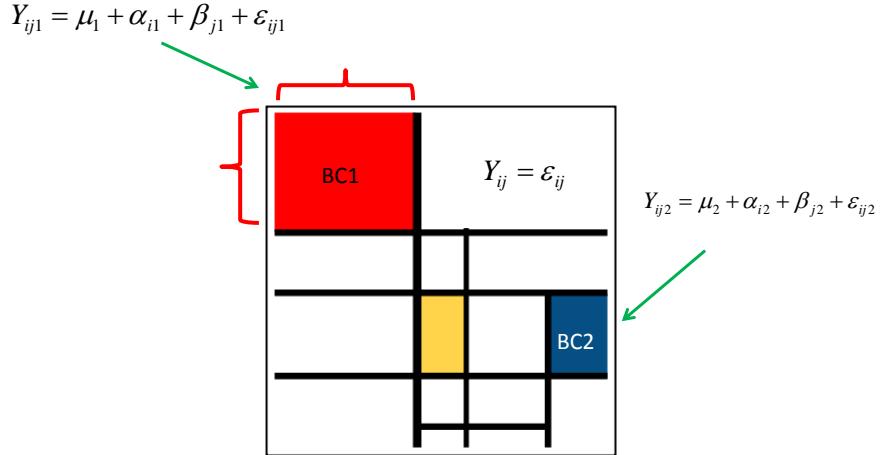
The plaid model

Turner, H., Bailey, T. and Krzanowski, W. (2005) Improved biclustering of microarray data demonstrated through systematic performance tests. *Computational Statistics and Data Analysis*, 48, 235–254.

Additive biclusters: the signal structure



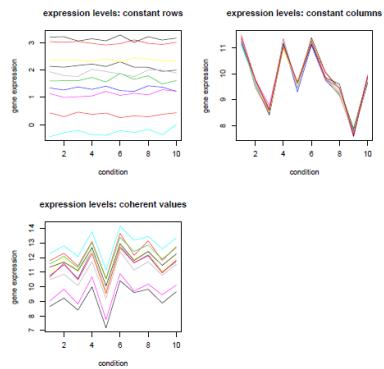
The plaid model



The plaid model: mean structure

$$Y_{ijk} = \mu_0 + \sum_{k=1}^K \theta_{ijk} \rho_{ik} K_{jk} + \varepsilon_{ijk}$$

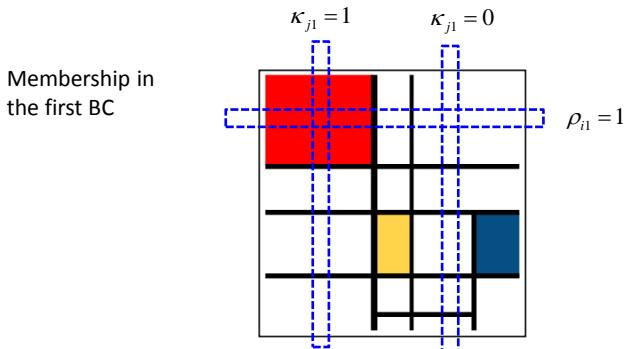
$$\theta_{ijk} = \begin{cases} \mu_k, & \text{Constant bicluster} \\ \mu_k + \alpha_{ik}, & \text{Constant rows} \\ \mu_k + \beta_{jk}, & \text{Constant cols.} \\ \mu_k + \alpha_{ik} + \beta_{jk}, & \text{Rows and cols. effects} \end{cases}$$



The plaid model: membership

$$Y_{ijk} = \mu_0 + \sum_{k=1}^K \theta_{ijk} \rho_{ik} \kappa_{jk} + \varepsilon_{ijk}$$

$$\kappa_{jk} = \begin{cases} 1 & \text{condition } j \text{ belongs to bicluster } k, \\ 0 & \text{otherwise.} \end{cases} \quad \rho_{ik} = \begin{cases} 1 & \text{gene } i \text{ belongs to bicluster } k, \\ 0 & \text{otherwise,} \end{cases}$$



Estimation the BC parameters

Given the membership in the k'th BC

$$Y_{ijk} = \mu_0 + \sum_{k=1}^K \theta_{ijk} \rho_{ik} \kappa_{jk} + \varepsilon_{ijk}$$

$$Y_{ijk} = (\mu_k + \alpha_{ik} + \beta_{jk} + \varepsilon_{ijk}) \rho_{ik} \times \kappa_{jk} + \varepsilon_{ijk}$$

$$\rho_{ik} = \kappa_{jk} = 1$$

$$Y_{ijk} = \underbrace{\mu_k + \alpha_{ik} + \beta_{jk}}_{\theta_{ijk}} + \varepsilon_{ijk}$$

A two way ANOVA with one observation per cell.

Estimation the BC parameters

Minimize the sum of squares for the k'th
BC

$$Q_k = \sum_{ij} (Y_{ijk} - \mu_k + \alpha_{ik} + \beta_{jk})^2$$

For all BCs

$$Q = \sum_{k=1}^K \sum_{ij} (Y_{ijk} - \mu_k + \alpha_{ik} + \beta_{jk})^2$$

In practice, per BC, two-way ANOVA with one observation per cell.

Estimation the membership parameters (rows)

Given the rows and columns effects
and the membership for the
columns.

$$Y_{ijk} = (\mu_k + \alpha_{ik} + \beta_{jk} + \varepsilon_{ijk}) \rho_{ik} \times \kappa_{jk} + \varepsilon_{ijk}$$

$\kappa_{jk} = 1$

$$Y_{ijk} = \rho_{ik} [(\mu_k + \alpha_{ik} + \beta_{jk}) \times \kappa_{jk}] + \varepsilon_{ijk}$$

“Known”

The only unknown is the row
membership.

Condition on the parameter
Estimates for BC effects and
membership (columns).

$$Y_{ijk} = \rho_{ik} [\hat{\mu}_k + \hat{\alpha}_{ik} + \hat{\beta}_{jk}] \times \hat{\kappa}_{jk} + \varepsilon_{ijk}$$

Minimize the residuals sum of
squares:

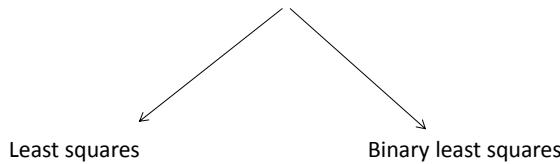
$$Q = \sum (Y_{ijk} - \rho_{ik} [\hat{\mu}_k + \hat{\alpha}_{ik} + \hat{\beta}_{jk}] \times \hat{\kappa}_{jk})^2$$

Estimation the membership parameters (rows)

Condition on the parameter estimates, linear regression model with one parameter

$$Y_{ijk} = \rho_{ik} [(\hat{\mu}_k + \hat{\alpha}_{ik} + \hat{\beta}_{jk}) \times \hat{\kappa}_{jk}] + \varepsilon_{ijk}$$

$$Q = \sum (Y_{ijk} - \rho_{ik} [(\hat{\mu}_k + \hat{\alpha}_{ik} + \hat{\beta}_{jk}) \times \hat{\kappa}_{jk}])^2$$



See later

Estimation the membership parameters (columns)

Given the rows and columns effects and the membership for the rows

$$Y_{ijk} = (\mu_k + \alpha_{ik} + \beta_{jk} + \varepsilon_{ijk}) \underbrace{\rho_{ik} \times \kappa_{jk}}_{\rho_{ik}=1} + \varepsilon_{ijk}$$

$$Y_{ijk} = \kappa_{jk} [(\mu_k + \alpha_{ik} + \beta_{jk}) \times \rho_{ik}] + \varepsilon_{ijk}$$

“Known”

The only unknown is the columns membership

Condition on the parameter estimates and membership (rows):

$$Y_{ijk} = \kappa_{ik} [(\hat{\mu}_k + \hat{\alpha}_{ik} + \hat{\beta}_{jk}) \times \hat{\rho}_{jk}] + \varepsilon_{ijk}$$

Estimation the membership parameters (columns)

Condition on the parameter estimates, linear regression model with one parameter

$$Y_{ijk} = \kappa_{ik} [(\hat{\mu}_k + \hat{\alpha}_{ik} + \hat{\beta}_{jk}) \times \hat{\rho}_{jk}] + \varepsilon_{ijk}$$

$$Q = \sum (Y_{ijk} - \kappa_{ik} [(\hat{\mu}_k + \hat{\alpha}_{ik} + \hat{\beta}_{jk}) \times \hat{\rho}_{jk}])^2$$

Search algorithm

Data structure for K BCs

$$Y_{ijk} = \mu_0 + \sum_{k=1}^K \theta_{ijk} \rho_{ik} \kappa_{jk} + \varepsilon_{ijk}$$

Residuals

$$Z_{ijk} = Y_{ijk} - \left(\mu_0 + \sum_{k=1}^K \hat{\theta}_{ijk} \hat{\rho}_{ik} \hat{\kappa}_{jk} \right)$$

Observed data Estimated BC and membership parameters

Search algorithm

Let us assume that L-1 BCs were found and we are looking for the L'th BC

$$Y_{ijk} = \mu_0 + \sum_{k=1}^{L-1} \theta_{ijk} \rho_{ik} \kappa_{jk} + \varepsilon_{ijk}$$

Residuals:

$$\hat{Z}_{ijk} = Y_{ijk} - \left(\mu_0 + \sum_{k=1}^{L-1} \theta_{ijk} \rho_{ik} \kappa_{jk} \right)$$

Residuals matrix:

\hat{Z} The input matrix for the next BC (the L'th BC)

Search algorithm

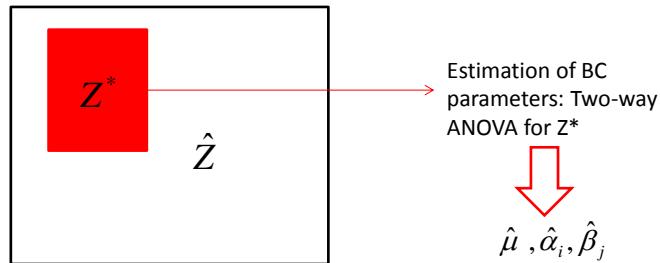
Input for the analysis of the L'th BC \hat{Z}

1. Compute \hat{Z} : matrix of residuals from the current model.
2. Compute starting values or initial memberships $\hat{\rho}_i^0$ and $\hat{\kappa}_i^0$.
3. Set $s=1$.
4. Update the layer effects using Z^* : submatrix of \hat{Z} indicated by $\hat{\rho}_i^{(s-1)}$ and $\hat{\kappa}_j^{(s-1)}$; $\hat{\mu}^s$, $\hat{\alpha}_i^s$ and $\hat{\beta}_j^s$.
5. Update cluster membership parameters: $\hat{\rho}_i^s$ and $\hat{\kappa}_j^s$
6. Repeat steps 4 and 5 for $s = 2, \dots, S$ iterations.
7. Compute $\hat{\mu}^{s+1}$, $\hat{\alpha}^{s+1}$, and $\hat{\beta}^{s+1}$ as in step 4.
8. Prune the bicluster to remove poor fitting rows and columns (see below).
9. Calculate layer sum of squares (LSS)
10. Permute \hat{Z} B times and follow steps 2 to 9 for each permutation.
11. Accept the bicluster if its LSS is greater than all permuted runs, otherwise stop.
12. sequentially, refit all layers in the model R times, then search for the next layer.

Search algorithm

Input for the analysis of the L'th BC \hat{Z}

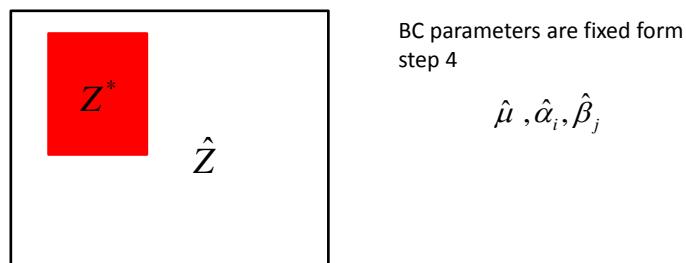
4. Update the layer effects using Z^* : submatrix of \hat{Z} indicated by $\hat{\rho}_i^{(s-1)}$ and $\hat{\kappa}_j^{(s-1)}$. $\hat{\mu}^s$, $\hat{\alpha}_i^s$ and $\hat{\beta}_j^s$.



Search algorithm

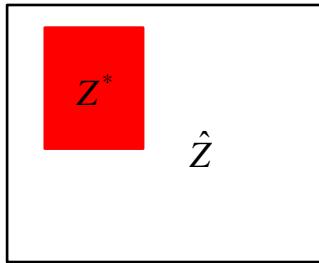
Input for the analysis of the L'th BC \hat{Z}

5. Update cluster membership parameters: $\hat{\rho}_i^s$ and $\hat{\kappa}_j^s$



Estimation the membership parameters: least squares (rows)

For the current BC:



$$Z_{ijk} = \rho_{ik} \left[(\hat{\mu}_k + \hat{\alpha}_{ik} + \hat{\beta}_{jk}) \times \hat{\kappa}_{jk} \right] + \varepsilon_{ijk}$$

$\hat{\theta}_{ijk}$ fixed form step 4

$$Z_{ijk} = \rho_{ik} \left[\hat{\theta}_{ijk} \times \hat{\kappa}_{jk} \right] + \varepsilon_{ijk}$$

$$Q = \sum (Z_{ijk} - \rho_{ik} [\hat{\theta}_{ijk} \times \hat{\kappa}_{jk}])^2$$

Least squares solution

$$Q = \sum (Z_{ijk} - \rho_{ik} [\hat{\theta}_{ijk} \times \hat{\kappa}_{jk}])^2 \quad \rightarrow \quad \rho_i = \frac{\sum_j \kappa_j \theta_{ij} Z_{ij}}{\sum_j \kappa_j^2 \theta_{ij}^2}$$

Estimation the membership parameters: least squares

rows

$$Q = \sum (Z_{ijk} - \rho_{ik} [\hat{\theta}_{ijk} \times \hat{\kappa}_{jk}])^2 \quad \rightarrow \quad \rho_i = \frac{\sum_j \kappa_j \theta_{ij} Z_{ij}}{\sum_j \kappa_j^2 \theta_{ij}^2}$$

columns

$$Q = \sum (Z_{ijk} - \kappa_{ik} [\hat{\theta}_{ijk} \times \hat{\rho}_{jk}])^2 \quad \rightarrow \quad \kappa_j = \frac{\sum_i \rho_i \theta_{ij} Z_{ij}}{\sum_i \rho_i^2 \theta_{ij}^2}$$

Search algorithm

9. Calculate layer sum of squares (LSS)
10. Permute $\hat{\mathbf{Z}}$ B times and follow steps 2 to 9 for each permutation.
11. Accept the bicluster if its LSS is greater than all permuted runs, otherwise stop.

Search algorithm

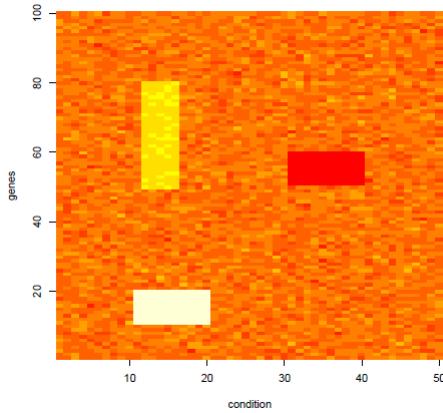
8. Prune the bicluster to remove poor fitting rows and columns

$$\hat{\rho}_i^s = \begin{cases} 1 & \text{if } \sum_j [Z_{ij} - \hat{\kappa}_j^{(s-1)}(\hat{\mu}^s + \hat{\alpha}_i^s + \hat{\beta}_j^s)]^2 < (1 - \tau_1) \sum_j Z_{ij}^2, \\ 0 & \text{otherwise.} \end{cases}$$


 $0 \leq \tau_1 \leq 1$

This means: a row is included if it leads to a reduction of τ_1 in the residuals sum of squares.

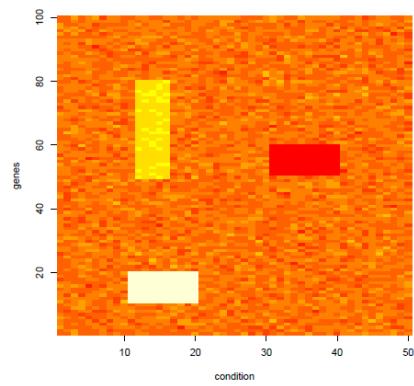
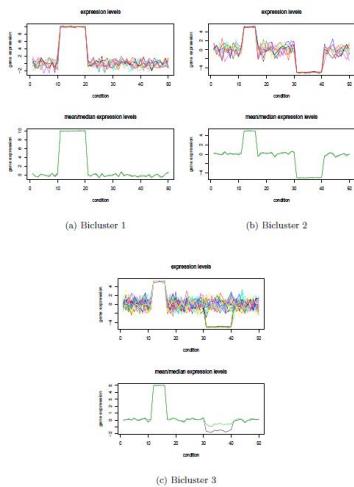
Example: the test data



A 100×50 data matrix with
3 BCs.

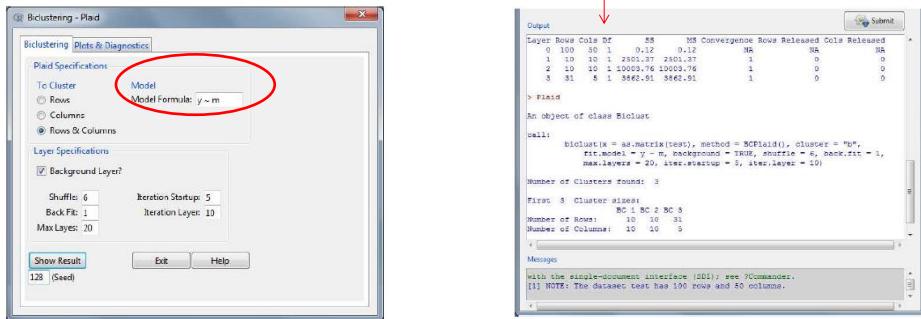
A group of rows are members
in two BCs.

Response profiles in the three BCs



The plaid model in R (I)

Constant BC

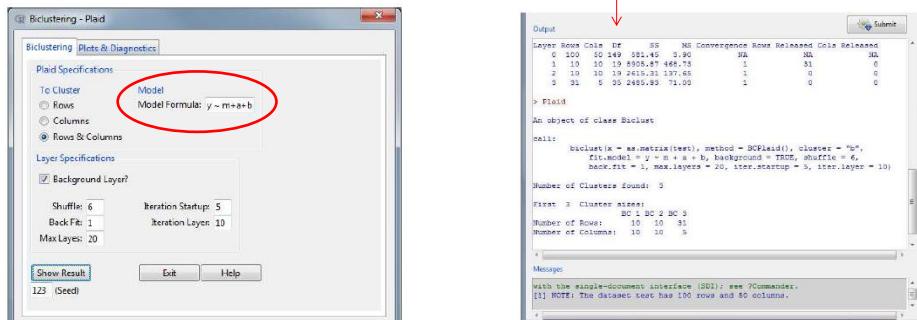


$$Y_{ijk} = \mu_0 + \sum_{k=1}^K \theta_{ijk} \rho_{ik} \kappa_{jk} + \varepsilon_{ijk}$$

$$\theta_{ijk} = \mu_k$$

The plaid model in R (II)

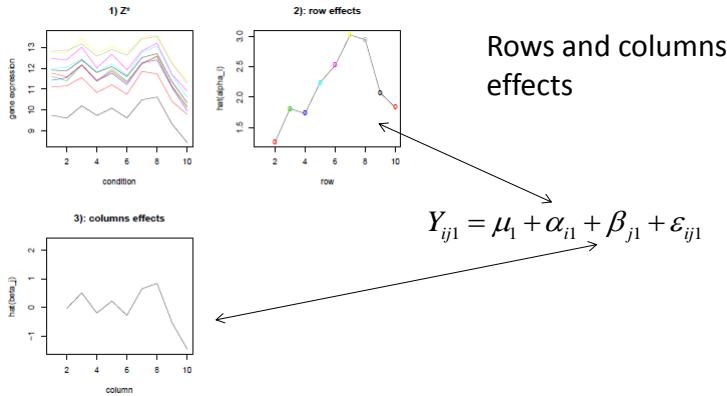
Rows and columns effects



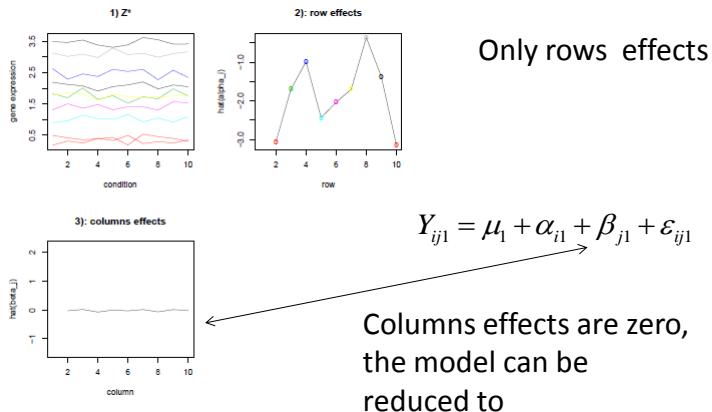
$$Y_{ijk} = \mu_0 + \sum_{k=1}^K \theta_{ijk} \rho_{ik} \kappa_{jk} + \varepsilon_{ijk}$$

$$\theta_{ijk} = \mu_k + \alpha_{ik} + \beta_{jk}$$

The mean structure within a BC (I)



The mean structure within a BC (II)



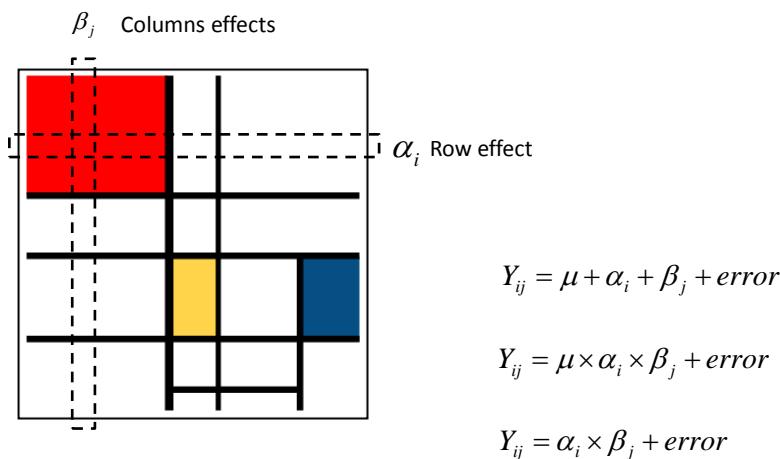
$$Y_{ij1} = \mu_1 + \alpha_{i1} + \varepsilon_{ij1}$$

Part 3.3

FABIA

Hochreiter, S., Bodenhofer, U., Heusel, M., Mayr, A., Mitterecker, A., Kasim, A., Khamiakova, T., Van Sanden, S., Lin, D., Talloen, W., Bijnens, L., Góohlmann, H. W. H., Shkedy, Z. and Clevert, D.-A. (2010a) Fabia: factor analysis for bicluster acquisition. *Bioinformatics*, 26, 1520–1527.

Multiplicative versus additive biclusters



Multiplicative bicluster: signal structure

$$\beta_j = \begin{cases} \beta_j & C_j \in BC \\ 0 & C_j \notin BC \end{cases}$$

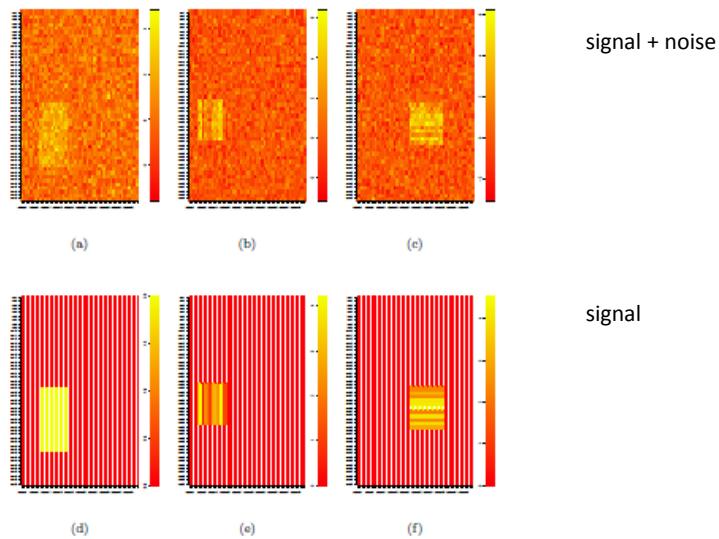
The diagram shows three components:

- α : A vertical vector of length 10, with values 0, 0, 0, 1, 2, 3, 4, 0, 0, 0.
- β^T : A horizontal vector of length 10, with values 0, 0, 0, 0, 0, 1, 2, 3, 4, 5, 0, 0, 0.
- $\alpha * \beta^T$: A 10x10 matrix where each element is the product of the corresponding elements from α and β^T . The resulting matrix has non-zero entries at positions (1,1), (2,2), (3,3), (4,4), (1,5), (2,6), (3,7), (4,8), (5,9), and (6,10).

 An arrow labeled "Membership vectors" points from the vectors to the matrix.

$$\alpha_i = \begin{cases} \alpha_i & R_i \in BC \\ 0 & R_i \notin BC \end{cases} \quad signal_{ij} = \alpha_i \times \beta_j$$

Examples



Multiplicative model

$$\alpha \begin{array}{c} * \\ \downarrow \end{array} \beta^T = \alpha * \beta^T$$

Observed data

$$Y_{ij} = signal_{ij} + error = \alpha_i \times \beta_j + error$$

A factor analysis model
in which a BC is a factor.

Signal structure:

$$Y = BC_1 + BC_2 + \dots + BC_K + error$$

A factor analysis model
with K factors.

FABIA: model formulation

$$Y = \sum_{k=1}^K \alpha_k \times \beta_k^T + Z$$

↓ ↓ ↓

Rows scores Columns scores error

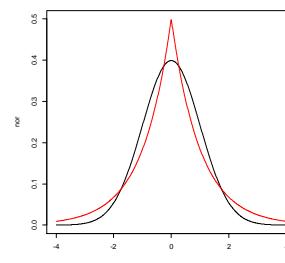
Model formulation of a
factor analysis model
with K factors.

Multiplicative signal.

For FABIA, priors for
factor loadings and
scores:

$$P(\alpha_k) \quad P(\beta_k)$$

Laplace distribution



FABIA: model formulation

$$Y = \sum_{k=1}^K \alpha_k \times \beta_k^T + Z$$

↓ ↓ ↓

Rows scores Columns scores error

$P(\alpha_k)$ $P(\beta_k)$ $N(0, \sigma^2)$

{ { { } } }

Laplace distribution

A factor analysis model:
 Rows scores (membership vector for rows): factor loadings.
 Columns scores (membership vector for columns): factor scores.

FABIA: example – a data matrix with one BC

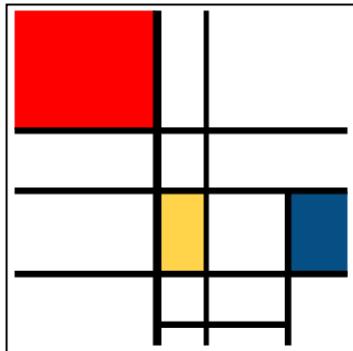
$$\begin{matrix} \alpha \\ * \\ \beta^T \end{matrix} = \begin{matrix} \alpha * \beta^T \end{matrix}$$

$Y = \sum_{k=1}^K \alpha_k \times \beta_k^T + Z$

↓ ↓ ↓

Rows scores Columns scores error

FABIA: example – a data matrix with three BCs

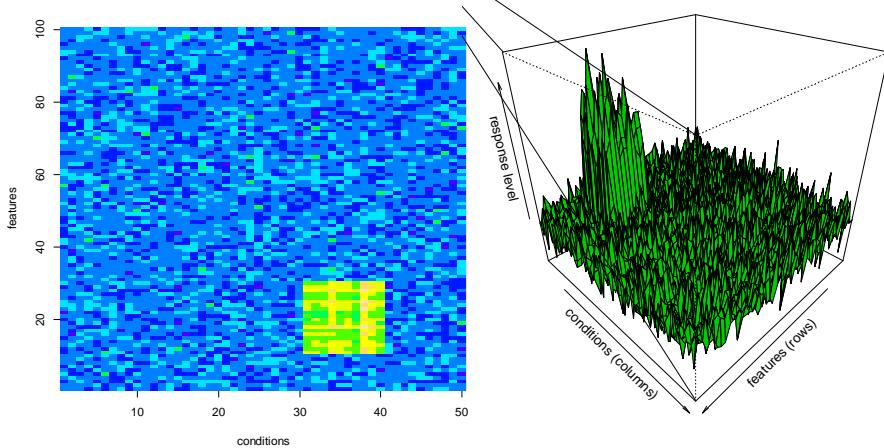


$$\begin{aligned}
 Y &= \underbrace{\alpha_1 \times \beta_1^T}_{\text{Factor 1}} + \underbrace{\alpha_2 \times \beta_2^T}_{\text{Factor 2}} + \underbrace{\alpha_3 \times \beta_3^T}_{\text{Factor 3}} + Z \\
 &= \text{Factor 1} + \text{Factor 2} + \text{Factor 3} + Z
 \end{aligned}$$

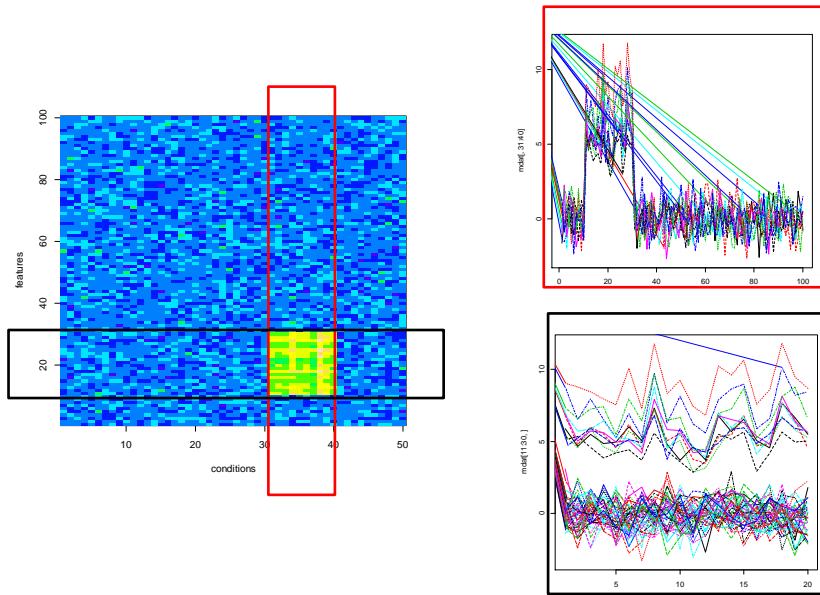
Factor 1 Factor 2 Factor 3
 First BC Second BC Third BC

Example: one BC

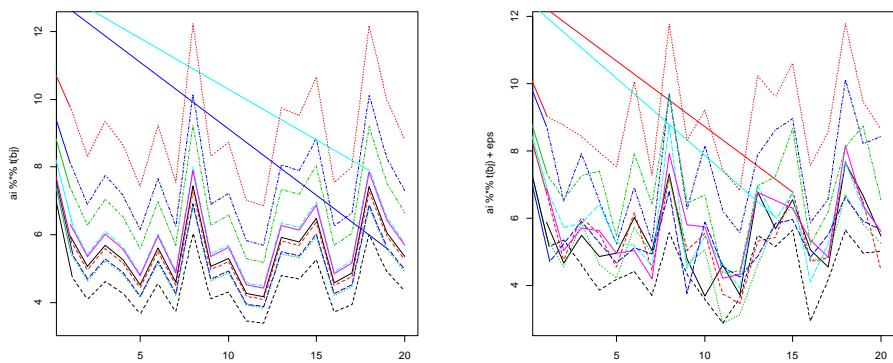
A 100X50 data matrix with one BC



Example: one BC – rows and columns



Example: signal, and signal+noise



Data analysis

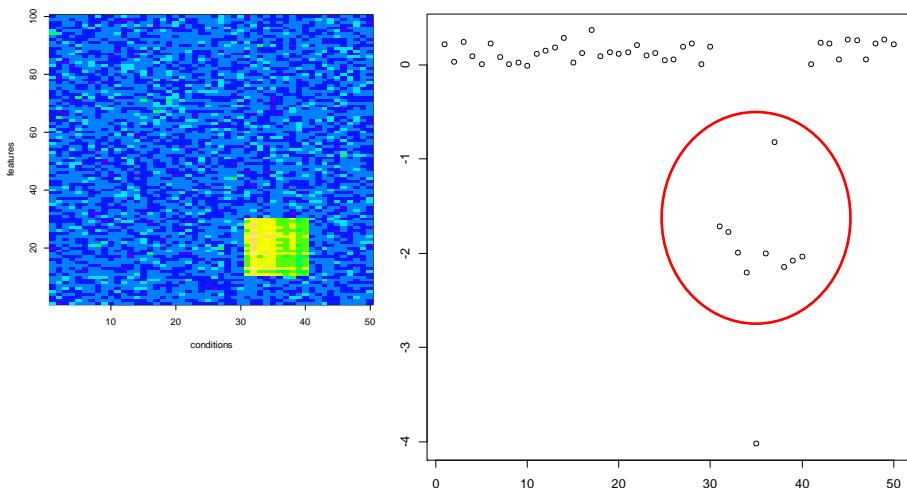
A factor analysis model
with one factor:

$$Y = \alpha_1 \times \beta_2^T + Z$$

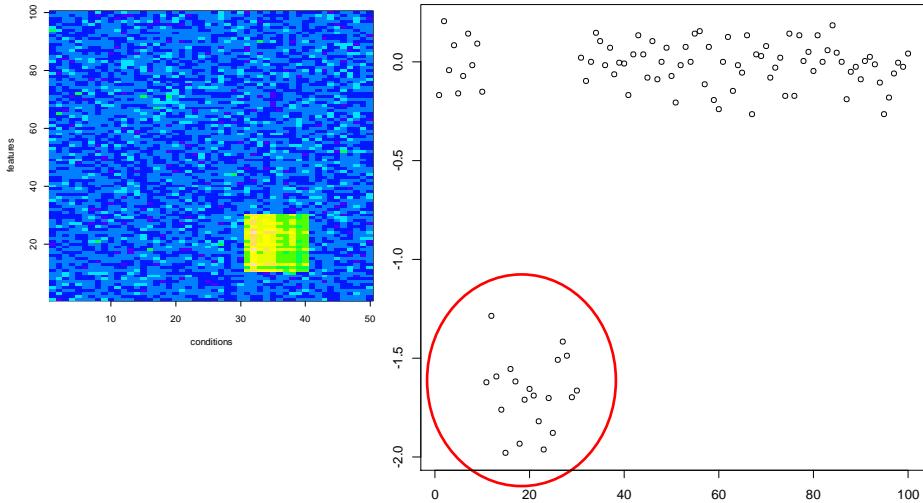
In R:

```
> fabRes <- fabia(mdat,p=1)
```

Results: factor scores (columns)



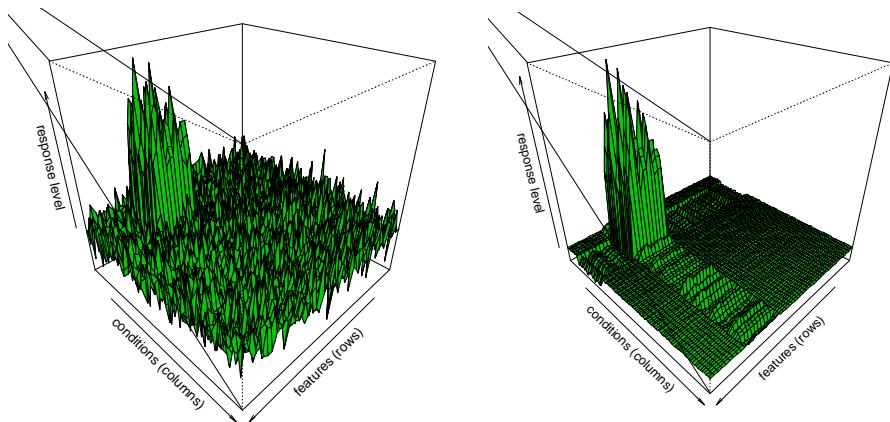
Results: factor loadings (rows)



Observed and predicted data

$$Y = \sum_{k=1}^K \alpha_k \beta_k^T + Z$$

$$\hat{Y} = \sum_{k=1}^K \hat{\alpha}_k \hat{\beta}_k^T$$



Short summary: methods

- Many other methods were developed.
- Local patterns.
- Trying to discover the signal in a noisy data.
- Best method ? No, completely data dependent.
- For all method: subjective selection of parameter settings !
- For most of the methods, multiple runs leads to multiple results !!
- Robust analysis should be performed !!!

Short summary: software

- Method specific (many methods and packages are available).
- Genreal:
 - biclust.
 - biclustGUI.
 - biclust shiny App.
 - online and cloud products.

Part 4

Case Studies

Part 4.1

Biclustering for Market segmentation

Market segmentation

- Market segmentation is essential for marketing success.
- The most successful firms drive their businesses based on segmentation.
- In tourism:
 - identify groups of tourists who share common characteristics.
 - Make it possible to develop a tailored marketing mix to most successfully attract such subgroups of the market.
 - Focusing on subgroups increases the chances of success within the subgroup.

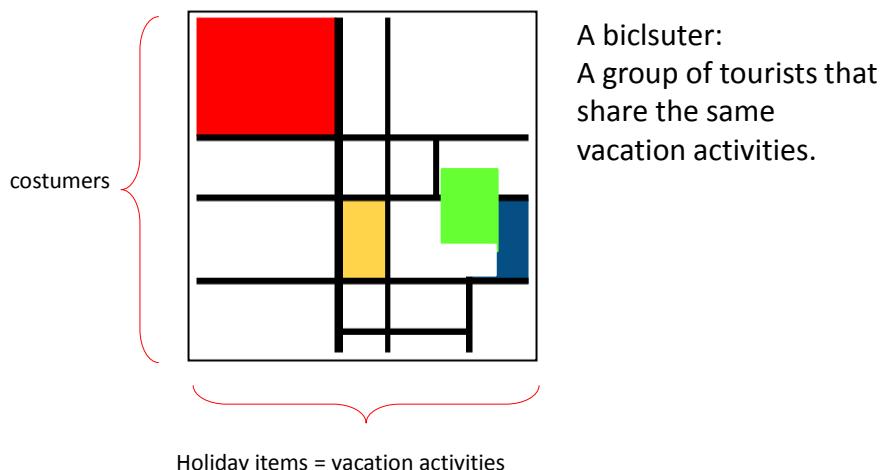
Dimensionality problem

- One of the typical methodological challenge:
 - large amount of information (responses to many survey questions) is available from tourists....
 - But typically the sample sizes are too low given the number of variables used to conduct segmentation analysis.
- Solution: collect large samples that allow segmentation with a large number of variables.

The tourism survey

- The data set used for this illustration is a tourism survey of adult Australians (internet based survey).
- Participants were asked questions about their **general travel behavior**, their travel behavior on their last Australian vacation, benefits they perceive of undertaking travel, and image perceptions of their ideal tourism destination.
- Information was also collected about the participants age, gender, annual household income, marital status, education level, occupation, family structure, and media consumption.

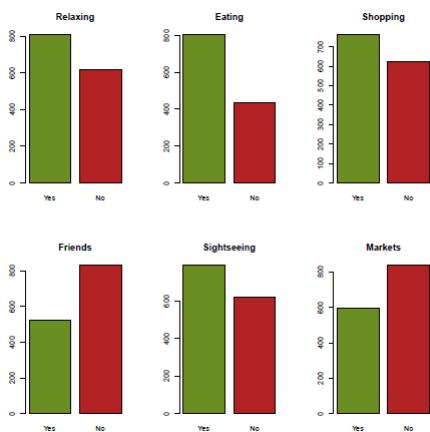
Data structure



Consumption of holiday activities

- In the present data set 1,003 respondents were asked to state for 44 vacation activities whether they engaged in them during their last vacation.
- Activities includes: relaxing, eating in reasonably priced eateries, shopping, sightseeing, visiting industrial attractions (such as wineries, breweries, mines, etc.), going to markets, scenic walks, visiting museums and monuments, botanic and public gardens, and the countryside/farms.

Consumption of holiday activities



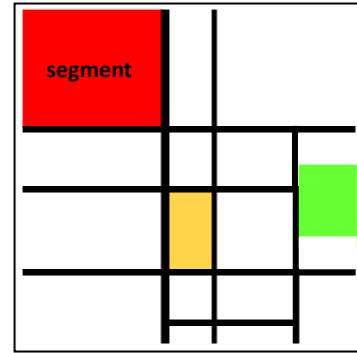
Distribution per item gives information how popular is an item among the costumers **but...**

..we do not know which items are consumed together.

Bicluster configuration: market segmentation

vacation activities

costumers	vacation activities
1 1 1	0 0 0 0 0 0 0 0 0
1 1 1	0 0 0 0 0 0 0 0 0
1 1 1	0 0 1 0 0 0 0 0 0
1 1 1	0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0	1 1 1
0 0 1 0 0 0 0 0 0	1 1 1
0 0 0 0 1 1 0 0 0	1 1 1
0 0 0 0 1 1 0 0 0	0 0 0 0
0 0 0 0 1 1 0 0 0	0 0 0 0
0 1 0 0 0 1 0 0 0	0 0 0 0

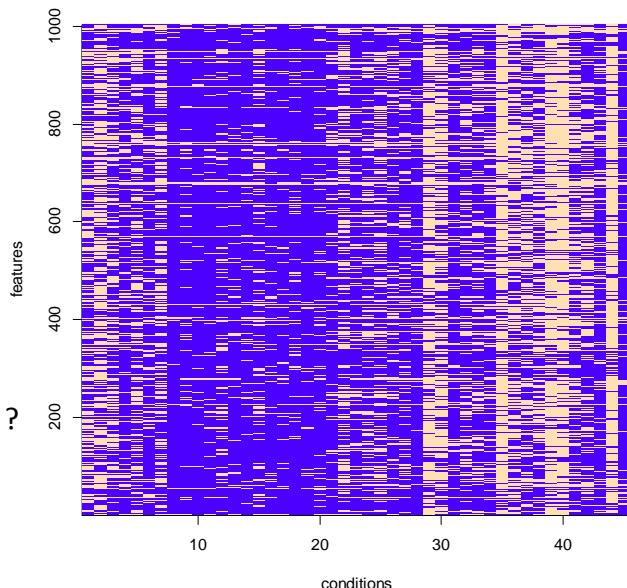


Observed data

```
[1 1 1 0 0 0 0 0 0 0 0 0
1 1 1 0 0 0 0 0 0 0 0 0
1 1 1 0 0 0 1 0 0 0 0 0
1 1 1 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 1 1 1
0 0 1 0 0 0 0 0 1 1 1
0 0 0 0 1 1 0 1 1 1
0 0 0 0 1 1 0 0 0 0
0 0 0 0 1 1 0 0 0 0
0 1 0 0 0 1 0 0 0 0]
```

A 1003X45 binary data.

Observed patterns ?



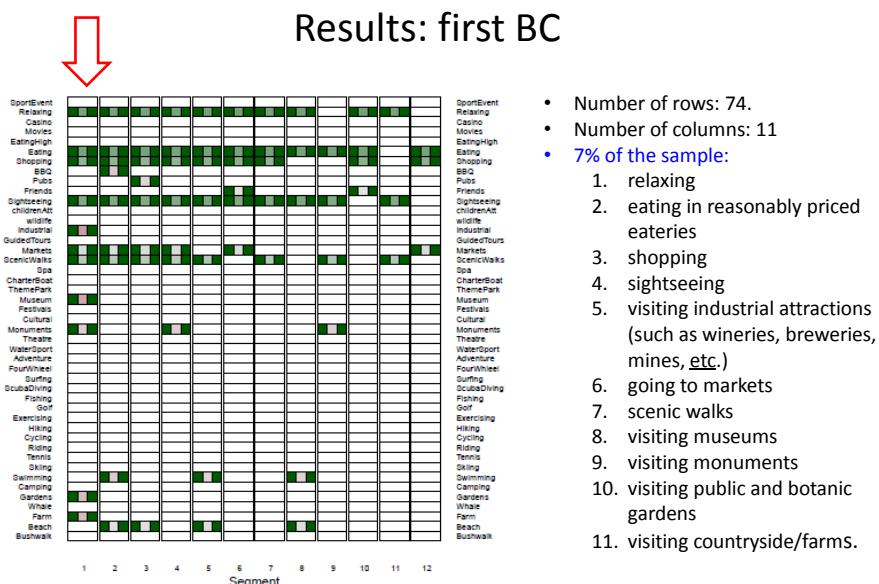
Data analysis using Bimax

1	1	1	1	0	0	0	0	0	0	0
1	1	1	0	0	0	0	0	0	0	0
1	1	1	1	0	1	0	0	0	0	0
1	1	1	1	0	0	0	0	0	0	0
0	0	0	0	0	0	0	1	1	1	
0	0	1	0	0	0	0	1	1	1	
0	0	0	0	1	1	0	1	1	1	
0	0	0	0	1	1	0	0	0	0	
0	0	0	0	1	1	0	0	0	0	
0	1	0	0	0	1	0	0	0	0	

Find group od subjects with the same sequence of 1s.

Minimal size of BC ?

Results: first BC

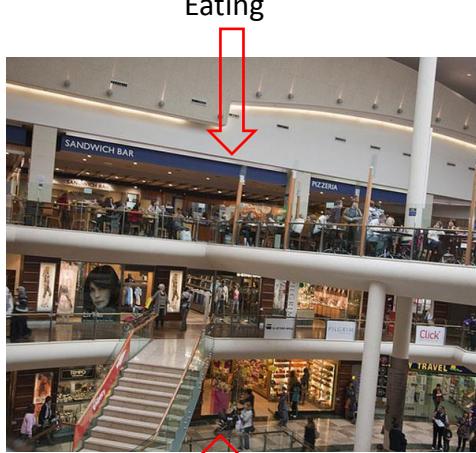


Results: first BC



Segment	1	2	3	4	5	6	7	8	9	10	11	12
Relaxing	■	■	■	■	■	■	■	■	■	■	■	■
Entertaining	■	■	■	■	■	■	■	■	■	■	■	■
Cooking	■	■	■	■	■	■	■	■	■	■	■	■
Swimming	■	■	■	■	■	■	■	■	■	■	■	■
Dining	■	■	■	■	■	■	■	■	■	■	■	■
Entertainment	■	■	■	■	■	■	■	■	■	■	■	■
Reading	■	■	■	■	■	■	■	■	■	■	■	■
Walking	■	■	■	■	■	■	■	■	■	■	■	■
Parks	■	■	■	■	■	■	■	■	■	■	■	■
Photographing	■	■	■	■	■	■	■	■	■	■	■	■
Sightseeing	■	■	■	■	■	■	■	■	■	■	■	■
Industrial	■	■	■	■	■	■	■	■	■	■	■	■
Outdoors	■	■	■	■	■	■	■	■	■	■	■	■
Markets	■	■	■	■	■	■	■	■	■	■	■	■
Scenic walks	■	■	■	■	■	■	■	■	■	■	■	■
Chocolatier	■	■	■	■	■	■	■	■	■	■	■	■
Monuments	■	■	■	■	■	■	■	■	■	■	■	■
Waterport	■	■	■	■	■	■	■	■	■	■	■	■
Adventure	■	■	■	■	■	■	■	■	■	■	■	■
Hobbies	■	■	■	■	■	■	■	■	■	■	■	■
Culture	■	■	■	■	■	■	■	■	■	■	■	■
Museums	■	■	■	■	■	■	■	■	■	■	■	■
Festivals	■	■	■	■	■	■	■	■	■	■	■	■
Events	■	■	■	■	■	■	■	■	■	■	■	■
Evening	■	■	■	■	■	■	■	■	■	■	■	■
Shopping	■	■	■	■	■	■	■	■	■	■	■	■
Swimming	■	■	■	■	■	■	■	■	■	■	■	■
Gardening	■	■	■	■	■	■	■	■	■	■	■	■
Entertaining	■	■	■	■	■	■	■	■	■	■	■	■
Whale	■	■	■	■	■	■	■	■	■	■	■	■
Beach	■	■	■	■	■	■	■	■	■	■	■	■

- Number of rows: 74.
- Number of columns: 11
- 7% of the sample:
 1. relaxing
 2. **eating in reasonably priced eateries**
 3. shopping
 4. sightseeing
 5. visiting industrial attractions (such as wineries, breweries, mines, etc.)
 6. going to markets
 7. scenic walks
 8. visiting museums
 9. visiting monuments
 10. visiting public and botanic gardens
 11. visiting countryside/farms.



Shopping



Results: second BC

Segment	1	2	3	4	5	6	7	8	9	10	11	12
Relaxing	■	■	■	■	■	■	■	■	■	■	■	■
Entertaining	■	■	■	■	■	■	■	■	■	■	■	■
Cooking	■	■	■	■	■	■	■	■	■	■	■	■
Swimming	■	■	■	■	■	■	■	■	■	■	■	■
Dining	■	■	■	■	■	■	■	■	■	■	■	■
Entertainment	■	■	■	■	■	■	■	■	■	■	■	■
Reading	■	■	■	■	■	■	■	■	■	■	■	■
Walking	■	■	■	■	■	■	■	■	■	■	■	■
Parks	■	■	■	■	■	■	■	■	■	■	■	■
Photographing	■	■	■	■	■	■	■	■	■	■	■	■
Sightseeing	■	■	■	■	■	■	■	■	■	■	■	■
Industrial	■	■	■	■	■	■	■	■	■	■	■	■
Outdoors	■	■	■	■	■	■	■	■	■	■	■	■
Markets	■	■	■	■	■	■	■	■	■	■	■	■
Scenic walks	■	■	■	■	■	■	■	■	■	■	■	■
Chocolatier	■	■	■	■	■	■	■	■	■	■	■	■
Monuments	■	■	■	■	■	■	■	■	■	■	■	■
Waterport	■	■	■	■	■	■	■	■	■	■	■	■
Adventure	■	■	■	■	■	■	■	■	■	■	■	■
Hobbies	■	■	■	■	■	■	■	■	■	■	■	■
Culture	■	■	■	■	■	■	■	■	■	■	■	■
Museums	■	■	■	■	■	■	■	■	■	■	■	■
Festivals	■	■	■	■	■	■	■	■	■	■	■	■
Events	■	■	■	■	■	■	■	■	■	■	■	■
Evening	■	■	■	■	■	■	■	■	■	■	■	■
Shopping	■	■	■	■	■	■	■	■	■	■	■	■
Swimming	■	■	■	■	■	■	■	■	■	■	■	■
Gardening	■	■	■	■	■	■	■	■	■	■	■	■
Entertaining	■	■	■	■	■	■	■	■	■	■	■	■
Whale	■	■	■	■	■	■	■	■	■	■	■	■
Beach	■	■	■	■	■	■	■	■	■	■	■	■

- Number of rows: 87.
- Number of columns: 9
- 8.6% of the sample:
 1. relaxing
 2. **eating in reasonably priced eateries**
 3. Shopping
 4. **BBQ**
 5. sightseeing
 6. going to markets
 7. scenic walks
 8. Swimming
 9. Beach

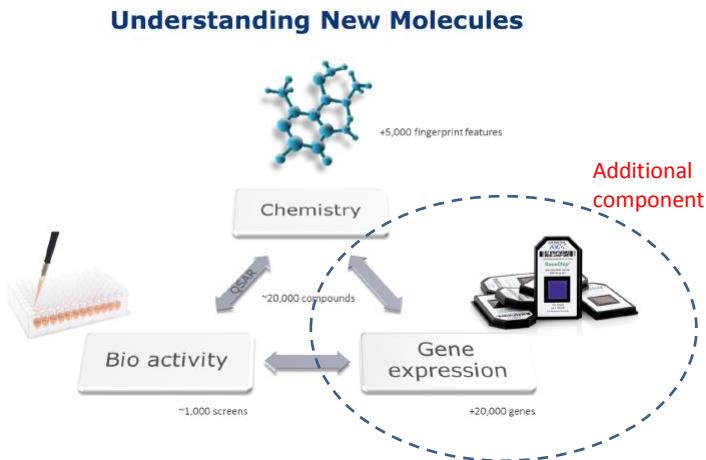


Software

Part 4.2

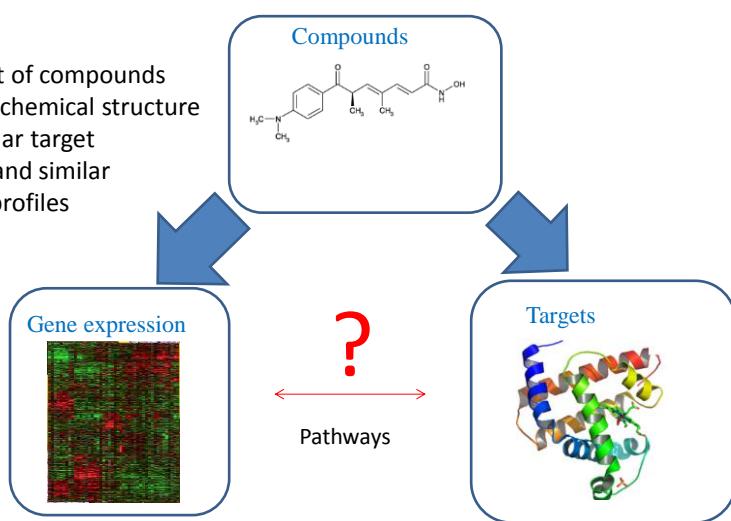
Drug Discovery (I):
Biclustering methods for chemoinformatics

Quantifying S-Transcription-A-R



Relating gene expression profiles to Protein targets via compounds

Aim:
find a subset of compounds
with similar chemical structure
(= have similar target
prediction) and similar
expression profiles

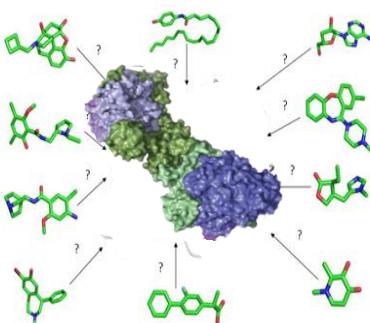


110

Protein Targets and Target prediction

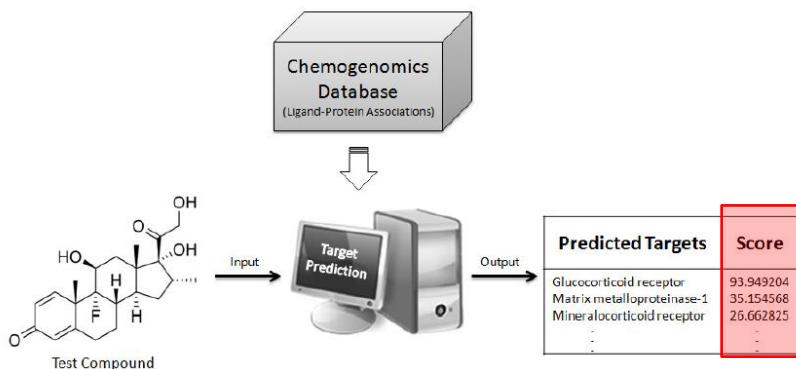
- Many candidate molecules
 - With unknown mechanism of action
- One drug – many targets
- One target – many active sites (for binding)
- Difficult (expensive) to measure activity of a molecule in all assays
- Predict if drugs will bind to a target given its chemical structure and already known drug-target associations?
- Target prediction

e.g. Histone deacetylase enzyme



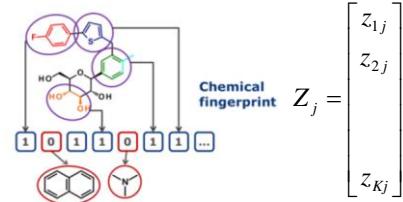
Target Prediction Score

- Likelihood of binding of a compound to every protein target (Koutsoukas, 2011)



Protein Targets and Target prediction

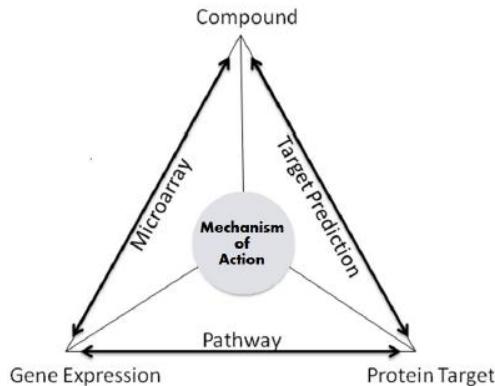
- Predict if drugs will bind to a target given its chemical structure and already known drug-target associations?
- Target prediction



$$P(\text{TARGET}_j) = f(FP_1, FP_2, \dots, FP_K)$$

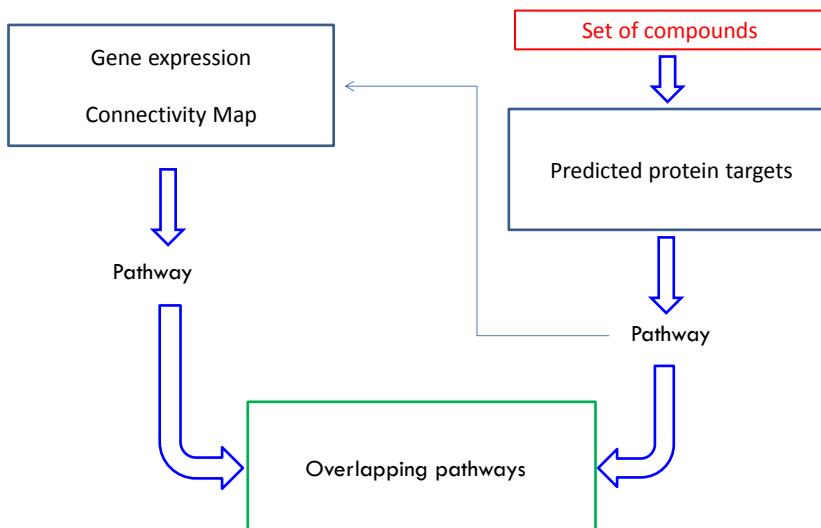
The setting

Mechanism of Action of compound



- Drugs regulating similar protein targets (similar structure) affects similar set of genes

Genes and protein targets pathway overlap



116

Gene expression profiles

- Connectivity map data :
 - ❖ 4 cell lines(MCF7,PC3,HL60 and SKMEL5)
 - ❖ After pre-processing ~2400 genes
 - ❖ 1309 drug like compounds
 - ❖ Similar concentration and time of compound exposure

X = Gene Expression
J = 2340 genes
I = 36 compounds
(MCF7 cell line, 6 hours,
10micromolars)

$$X = \begin{bmatrix} x_{11} & x_{21} & \dots & x_{n1} \\ x_{12} & x_{22} & \dots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1m} & x_{2m} & \dots & x_{nm} \end{bmatrix}$$

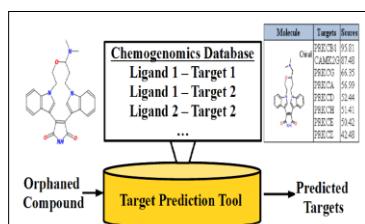
genes

compounds

117

Target prediction: binary scores

- Biosar (Naïve Bayes) used to predict targets
- Individual cut off used for each target



$$t_{ip} = \begin{cases} 1 & \text{Comp } i \text{ hit on target } p \\ 0 & \text{otherwise} \end{cases}$$

$$T_{C_i} = (0, 1, 1, 0, 0, 0, \dots, 1, 0)$$

Target scores matrix

$$T = \begin{bmatrix} t_{11} & t_{21} & \dots & t_{n1} \\ t_{12} & t_{22} & \dots & t_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ t_{1m} & t_{2m} & \dots & t_{nm} \end{bmatrix}$$

targets

compounds

118

Pathways

- Specific group of compounds
- A group of genes and targets that share:
 - a biological pathway
 - a statistical pathway

119

Biological pathways

- KEGG (Kyoto Encyclopedia of Genes and Genomes)- is a freely available information repository of the network of genes and molecules for practical analysis of the gene functions
- GO (Gene Ontology)- is a bioinformatics project that is the largest repository for catalogue gene function that unifies the representation of gene and gene product attribute across all the species.
- KEGG and Go pathways were annotated to proteins and genes.

120

Part II Data analysis

121

Data analysis steps

- Target based clustering.
 similarity matrix based on target prediction scores.
- Gene expression profiling.
- Enrichment of the gene set.
- Pathway identification.

122

Correlation between compounds

Target scores matrix

$$T = \begin{bmatrix} t_{11} & t_{21} & \dots & t_{n1} \\ t_{12} & t_{22} & \dots & t_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ t_{1m} & t_{2m} & \dots & t_{nm} \end{bmatrix}$$

compounds

targets

- N_{c1} and N_{c2} are the number of fingerprint features present in compound 1 and compound 2.
- N_{c12} is the number of features common to both compounds.

Tanimoto scores

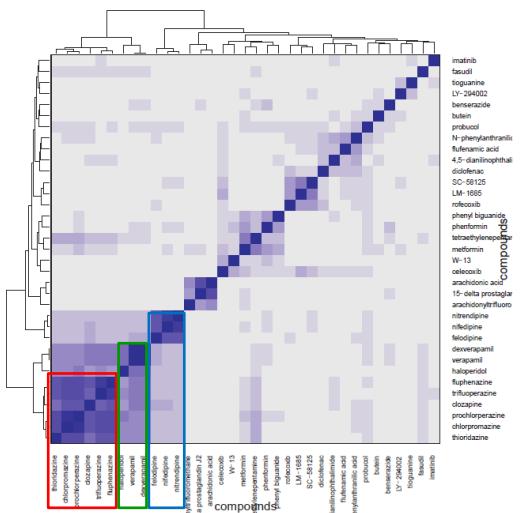
$$TC = \frac{N_{c12}}{N_{c1} + N_{c2} - N_{c12}}$$

$TC = 1$ 2 compounds are identical given the set of chemical structures

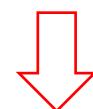
$TC = 0$ 2 compounds do not share any chemical structure

123

Target similarity matrix



$$TC = \frac{N_{c12}}{N_{c1} + N_{c2} - N_{c12}}$$



Cluster compound based on similarity scores

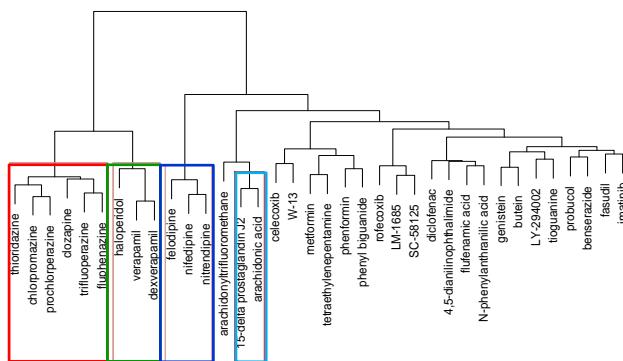
124

Hierarchical clustering

- Input: Similarity matrix
- Start: each compound is a cluster
- Merge compounds according to a criterion
- Ward's distance
- End: single cluster of all compounds

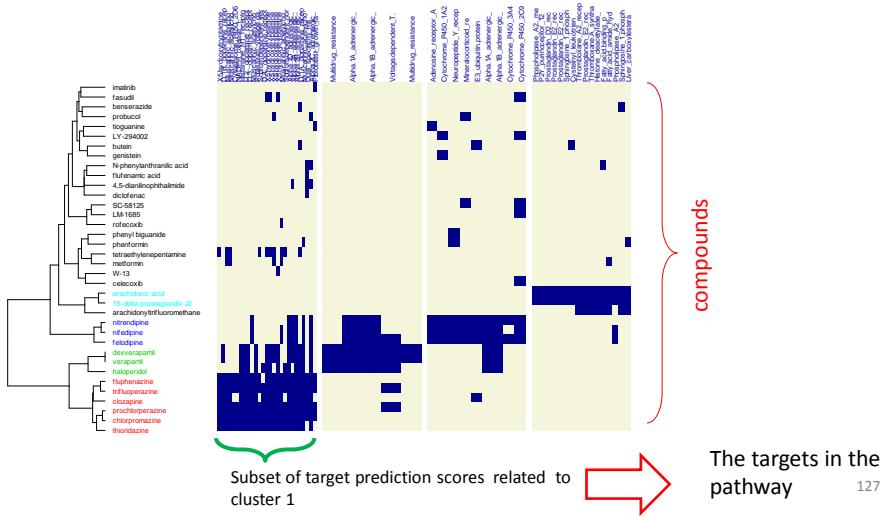
Target prediction based clustering

For each cluster, identify target scores in common for all compounds in the cluster.



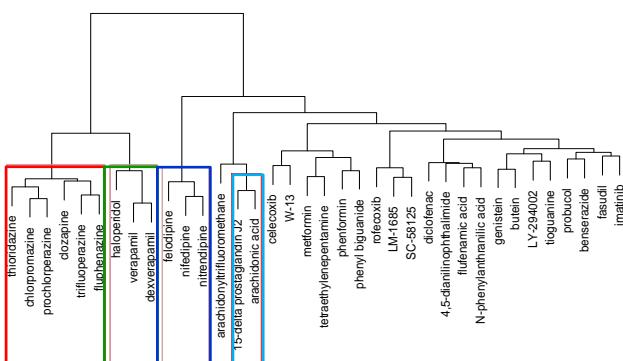
Target prediction based clustering

Identification of target prediction scores which are in common for the compounds in a cluster.



Target-based clustering

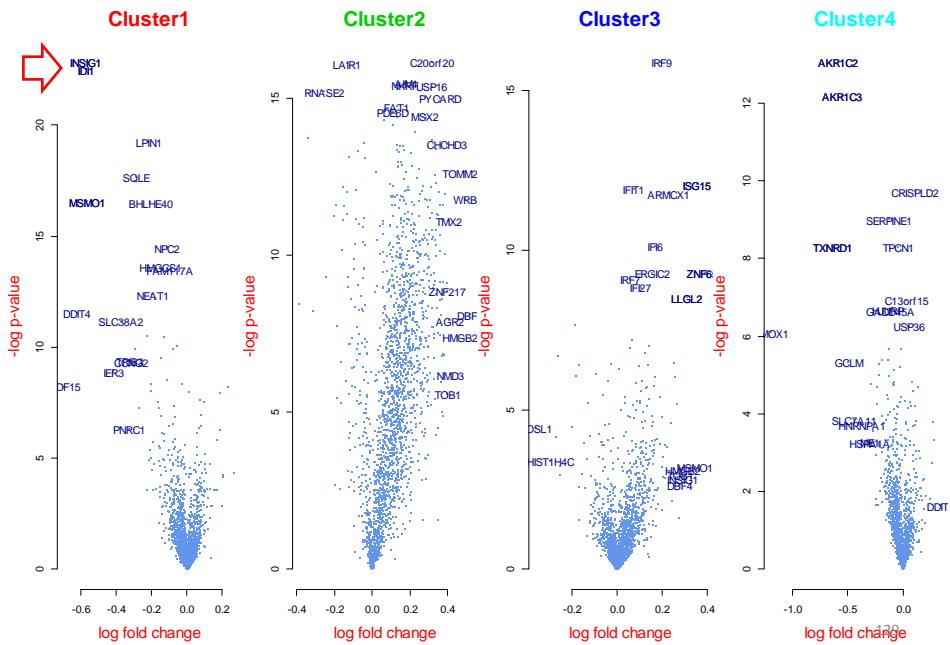
Identify genes which have different expression profile between a cluster of interest and the rest of the compounds.



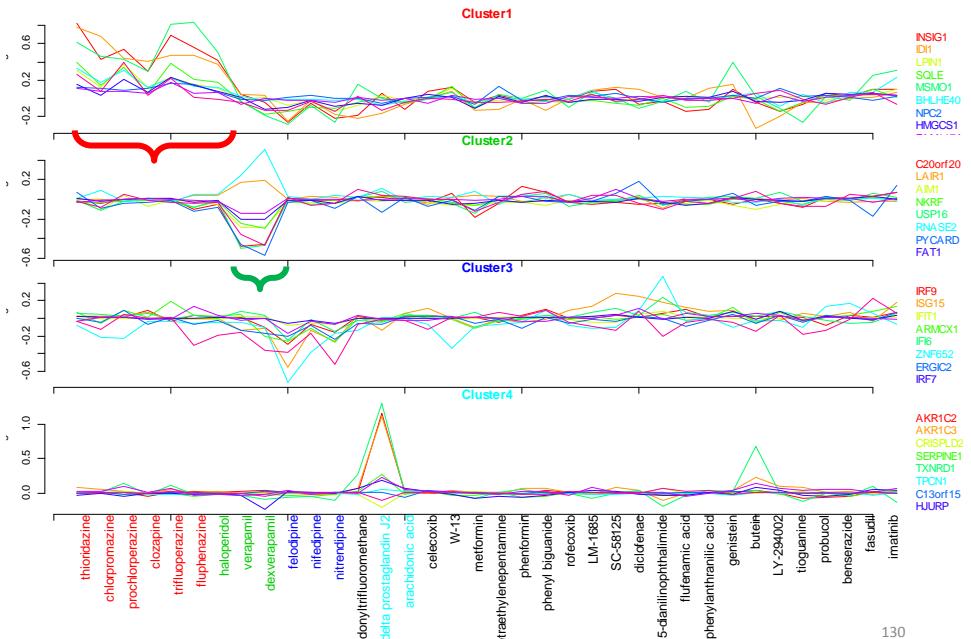
Which genes are related to this specific cluster ?

128

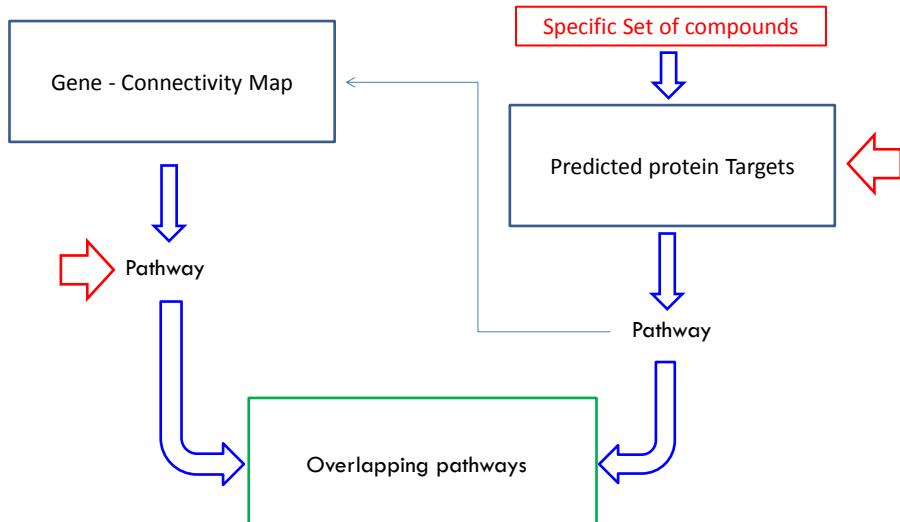
Differentially expressed genes



Profiles plots for top 8 genes by cluster



Genes and protein targets pathway overlap



131

Biological pathways: cluster 1

Use:

top K genes.

Target scores which are in common among compounds in cluster 1 (search in KEGG , GO)

Gene set analysis with MLP was done as well (to discover more genes).

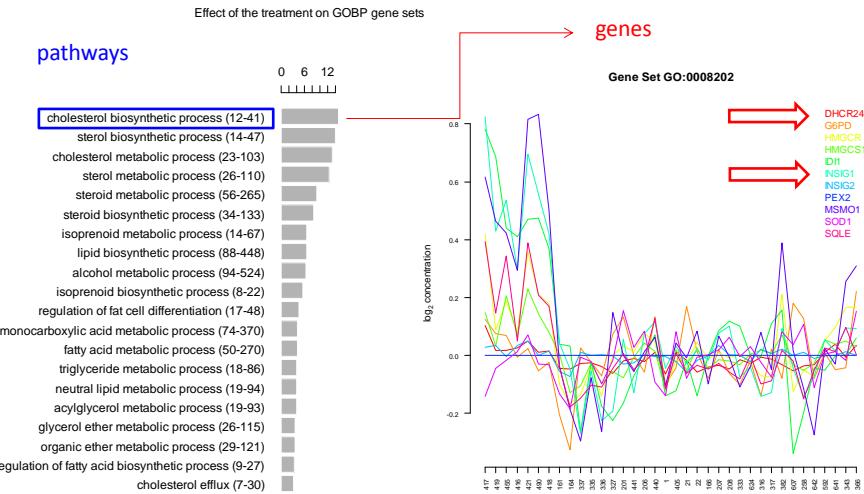
Identify :

biological pathways

Compound Clusters	Compound Names	Targets	Genes	Pathways
antipsychotic	"clozapine"	CytochromeP 4502D6	INSIG1; LDLR	GO:0008202; P:steroid metabolic process; IMP:BHF-UCL
	"thioridazine"	Dual specificity mitogen: activated protein kinase kinase1	DUSP4	hsa04010: MAPKsignalingpathway
	"chlorpromazine"	Fibroblast growth f actor receptor1		
	"trifluoperazine"	Dual specificity mitogen: activated protein kinase kinase1	LAMA3	hsa04510: Focal adhesion
	"prochlorperazine"	Dual specificity mitogen: activated protein kinase kinase1	TUBA1A	hsa04540: Gapjunction

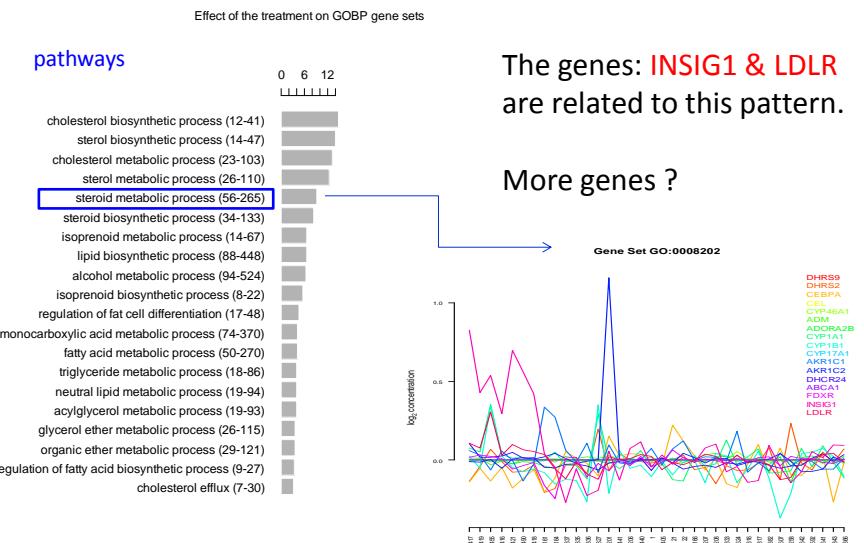
132

MLP: cluster 1



133

MLP: cluster 1

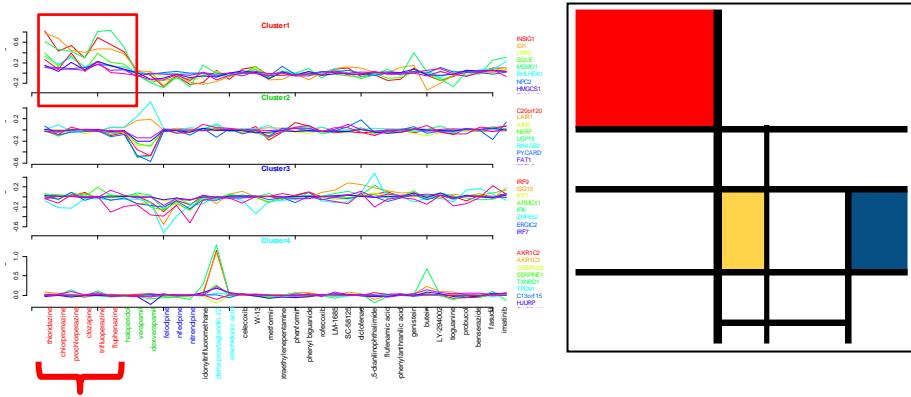


The genes: **INSIG1 & LDLR**
are related to this pattern.

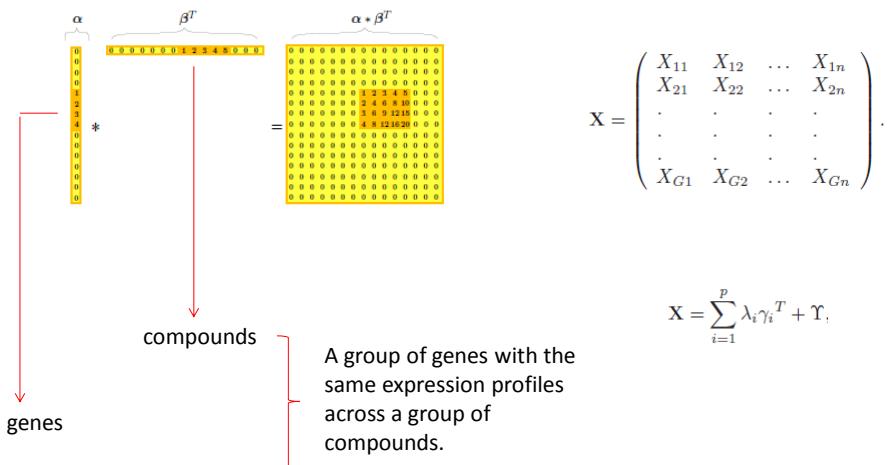
More genes ?

134

Why this is a bicluster ?



Applying FABIA: data structure + model



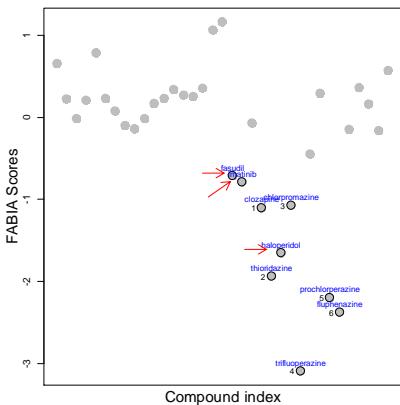
Software: biclustering using FABIA

$$\text{gMat} = \mathbf{X} = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1B} \\ X_{21} & X_{22} & \dots & X_{2B} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ X_{(B)} & X_{B2} & \dots & X_{GB} \end{pmatrix}$$

G

```
>fabRes <- fabia(gMat, alpha=0.1, p=20, cyc=1000, spl=1,
                    spz=0.5)
>rb <- extractBic(fabRes)
```

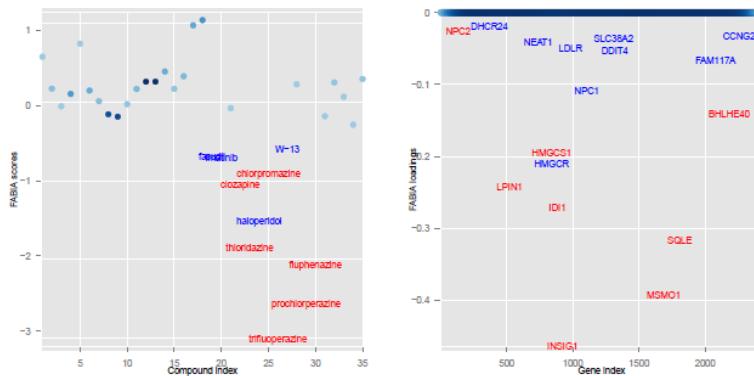
Results: biclustering using FABIA- compound scores



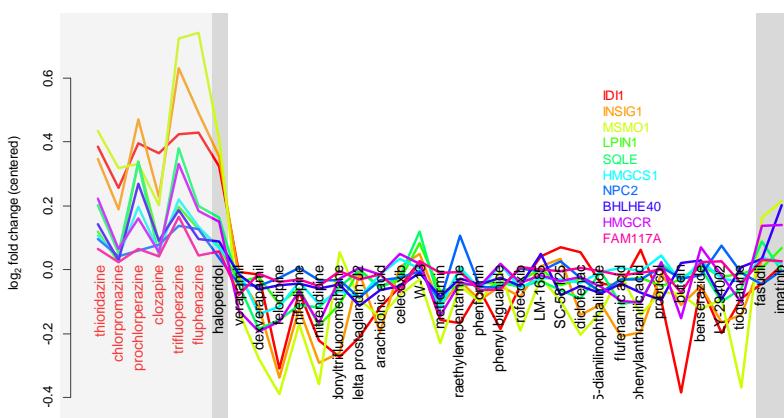
- Bicluster 1 is similar to cluster 1 compound set
- With 3 extra compounds

```
> str(bicList[[1]])
List of 2
$ compounds: chr [1:9] "trifluoperazine" "fluphenazine" ...
$ genes      : chr [1:13] "MSMO1" "INSIG1" "IDI1" "SQLE" ...
```

Factor scores (compounds) and factor loadings (genes)



A bicluster (FABIA)



Discussion

- An exploratory tool for discovering subgroups with aligned multiple properties
- Could be applicable in other research fields
- One of the integrative clustering approaches included in the package *IntClust*.

141

Part 4.3

Sport:

Using biclustering method to detection
of local patterns in NBA data



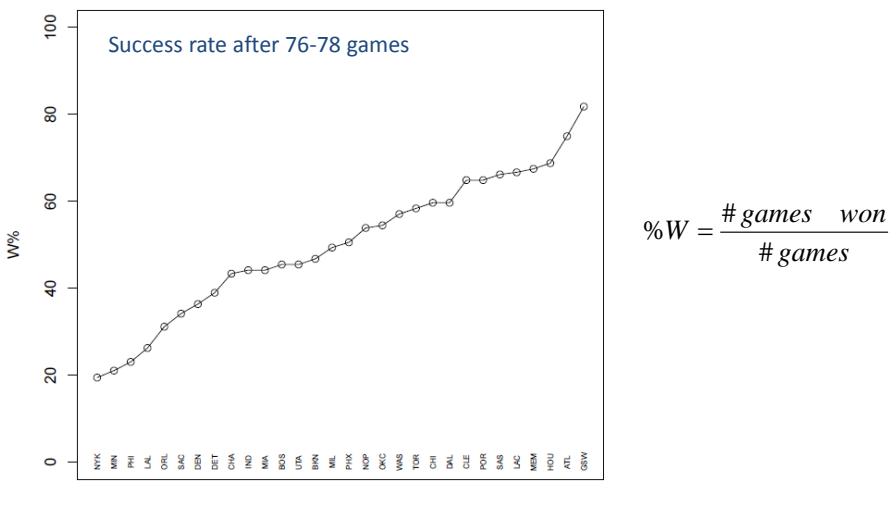
NBA



- 30 teams
- Regular season : 82 games per team
- 16 teams go to the Play-offs
- Performance Statistics (teams and individuals) is well developed

Our aim: to develop a multivariate performance indicator

Success rate of NBA teams in the regular season of 2014/2015



Commonly used performance indicators in NBA

- 2-pt / 3-pt Successful
- 2-pt / 3-pt Unsuccessful
- Free Throw Successful / Unsuccessful
- Defensive / Offensive Rebounds
- Assists
- Turnovers
- Steals
- Dunks
- Blocks Committed / Received
- Fouls Committed / Received

Garcia et al (2013) showed that these variables are good performance indicators for Regular Season as well as for Playoff Games

Data Structure

- 7 online databases in the NBA website:
 - Traditional Stats
 - Advanced Stats
 - Four Factors
 - Misc. Stats
 - Scoring
 - Opponent
 - Shooting

Updated after each game

Data structure

$$X = [X_1, X_2, X_3, X_4, X_5, X_6, X_7]$$

1. Advanced Stat
2. Four Factors
3. Misc. Stats
4. Scoring
5. Opponent
6. Shooting

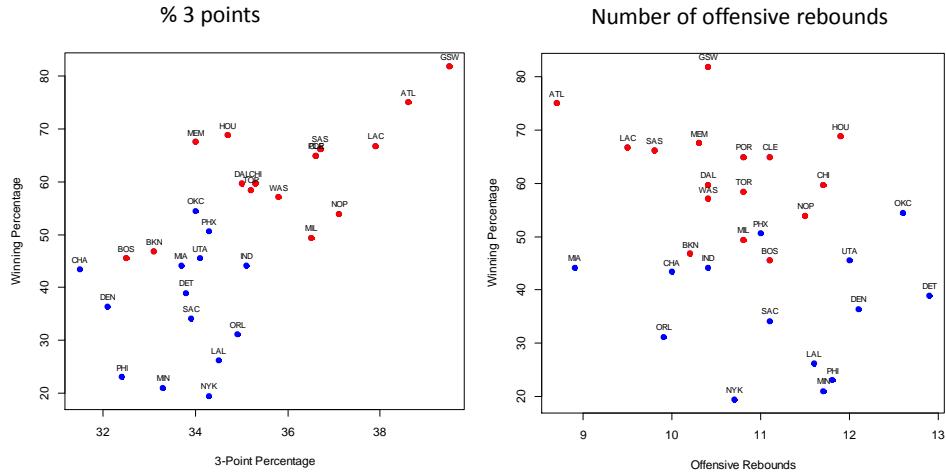
- Each matrix is a $30 \times n_i$
 Example :
- 2-pt / 3-pt Successful
 - 2-pt / 3-pt Unsuccessful
 - Free Throw Successful / Unsuccessful
 - Defensive / Offensive Rebounds
 - Assists
 - Turnovers
 - Steals
 - Dunks
 - Blocks Committed / Received
 - Fouls Committed / Received

Analysis plan

- Step 1: PCA for the Traditional Stats:
 - 2-pt / 3-pt Successful
 - 2-pt / 3-pt Unsuccessful
 - Free Throw Successful / Unsuccessful
 - Defensive / Offensive Rebounds
 - Assists
 - Turnovers
 - Steals
 - Dunks
 - Blocks Committed / Received
 - Fouls Committed / Received
- Step 2: Multiple factor analysis

Can we find patterns among these indicators ??

Example: traditional performance indicator in NBA



What do we want to do ?

1. Develop a performance score that will tell us who are the “best teams” (in terms of performance, i.e % Win).
2. Multivariate performance score.

3. Analysis in two steps:
 1. First step: PCA for Traditional Stats
 2. Second step: multiple factor analysis for all data (data integration).

MFA: data structure

$$X = [X_1 | X_2 X_3 X_4 X_5 X_6 X_7]$$

Set of the traditional stats indicators

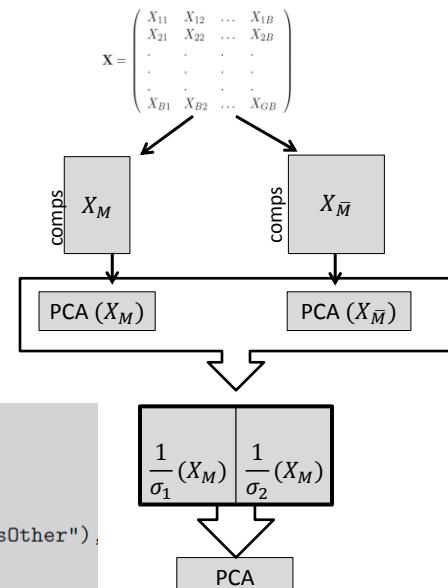
- 3-pt Percentage
- Free Throws Percentage
- Defensive Rebounds
- Offensive Rebounds
- Assists
- Turnovers
- Steals
- Field Goals Percentage
- Blocks Committed / Received
- Fouls Committed / Received

- Advanced Stats
- Four Factors
- Misc. Stats
- Scoring
- Opponent
- Shooting

Multiple factor analysis

- All variables standardized
- Normalize each data matrix
- Concatenate all normalized datasets and perform PCA on the combined weighted data
 - Factor scores describe compounds
 - Factor loadings describe variables

```
>resMFA <- MFA(dataMFA,
group = c(ncol(Mat1), ncol(Mat2)),
type = c("c", "c"),
ncp = 2,
name.group = c("genesInitial", "genesOther"),
graph=FALSE
)
```



Step 1:

PCA for the leading performance indicator

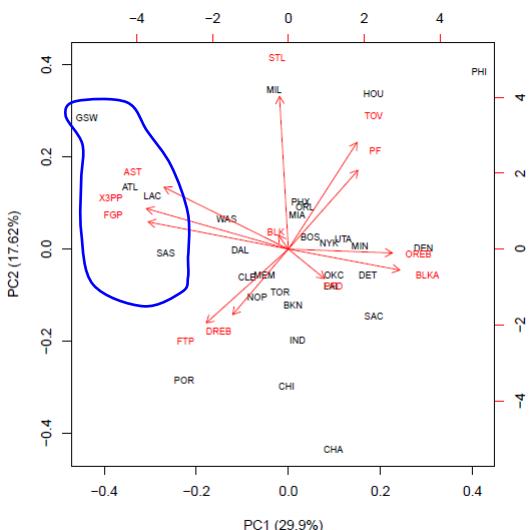
Leading Performance Indicators:

- 3-pt Percentage
- Free Throws Percentage
- Defensive Rebounds
- Offensive Rebounds
- Assists
- Turnovers
- Steals
- Field Goals Percentage
- Blocks Committed / Received
- Fouls Committed / Received

PCA:

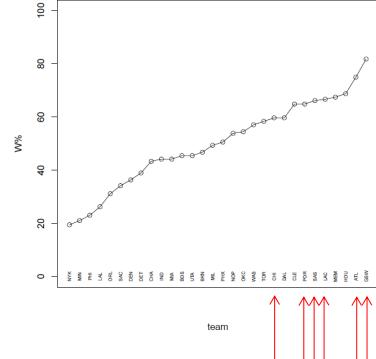
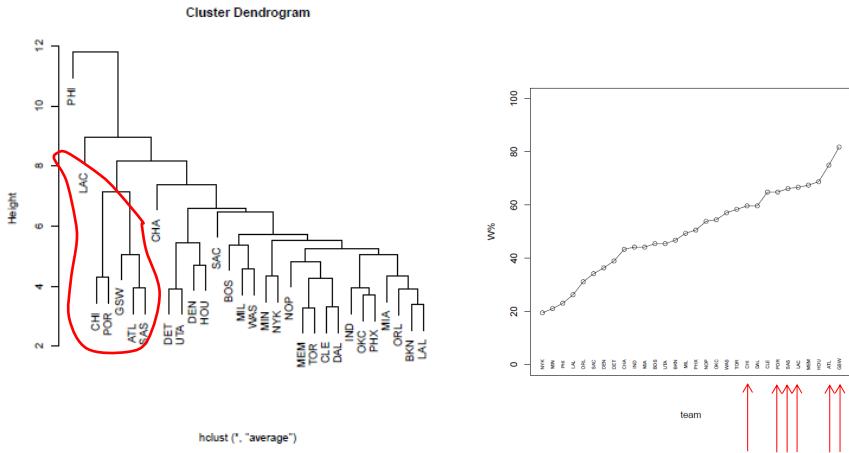
To convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables

PCA for the leading performance indicator

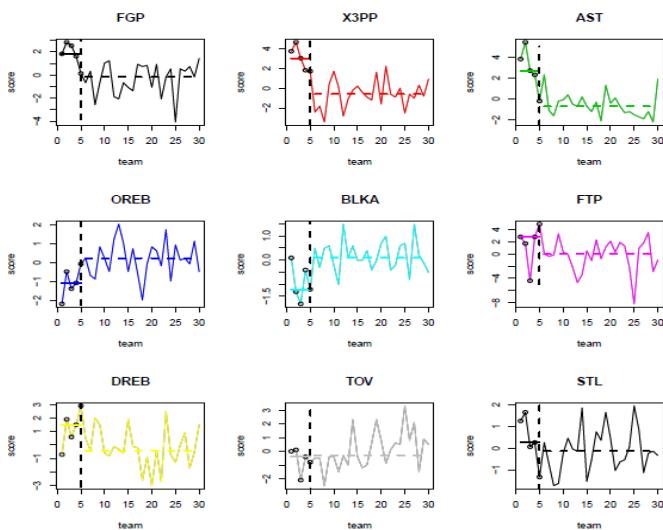


- Clear cluster of teams based on PC1
- Similar pattern was detected by Hierarchical clustering
- Indicators:
 - AST
 - X3PP
 - FGP

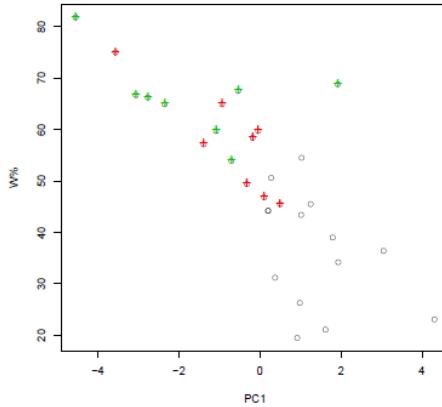
Hierarchical clustering



PCA for the leading performance indicator



Performance score



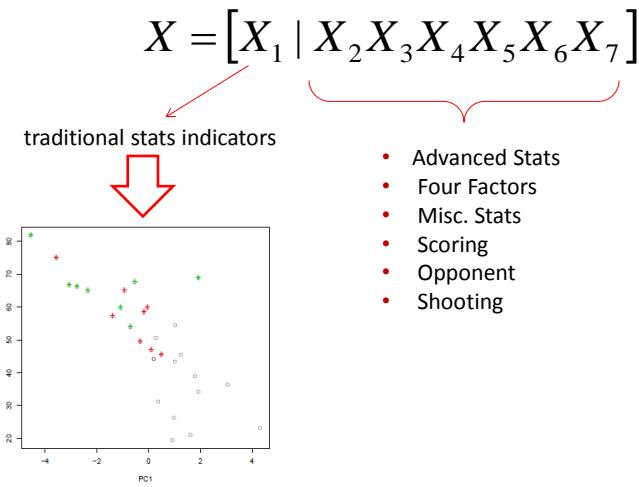
First PC versus % win.

Correlation.

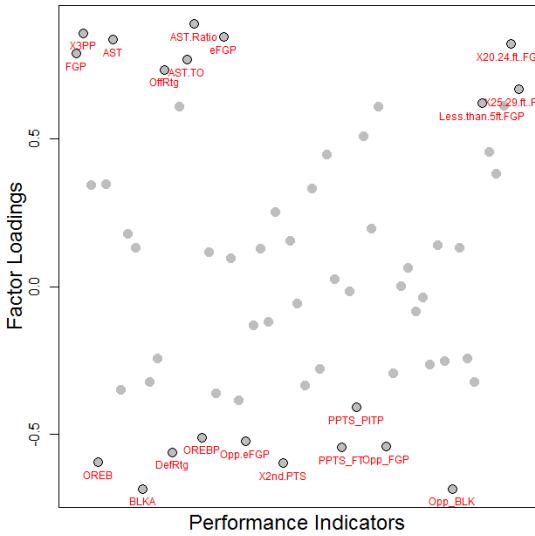
Performance score:

$$PC_1 = \sum \ell_j X_j$$

Data integration: MFA

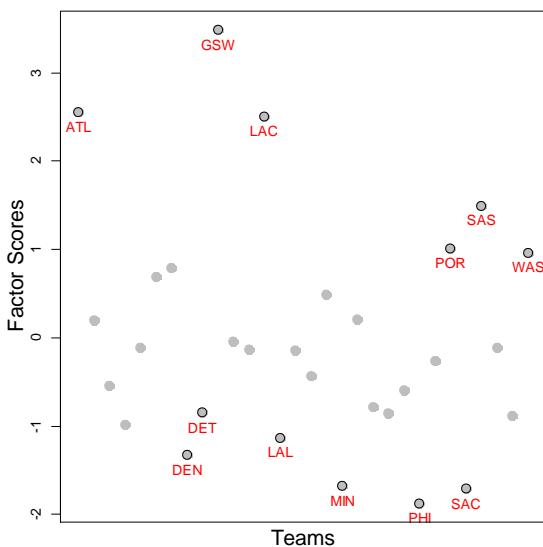


MFA: factor loadings (variables)



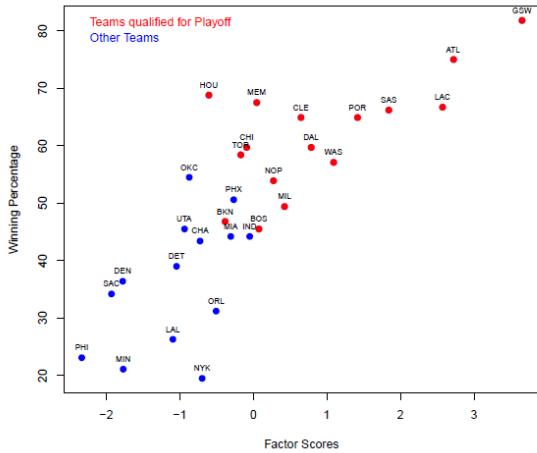
- Traditional Stats with high loadings
- Some new indicators were discovered

MFA: factor scores (teams)



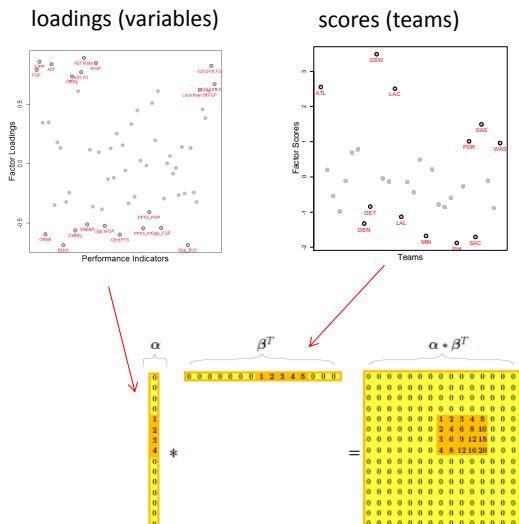
- Similar set of teams with high factor scores

Multivariate performance score



$$MPS(team_k) = \sum \ell_j X_j$$

Biclustering using FABIA



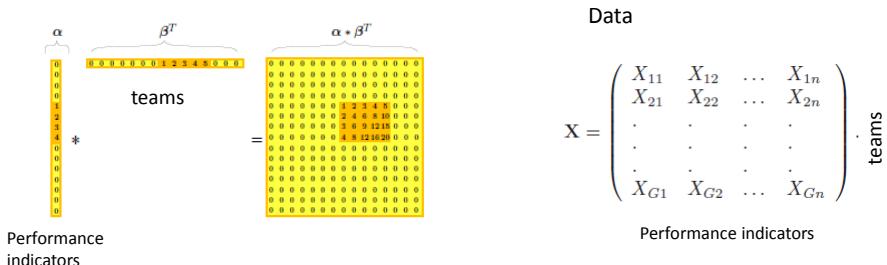
MFA:
After normalization: FA
with one factor

$$X = [X_1 | X_2 X_3 X_4 X_5 X_6 X_7]$$

$$\mathbf{X} = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1n} \\ X_{21} & X_{22} & \dots & X_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ X_{G1} & X_{G2} & \dots & X_{Gn} \end{pmatrix}.$$

FABIA with one factor (BC)

Applying FABIA: data structure + model



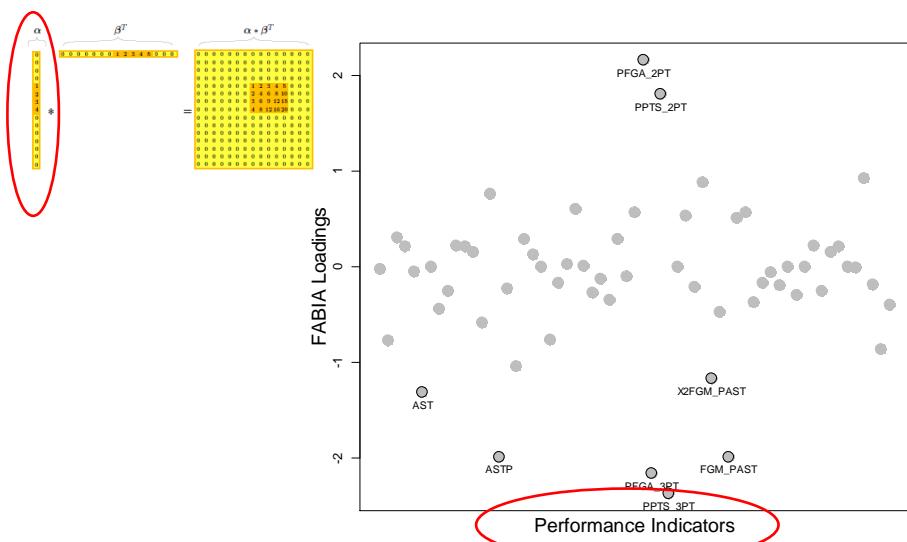
Aim:

1. Find a group of teams that share patterns in performance indicators.
2. Correlation to overall performance.

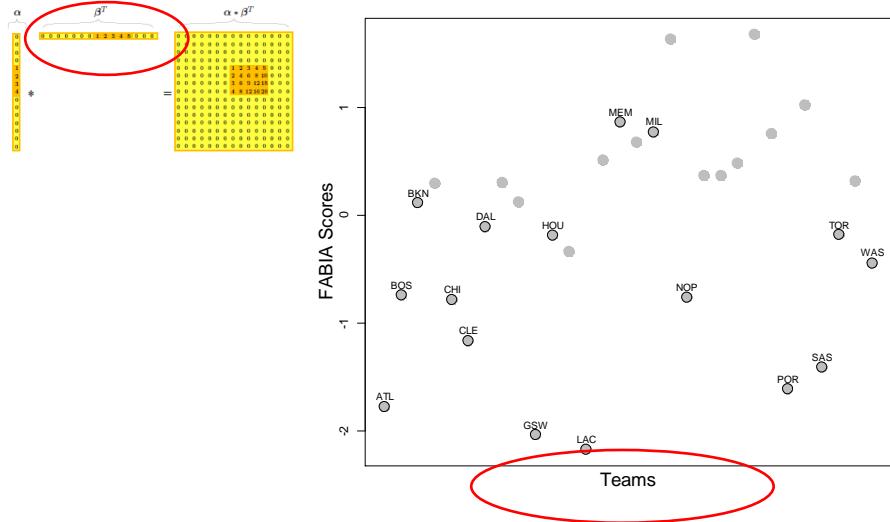
The model

$$X = \sum_{i=1}^p \lambda_i \gamma_i^T + \Upsilon,$$

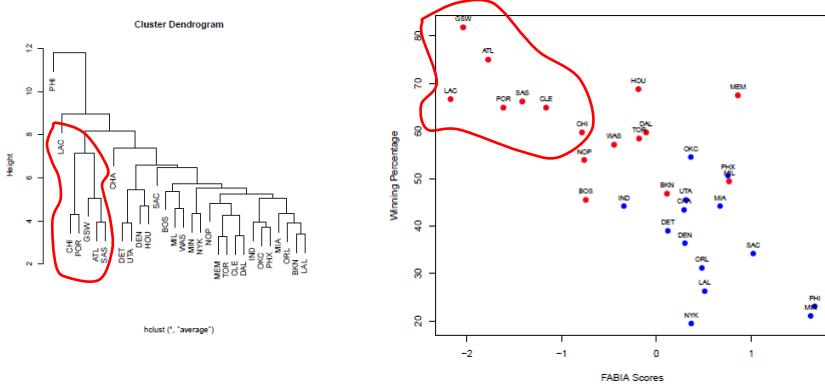
BC1 FABIA: factor loadings



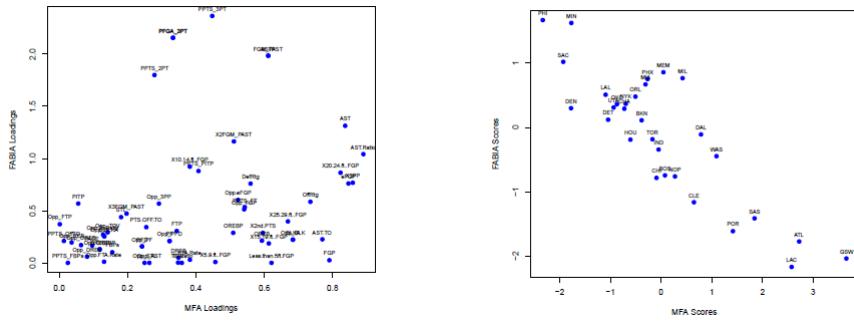
BC1 FABIA: factor scores



Overall performance score



FABIA and MFA



Software

Part 4.4

Enrichment of Gene Expression Modules using Multiple Factor Analysis and Biclustering

Motivation

- A bicluster contains
genes that are
coordinately regulated
under a subset of conditions

Gene module

- Summarized expression profiles of these genes
- Many biclusters -> many possible gene modules

Motivation

- A bicluster contains
 - genes that are **not only**
 - coordinately regulated
 - under a subset of conditions
 - but are also mostly functionally coherent.**

Gene module

- Summarized expression profiles of these genes that act in concert to carry out a specific function

Motivation

- A bicluster contains
 - genes that are **not only**
 - coordinately regulated
 - under a subset of conditions
 - but are also mostly functionally coherent.**
- Availability of a subset of “lead” genes/compounds
 - Genes related to a phenotype of interest
 - Genes that are known to be part of a biological pathway
 - Some hypothesis generated from previous experiments

Aim

- To enrich this set of M “lead” genes

$$\mathbf{X}_M = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1n} \\ X_{21} & X_{22} & \dots & X_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ X_{M1} & X_{M2} & \dots & X_{Mn} \end{pmatrix}$$

Idea

- Run biclustering algorithm
 - Search for the bicluster that contains most of the genes in the list of “lead” genes
 - Not necessarily the first (ranking) bicluster
 - Dependent on the sparsity parameter, etc.
- MFA
 - find links between datasets (presence of common

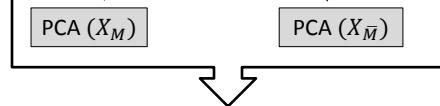
```
> install.packages("FactoMineR")
> library(FactoMineR)
```

MFA

- All variables standardized
- Normalize each data matrix
- Concatenate all normalized datasets and perform PCA on the combined weighted data
 - Factor scores describe compounds
 - Factor loadings describe variables

```
> resMFA <- MFA(dataMFA,
group = c(ncol(Mat1), ncol(Mat2)),
type = c("c", "c"),
ncp = 2,
name.group = c("genesInitial", "genesOther"),
graph=FALSE
)
```

$$\mathbf{X} = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1B} \\ X_{21} & X_{22} & \dots & X_{2B} \\ \vdots & \vdots & \ddots & \vdots \\ X_{B1} & X_{B2} & \dots & X_{GB} \end{pmatrix}$$



$$\frac{1}{\sigma_1}(X_M) \quad \frac{1}{\sigma_2}(X_{\bar{M}})$$

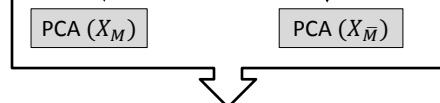
PCA

MFA

- All variables standardized
- Normalize each data matrix
- Concatenate all normalized datasets and perform PCA on the combined weighted data
 - Compound scores
 - Gene loadings

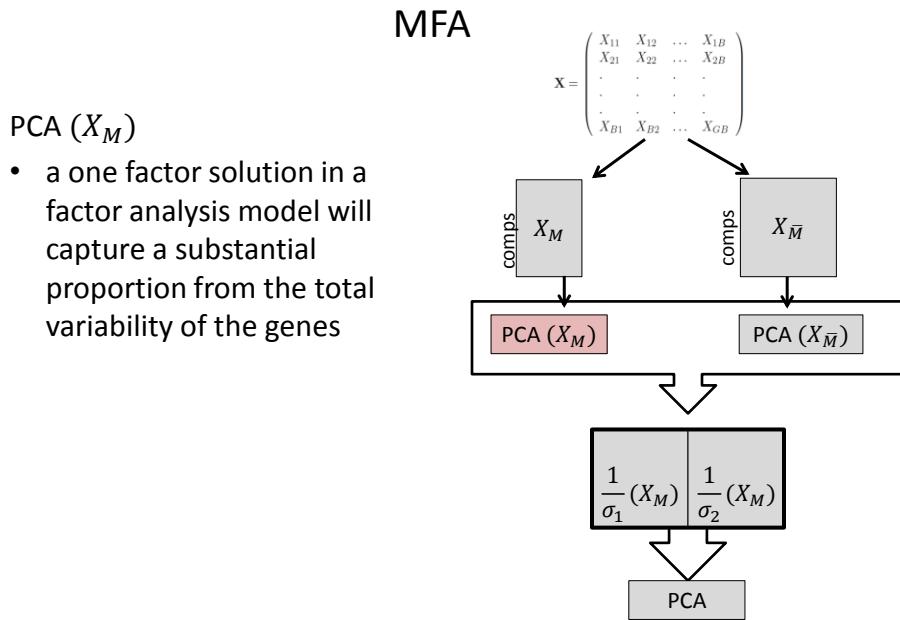
```
> loadings1 <- resMFA$quanti.var$coord[,1]
> scores1 <- resMFA$ind$coord[,1]
```

$$\mathbf{X} = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1B} \\ X_{21} & X_{22} & \dots & X_{2B} \\ \vdots & \vdots & \ddots & \vdots \\ X_{B1} & X_{B2} & \dots & X_{GB} \end{pmatrix}$$



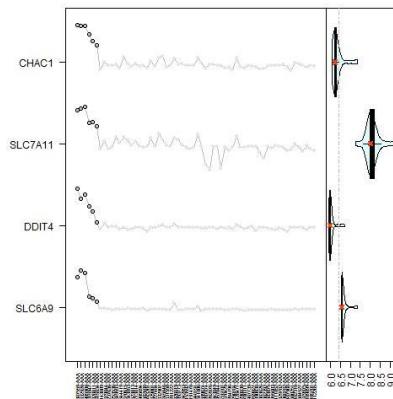
$$\frac{1}{\sigma_1}(X_M) \quad \frac{1}{\sigma_2}(X_{\bar{M}})$$

PCA



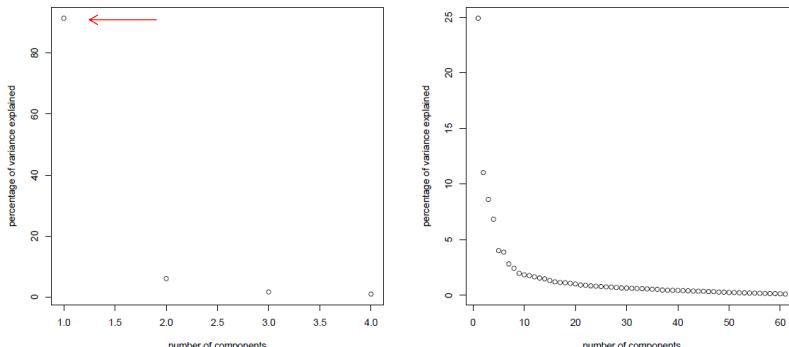
Motivating Data: mGlu2 project

- n= 62 compounds
- G = 566 genes
- M=4 genes that are known to be biologically related and are linked to the phenotype of interest.



Scree plots

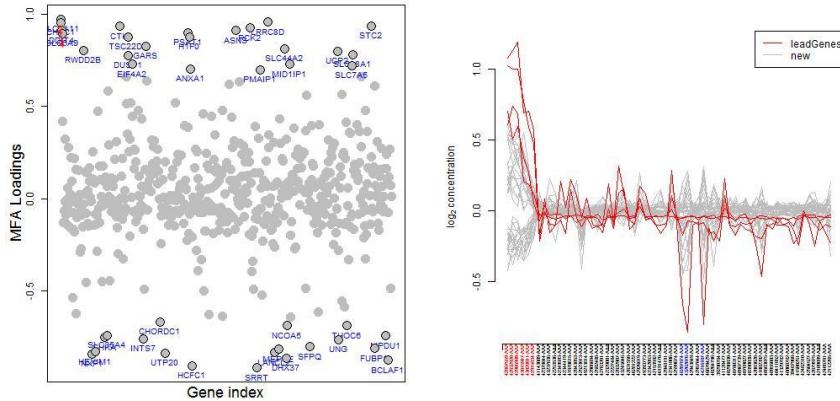
- One known structure in X_M
- Escoufier's Rv coefficient = 26%

(a) scree plot for \mathbf{X}_M (b) scree plot for $\mathbf{X}_{\bar{M}}$

Data Contribution to the main factors

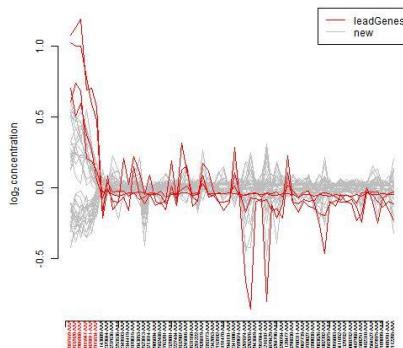
Data	Factor 1	Factor 2
\mathbf{X}_M	76.48	0.68
$\mathbf{X}_{\bar{M}}$	23.52	99.32

MFA 1 Gene loadings

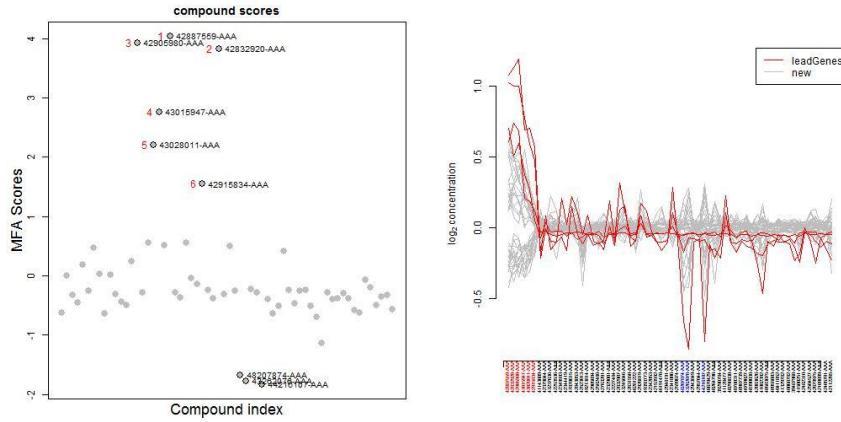


Fabia Bicluster 2

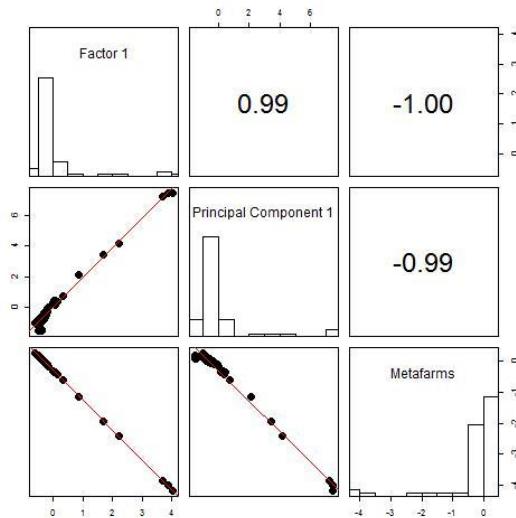
- Absence of lead genes
- Fabia searches only for correlated profiles across a subset of samples
- MFA uses the similarity of gene profiles across all compounds.
- As a result, some genes discovered by MFA are not part of the fabia bicluster and vice versa



MFA 1 Compound Scores



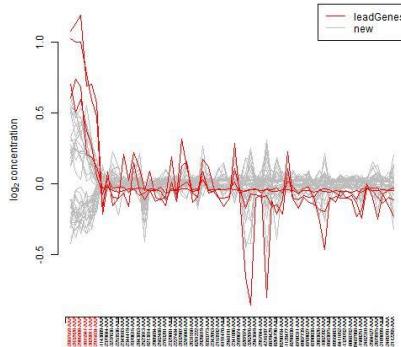
Gene Module Summarization – one



Fabia Bicluster 2

FABIA BC 2 genes

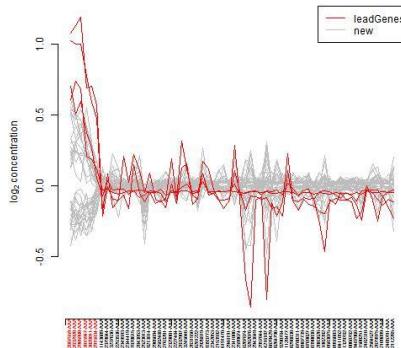
- Absence of lead genes
- Run biclustering
- explore interesting biclusters
- Lead genes are in Fabia bicluster 2



Fabia Bicluster 2

FABIA BC 2 genes

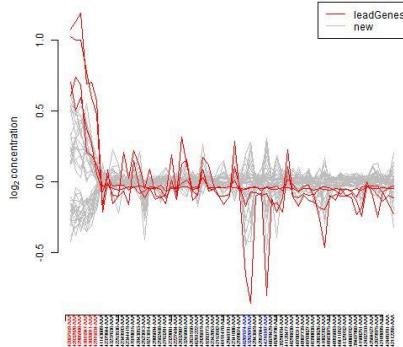
- Absence of lead genes
- Fabia → correlated profiles across a subset of samples
- MFA → similarity of gene profiles across all compounds
 - some genes discovered by MFA are not part of the fabia bicluster and vice versa



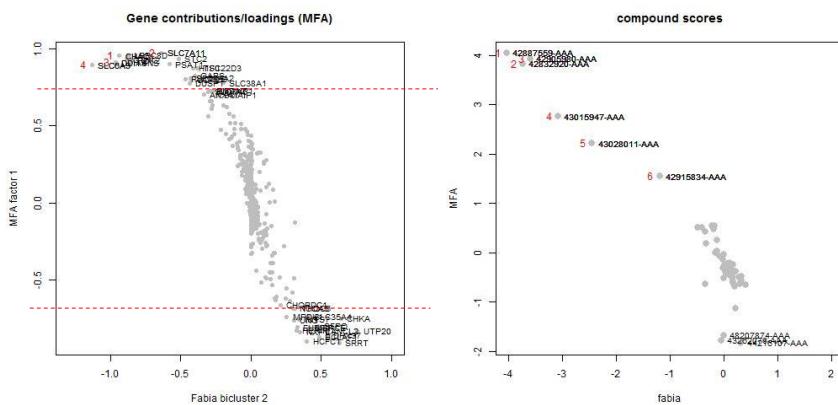
MFA 1

MFA 1 genes

- Absence of lead genes
- Fabia → correlated profiles across a subset of samples
- MFA → similarity of gene profiles across all compounds
 - some genes discovered by MFA are not part of the fabia bicluster and vice versa



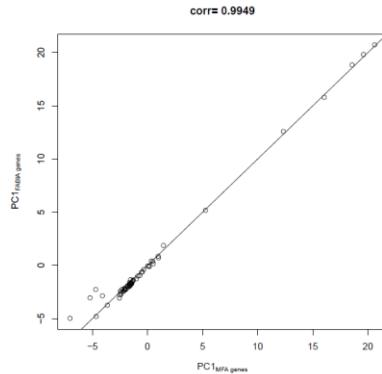
Gene loadings and compound scores



Gene Module

MFA 1 genes

- Absence of lead genes
- Fabia → correlated profiles across a subset of samples
- MFA → similarity of gene profiles across all compounds
 - some genes discovered by MFA are not part of the fabia bicluster and vice versa
- The underlying latent structure is almost identical



Part 4.5

Drug Discovery (II)
Ranking of BCs

Motivation

- how to determine which biclusters are most informative and rank them on the basis of their importance?
 - Data-driven, statistical measure (information content (FABIA))
 - biological context - gene ontology annotations or other literature-based enrichment analysis

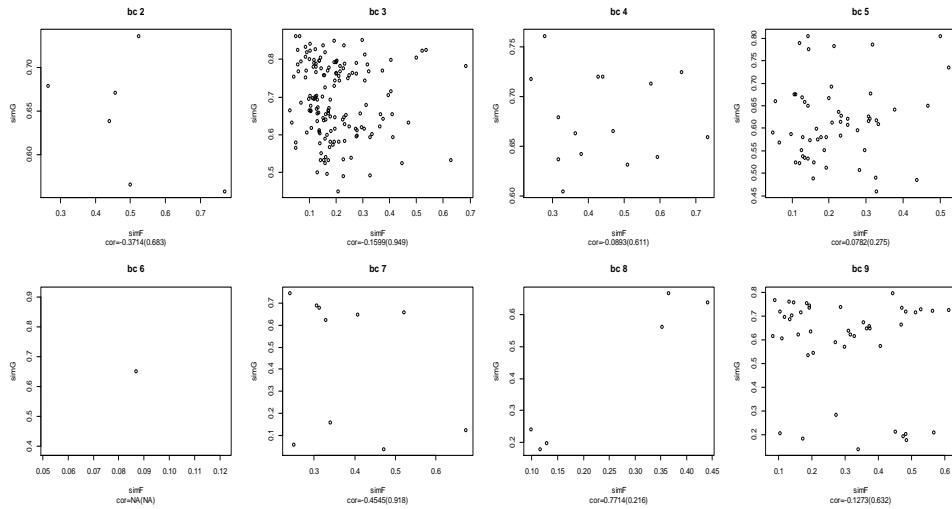
Idea for early drug discovery

- rank based on another source of information, (e.g. the chemical structure, target predictions, HCS, etc)
- investigate whether compounds in a bicluster are also structurally similar
- Similar activity and similar structure → desirable compound set!

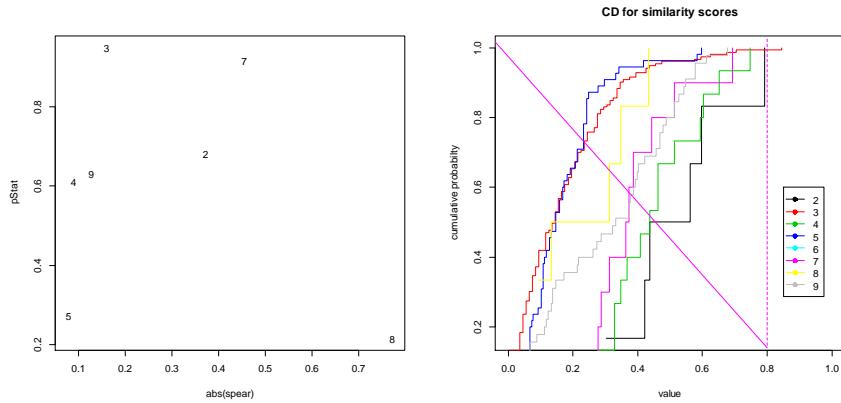
Biclustering Results

	nCompounds	nGenes
1	1	63
2	4	50
3	18	41
4	6	53
5	11	28
6	2	12
7	5	26
8	4	10
9	10	2
10	21	1

Spearman correlation of similarity scores



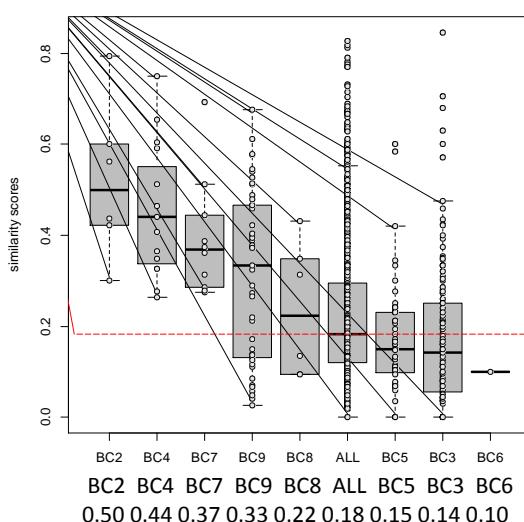
Ranking statistics



BC Ranking based on median similarity scores

(C)

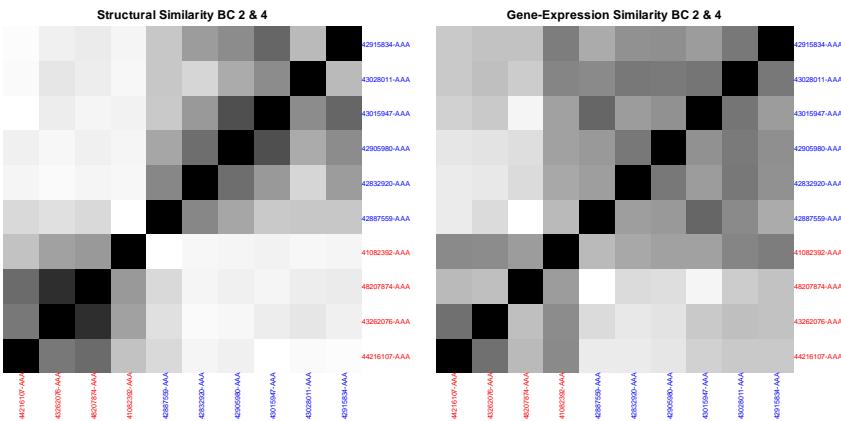
Boxplot of Compound Similarity Scores



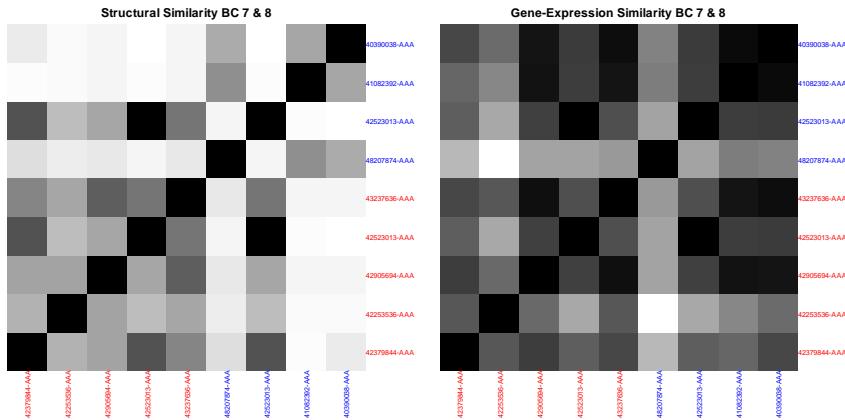
other statistics

BC	mean	median	sd	range	CV
2	0.52	0.50	0.17	0.49	0.33
3	0.17	0.14	0.16	0.85	0.91
4	0.45	0.44	0.14	0.49	0.32
5	0.17	0.15	0.12	0.60	0.74
6	0.10	0.10		0.00	
7	0.39	0.37	0.13	0.42	0.33
8	0.24	0.22	0.15	0.34	0.62
9	0.31	0.33	0.19	0.65	0.61

Similarity Scores of BC2 and BC4



Similarity Scores of BC7 and BC8



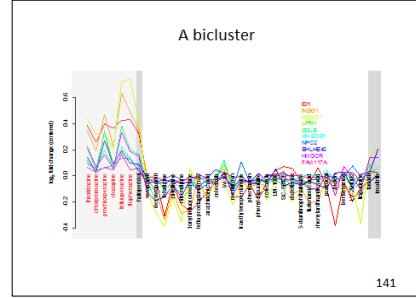
Discussion

- Not ranking per se but to prioritize more interesting biclusters using extra information available
 - Software: *bcRank*

Summary



Piet Mondrian, Tate modern.



Perualila et al, 2016

Biclustering: local patterns to understand the big picture.
Many areas of applications.
Many methods.
Software.