

# Sparse Multiple Correspondence Analysis for selection of Single Nucleotide Polymorphisms

**Anne BERNARD**

QFAB Bioinformatics, University of Queensland, Brisbane  
December 2<sup>nd</sup>, 2015





## High-dimensional data

Analysis of the data structure and observation of a possible natural separation between individuals depending on their human genetic heritage.

One data set  $\mathbf{X}$  ( $I \times J$ )  $\Rightarrow$  Unsupervised multivariate analysis



Continuous data: PCA

Categorical data: MCA

In case of high dimensional data ( $I \gg J$ ):  
results difficult to interpret.

**Solution:** Use/Develop appropriate statistical methods  
to **select relevant variables** and **facilitate interpretation** of the  
results.





















## SNP Coding

ID	SNP.1	SNP.2
AK	GG	AA
JD	GT	CC
MR	GT	AC
GB	GG	CC
NH	TT	AC
AL	GG	CC
DO	TT	AC
JM	TT	CC
ED	TT	AC
CB	TT	CC
CF	GG	AC
OD	TT	CC
DM	TT	AC
NS	GG	CC
JR	GT	AA

Original Coding

GG	GT	TT	CC	AC	AA	ID
1	0	0	0	0	1	AK
0	1	0	1	0	0	JD
0	1	0	0	1	0	MR
1	0	0	1	0	0	GB
0	0	1	0	1	0	NH
1	0	0	1	0	0	AL
0	0	1	0	1	0	DO
0	0	1	1	0	0	JM
0	0	1	0	1	0	ED
0	0	1	1	0	0	CB
1	0	0	0	1	0	CF
0	0	1	1	0	0	OD
0	0	1	0	1	0	DM
1	0	0	1	0	0	NS
0	1	0	0	0	1	JR

Nominal Coding

## SNP Coding

ID	SNP.1	SNP.2
AK	GG	AA
JD	GT	CC
MR	GT	AC
GB	GG	CC
NH	TT	AC
AL	GG	CC
DO	TT	AC
JM	TT	CC
ED	TT	AC
CB	TT	CC
CF	GG	AC
OD	TT	CC
DM	TT	AC
NS	GG	CC
JR	GT	AA

Original Coding

GG	GT	TT	CC	AC	AA	ID
1	0	0	0	0	1	AK
0	1	0	1	0	0	JD
0	1	0	0	1	0	MR
1	0	0	1	0	0	GB
0	0	1	0	1	0	NH
1	0	0	1	0	0	AL
0	0	1	0	1	0	DO
0	0	1	1	0	0	JM
0	0	1	0	1	0	ED
0	0	1	1	0	0	CB
1	0	0	0	1	0	CF
0	0	1	1	0	0	OD
0	0	1	0	1	0	DM
1	0	0	1	0	0	NS
0	1	0	0	0	1	JR

SNP1      SNP2

## SNP Coding

ID	SNP.1	SNP.2
AK	GG	AA
JD	GT	CC
MR	GT	AC
GB	GG	CC
NH	TT	AC
AL	GG	CC
DO	TT	AC
JM	TT	CC
ED	TT	AC
CB	TT	CC
CF	GG	AC
OD	TT	CC
DM	TT	AC
NS	GG	CC
JR	GT	AA

Original Coding

GG	GT	TT	CC	AC	AA	ID
1	0	0	0	0	1	AK
0	1	0	1	0	0	JD
0	1	0	0	1	0	MR
1	0	0	1	0	0	GB
0	0	1	0	1	0	NH
1	0	0	1	0	0	AL
0	0	1	0	1	0	DO
0	0	1	1	0	0	JM
0	0	1	0	1	0	ED
0	0	1	1	0	0	CB
1	0	0	0	1	0	CF
0	0	1	1	0	0	OD
0	0	1	0	1	0	DM
1	0	0	1	0	0	NS
0	1	0	0	0	1	JR

Nominal Coding

## SNP Coding

ID	SNP.1	SNP.2
AK	GG	AA
JD	GT	CC
MR	GT	AC
GB	GG	CC
NH	TT	AC
AL	GG	CC
DO	TT	AC
JM	TT	CC
ED	TT	AC
CB	TT	CC
CF	GG	AC
OD	TT	CC
DM	TT	AC
NS	GG	CC
JR	GT	AA

Original Coding

GG	GT	TT	CC	AC	AA	ID
1	0	0	0	0	1	AK
0	1	0	1	0	0	JD
0	1	0	0	1	0	MR
1	0	0	1	0	0	GB
0	0	1	0	1	0	NH
1	0	0	1	0	0	AL
0	0	1	0	1	0	DO
0	0	1	1	0	0	JM
0	0	1	0	1	0	ED
0	0	1	1	0	0	CB
1	0	0	0	1	0	CF
0	0	1	1	0	0	OD
0	0	1	0	1	0	DM
1	0	0	1	0	0	NS
0	1	0	0	0	1	JR

Nominal Coding



## SNP Coding

ID	SNP.1	SNP.2
AK	GG	AA
JD	GT	CC
MR	GT	AC
GB	GG	CC
NH	TT	AC
AL	GG	CC
DO	TT	AC
JM	TT	CC
ED	TT	AC
CB	TT	CC
CF	GG	AC
OD	TT	CC
DM	TT	AC
NS	GG	CC
JR	GT	AA

Original Coding

GG	GT	TT	CC	AC	AA	ID
1	0	0	0	0	1	AK
0	1	0	1	0	0	JD
0	1	0	0	1	0	MR
1	0	0	1	0	0	GB
0	0	1	0	1	0	NH
1	0	0	1	0	0	AL
0	0	1	0	1	0	DO
0	0	1	1	0	0	JM
0	0	1	0	1	0	ED
0	0	1	1	0	0	CB
1	0	0	0	1	0	CF
0	0	1	0	0	0	OD
0	0	1	0	1	0	DM
1	0	0	1	0	0	NS
0	1	0	0	0	1	JR

Nominal Coding

## SNP Coding

ID	SNP.1	SNP.2
AK	GG	AA
JD	GT	CC
MR	GT	AC
GB	GG	CC
NH	TT	AC
AL	GG	CC
DO	TT	AC
JM	TT	CC
ED	TT	AC
CB	TT	CC
CF	GG	AA
OD	TT	CC
DM	TT	AC
NS	GG	CC
JR	GT	AA

Original Coding

GG	GT	TT	CC	AC	AA	ID
1	0	0	0	0	1	AK
0	1	0	1	0	0	JD
0	1	0	0	1	0	MR
1	0	0	1	0	0	GB
0	0	1	0	1	0	NH
1	0	0	1	0	0	AL
0	0	1	0	1	0	DO
0	0	1	1	0	0	JM
0	0	1	0	1	0	ED
0	0	1	1	0	0	CB
1	0	0	0	1	0	CF
0	0	0	1	0	0	OD
0	0	1	0	1	0	DM
1	0	0	1	0	0	NS
0	1	0	0	0	1	JR

Nominal Coding

# Case study: Material

## Variables to be explained

Phenotypes

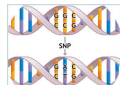


- Sagging score [0-10]

and others (global photoageing, wrinkling score, lentigines score)

## Explained variables

Genetic data



795 063 SNPs analysed

Targeted set of 537 SNPs

- "Candidate GWAS"
- 1611 disjunctive columns
- 537 blocks

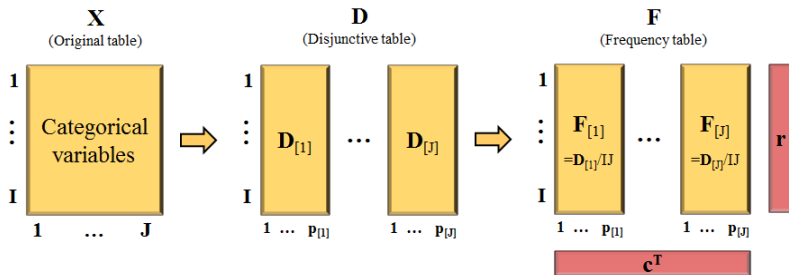
## Application on SNPs data: Exploratory analysis

**Step 1** Visualization of links between SNPs and between samples using **Multiple Correspondence Analysis (MCA)**

**Step 2** Selection of the most important SNPs for a component using the sparse extension of MCA to select variables:  
**Sparse MCA**

## Multivariate Exploratory Methods

When  $\mathbf{X}$  matrix of categorical variables



$$\mathbf{r} = \mathbf{F}\mathbf{1}$$

$$\mathbf{c} = \mathbf{F}^T \mathbf{1}$$

$p_{[j]} = \text{Nb of modalities of variable } j$

## Multiple Correspondence Analysis

### MCA via Generalized SVD of $\mathbf{F}$

$$\mathbf{F} = \mathbf{P}\mathbf{\Delta}\mathbf{Q}^T \quad \text{with } \mathbf{P}^T\mathbf{M}\mathbf{P} = \mathbf{Q}^T\mathbf{W}\mathbf{Q} = \mathbf{I}$$

where  $\mathbf{F} = [\mathbf{F}_{[1]} | \dots | \mathbf{F}_{[j]} | \dots | \mathbf{F}_{[J]}]$  and  $\mathbf{Q} = [\mathbf{Q}_{[1]} | \dots | \mathbf{Q}_{[j]} | \dots | \mathbf{Q}_{[J]}]$

In the case of PCA:  $\mathbf{M} = \mathbf{W} = \mathbf{I}$

In the case of MCA:  $\mathbf{M} = \mathbf{D}_r^{-1}$        $\mathbf{W} = \mathbf{D}_c^{-1}$

### GSVD as low rank approximation of matrices

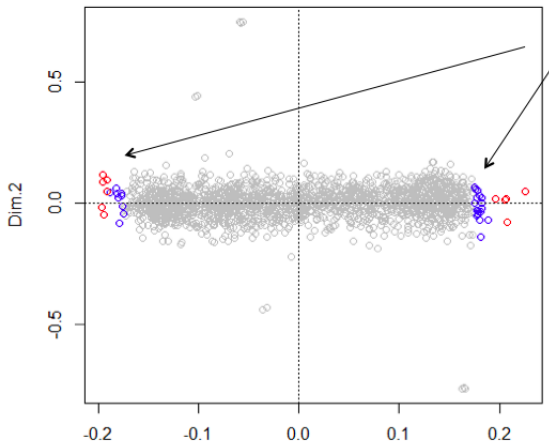
MCA can be seen as the solution of

$$\min_{\tilde{\mathbf{p}}, \tilde{\mathbf{q}}} \|\mathbf{F} - \tilde{\mathbf{p}}\tilde{\mathbf{q}}^T\|_{\mathbf{W}}^2 \quad \tilde{\mathbf{p}}^T\mathbf{M}\tilde{\mathbf{p}} = \tilde{\mathbf{q}}^t\mathbf{W}\tilde{\mathbf{q}} = 1 \quad (1)$$

with  $\mathbf{F}^{(1)} = \tilde{\mathbf{p}}\tilde{\mathbf{q}}$  the best rank-one matrix approximation of  $\mathbf{F}$

## Application on SNPs data: MCA analysis

## Visualization of links between SNPs



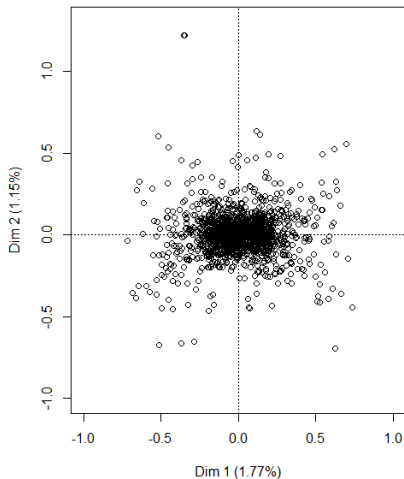
- The SNPs the most contributing to the first axis

- 2 SNPs are close if individuals have the same genetic

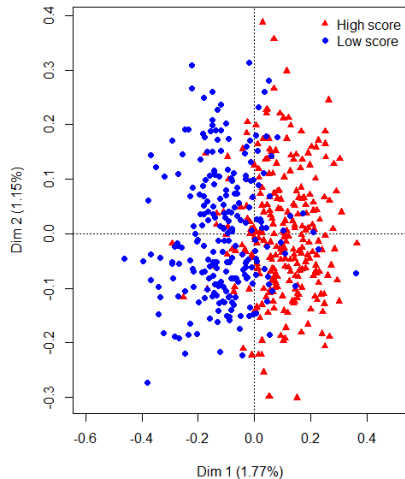
- Too many SNPs  
→ we want to select the most relevant ones

# Application on SNPs data: sparse MCA

Variables Map  
MCA



Individual Map  
MCA





## Application on SNPs data

**Step 1** Visualization of links between SNPs and between samples using **Multiple Correspondence Analysis (MCA)**

**Step 2** Selection of the most important SNPs for a component using the sparse extension of MCA to select variables:  
**Sparse MCA**

## From MCA to sparse MCA

### Challenge

To facilitate interpretation of MCA results

⇒ Select the most contributing SNPs on each axis (easier visualisation of relationship between SNPs and phenotype)

### How?

Constraints imposed in the MCA problem to set coefficients to zero

# From MCA to sparse MCA

ID	SNP.1	SNP.2
AK	GG	AA
JD	GT	CC
MR	GT	AC
GB	GG	CC
NH	TT	AC
AL	GG	CC



GG	GT	TT	CC	AC	AA	ID
1	0	0	0	0	1	AK
0	1	0	1	0	0	JD
0	1	0	0	1	0	MR
1	0	0	1	0	0	GB

SNP1      SNP2

Selection of **1 column** in the original table  
(categorical variable **X**)

=

Selection of **a block of indicator variables** in the complete disjunctive table

## From MCA to sparse MCA

## Sparse MCA via GSVD

$$\min_{\tilde{\mathbf{p}}, \tilde{\mathbf{q}}} \|\mathbf{F} - \tilde{\mathbf{p}}\tilde{\mathbf{q}}^T\|_{\mathbf{W}}^2 \quad \tilde{\mathbf{p}}^T \mathbf{M} \tilde{\mathbf{p}} = \tilde{\mathbf{q}}^t \mathbf{W} \tilde{\mathbf{q}} = 1 \quad (2)$$

## From MCA to sparse MCA

## Sparse MCA via GSVD

+regularization penalty function applied on  $\tilde{\mathbf{q}}$ 

$$\min_{\tilde{\mathbf{p}}, \tilde{\mathbf{q}}} \|\mathbf{F} - \tilde{\mathbf{p}}\tilde{\mathbf{q}}^T\|_{\mathbf{W}}^2 + P_{\lambda}(\tilde{\mathbf{q}}) \quad \tilde{\mathbf{p}}^T \mathbf{M} \tilde{\mathbf{p}} = \tilde{\mathbf{q}}^t \mathbf{W} \tilde{\mathbf{q}} = \mathbf{1} \quad (2)$$

 $P_{\lambda}$  is a penalty function with tuning regularization parameter  $\lambda$

## From MCA to sparse MCA

## Sparse MCA via GSVD

+regularization penalty function applied on  $\tilde{\mathbf{q}}$ 

$$\min_{\tilde{\mathbf{p}}, \tilde{\mathbf{q}}} \|\mathbf{F} - \tilde{\mathbf{p}}\tilde{\mathbf{q}}^T\|_{\mathbf{W}}^2 + P_{\lambda}(\tilde{\mathbf{q}}) \quad \tilde{\mathbf{p}}^T \mathbf{M} \tilde{\mathbf{p}} = \tilde{\mathbf{q}}^t \mathbf{W} \tilde{\mathbf{q}} = \mathbf{1} \quad (2)$$

 $P_{\lambda}$  is a penalty function with tuning regularization parameter  $\lambda$  $\Rightarrow$  Use the **Group LASSO penalization**

$$P_{\lambda}(\boldsymbol{\beta}) = \lambda \sum_{k=1}^K \sqrt{J_{[k]}} \|\boldsymbol{\beta}_{[k]}\|_2$$

 $J_{[k]}$ : number of variables in block  $k$  $\lambda$ : penalty parameter to determine (cross validation, ad hoc approach,...)

## From MCA to sparse MCA

## Sparse MCA via GSVD

+regularization penalty function applied on  $\tilde{\mathbf{q}}$ 

$$\min_{\tilde{\mathbf{p}}, \tilde{\mathbf{q}}} \|\mathbf{F} - \tilde{\mathbf{p}}\tilde{\mathbf{q}}^T\|_{\mathbf{W}}^2 + P_{\lambda}(\tilde{\mathbf{q}}) \quad \tilde{\mathbf{p}}^T \mathbf{M} \tilde{\mathbf{p}} = \tilde{\mathbf{q}}^t \mathbf{W} \tilde{\mathbf{q}} = \mathbf{1} \quad (2)$$

 $P_{\lambda}$  is a penalty function with tuning regularization parameter  $\lambda$ ⇒ Use the **Group LASSO penalization**

$$P_{\lambda}(\boldsymbol{\beta}) = \lambda \sum_{k=1}^K \sqrt{J_{[k]}} \|\boldsymbol{\beta}_{[k]}\|_2$$

 $J_{[k]}$ : number of variables in block  $k$  $\lambda$ : penalty parameter to determine (cross validation, ad hoc approach,...)**Result:** Entire blocks of dummy variables are selected or removed

## Penalty parameter influence

Tuning parameter  $\lambda = 0 \Rightarrow$  sparse MCA=MCA

