

Summary

- 1 GWAS and Block of linkage disequilibrium
 - Genome Wide Association Studies
 - Blocks of linkage disequilibrium
 - Hierarchical Clustering with Adjacency Constraints
 - How to improve?
 - Some computation times
- 2 Epistasis
- 3 Method
 - The G-GEE modeling approach
 - Simulations
- 4 Application
 - Ankylosing Spondylitis
 - First results

High-dimension in Genomics

'Omics'

- Genomics, Transcriptomics, Proteomics, Metabolomics, Epigenomics, Metagenomics

Large p small n

- n of the order of 10 to 10^4 : Each statistical individual is a 'costly' experiment
- p of the order of 10 to 10^9 : Transcriptomics (10^4), Genomics (10^9) ...

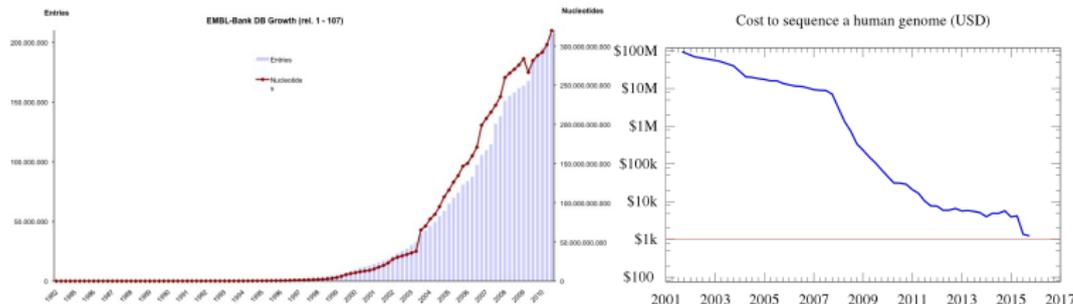
Need for dimensionality reduction

- Selection (of variables)
- Projection in low dimensional subspace
- Clustering

High-dimension in Genomics

Technology evolution and Genomic Genome

- Human Genome Project (1990-2003)
- 2002 launch of HapMap project (report in Nature 2005)
- 2008-2012 : the 1000 Genomes Project
- NGS (new generation sequencing)



Nature Biotechnology

Wikipedia

- 1 GWAS and Block of linkage disequilibrium
 - Genome Wide Association Studies
 - Blocks of linkage disequilibrium
 - Hierarchical Clustering with Adjacency Constraints
 - How to improve?
 - Some computation times
- 2 Epistasis
- 3 Method
 - The G-GEE modeling approach
 - Simulations
- 4 Application
 - Ankylosing Spondylitis
 - First results

Single-Nucleotide Polymorphism Data

- 90 % of human genetic variation,
- In human genom, SNP with allelic frequency greater than 1 % are present every 300 base pairs (in average)
- 2 SNP among 3 substitute cytosine with thymine

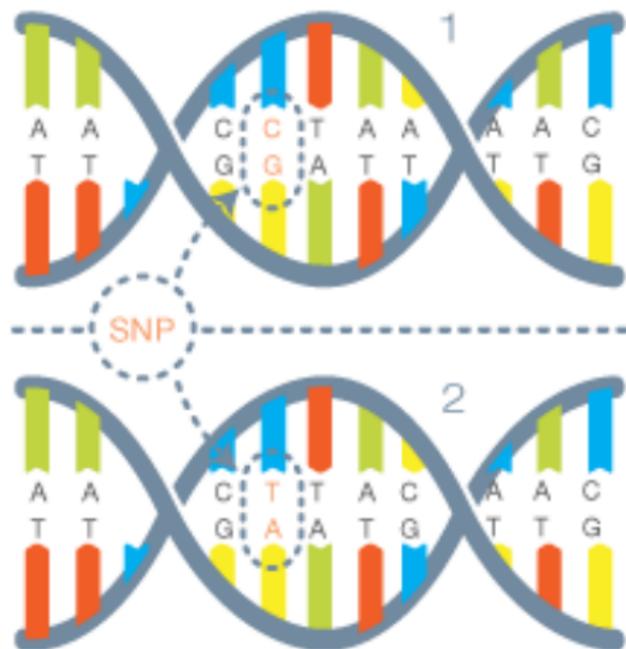


Figure: SNP (wikipedia)

SNP Data

Four possible C/T configurations

$X_{ij} \in \{0, 1, 2\}$ (Individual i at locus j)	Mother	Father	
		C	T
	C	0	1
	T	1	2

High-dimension

p markers

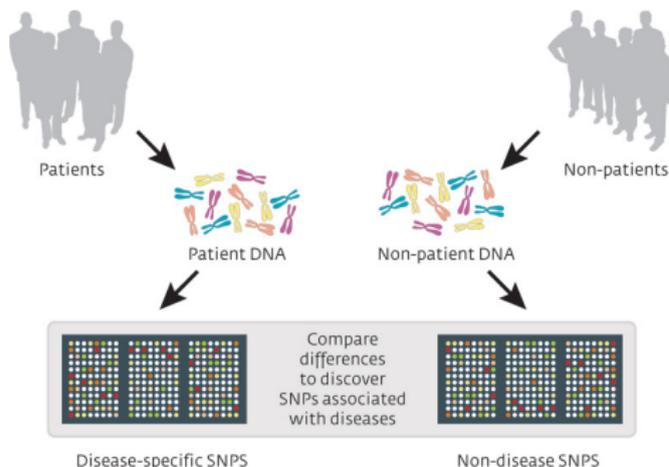
$$n \text{ individuals} \begin{pmatrix} \mathbf{X}_{11} & \cdots & \cdots & \mathbf{X}_{1p} \\ \vdots & & & \vdots \\ \mathbf{X}_{n1} & \cdots & \cdots & \mathbf{X}_{np} \end{pmatrix}$$

\approx up to few million of SNPs can be genotyped !
 \Rightarrow number of variables \gg number of individuals ($p \gg n$)

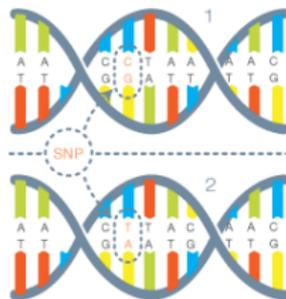
Genome-Wide Association Studies

GWAS characteristics :

- **Objective** : find associations between genetic markers ($SNP_{i,j} \in \{0, 1, 2\}$) and a phenotypic trait ($Y_i \in \{0, 1\}$ or $Y_i \in \mathbb{R}$)



Genetic markers \rightarrow SNP



<http://www.siriusgenomics.com/technology/>

- Generalized Linear Model

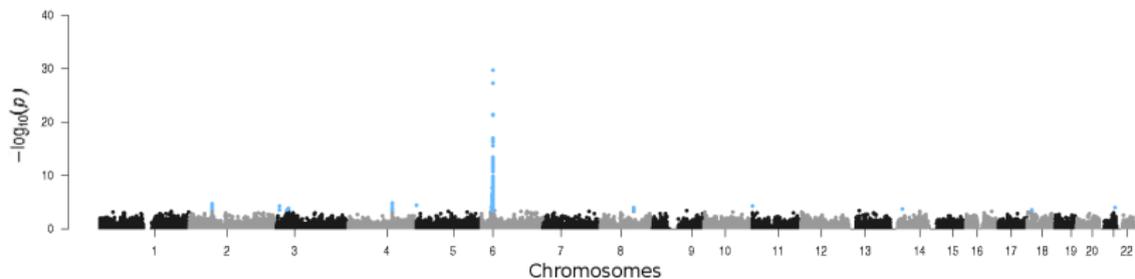
$$g(E[Y_i|x_i]) = \beta_0 + \sum_{j=1}^p \beta_j x_{ij} \quad , i = 1, \dots, n$$

- n : number of individuals
- p : number of covariates
- Y_i : response for the individual i
- x_j : observations for covariate j (coded in 0, 1 or 2)

Genome-Wide Association Studies

SNP analysis

Differences between cases and controls at a specific SNP



Where is the missing heritability?

Missing Heritability

- gene-gene interaction,
- gene-environment interactions,
- rare variant effects (mutation present in less than 1% of the population),
- additivity of numerous common variants (not detected using univariate strategies)....

Data characteristics

- High dimension ($p \gg n$)
- Spatial structure
- Small effects

The LD measures

Linkage Disequilibrium

- non-random association of alleles at two or more loci
- depends on the difference between observed allelic frequencies and those expected from a independent randomly distributed model.

Computation

- Z_j the indicator of the presence of minor allele for SNP j .
- $Z_j \sim \mathcal{B}(p_j)$

$$D(j, k) = p_{jk} - p_j p_k = E[Z_j Z_k] - p_j p_k = \text{cov}(Z_j, Z_k)$$

$$r^2(j, k) = \text{corr}(Z_j, Z_k)$$

ou

$$D'(j, k) = D(j, k) / D_{\max}(j, k)$$

How to estimate LD?

snp	vv	vV	VV
uu	a	b	c
uU	d	e	f
UU	g	h	i

snp	v	V
u	α	β
U	γ	δ



Only the genotype data table is observed

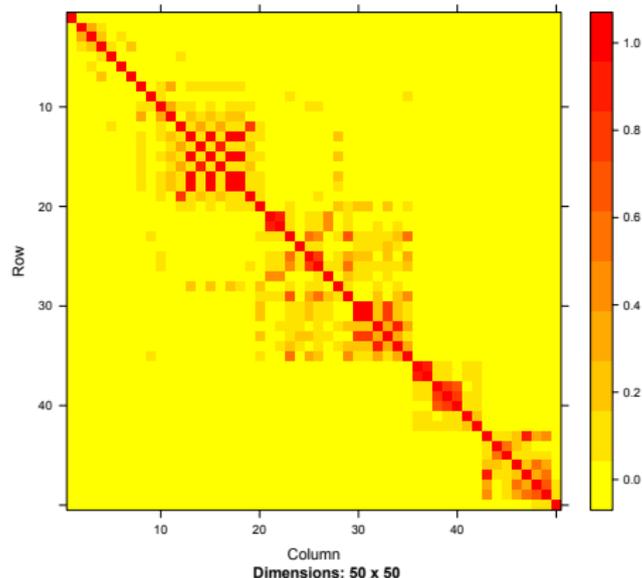
- $\alpha, \beta, \gamma, \delta$ are estimated
- a system of equations. e.g : $\alpha = 2a + b + d + pe$

with p the "probability" of the haplotype (uv, UV).

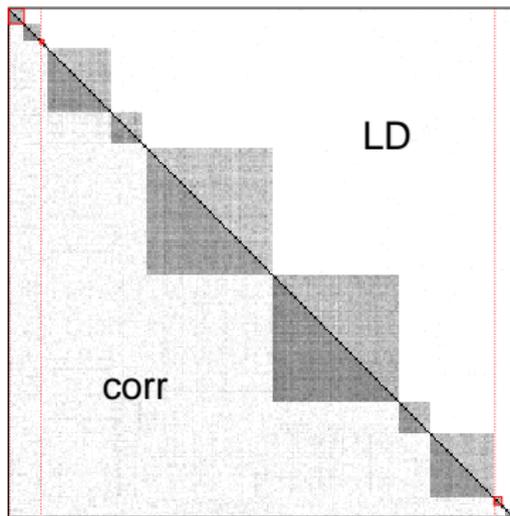
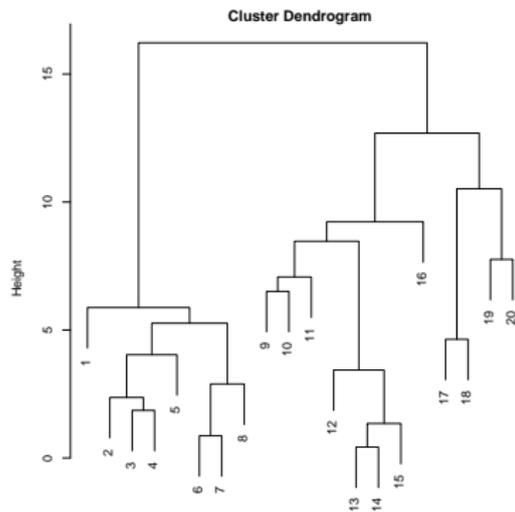
\Rightarrow estimating p , then $(\alpha, \beta, \gamma, \delta)$ and finally
 $D = p_{UV} - p_U p_V$.

The LD block structure

- the r^2 coefficients among the **50 first SNPs** of the Chromosome 22 (Dalmasso et al. 2008)
- LD structured in blocks



Hierarchical Clustering with Adjacency Constraints

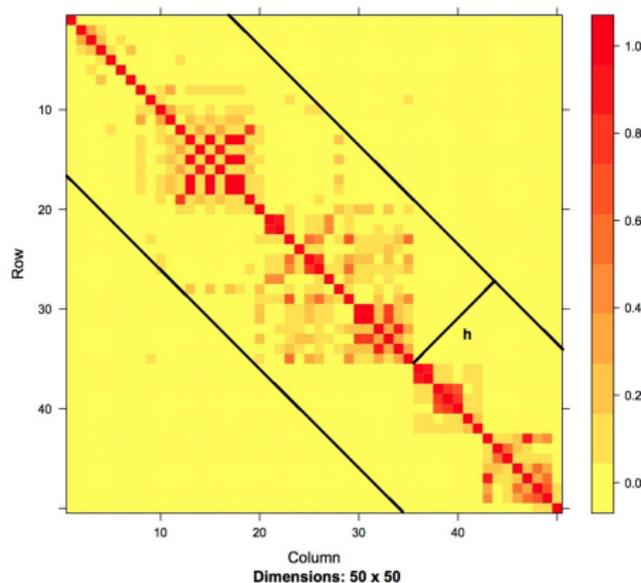


Block-Wise Approach using Linkage Disequilibrium (BALD)

- 1 Hierarchical clustering of the SNPs with adjacency constraint and using the LD similarity.
- 2 Estimation of the optimal number of groups using the Gap statistic (Tibshirani et. al., 2001).

The h-band

- All coefficients outside the band “h” are null
- a $p \times h$ similarity matrix



⇒ a hierarchical clustering with adjacency constraint

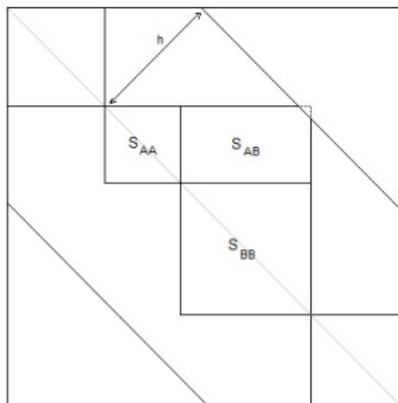
A pseudocode

```
Data:  $\mathbf{X} \in \{0, 1, 2\}^{n \times p}$ , Sim  
 $\mathcal{C} \leftarrow \{C_i = \{\mathbf{X}_{.i}\}, i \in 1, \dots, p\}$  /* clusters = singletons  
*/ ;  
 $D \leftarrow \{1 - \text{Sim}(\mathbf{X}_{.i}, \mathbf{X}_{.(i+1)}), i \in 1, \dots, p - 1\}$  ;  
for step = 1 to  $p - 1$  do  
     $i^* \leftarrow \operatorname{argmin}_{i \in \{1, \dots, p - \text{step}\}} D(C_i, C_{i+1})$  ;  
     $\mathcal{C} \leftarrow \mathcal{C} \setminus \{C_{i^*}, C_{i^*+1}\} \cup \{C_{i^*} \cup C_{i^*+1}\}$  ;  
     $d_1 \leftarrow D(C_{i^*-1}, C_{i^*} \cup C_{i^*+1})$  ;  
     $d_2 \leftarrow D(C_{i^*} \cup C_{i^*+1}, C_{i^*+2})$  ;  
     $D \leftarrow D \setminus \{D(C_{i^*-1}, C_{i^*}), D(C_{i^*}, C_{i^*+1})\} \cup \{d_1, d_2\}$  ;  
end
```

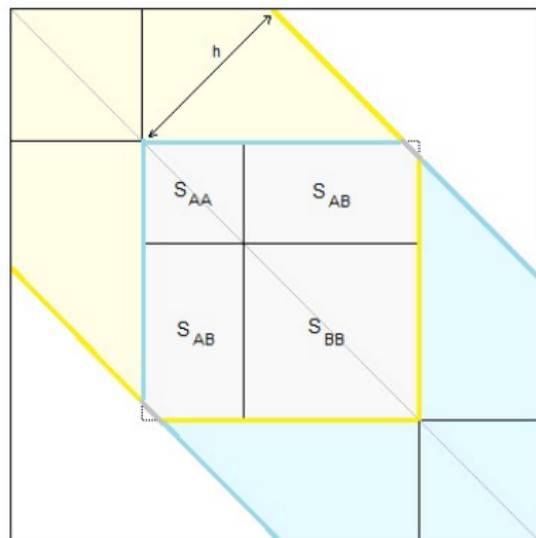
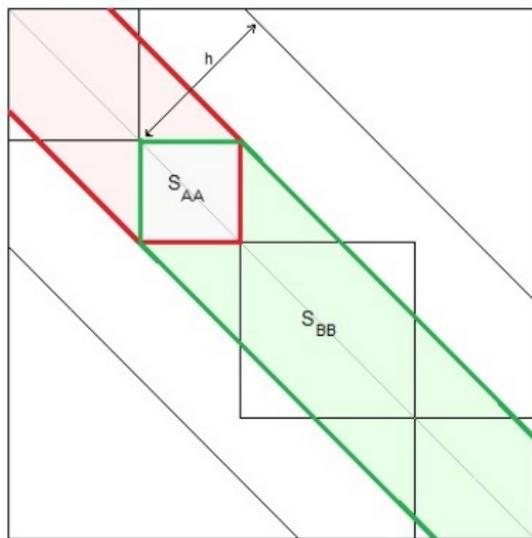
The Ward's distance

Ward Constrained Hierarchical Clustering

$$d(A, B) = \frac{n_A n_B}{n_A + n_B} \left(\frac{1}{n_A^2} S_{A,A} + \frac{1}{n_B^2} S_{B,B} - \frac{2}{n_A n_B} S_{A,B} \right)$$



The pencils' trick : Calculating S_{AA} and S_{AB}



The “pencils”

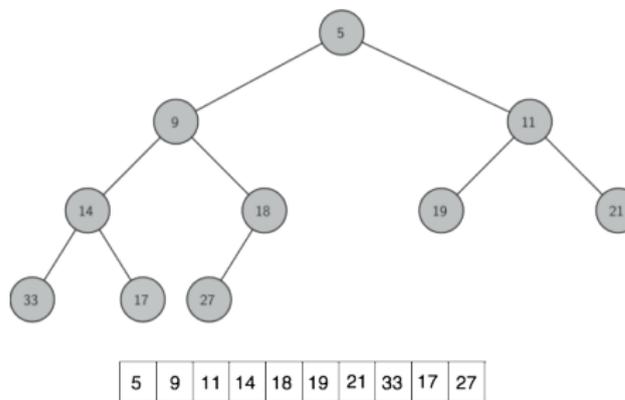
Assessing S_{AA} , S_{BB} and S_{AB} requires the calculation of sums of LD measures within *pencil-shaped areas* defined by :

- direction : right or left
- depth : hLoc
- end point : lim

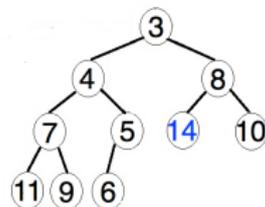
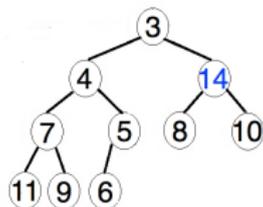
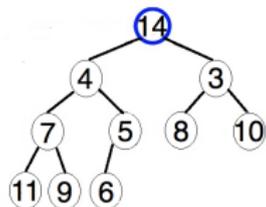
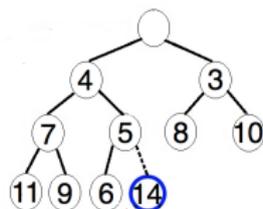
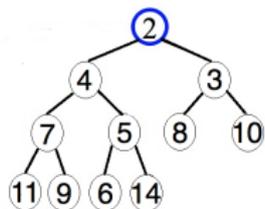
⇒ Two arrays of sizes $p \times h$ for storing the pencils sums.

The binary min-heap

- All nodes are either **less than or equal** to each of its children.
- Uniquely represented by storing its level order traversal in an array.
Given a position i :
 - $\text{Parent}(i) = \lfloor i/2 \rfloor$
 - $\text{Left}(i) = 2i$
 - $\text{Right}(i) = 2i + 1$

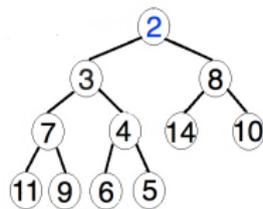
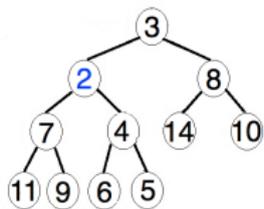
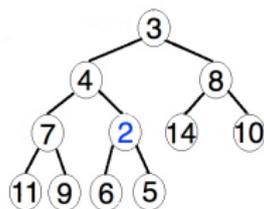
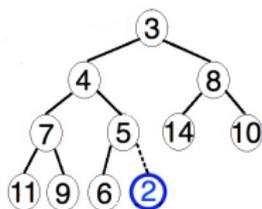
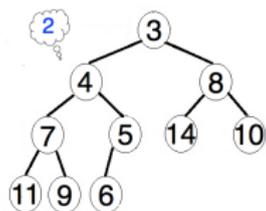


DeleteMin



Time complexity : $\mathcal{O}(\log(p))$

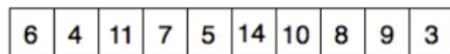
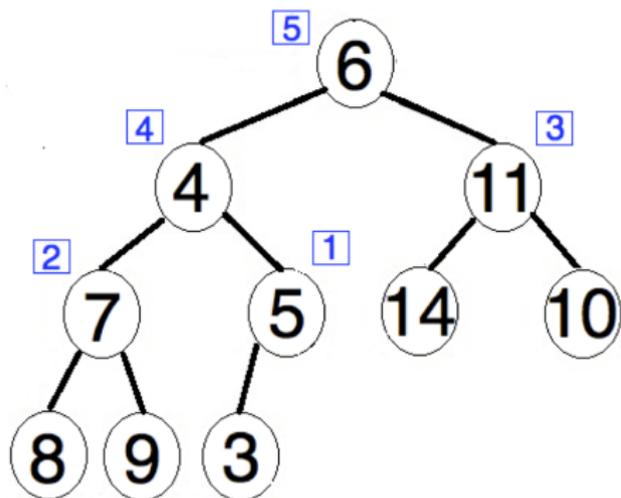
InsertHeap



Time complexity : $\mathcal{O}(\log(p))$

BuildHeap

Data: An array A
Result: A min-heap H
for $i = \lfloor \text{length}(A)/2 \rfloor$ down to 1
do
 | PercolDown(A, i);
end



Time complexity : $\mathcal{O}(p \log(p))$

Time complexity of some operations

	findMin	insert	deleteMin
unordered array	$\mathcal{O}(p)$	$\mathcal{O}(1)$	$\mathcal{O}(p)$
binary heap	$\mathcal{O}(1)$	$\mathcal{O}(\log(p))$	$\mathcal{O}(\log(p))$

cWard in seconds...

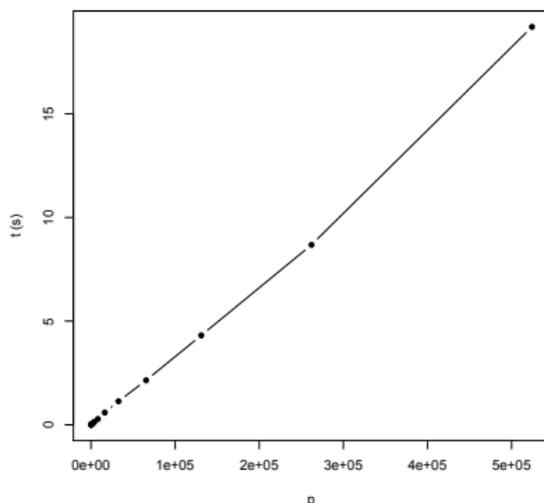


Figure: The mean computation time t versus the number of markers p for the cWard algorithm applied to randomly sampled SNP matrices. $N = 100$, $h = 30$ and t is averaged across 50 simulation runs.

Compared to a former implementation

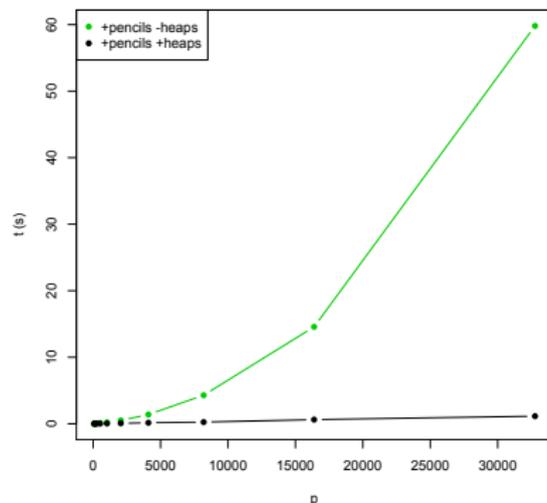


Figure: The mean computation time t versus the number of markers p for the cWard algorithm and an implementation without heaps. t is averaged across 20 simulation runs.

Scalable Hierarchical Clustering with pencils and binary heap

To sum up :

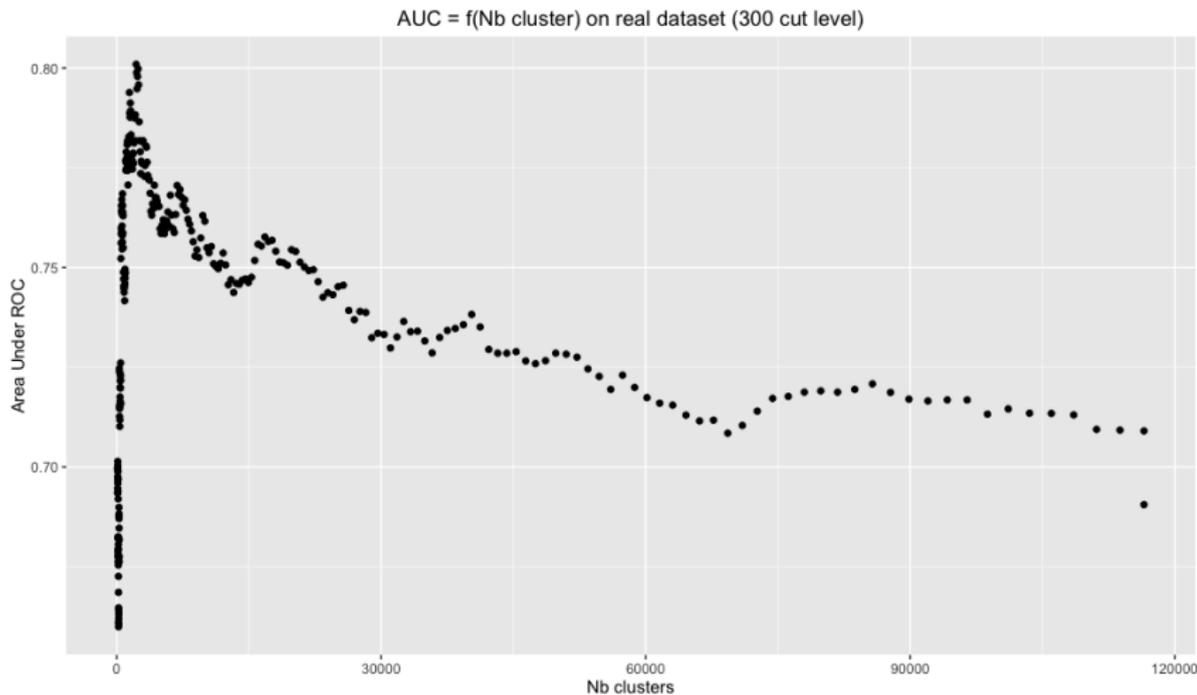
- 1 A $\mathcal{O}(p^2)$ algorithm does not scale for GWA studies.
- 2 The Ward distance written in a simple way.
- 3 Space complexity of $\mathcal{O}(ph)$ by using the pencils' trick.
- 4 Time complexity of :

$$\underbrace{\mathcal{O}(ph)}_{2\ p \times h \text{ arrays of pencils' sums}} + \underbrace{\mathcal{O}(p \log(p))}_{\text{building the heap and insert/delete heaps' operations within the loop}}$$

Ongoing work :

- Currently implemented with a genotype matrix as input.
⇒ can be generalized to any band similarity matrix.

What is the appropriate scale for association analysis?



- 1 GWAS and Block of linkage disequilibrium
 - Genome Wide Association Studies
 - Blocks of linkage disequilibrium
 - Hierarchical Clustering with Adjacency Constraints
 - How to improve ?
 - Some computation times
- 2 Epistasis
- 3 Method
 - The G-GEE modeling approach
 - Simulations
- 4 Application
 - Ankylosing Spondylitis
 - First results

Definition

Interaction of alleles effects from different markers

Existing methods

- mainly SNP \times SNP
- some at the block (gene) scale

Advantages of gene (or block) scale approaches

- results biologically interpretable
- genetic effects may be easier to detect
- reduce the number of variables

Epistasis - Gene scale methods

Existing gene scale methods :

Two or few genes

- PCA + logistic regression (*He et al. 2011, Li et al. 2009, Zhang et al. 2008*)
- PLS + logistic regression (*Wang T et al. 2009*)

For a larger number of genes

- PCA + LASSO (*D'Angelo et al. 2009*)
- PCA + pathway-guided penalized regression (*Wang X et al. 2014*)

Epistasis - Gene scale methods

Existing gene scale methods :

Two or few genes

- PCA + logistic regression (*He et al. 2011, Li et al. 2009, Zhang et al. 2008*)
- PLS + logistic regression (*Wang T et al. 2009*)

For a larger number of genes

- PCA + LASSO (*D'Angelo et al. 2009*)
- PCA + pathway-guided penalized regression (*Wang X et al. 2014*)

Objectives : To develop a new gene scale method

- considers interaction variables,
- takes into account the group structure,
- is applicable with many groups (genes)

- 1 GWAS and Block of linkage disequilibrium
 - Genome Wide Association Studies
 - Blocks of linkage disequilibrium
 - Hierarchical Clustering with Adjacency Constraints
 - How to improve ?
 - Some computation times
- 2 Epistasis
- 3 **Method**
 - The G-GEE modeling approach
 - Simulations
- 4 Application
 - Ankylosing Spondylitis
 - First results

Group modeling approach

	$SNP_{1,1}$..	SNP_{1,p_1}	..	$SNP_{G,1}$..	SNP_{G,p_G}	Pheno
Ind_1	1		0		0		1	y_1
Ind_2	0		0		2		1	y_2
.	2		1		1		2	.
.	0		1		0		0	.
Ind_i	0		2		1		0	y_i

⏟
⏟
gene₁
gene_G

We note
 $SNP_{1,1} =$
 $X_{1,1}$

model :

$$g(E[y|\mathbf{X}]) = \underbrace{\sum_g \sum_{p_g} \beta_{g,p_g} \mathbf{X}_{g,p_g}}_{\text{Main effects}}$$

$$\beta = \left(\underbrace{\beta_{1,1}, \beta_{1,2}, \dots, \beta_{1,p_1}}_{\text{gene 1}}, \dots, \underbrace{\beta_{G,1}, \dots, \beta_{G,p_G}}_{\text{gene G}} \right)^T$$

Group modeling approach

	$SNP_{1,1}$..	SNP_{1,p_1}	..	$SNP_{G,1}$..	SNP_{G,p_G}	Pheno
Ind_1	1		0		0		1	y_1
Ind_2	0		0		2		1	y_2
.	2		1		1		2	.
.	0		1		0		0	.
Ind_i	0		2		1		0	y_i

⏟
 $gene_1$
⏟
 $gene_G$

We note
 $SNP_{1,1} = X_{1,1}$
 r, s two genes

model :

$$g(E[y|\mathbf{X}]) = \underbrace{\sum_g \sum_{p_g} \beta_{g,p_g} \mathbf{X}_{g,p_g}}_{\text{Main effects}} + \underbrace{\sum_{r,s} \gamma_{r,s} \mathbf{Z}_{r,s}}_{\text{Interaction effects}}$$

$$\beta = \left(\underbrace{\beta_{1,1}, \beta_{1,2}, \dots, \beta_{1,p_1}}_{\text{gene 1}}, \dots, \underbrace{\beta_{G,1}, \dots, \beta_{G,p_G}}_{\text{gene G}} \right)^T \quad \gamma = \left(\gamma_{12}, \dots, \underbrace{\gamma_{1G}}_{\text{gene 1 \& G}}, \dots, \gamma_{(G-1)G} \right)$$

Interaction variables construction : Gene-Gene Eigen Epistasis (G-GEE)

$Z^{rs} = f_{\mathbf{u}}(\mathbf{X}^r, \mathbf{X}^s)$ the interaction between genes r, s .

$$\hat{\mathbf{u}} = \arg \max_{\mathbf{u}, \|\mathbf{u}\|=1} \text{cov}^2(\mathbf{y}, f_{\mathbf{u}}(\mathbf{X}^r, \mathbf{X}^s))$$

We set : $f_{\mathbf{u}}(\mathbf{X}^r, \mathbf{X}^s) = \mathbf{W}^{rs} \mathbf{u}$ with

$$\mathbf{W}^{rs} = \{X_{ij}^r X_{ik}^s\}_{i=1 \dots n}^{j=1 \dots, p_r; k=1 \dots, p_s}$$

$$\max_{\mathbf{u}, \|\mathbf{u}\|=1} \|\widehat{\text{cov}}[\mathbf{W}^{rs} \mathbf{u}, \mathbf{y}]\|^2 = \max_{\mathbf{u}, \|\mathbf{u}\|=1} \mathbf{u}^T \mathbf{W}^{rsT} \mathbf{y} \mathbf{y}^T \mathbf{W}^{rs} \mathbf{u}$$

\mathbf{u} : eigen vector associated to the largest eigenvalue of $\mathbf{W}^{rsT} \mathbf{y} \mathbf{y}^T \mathbf{W}^{rs}$

$$\mathbf{u} = \mathbf{W}^{rsT} \mathbf{y}$$

Coefficients estimation (linear model)

Group LASSO regression

$$(\hat{\beta}, \hat{\gamma}) = \underset{\beta, \gamma}{\operatorname{argmin}} \sum_i (y_i - \mathbf{x}_i \beta - \mathbf{z}_i \gamma)^2 + \lambda \left(\sum_g \sqrt{p_g} \|\beta^g\|_2 + \sum_{rs} \sqrt{p_r p_s} \|\gamma^{rs}\|_2 \right)$$

Limits of the groupLASSO regression :

- $P(S^* \subset \hat{S}) \xrightarrow[n \rightarrow +\infty]{} 1$ but $|\hat{S}| \gg |S^*|$
- Difficult to compute p-value or confidence interval

Coefficients estimation (linear model)

Group LASSO regression

$$(\hat{\beta}, \hat{\gamma}) = \underset{\beta, \gamma}{\operatorname{argmin}} \sum_i (y_i - \mathbf{x}_i \beta - \mathbf{z}_i \gamma)^2 + \lambda \left(\sum_g \sqrt{p_g} \|\beta^g\|_2 + \sum_{rs} \sqrt{p_r p_s} \|\gamma^{rs}\|_2 \right)$$

Limits of the groupLASSO regression :

- $P(S^* \subset \hat{S}) \xrightarrow{n \rightarrow +\infty} 1$ but $|\hat{S}| \gg |S^*|$
- Difficult to compute p-value or confidence interval

Adaptive-Ridge Cleaning *Becu JM, 2015*

- Use of a specific penalty for group LASSO
- Permutation test based on Fisher test approach for each group

Genotype :

From a real data set composed of 763 individuals and 63340 SNPs structured in 7216 genes.

Continuous phenotype simulated under two different schemes :

→ from Wang X et al., 2014 :

$$Y_i = \beta_0 + \sum_g \beta_g \left(\sum_{k \in \mathcal{C}} X_{ik}^g \right) + \sum_{rs} \gamma_{rs} \left(\sum_{(j,k) \in \mathcal{C}^2} X_{ij}^r X_{ik}^s \right) + \epsilon_i \quad (1)$$

→ PCA model :

$$Y_i = \beta_0 + \sum_g \beta_g \left(\sum_{k \in \mathcal{C}} X_{ik}^g \right) + \sum_{rs} \gamma_{rs} C_{i1}^r C_{i1}^s + \epsilon_i. \quad (2)$$

Scenarios :

For each setting we consider 6 genes.

→Five settings :

- same genes for main and interaction effects,
- different genes for main and interaction effects,
- only one interaction effect,
- only two main effects,
- no effects.

Simulations results - Interactions power

Settings →

Main effects :

gene 1
gene 2

Interaction effects :

gene 1 x gene 2

Main effects :

gene 1
gene 2

Interaction effects :

gene 3 x gene 4

Main effects :

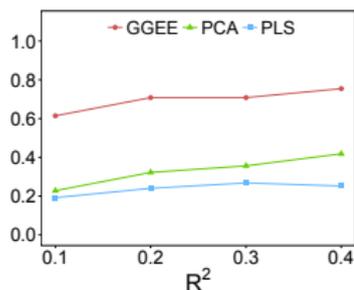
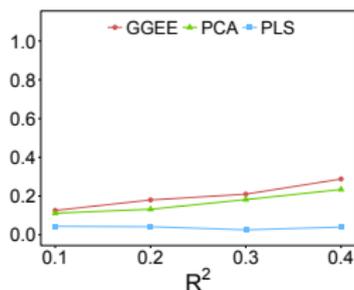
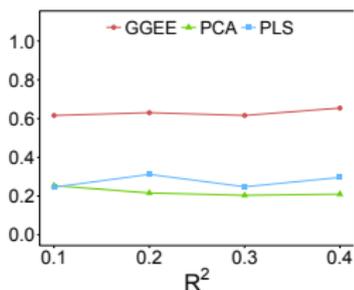
-

Interaction

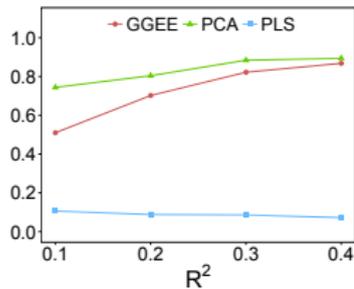
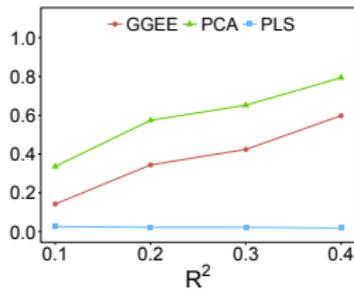
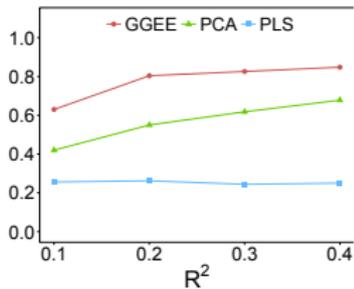
effects :

gene 1 x gene 2

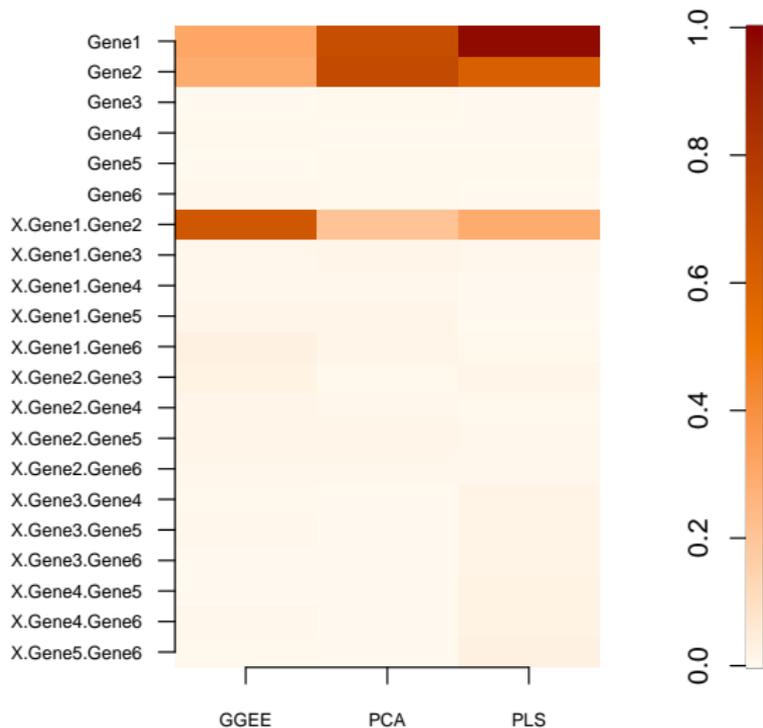
Wang si-
mulation
model



PCA si-
mulation
model



Simulations results - Discoveries matrix



Simulations results - $R^2 = 0.2$

Settings →

Main effects :

gene 1
gene 2

Interaction effects :

gene 1 x gene 2

Main effects :

gene 1
gene 2

Interaction effects :

gene 3 x gene 4

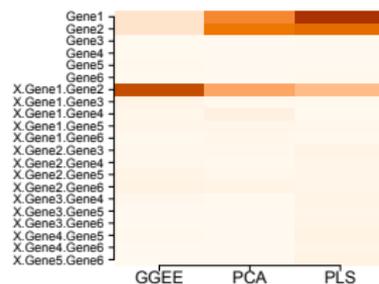
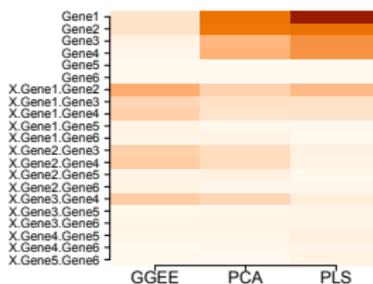
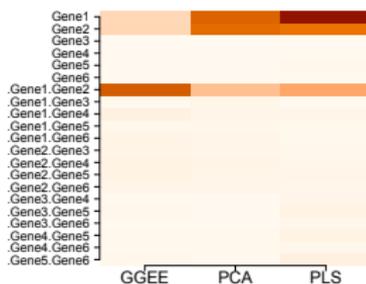
Main effects :

-

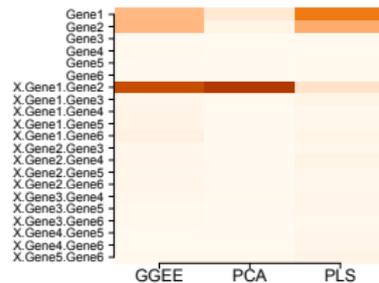
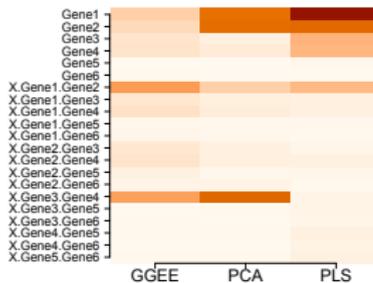
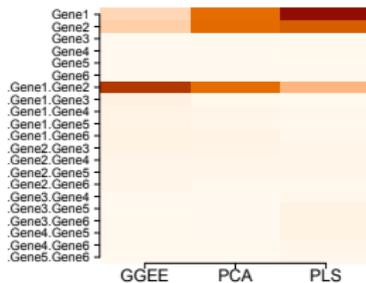
Interaction effects :

gene 1 x gene 2

Wang simulation model



PCA simulation model



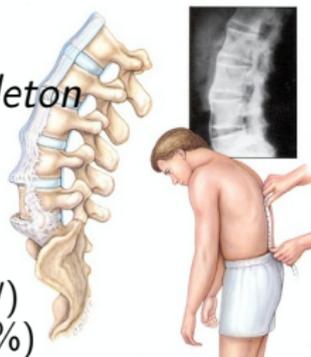
- 1 GWAS and Block of linkage disequilibrium
 - Genome Wide Association Studies
 - Blocks of linkage disequilibrium
 - Hierarchical Clustering with Adjacency Constraints
 - How to improve ?
 - Some computation times
- 2 Epistasis
- 3 Method
 - The G-GEE modeling approach
 - Simulations
- 4 Application
 - Ankylosing Spondylitis
 - First results

Ankylosing Spondylitis

Chronic inflammatory disease of the axial skeleton

Epidemiology :

- Age at first symptoms : 20 - 30 years
- Sexe : predominance for men (sex ratio 2M :1W)
- Prevalence : depend of populations (0.1% - 1.4%)



Right etiology unknown :

- Environmental factors ?
- Genetic factors ?
→ Importance of HLA complex

HLA complex :

- Localized on chromosome 6
- Regroup about 200 genes
- Coding the immunity system
- Antigen HLA-B27 :
associated to SPA

→ Effect from other gene in HLA group ?

→ Outside HLA group ?

Known genes

Table 1 Summary of ankylosing spondylitis-susceptibility genes identified by genome-wide association studies

<i>RUNX3</i>	Runt-related transcription factor 3
<i>IL23R</i>	Interleukin 23 receptor
<i>IL12Rβ2</i>	Interleukin 12 receptor, β2
<i>GRP25</i>	G-protein-coupled receptor 25
<i>KIF21B</i>	Kinesin family member 21B
<i>PTGER4</i>	Prostaglandin E receptor 4 (subtype EP ₄)
<i>ERAP1</i>	Endoplasmic reticulum aminopeptidase 1
<i>ERAP2</i>	Endoplasmic reticulum aminopeptidase 2
<i>LNPEP</i>	Leucyl/cystinyl aminopeptidase
<i>IL12B</i>	Interleukin 12B
<i>CARD9</i>	Caspase recruitment-domain family member 9
<i>LTβR</i>	Lymphotoxin β-receptor (TNFR superfamily, member 3)
<i>TNFRSF1A</i>	Tumor-necrosis factor-receptor superfamily member 1A
<i>NPEPPS</i>	Aminopeptidase puromycin-sensitive
<i>TBxBP1</i>	TNFR-associated factor family member-associated nuclear factor-κB-binding kinase 1-binding protein
<i>TBX21</i>	T-box 21

<i>IL6R</i>	Interleukin 6 receptor
<i>FCGR2A</i>	Fc fragment of immunoglobulin G, low-affinity IIa, receptor (CD32)
<i>UBE2E3</i>	Ubiquitin-conjugating enzyme E2E 3
<i>GPR35</i>	G-protein-coupled receptor 35
<i>NKX2-3</i>	NK2 homeobox 3
<i>ZMIZ1</i>	Zinc finger, MIZ type-containing 1
<i>SH2B3</i>	Src homology 2B adaptor protein 3
<i>GPR65</i>	G-protein-coupled receptor 65
<i>IL27</i>	Interleukin 27
<i>SULT1A1</i>	Sulfotransferase family cytosolic 1A
<i>TYK2</i>	Tyrosine kinase 2
<i>ICOSLG</i>	Inducible T-cell costimulator ligand
<i>EOMES</i>	Eomesodermin
<i>IL7R</i>	Interleukin 7 receptor
<i>BACH2</i>	BTB and CNC homology 1, basic leucine-zipper transcription-factor 2

Abbreviation: CD, classification determinant.

Tsui et al., 2014 : The genetic basis of ankylosing spondylitis : new insights into disease pathogenesis, The Application of Clinical Genetics :7 105-115

→ 29 susceptibility genes identified by GWAS

	Significant results
G-GEE	HLA-B x SULT1A1 IL23R x ERAP2
PLS	HLA-B EOMES x BACH2
PCA	HLA-B

Conclusions and perspectives

The G-GEE method

- Takes into account the gene structure of data
- Can be applied on a large number of genes
- Uses a specific interaction modeling approach

Ankylosing Spondylitis

- Identification of potential interactions to discuss with medical doctors
- HLA-B effect

Perspectives

- Explore new $f_u(\mathbf{X}_i^r, \mathbf{X}_i^s)$
- Definition Additional simulations on larger data set
- New applications on other data sets

Thank you for your attention !



Simulations results - Main effect power

Wang simu model

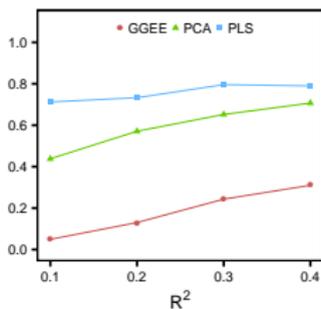
Main effects :

gene 1

gene 2

Interaction effects :

gene 1 x gene 2



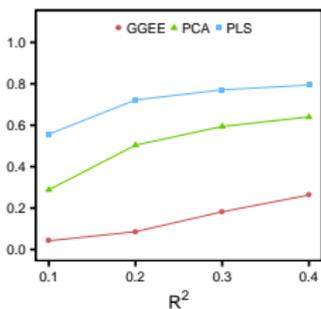
Main effects :

gene 1

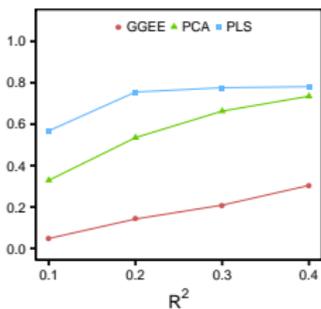
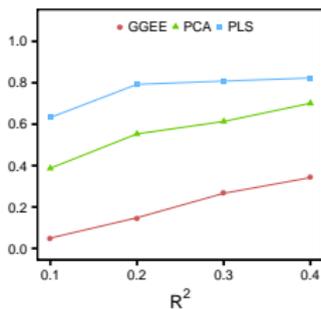
gene 2

Interaction effects :

gene 3 x gene 4



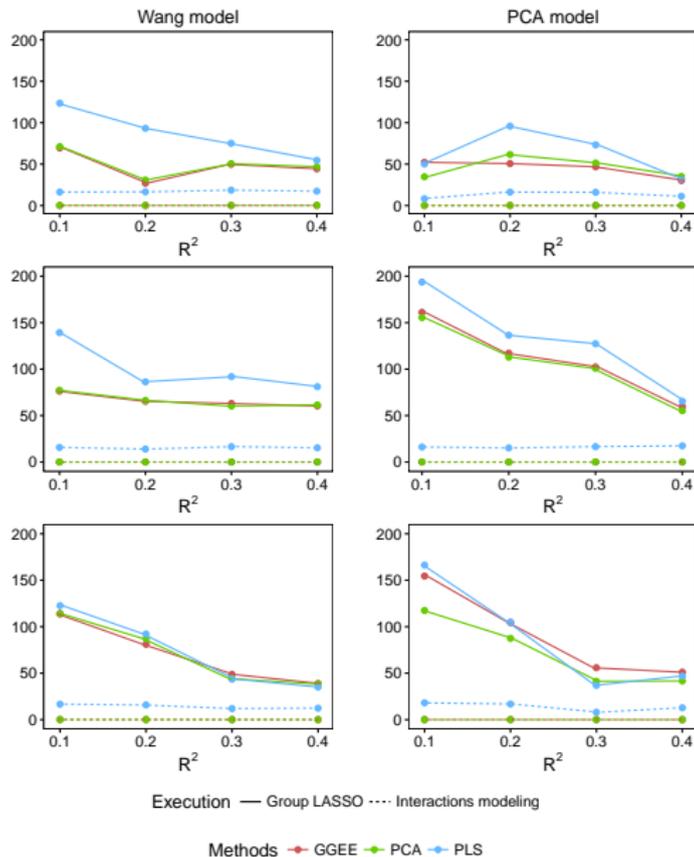
PCA simu model



Group LASSO regression

$$\hat{\theta} = (\hat{\beta}, \hat{\gamma}) = \underset{\beta, \gamma}{\operatorname{argmin}} \left(\sum_i -L(y_i; \mathbf{X}_i \beta + \mathbf{Z}_i \gamma) + \lambda \left[\sum_g \sqrt{p_g} \|\beta^g\|_2 + \sum_{rs} \sqrt{p_r p_s} \right] \right)$$

Execution median time (seconds)



Main effects :

gene 1
gene 2

Interaction effects :

gene 1 x gene 2

Main effects :

gene 1
gene 2

Interaction effects :

gene 3 x gene 4

Main effects :

Interaction effects :

gene 1 x gene 2

Simulations design

Genotype :

$\mathbf{X}_i \sim \mathcal{N}_p(\mathbf{0}, \mathbf{\Sigma})$ with $\mathbf{\Sigma}$ a block diagonal correlation matrix
($\rho = 0.8$ for two SNPs in the same gene)

$MAF_j \sim \mathcal{U}[0.05, 0.5]$ with $MAF_j = 0.2$ if j causal SNP

Scenarios :

We consider 600 subjects and 6 SNPs by gene

→ First scenario on 6 genes, two settings :

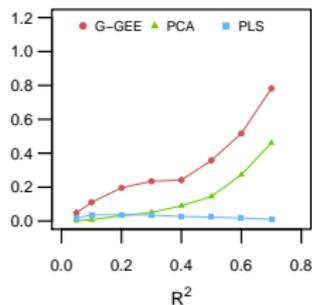
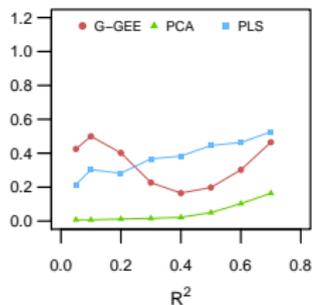
- same genes for main and interaction effects,
- different genes for main and interaction effects.

→ Second scenario on 25 genes, one setting :

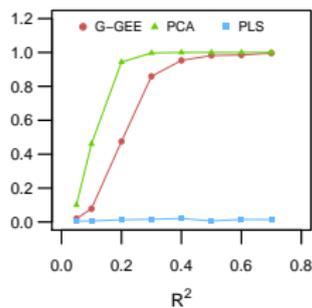
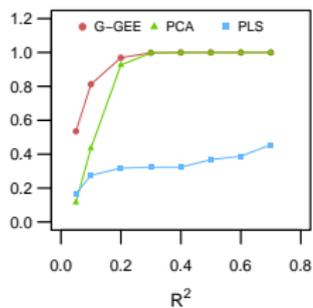
- different genes for main and interaction effects.

Simulations results - Interactions power ; First scenario on 6 genes

Wang model



PCA model



→ Main effects :

gene 1

gene 2

→ Interaction

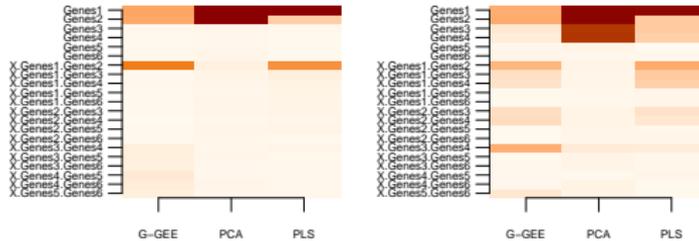
→ Main effects :

gene 1

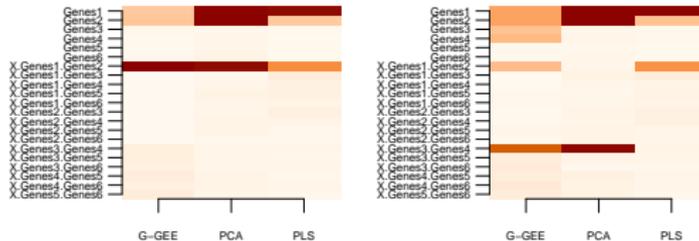
gene 2

Simulations results - $R^2 = 0.2$; First scenario on 6 genes

Wang model



PCA model



→ Main effects :

gene 1

gene 2

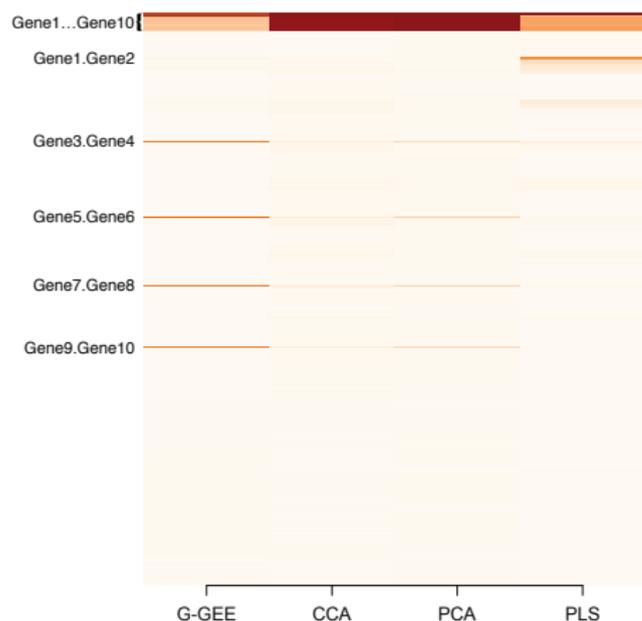
→ Interaction

→ Main effects :

gene 1

gene 2

Simulations results - Second scenario on 25 genes



→ Main effects :

gene 1

gene 2

→ Interaction effects :

gene 3 x gene 4

gene 5 x gene 6

gene 7 x gene 8

gene 9 x gene 10

Figure: Wang X et al. model, $R^2 = 0.7$

Adaptive-Ridge Cleaning

specific penalty for group LASSO : $\frac{\lambda}{\sqrt{|k(j)| \sum_{m \in k(j)} \hat{\theta}_m^2}}$