

Régression en grande dimension et épistasie par blocs pour les études d'association

V. Stanislas, C. Dalmasso, C. Ambroise

Laboratoire de Mathématiques et Modélisation d'Évry "Statistique et Génome"



Summary

GWAS and Block of linkage desequilibrium

- Genome Wide Association Studies
- Blocks of linkage desiquilibrium
- Hierachical Clustering with Adjacency Constraints
- How to improve?
- Some computation times
- 2 Epistasis
 - The Gene-Gene Eigen Epistasis Modeling approach
 - Simulations
- 3 Application
 - Ankylosing Spondylitis
 - First results

Sommaire

GWAS and Block of linkage desequilibrium

- Genome Wide Association Studies
- Blocks of linkage desiguilibrium
- Hierachical Clustering with Adjacency Constraints
- How to improve?
- Some computation times

2 Epistasis

• The Gene-Gene Eigen Epistasis Modeling approach

Simulations

3 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



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 $Y_i \in \mathbb{R}$



©Pasieka, Science Photo Library

• Generalized Linear Model

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

- *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



- GWAS limits :
 - Reproductibility
 - Heritability

- Data particularities :
 - Structuration
 - High dimension (p » n)
 - Small effects

The LD measures

Linkage Desiquilibrium

- 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_i ∼ B(p_i)

$$D(j,k) = p_{jk} - p_j p_k = E[Z_j Zk] - p_j p_k = cov(Z_j, Z_k)$$
$$r^2(j,k) = corr(Z_j, Z_k)$$

ou

$$D'(j,k) = D(j,k)/Dmax(j,k)$$

How to estimate LD?

snp	vv	vV	VV	snn	V	V
uu	а	b	С	sip	V	
uU	d	е	f	<u> </u>	α	β
	~	h		U U	γ	δ
00	g	п	I			

Only the genotype data table is observed

- α , β , γ , δ are estimated
- a system of equations. e.g : $\alpha = 2a + b + d + pe$

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

⇒ estimating p, then (α , β , γ , δ) 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



Hierachical Clustering with Adjacency Constraints



Block-Wise Approach using Linkage Disequilibrium (BALD)

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

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



 \Rightarrow a hierarchical clustering with adjacency constraint

A pseudocode

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

The Ward's distance

Ward Constrained Hierarchical Clustering

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



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





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

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

- 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* :
 - $Parent(i) = \lfloor i/2 \rfloor$
 - Left(i) = 2i
 - $\operatorname{Right}(i) = 2i + 1$



DeleteMin











Time complexity : O(log(p))

InsertHeap











Time complexity : O(log(p))

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



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

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 :

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



 $\mathcal{O}(plog(p))$

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.

Sommaire

1 GWAS and Block of linkage desequilibrium

- Genome Wide Association Studies
- Blocks of linkage desiguilibrium
- Hierachical Clustering with Adjacency Constraints
- How to improve?
- Some computation times

2 Epistasis

- The Gene-Gene Eigen Epistasis Modeling approach
- Simulations

3 Application

- Ankylosing Spondylitis
- First results

Definition

Interaction of alleles effects from different markers

Existing methods

- mainly SNP x 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

- \rightarrow considers a more accurate definition of interaction variables,
- \rightarrow is applicable with with genes,
- \rightarrow takes into account the group structure





Model $\boldsymbol{Y_i} = \beta_0 + \sum_{g} \beta_g \left(\sum_{k \in \mathcal{C}} \boldsymbol{X_{ik}^g} \right)$ $+\epsilon_i$ Main effects $\boldsymbol{\beta} = \left(\underbrace{\beta_{1,1}, \beta_{1,2}, \cdots, \beta_{1,p_1}}_{\text{gener}}, \cdots, \underbrace{\beta_{G,1}, \cdots, \beta_{G,p_G}}_{\text{gener}}\right)^{\prime}$

Model

$$Y_{i} = \beta_{0} + \sum_{g} \beta_{g} \left(\sum_{k \in \mathcal{C}} X_{ik}^{g} \right) + \sum_{\substack{r,s \\ \text{Main effects}}} \gamma_{r,s} Z_{i}^{r,s} + \epsilon_{i}$$

$$\beta = \left(\underbrace{\beta_{1,1}, \beta_{1,2}, \dots, \beta_{1,p_{1}}, \dots, \underbrace{\beta_{G,1}, \dots, \beta_{G,p_{G}}}_{gene_{G}} \right)^{T} \qquad \gamma = \left(\gamma_{12}, \dots, \underbrace{\gamma_{1G}}_{\gamma_{1G,1}, \dots, \gamma_{1G,q}}, \dots, \gamma_{(G-1)G} \right)$$

$$q : \# \text{ of interaction variables for a}$$

Interaction variable : Gene-Gene Eigen Epistasis (G-GEE)

We consider $f_{\boldsymbol{u}}(\boldsymbol{X}_{i}^{r}, \boldsymbol{X}_{i}^{s})$ to represent the interaction between genes r, s.

$$\hat{\boldsymbol{u}} = arg \max_{\boldsymbol{u}, \|\boldsymbol{u}\|=1} cor(\boldsymbol{y}, f_{\boldsymbol{u}}(\boldsymbol{X}^r, \boldsymbol{X}^s))$$

Eigen Epistasis

$$f_{\boldsymbol{u}}(\boldsymbol{X}^r, \boldsymbol{X}^s) = \boldsymbol{W}^{rs} \boldsymbol{u}$$
 with $\boldsymbol{W}^{rs} = \{X_{ij}^r X_{ik}^s\}_{i=1\cdots n}^{j=1\cdots, p_r; k=1, \cdots, p_s}$

$$\max_{\boldsymbol{u},\|\boldsymbol{u}\|=1} ||\hat{cor}[\boldsymbol{W}^{rs}\boldsymbol{u},\boldsymbol{y}]||^2 == \max_{\boldsymbol{u},\|\boldsymbol{u}\|=1} \boldsymbol{u}^T \boldsymbol{W}^{rsT} \boldsymbol{y} \boldsymbol{y}^T \boldsymbol{W}^{rs} \boldsymbol{u}$$

u: only eigen vector associated of $W^{rsT}yy^TW^{rs}$ For each couple $(r, s) \rightarrow Z^{rs} = W^{rs}u$

Coefficients estimation

Group LASSO regression

$$(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}}) = \underset{\boldsymbol{\beta}, \boldsymbol{\gamma}}{\operatorname{argmin}} \sum_{i} (y_{i} - \boldsymbol{X}_{i} \boldsymbol{\beta} - \boldsymbol{Z}_{i} \boldsymbol{\gamma})^{2} + \\ \lambda \left(\sum_{g} \sqrt{p_{g}} || \boldsymbol{\beta}^{g} ||_{2} + \sum_{rs} \sqrt{p_{r} p_{s}} || \boldsymbol{\gamma}^{rs} ||_{2} \right)$$

Limits of the groupLASSO regression :

• Difficult to compute p-value or confidence interval

Coefficients estimation

Group LASSO regression

$$(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}}) = \underset{\boldsymbol{\beta}, \boldsymbol{\gamma}}{\operatorname{argmin}} \sum_{i} (y_{i} - \boldsymbol{X}_{i}\boldsymbol{\beta} - \boldsymbol{Z}_{i}\boldsymbol{\gamma})^{2} + \\ \lambda \left(\sum_{g} \sqrt{p_{g}} ||\boldsymbol{\beta}^{g}||_{2} + \sum_{rs} \sqrt{p_{r}p_{s}} ||\boldsymbol{\gamma}^{rs}||_{2} \right)$$

Limits of the groupLASSO regression :

• 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 $P_k = \frac{1}{B} \# \{F_k^* \ge F_k\}$

Interaction variable modeling approaches comparison

methods	criteria	$Z_i^{rsT}\gamma^{rs}$
G-GEE	$cor(\boldsymbol{Y}, f_{\boldsymbol{u}}(\boldsymbol{X}^r, \boldsymbol{X}^s))$	$W^{rs}u\gamma^{rs}$
PCA	$var(G_r v)$ and $var(G_s v)$	$\sum_{j=1}^{q}\sum_{k=1}^{q}\gamma_{jk}^{rs}C_{j}^{r}C_{k}^{s}$
ССА	$cor(\boldsymbol{G}_r \boldsymbol{a}, \boldsymbol{G}_s \boldsymbol{b})$	$\sum_{j=1}^q \gamma_j^{ m rs} oldsymbol{A}_j^r oldsymbol{B}_j^s$
PLS	$cov(\mathbf{Y}\mathbf{G}_{r}\mathbf{c},\mathbf{G}_{s}\mathbf{w})$	$\sum_{j=1}^q \gamma_j^{ m rs} {m au}_j^{ m rs}$

Simulation design

$\begin{array}{l} \textbf{Genotype:} \\ \textbf{\textit{X}}_i \sim \mathcal{N}_{\rho}(\textbf{0}, \textbf{\Sigma}) \text{ with } \textbf{\textit{\Sigma}} \text{ a block diagonal correlation matrix} \\ (\rho = 0.8 \text{ for two SNPs in the same gene}) \end{array}$

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

Continuous phenotype simulated under two different schemes :

 \rightarrow 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 :

$$Y_{i} = \beta_{0} + \sum_{g} \beta_{g} \left(\sum_{k \in \mathcal{C}} X_{ik}^{g} \right) + \sum_{rs} \gamma_{rs} C_{i1}^{r} C_{i1}^{s} + \epsilon_{i}.$$
(2)

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.

 \rightarrow Second scenario on 25 genes, one setting :

• different genes for main and interaction effects.

Simulations results - First scenario on 6 genes



- → 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

Simulations results - First scenario on 6 genes



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, $r^2 = 0.7$

Sommaire

1 GWAS and Block of linkage desequilibrium

- Genome Wide Association Studies
- Blocks of linkage desiguilibrium
- Hierachical Clustering with Adjacency Constraints
- How to improve?
- Some computation times

2 Epistasis

- The Gene-Gene Eigen Epistasis Modeling approach
- Simulations

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?

41 / 45

Known genes

Table I Summary of ankylosing spondylitis-susceptibility genes

identified b	ferrome-mide association scudies	ILOR	Interleukin 6 receptor	
RUNX3 IL23R	Runt-related transcription factor 3 Interleukin 23 receptor	FCGR2A	Fc fragment of immunoglobulin G, low-affinity IIa, receptor (CD32)	
ILI 2RB2	Interleukin 12 receptor, 82	UBE2E3	Ubiquitin-conjugating enzyme E2E 3	
GRP25	G-protein-coupled receptor 25	GPR35	G-protein-coupled receptor 35	
KIF21B	Kinesin family member 21B	NKX2-3	NK2 homeobox 3	
PTGER4	Prostaglandin E receptor 4 (subtype EP.)	ZMIZI	Zinc finger, MIZ type-containing I	
ERAPI	Endoplasmic reticulum aminopeptidase I	SH2B3	Src homology 2B adaptor protein 3	
ERAP2	Endoplasmic reticulum aminopeptidase 2	GPR65	G-protein-coupled receptor 65	
LNPEP	Leucyl/cystinyl aminopeptidase	IL27	Interleukin 27	
ILI 2B	Interleukin 12B	SULTIAI	Sulfotransferase family cytosolic IA	
CARD9	Caspase recruitment-domain family member 9	TYK2	Tyrosine kinase 2	
LTβR	Lymphotoxin β-receptor (TNFR superfamily, member 3)	ICOSLG	Inducible T-cell costimulator ligand	
TNFRSFIA	Tumor-necrosis factor-receptor superfamily member IA	EOMES	Eomesodermin	
NPEPPS	Aminopeptidase puromycin-sensitive	IL7R	Interleukin 7 receptor	
TBKBPI	TNFR-associated factor family member-associated	BACH2	BTB and CNC homology I, basic leucine-zipper	
	nuclear factor-xB-binding kinase 1-binding protein		transcription-factor 2	
TBX21	T-box 21	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

\rightarrow 29 susceptibility genes identified by GWAS

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

Conclusions and perspectives

The G-GEE method :

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

Ankylosing Spondylitis :

 $\boldsymbol{\rightarrow}$ Identification of potential interactions to discuss with doctors

→ HLA-B effect

Perspectives :

- \rightarrow Explore new $f_{u}(X_{i}^{r}, X_{i}^{s})$ definition
- \rightarrow Additional simulations on larger data set
- → New applications on other data set

Thank you for your attention !



→ Adaptive-Ridge Cleaning

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