## **GBIO0009**

# Effects of confounding factors on linkage disequilibrium and GWAS

17.10.2018 Sandra Negro

#### PLAN

#### I. Introduction

- Reminder on genetic terms
- $\rightarrow$  genotype, haplotype, homozygote, etc.
- Reminder on the principal of GWAS and linkage disequilibrium

II. Confounding factors

- Effects on LD
- Effects on GWAS

III. Existing approaches to resolve the confounder effects

#### Introduction Reminder on genetic terms



								Genotype (forward)
$\left\{ \right.$	Forward strand	A	A	A	<b>T</b>		С	of a diploïde organism:
	Reverse strand	A	T	A	T	-	G	
								homozygote heterozygote
	Forward strand	<u> </u>	G	A	<b>T</b>	-	G	Δn <b>hanlotyne</b> is a
	Reverse strand	A	С	A	T	-	G	sequence of alleles:



#### Introduction Reminder on genetic terms

Example of the type of data in plants :

	SNP 1	SNP 2	SNP 3	SNP 4	SNP 5	SNP 6	SNP 7	SNP 8	SNP 9
Ind1	GG	GG	AG	TT	TT	TT	CC	AA	AA
Ind2	AA	GG	GG	TT			CC	AA	СС
Ind3	GG	CC		TT	CC	CC	CC	TT	AA
Ind4		CG	AA	GG	TT	CC	TT		AA

Rate of heterozygosity in human ~= 95 % Rate of heterozygosity in inbred lines (plants) ~= 0.3 %

•

Alleles of this bi-allelic marker is « G » and « A »

#### NB :

indels (=insertion-deletion) can be found at certain SNP postions, but it remains a minority Rarely, SNPs can be tri-allelic or more

In general SNPs are bi-allelic and co-dominant

Introduction Reminder on the principal of GWAS

Quantitative Trait Loci (QTL) = genomic regions involved in quantitative effects on the phenotype



GWAS approach rely on the fact that molecular markers can capture QTL effect thanks to linkage disequilibrium

Linkage Disequilibrium (LD) = non random association between alleles of different loci

Two markers in LD B Ind 1 Α Β Ind 2 Ind 3 B Α Ind 4 b Α Ind 5 b a Ind 6 Ind 7 b a Ind 8 b a Ind 9 B a Ind 10

Considering these 2 SNPs, the possible haplotypes are: AB, Ab, aB and ab

- Linkage equilibrium :

Haplotype frequency = product of the corresponding allelic frequency

- Linkage disequilibrium refers to the deviation from this equilibrium :

$$\Rightarrow D_{AB} = p_{AB} - p_A p_B$$
$$\Rightarrow r^2 = \frac{D_{AB}^2}{p_A (1 - p_A) p_B (1 - p_B)}$$

One of the LD estimator

#### Introduction Reminder on Linkage disequilirium

The GWAS approach rely on the fact that SNPs can capture QTL effect thanks to linkage diseguilibrium



LD can be due to:

- 1) physical linkage
- 2) genetic structure in the population
- 3) kinship between individuals

In genome-wide association studies (GWAS), we are only interested by physical linkage, Population structure and kinship have to be controlled to prevent false positives



## Introduction Reminder on the principal of GWAS and linkage disequilibrium



Is there an association?:

If there is an association between the marker and the phenotype, then the marker is correlated (in LD) to a QTL



#### Confounding factors Relatedness



Two estimators (Identity-By-Descent, IBD and Identity-By-State, IBS) can be used to estimate relatedness amongst several individuals.



#### Confounding factors Population structure

#### Example of population struture in maize using ADMIXTURE

6 groups identified by Admixture (6 ancestral fractions per ind.) It is important to understand the organization of genetic diversity within a panel used in association genetics to define statistical models, to analyze the relationship between genetic polymorphism and the variation of traits, and to define the density of markers that is suitable to localize causal polymorphisms.





### **Confounding factors**



#### PCoA on the distance matrix (IBD, PANZEA) using 6 groups identified by admixture

1st axis: 8.9 % of the variance

- The two first PCoA axes explain an important part of the variance (15%)
- Major genetic groups are visualized
- Genetic groups and PCoA estimates are in accordance with breeder's knowledge

#### Confounding factors Effects on the LD and GWAS

The presence of ind. from different pop with different genetic origins within a sample can produce **LD between unlinked** loci, simply because of differences of allelic frequencies  $\rightarrow$  Hence, such structured sample can lead to a bias estimate of LD  $\Rightarrow$  which may increase the rate of false positive in GWAS  $\Rightarrow$  which could lead to inappropriate choice of marker density and therefore to a decreased power

A biased estimate of LD is also obtenained when genotyped ind. are not independent

III All the analyses using r<sup>2</sup> rely on the assumption that the extent of r<sup>2</sup> around the causal polymporphism depends only a drift-recombination process in a random mating pop without selection III  $\rightarrow$  But it is not always the case in real life data

#### Confounding factors Existing approaches to resolve the confounder effects

Structure and relatedness can be corrected for in diversity analyses and association mapping (GWAS) to take into account the non-indepence of loci due to both population differenciation and uneven levels of relatedness.

 $\Rightarrow$  Read paper provided (Mangin et al. 2012) to have details on the new adjusted r<sup>2</sup>

Real life data and simulation analyses highlighted:

1) in a two-population structured sample; r<sup>2</sup> bias increases with the differentiation of loci and with the decrease of LD

2) in a highly related sample; r<sup>2</sup> overestimate the true LD value

#### Confounding factors genome-wide linkage desequilibrium between all loci within and between chromosomes : r<sup>2</sup> versus adjusted r<sup>2</sup>

1.0

- 0.8

0.6

0.4

0.2



#### Confounding factors Existing approaches to resolve the confounder effects

#### Statistical approaches to resolve the confounder effects :

Four statistical models can be tested to limit false positive

```
Model 1: without correction; Y \sim X\beta + E
```

```
Model 2: takes into account group structure; Y \sim Qs + X\beta + E
```

Model 3: takes into account kinship between individuals;  $Y \sim X\beta + Ku + E$ 

Model 4: takes into account both group structure and kinship;  $Y \sim Qs + X\beta + Ku + E$ 

Home-made scripts or different informatic tools (EMMAX, ASReml, FASTLMM, etc...) can incorporate these models.