

Problems with measuring LD

All measures of LD depend on the allele frequencies of the pair of SNPs measured.

In different parts of the genome SNPs have different allele frequencies

'measure ?'

Dependence of the alleles frequencies of LD measures is a real bottleneck!

→ Doing better, more robust with
→ D:

$$P_{11}, P_{12}, P_{21}, P_{22} \geq 0$$

$$\begin{aligned} - P_{11} \geq 0 &\Rightarrow D = P_{11} - p_1 q_1, \\ , P_{22} \geq 0 &\Rightarrow D \geq -p_2 q_2 \end{aligned} \quad \left. \begin{array}{l} \\ D \geq -p_1 q_1 \end{array} \right\} \min$$

$$D_{\min} = \max \{-p_1 q_1, -p_2 q_2\}$$

$$\begin{aligned} P_{12} \geq 0 &\Rightarrow D \leq p_1 q_2, \\ P_{21} \geq 0 &\Rightarrow D \leq p_2 q_1 \end{aligned} \quad \left. \begin{array}{l} \\ D \leq p_2 q_1 \end{array} \right\} \max$$

$$D_{\max} = \min \{p_1 q_2, p_2 q_1\}$$

We introduce : D'

$$D' = \begin{cases} \frac{D}{D_{\max}} & \text{if } D > 0 \\ \frac{D}{D_{\min}} & \text{if } D < 0 \end{cases}$$

Another measure of LD: σ^2

$$\sigma^2 = \frac{D^2}{P_1 P_2 q_1 q_2}$$

correlation coefficient measure

The two measures σ^2 and D' are used together: they have desirable complementary type of information they measure

Both D' and σ^2 are having values in $[0, 1]$

r^2 has the property that
the intersection values

$$0 < \rho < 1$$

can be interpreted !!

HAPMAP PROJECT

Based on this "interpretability"

D' does not have the
interpretability property

The axioms / desiderata from the
decades of ~~top~~ research in
Human Genetics:

Ax 0. Independence of allele
frequencies

Ax1. The Curse of the pairwise:
we want LD measures
to be extensible to
multi markers/SNP in
a conservative way

Ax3. Interpretability of
intermediate values

BIG OPEN PB: To find
an LD measure satisfying
the above three basic/other
axioms simultaneously.

Interpretability:

- ① If two SNPs are in equal
frequency and create

between them only

2 haplotypes, then

$$r^2 = 1, D' = 1$$

A C

$$\begin{matrix} A & C \end{matrix} \quad r^2 = 1$$

$$\begin{matrix} T & G \end{matrix} \quad D' = 1$$

$$\begin{matrix} T & G \end{matrix}$$

-
- ② If the two SNPs are in unequal frequency, and at least three haplotypes are present then:

Exactly three haplotypes are present, then $D' = 1, r^2 < 1$

A C

$$\begin{matrix} A & G \end{matrix} \quad r^2 < 1$$

T G

$$\begin{matrix} T & G \end{matrix} \quad D' = 1$$

T G

③ when all four haplotypes
are present then $D' \leq 1$, $\pi^2 \leq 1$

$$\begin{array}{ll} AC & \pi^2 \leq 1 \\ TC & \\ AG & D' \leq 1 \\ TG & \end{array}$$

Population

- Sample of individuals : a random sample from the population we are studying
- Computing LD for the sample we think of it as an approximation of LD for the population.

Testing for LD

is a Statistical process

LD is an empirical pattern
Statistical Significance of LD

Are the two sites significantly associated? I.e. whether the LD detected is statistically significant.

The Null Hypothesis:

H_0 : they are in Linkage Equilibrium

$$P_{ij} = P_i q_j$$

We use the counts of each of

the alleles occurring together

There are four haplotypes
at the two sites/SNP.

The Null hypo. states
that the two sites are
independent.

The natural statistical
test is a test for independence
of 2×2 contingency table
with elements n_{ij} , where:

n_{ij} = the number of haplo-
carrying allele i at
position 1 and
allele j at position 2

$$\sum_{j=1}^2 \sum_{i=1}^2 n_{ij} = n = \text{total sample size}$$

The most precise test is

Fisher's exact test of independence.

However, Fisher's exact test is complicated to calculate and when n is large, we may calculate instead

$$X = \sum \frac{(n_{ij} - e_{ij})^2}{e_{ij}}$$

e_{ij} = the estimated ~~expected~~ of n_{ij} under the Null Hypothesis

i.e. assumption of no association

If e_{ij} are sufficiently large (> 5) then χ^2 is approximately χ^2 distributed with one degree of freedom.

An example: NM blood group in a British population

A sample of 1000 people with genotype frequencies as follows:

Gene NM

Two alleles: N, M

genotypes	$\left\{ \begin{array}{l} MM \\ MN \\ NN \end{array} \right.$	298 "Blood group M"
		489 "
		- - - MN"
$\overline{sum = 1000}$		

Estimate allele frequencies:

M: ♂ , N: ♀

$$\hat{p} = \frac{1085}{2000} = 0.54 \quad \left. \right\} 1085 = 2 \times 298 + 489$$

$$\hat{q} = \frac{915}{2000} = 0.45$$

Gene S, A

$$S: \hat{q}_1 = 0.30$$

$$A: \hat{q}_2 = 0.69$$

SS	99
SA	418
AA	483

2 Genes : haplotypes

Haplotypes	Observed	Expected
M S	474	$0.54 \times 0.30 \times 2000 = 334$
M S	611	
N S	142	$0.54 \times 0.69 \times 2000 = 750$
N S	773	
	281	
	633	

$$\chi^2 = 184.7, \text{ with}$$

$$df = 1$$

$$P_1 P_2 P_3 P_4 : 4 - 3 = 1$$

$$-1, -1 \\ P_1 + q_1 = 1) \quad P_2 + q_2 = 1$$

P-value assoc with $\chi^2 = 184.7$

very small P-value < 0.001

Reject Null Hypothesis

There is LD.

How much LD:

To quantify the amount of LD, we use D'.

$$P_{11}, P_{12}, P_{21}, P_{22}$$

$$MS: \hat{P}_{11} = \frac{474}{2000}$$

$$MS: \hat{P}_{12} = \frac{611}{2000}$$

$$NS: \hat{P}_{21} = \frac{142}{2000}$$

$$NS: \hat{P}_{22} = \frac{773}{2000}$$

$$\hat{D} = \hat{P}_{11} \hat{P}_{22} - \hat{P}_{12} \hat{P}_{21} = 0.07$$

$$\hat{D}_{MAX} = \min \left\{ \begin{array}{l} \hat{P}_1 \hat{q}_2 = 0.38 \\ \hat{P}_2 \hat{q}_1 = 0.14 \end{array} \right\} = 0.14$$

$$\hat{D}' = \frac{\hat{D}}{\hat{D}_{MAX}} = \frac{0.07}{0.14} = 50\% \text{ of}$$

if theoretical max