

Linkage Disequilibrium

Typed up notes by Derek Aguiar

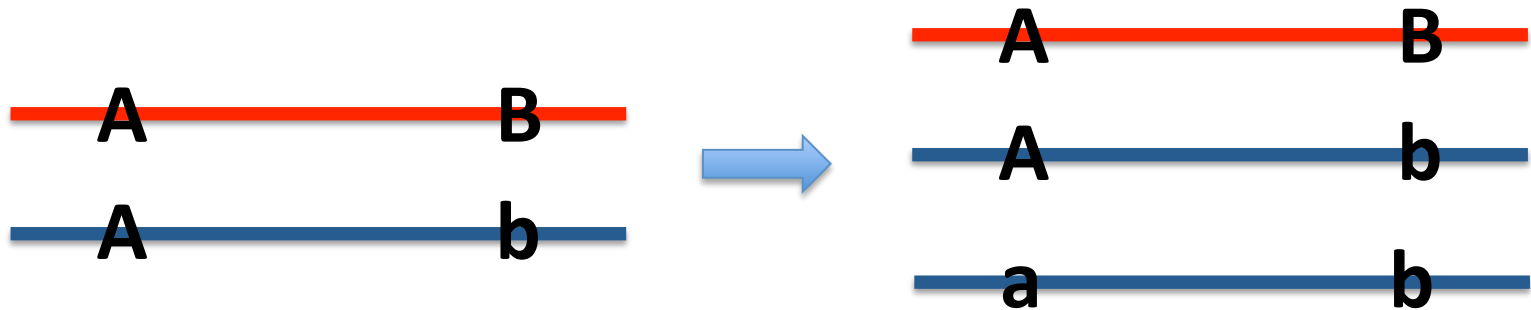
Linkage Disequilibrium

- Abbreviated LD, represents a statistical linkage between alleles on the same chromosome
- Recombination breaks linkage disequilibrium
- For example, we have one ancestral haplotype that gains a mutation, so now two haplotypes exist in the population.



Linkage disequilibrium

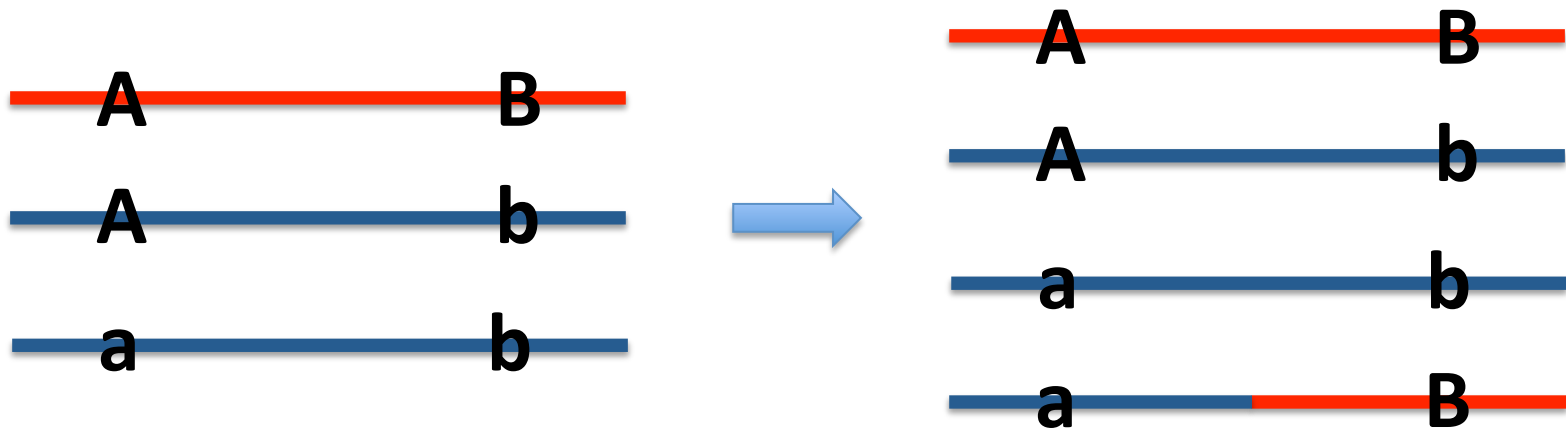
If another site is mutated, three haplotypes exist



But the *a* allele completely determines the *b* allele so the sites are still in strong linkage disequilibrium

Linkage disequilibrium

- Because the probability of mutating either the A/a or B/b variant is very low, linkage disequilibrium is more often broken by recombination



An example in drosophila

Three alleles (genes or variants)

Haplotype	Count
ABC	264
ABc	13
AbC	29
Abc	8
aBC	152
aBc	7
abC	15
abc	1

In the case of two alleles (pairwise measurements)

Haplotype	Frequency
AB	p1q1
Ab	p1q2
aB	p2q1
ab	p2q2

Allele	Frequency
A	p1
a	p2
B	q1
b	q2

Linkage equilibrium

Allele	Frequency
A	p1
a	p2
B	q1
b	q2

Haplotype	Frequency
AB	p1q1
Ab	p1q2
aB	p2q1
ab	p2q2

Real data frequency
P11=frequency of haplotype AB in population
P12=frequency of haplotype Ab in population
P21=frequency of haplotype aB in population
P22=frequency of haplotype ab in population

So we can measure the frequency of haplotypes and subtract of the expected frequency given from the allele frequencies.

This is the fundamental concept for linkage disequilibrium, D

$$D = P_{11} - p_1q_1$$

If $D=0$, we have linkage equilibrium,

If $LD \neq 0$, we have some level of linkage disequilibrium

D

Linkage disequilibrium D

$$P_{11} = p_1q_1 + D$$

$$P_{12} = p_1q_2 - D$$

$$P_{21} = p_2q_1 - D$$

$$P_{22} = p_2q_2 + D$$

Lemma

$$D = P_{11}P_{22} - P_{12}P_{21}$$

Proof

$$P_{11}P_{22} = (p_1q_1 + D)(p_2q_2 + D) = p_1q_1p_2q_2 + p_1q_1D + p_2q_2D + D^2$$

$$P_{12}P_{21} = (p_1q_2 - D)(p_2q_1 - D) = p_1p_2q_1q_2 - Dp_2q_1 - Dp_1q_2 + D^2$$

$$P_{11}P_{22} - P_{12}P_{21} = D$$

So, D is indeed a measure that takes into account all of the haplotype frequencies.

$$D'$$

$$D' = D/D_{\max} \text{ when } D > 0 \text{ and} \\ D/D_{\min} \text{ when } D < 0$$

where

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

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

$$D'$$

- We know $P_{11}, P_{12}, P_{21}, P_{22} \geq 0$
- If $D < 0$ then D_{\min} is negative and D' is positive.
- If $D > 0$ then D_{\max} is positive and D' is positive.
- If $D < 0$ then $P_{12}P_{21} > P_{11}P_{22}$
- If $D > 0$ then $P_{11}P_{22} > P_{12}P_{21}$

$$r^2$$

- $r^2 = D^2 / p_1 p_2 q_1 q_2$
- $r = D / \sqrt{p_1 p_2 q_1 q_2}$ is the correlation coefficient between pairs of variants

Statistical testing for LD

- Given a pair of sites, are they significantly associated?
- Hypothesis testing
 - LD detected is statistically significant
- Null hypothesis
 - Linkage equilibrium, $P_{ij}=p_iq_j$, the frequency of the haplotype is the product of the allele frequencies
- Alternative hypothesis

Hypothesis testing

- 2x2 contingency table, n_{ij} = the number of alleles with i in the first position and j in the second
- Sum of the $n_{ij} = n$ over all i 's and j 's
- Apply Fisher's exact test or X^2 test
 - If n is large Fisher is difficult to compute but X^2 is a good approximation
 - $X = \sum [(n_{ij} - p_{ij})^2 / p_{ij}]$ where this is approximately X^2 with 1 degree of freedom distributed for large n

Example in drosophila

- Computing χ^2 for two sites (ignore C/c allele)

Haplotype	Count	Frequencies
AB	277	0.5656
Ab	37	0.0757
aB	159	0.3251
ab	16	0.03272

– $D=0.0061$ and $r^2=0.001659$

Example in drosophila

- $\chi^2 = r^2 * n = 0.81$
- χ^2 with 1 degree of freedom gives a probability of 0.37 which is not significant so we cannot reject the null hypothesis

Looking at three sites

- $X=13.0$
- X^2 prob of 0.0003
- $D_{\max}=\{0.053,0.102\}$
- $D'=0.012/0.053=.226$