Estimating experiment-wise error rates by permutation

ESP 2009 Bertram Müller-Myhsok

Multipe testing problem in Genome-wide Analysis

- Many tests performed
- 100,000 SNPs
- 300,000 SNPs
- 500,000 SNPs
- 1,000,000 SNPs

Multiple testing in GWA cont'd

- · Often several phenotypes
- · Often several genetic models
- · Example:
 - 10 phenotypes
 - -4 models
 - 500,000 SNPs

Multiple testing in GWA cont'd

- Number of tests then:
 10 * 4 * 500,000 = 20 million = 2e+07
- Necessary significance level for experiment-wise p of 0.05 using Bonferroni: - 0.05 / 2e+07 = 2.5e-09

But:

- · Tests are not necessarily independent
 - LD between SNPs
 - Correlation between phenotypes
 - Correlation between genotypes
- The effective number of tests is less
 than the number of tests performed

LD: From Genotypes to Haplotypes

	BB	Bb	bb
AA	AABB	AABb	Aabb
Aa	AaBB	AaBb	Aabb
аа	aaBB	aaBb	aabb

Joint genotype AaBb

A++a or A++a B++b b++B

Haplotypes

- For unrelated individuals Haplotypes can be reconstructed with 100% accuracy if:
 - All loci are homozygous
 - No more than 1 locus is heterozygous
- If more 2 or more loci are heterozygous then
 - The haplotypes are ambiguous
 - However usually one set (pair) of
 - haplotypes is more likely

Haplotype Reconstruction

- Haplotype construction is not robust over large genetic distances
- · Beware of using "best" haplotype configuration
 - It is important to take account of uncertainty in phase assignment if haplotypes are to be used in subsequent analysis
- The ability to estimate accurate haplotype assignments is dependent on the size of data set.
 - Generally smaller data sets give less accurate results









Another measure of the strength of association is r.

 $r=D/(p_1p_2q_1q_2)^{1/2}$

X²=r²N (1 df)

For this example $r = 0.06/(0.2 \cdot 0.8 \cdot 0.4 \cdot 0.6)^{1/2} = 0.3062$

X²=0.093746•100=9.376 p=0.0022

LD

• If one haplotype is not observed |D'|=1

- Complete LD There is not a 100% correlation between the allele at one locus and the allele at the second locus

If two haplotypes are not observed r²=1.0

- Perfect LD

- There is 100% correlation between the allele at one locus and the allele at the second locus
 If r²=1.0 for two loci genotyping one locus provides as much information as genotyping both loci

LD

- · D' provides information on historic recombination events
- r² provides information on the correlation of two loci - Better measure for association studies



Linkage Disequilibrium

- With random mating (assumptions: large population with no mutation, migration or selection) linkage equilibrium is eventually obtained
- · The rate of decay will depend on the recombination fraction between loci.
- The greater the rate of recombination the quicker the decay.





Three examples

Ex2: Two SNPs, SNP1 and SNP2, 1000 people genotyped

```
        A/A
        A/B
        B/B

        A/A
        375
        0
        0

        A/B
        176
        214
        106

        B/B
        0
        0
        129
```

Measures of LD: D' = 0.7357 r2 = 0.4650

Three examples						
Ex3: Two SNPs, SNP1 and SNP2, 1000 people genotyped						
A/A A/B B/B A/A 551 0 0 A/B 0 214 0 B/B 0 0 235						
Measures of LD: D' = 1.0000 r2 = 1.0000						



Permutation in principle							
Predictor variables, e.g. SNPsOutcome variable, e.g. affection status							
Affection status 1 0 0 1 1 0	SNP1 A/B A/A B/B B/B A/B	SNP2 A/B A/B B/B A/A A/B A/A					



Original data and 9 replicates (permutations)									
1	1	1	1	1	1	0	1	0	0
0	0	1	1	0	0	1	1	0	1
0	1	0	1	0	1	1	1	1	1
1	0	1	0	1	0	0	1	1	1
1	1	0	0	1	1	1	0	1	0
1	0	1	1	1	1	1	0	0	1
0	1	0	0	0	0	0	0	1	0



- Compute and store minimum of replicate-wise p-values
 -> test distribution of test statistic
 Compare p-value found against
- distribution of minimum p-values (equivalent to maximum of test statistic) – > correct for multiple testing















Quantile of distributions						
1%	5%	10%				
0.006	0.027	0.055	ex1			
0.006	0.029	0.060	ex2			
0.009	0.049	0.100	ex3			



Papers on the way

Am. J. Hum. Genet. 74:765–769, 2004

A Simple Correction for Multiple Testing for Single-Nucleotide Polymorphisms in Linkage Disequilibrium with Each Other

Dale R. Nyholt

Papers on the way

Am. J. Hum. Genet. 75:424-435, 2004

Efficient Computation of Significance Levels for Multiple Associations in Large Studies of Correlated Data, Including Genomewide Association Studies

Frank Dudbridge and Bobby P. C. Koeleman

Papers on the way

Hum Hered 2005;60:19-25

Evaluation of Nyholt's Procedure for Multiple Testing Correction

Daria Salyakina Shaun R. Seaman Brian L. Browning Frank Dudbridge Bertram Müller-Myhsok Genetic Epidemiology 32: 381-385 (2008)

Brief Report

Estimation of the Multiple Testing Burden for Genomewide Association Studies of Nearly All Common Variants

Itsik Pe'er,¹ Roman Yelensky,^{2–4} David Altshuler,^{23,3–7} and Mark J. Daly^{23,89} ¹Opartment of Computer Science, Colombia University, New York ²Centre for Human Geneir Research, Massachustetti Geneir Heingli, Bastu, Massachustts ¹Opartment of Madeault Bidegu, Massachustte Geneir Heingli, Bastu, Massachustts ¹Pitterna VII. Diesisier of Hehnik Science and Technology, Camirige, Massachustts ¹Diabets Unit, Massachustts Geneir Heingli, Bastu, Massachustts ¹Diabets Unit, Massachustts Geneira Heingli, Massachustts ²Dapartment of Geneiric, Harrard Maltal Schol, Seston, Massachusetts ¹Dapartment of Malcine, Harrard Maltal Schol, Seston, Massachusetts







