

GWA analysis of family studies

Yurii Aulchenko
Erasmus MC Rotterdam
27.08.2009

ESP29, 27.08.2009

Yurii Aulchenko

Outline

- Reasons for genetic association
- (Genome-Wide) Association analysis in pedigrees
- Conclusions

ESP29, 27.08.2009

Yurii Aulchenko

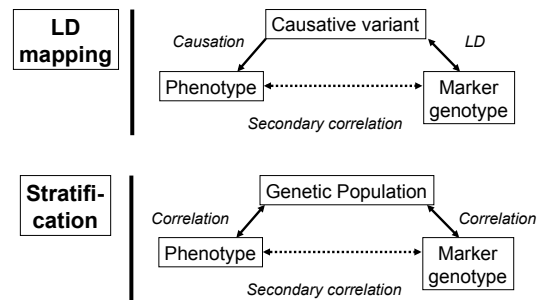
Benefits of studies of genetically isolated populations / family data

- Power to detect genes with rare variation
 - Variants with $0.0001 < \text{MAF} < 0.01$
 - May become more frequent compared to general population, providing high power
 - Replication issue: consortia of genetically isolated populations
 - “Singleton variants”
 - ... in outbred
 - Isolated populations: few copies in close relatives
- Ability to investigate complex genetic models
 - Parent-of-origin effects: imprinting, maternal
 - Parent x Child genotype interactions

ESP29, 27.08.2009

Yurii Aulchenko

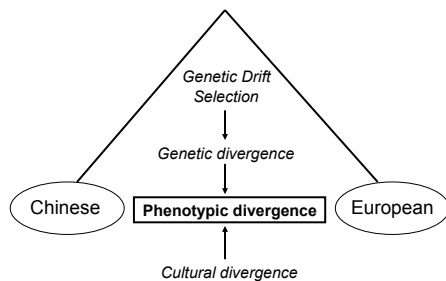
Confounding in genetic studies



ESP29, 27.08.2009

Yurii Aulchenko

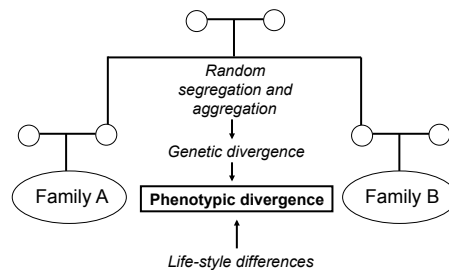
Genetic origin is a major confounder



ESP29, 27.08.2009

Yurii Aulchenko

Pedigree is a major confounder



ESP29, 27.08.2009

Yurii Aulchenko

The Importance of Genealogy in Determining Genetic Associations with Complex TraitsDINA L. NEWMAN,¹ MARK ARNEY,^{1,2}
MARY SARA McPECK,^{1,2} CAROLE OBER,¹
AND NANCY J. COX¹

- >750 Hutterites. Association tested between 3 quantitative traits (IgE level, LDL, BMI) and >500 markers with and without modeling the relatedness

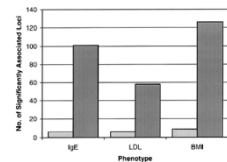


Figure 2 Number of significantly associated ($P < .01$) loci when pedigree structure is included (lighter bars) and when pedigree structure is not included (darker bars).

→ High level of false positive signals

Ignore relationship, apply GC?Vector of quantitative phenotype Y

$$Y = \mu + Bg + e$$

Score test for association:

$$T^2 = \frac{(g \cdot Y)^2}{g \cdot g} \sim \chi^2_1 \cdot \lambda$$

Lambda is estimated using genomic control (GC):

$$\lambda = \frac{\text{Median}(T_1^2, T_2^2, T_3^2, \dots, T_M^2)}{0.455}, \quad \lambda \geq 1$$

Computation time ~ N

ESP29, 27.08.2009

Yuri Aulchenko

When GC does not work (well)?

When stratification is large (say, $\lambda_{1000} > 1.1$) other, more powerful methods are to be used

GC assumes that stratification acts in the same manner across all loci

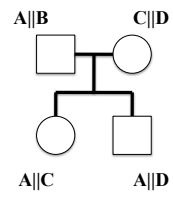
- This is not true for loci differentiated between population e.g. because of selection
- Such loci will still be falsely detected after GC correction

ESP29, 26.08.2009

Yuri Aulchenko

Relationship coefficient

- Two chromosomal regions are called Identical-by-descent (IBD) if they are copies of the same ancestral chromosomal region



- For a pair of people **relationship coefficient** r is the (expected) proportion of genome shared IBD
- Kinship coefficient = $\frac{1}{2}$ relationship coefficient

ESP29, 27.08.2009

Yuri Aulchenko

Correlation between relatives

r	relationship
0.5 (1/2)	parent-offspring
0.25 (1/4)	grandparent-grandchild
0.125 (1/8)	great grandparent-great grandchild
1	identical twins
0.5 (1/2)	full siblings
0.25 (1/4)	half siblings
0.125 (1/8)	first cousins
0.03125 (1/32)	second cousins ^[2]

- If a trait is controlled by genes only (heritability, $h^2=1$)
 - Identical twins would have exactly the same trait value ($\text{cor} = 1$)
 - Correlation between the phenotypes of sibs would be 0.5
 - For arbitrary relatives correlation would be r

If the proportion of trait's variance explained by genes is h^2

- Correlation between phenotypes of identical twins would be h^2
- Correlation between the phenotypes of sibs would be $0.5 h^2$
- For arbitrary relatives correlation would be $r h^2$

ESP29, 27.08.2009

Yuri Aulchenko

Mixed (polygenic) model

Linear Mixed Model (LMM) where the vector of quantitative phenotype Y is modeled as

$$Y = \mu + Bg + G + e$$

g : genotype indicator vector g_i in $\{0, 1, 2\}$

B : additive affect of the allele

e : is random residual effect $\sim \text{MVN}(\mathbf{0}, I\sigma_e^2)$

I : identity matrix

G : is random polygenic effect $\sim \text{MVN}(\mathbf{0}, \Phi \sigma_G^2)$

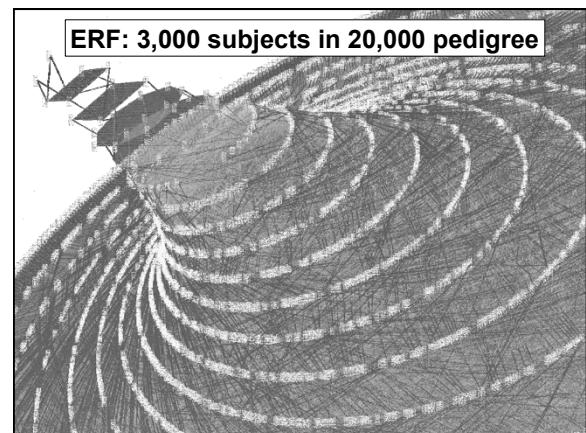
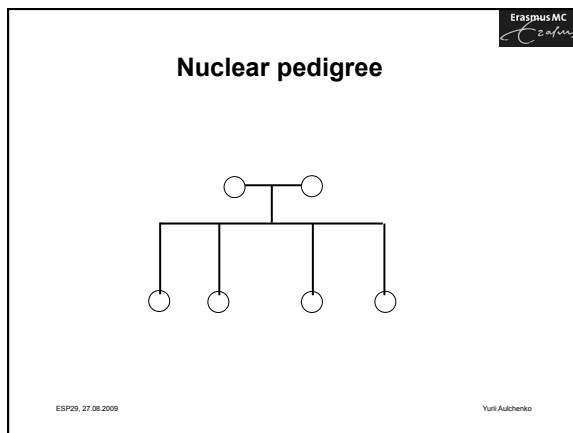
Φ : relationship matrix

Maximum Likelihood (ML) or Restricted ML (REML)

- Software packages available (SOLAR, MERLIN, QTDT, ASReml)

ESP29, 27.08.2009

Yuri Aulchenko



Erasmus MC
Erasmus

Analysis of large complex pedigrees

- Time required for GWA scan (2.5 million SNPs in 3,000 people)
- ML: 20 minutes per test => 95 years
- REML: 3-5 times faster

ESP29, 27.08.2009

Yuri Aulchenko

Erasmus MC
Erasmus

FASTA

Family Score Test for Association
Based on the mixed model $Y = \mu + Bg + G + e$

FASTA test for association:

- Estimate polygenic model $Y = \mu + G + e$
- Compute FASTA test

$$T^2 = \frac{g \cdot (\Phi \hat{\sigma}_G^2 + I \hat{\sigma}_e^2)^{-1} \cdot Y}{g \cdot (\Phi \hat{\sigma}_G^2 + I \hat{\sigma}_e^2)^{-1} \cdot g} \sim \chi_1^2$$

- Apply GC afterwards if $\lambda > 1$

Computation time $\sim N^2 + N$; no permutation testing

ESP29, 27.08.2009

Yuri Aulchenko
Chen & Abecasis, 2007

Erasmus MC
Erasmus

GRAMMAS

GW Rapid Association using Mixed Model And Score test
Based on the mixed model $Y = \mu + Bg + G + e$

GRAMMAS test for association:

- Estimate polygenic model $Y = \mu + G + e$
- Compute environmental residuals $Y^* = Y - (\hat{\mu} + \hat{G}) = \hat{e}$
- Runs score test on residuals

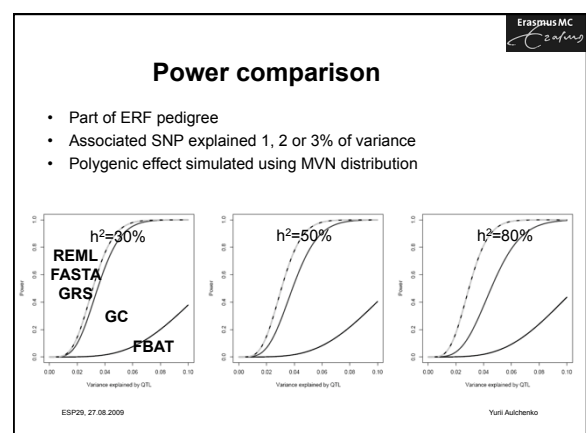
$$T^2 = \frac{(g \cdot Y^*)}{g \cdot g} = \frac{(g \cdot \hat{\sigma}_e^2 \cdot (\Phi \hat{\sigma}_G^2 + I \hat{\sigma}_e^2)^{-1} \cdot Y)}{g \cdot g}$$

- Apply GC (λ expected to be < 1)

Computation time $\sim N$; permutation testing possible

ESP29, 27.08.2009

Yuri Aulchenko
Aulchenko et al, 2007; Amin et al., 2007



Relationship between genomes

The estimate of kinship between i and j may be obtained from genomic data:

$$f_{ij} = \frac{1}{n} \sum_{k=1}^n \frac{(g_{ik} - p_k)(g_{jk} - p_k)}{p_k(1 - p_k)}$$

g_{ik} is the genotype (0, 0.5, 1) of the i -th person at k -th SNP

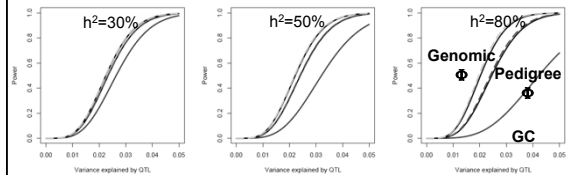
p_k is the frequency of "1" allele

ESP29, 27.08.2009

Yuri Aulchenko

Genomic vs. Pedigree kinship

- 1,400 ERF people genotyped for 6K Illumina Array
- Trait values simulated based on observed genotypes
- Associated SNPs explained from 0.3 to 4% of variance

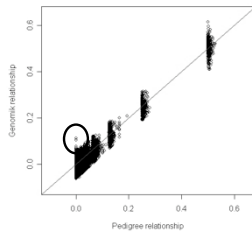


ESP29, 27.08.2009

Yuri Aulchenko

Genomic Φ is better than pedigree Φ

- Pedigree is not guaranteed to be correct
 - Missing links => increased type 1 error
- Pedigree relationship coefficient is the expected proportion of genome shared
 - Genomic relationship may better estimate true sharing



ESP29, 27.08.2009

Yuri Aulchenko

Testing association with binary trait in a general pedigree

- Bourgain et al., AJHG, 2003
 - Comparison of allele frequency between cases and controls
 - $\chi^2_{\text{corr}} \sim$ Genomic Control
 - CC-QLS (case-control quasi-likelihood score test)
 - Frequency is estimated under the null, using BLUP
 - Score test is performed
- William Astle, David Balding
 - Faster, more flexible methods based on Mixed Model

ESP29, 27.08.2009

Yuri Aulchenko