

GWA analysis of family studies

Yurii Aulchenko Erasmus MC Rotterdam 27.08.2009

ESP29, 27.08.2009

Outline

Erasmus MC

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

Erasmus MC 2 afms

Benefits of studies of genetically isolated populations / family data

- Power to detect genes with rare variation
 - Variants with 0.0001 < MAF < 0.01</p>
 - 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

Confounding in genetic studies







ESP29, 27.08.2009



Pedigree is a major confounder



Am. J. Hum. Genet. 69:1146-1148, 2001

The Importance of Genealogy in Determining Genetic Associations with Complex Traits Erasmus MC

DINA L. NEWMAN,¹ MARK ABNEY,^{1,2} MARY SARA MCPEEK,^{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



→ High level of false positive signals

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*).



Ignore relationship, apply GC?

Vector of quantitative phenotype Y

 $Y = \mu + B g + e$

Score test for association:

$$T^{2} = \frac{(g \cdot Y)^{2}}{g \cdot g} \sim \chi_{1}^{2} \cdot \lambda$$

Lambda is estimated using genomic control (GC):

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

Computation time ~ N

ESP29, 27.08.2009



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



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 = ½ relationship coefficient





Correlation between relatives

r	relationship
0.5 (1/2)	parent-offspring
0.25 (¼)	grandparent-grandchild
0.125 (1/8)	great grandparent-great grandchi
1	identical twins
0.5 (1/2)	full siblings
0.25 (¼)	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 (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 h²
- Correlation between the phenotypes of sibs would be 0.5 h²
- For arbitrary relatives correlation would be r h²
 Yurii 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 ~ MVN($\mathbf{0}, I\sigma_e^2$) *I:* identity matrix *G:* is random polygenic effect ~ 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



Nuclear pedigree



ERF: 3,000 subjects in 20,000 pedigree 187.5



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

Erasmus MC

FASTA

<u>FA</u>mily <u>S</u>core <u>Test</u> for <u>A</u>ssociation Based on the mixed model $Y = \mu + Bg + G + e$

FASTA test for association:

(a) Estimate polygenic model $Y = \mu + G + e$

(b) Compute FASTA test

$$T^{2} = \frac{\left(g \cdot \left(\Phi \hat{\sigma}_{G}^{2} + \mathrm{I} \hat{\sigma}_{e}^{2}\right)^{1} \cdot Y\right)}{g \cdot \left(\Phi \hat{\sigma}_{G}^{2} + \mathrm{I} \hat{\sigma}_{e}^{2}\right)^{1} \cdot g} \sim \chi_{1}^{2}$$

(c) Apply GC afterwards if $\lambda > 1$

Computation time ~ N²+N; no permutation testing

Chen & Abecasis, 2007

GRAMMAS

<u>GW Rapid Association using Mixed Model And Score test</u> Based on the mixed model $Y = \mu + Bg + G + e$ **GRAMMAS** test for association:

(a) Estimate polygenic model $Y = \mu + G + e$

(b) Compute environmental residuals $Y^* = Y - (\hat{\mu} + \hat{G}) = \hat{e}$

(c) Runs score test on residuals

$$T^{2} = \frac{\left(g \cdot Y^{*}\right)}{g \cdot g} = \frac{\left(g \cdot \hat{\sigma}_{e}^{2} \cdot \left(\Phi \hat{\sigma}_{G}^{2} + \mathrm{I} \hat{\sigma}_{e}^{2}\right)^{1} \cdot Y\right)}{g \cdot g}$$

(d) Apply GC (λ expected to be < 1)

Computation time ~ N; permutation testing possible

Yurii Aulchenko

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

ESP29. 27.08.2009



Power comparison

- Part of ERF pedigree
- Associated SNP explained 1, 2 or 3% of variance
- Polygenic effect simulated using MVN distribution



ESP29, 27.08.2009



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



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



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 <u>expected</u> proportion of genome shared
 - Genomic relationship may better estimate true sharing





Testing association with binary trait in a general pedigree

Bourgain et al., AJHG, 2003

- Comparison of allele frequency between cases and controls
- $-\chi^2_{corr}$ ~ 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