# ANALYSIS OF IMPUTED DATA

YURII ÁULCHENKO INDEPENDENT RESEARCHER & CONSULTANT YURII [DOT] AULCHENKO [AT] GMAIL [DOT] COM

### **ÅREAS OF EXPERTISE**

- (infra)Structure and project evaluation and advise
- Methodological advise (study design, planning of analyses, methods, software)
- Methods, algorithms and software development
- Data analyses
- Teaching and training

## ANALYSIS OF IMPUTED GENOTYPIC DATA IN GWAS

- Short review of standard methods
- Methods for analysis of imputed data
- Open questions and problems

### SINGE SNP ANALYSIS

- Analysis of each SNP in turn independent of others
- For each SNP, regression is performed, resulting in estimates of regression coefficients, their standard errors and the p-value

### LINEAR REGRESSION MODEL

The value of the trait in *i*-th individual is assumed to follow linear model

$$Y_i = m + b_g g_i + e_i$$

where *m* is intercept,  $g_i$  is the genotypic value (coded as 'B' allele dose - 0, 1 or 2), and  $e_i$  is random residual error











### SIGNIFICANCE TESTING

The estimate of  $b_g$  and its standard error  $s_g$  are computed using standard methods Under the null hypothesis, the test statistic  $T^2 = (b_g/s_g)^2$  is distributed as chisquared with 1 degree of freedom (1df)

For specific  $T^2$  we know the *p*-value: the probability that  $b_g$  deviates from zero purely by chance

### AS A RESULT...



### MULTIPLE TESTING

- Hence nominal single test p-value corresponding to experiment-wise type 1 error rate of 5% is << 0.05</li>
- Usual practice is to use a fixed threshold of 5x10<sup>-8</sup>
- Note: this threshold is defined for GWAS of *common variants* in a population of *European ancestry*

#### NON-ADDITIVE MODELS

... can be specified using linear model  $Y_i = m + b_1 I(g_i=1) + b_2 I(g_i=2) + e_i$ where  $I(g_i=k)$  is an indicator variable taking the value of 1 if  $g_i$  is equal to k and zero otherwise

### NON-ADDITIVE MODELS

In other words, the expected value of the trait for the genotype

- AA  $(g_i=0)$  is m
- AB ( $g_i=1$ ) is  $m + b_1$
- BB ( $g_i=2$ ) is  $m + b_2$

By varying m,  $b_1$  and  $b_2$ , any three genotypic means can be fit











### OTHER 1DF MODELS

Dominant B:  $b_1 = b_2$ Recessive B:  $b_1 = 0$ Over-dominant:  $b_2 = 0$ 

... and additive:  $b_2=2 b_1$ 

#### **INTERACTION MODELS**

The value of the trait in *i*-th individual is assumed to follow linear model

 $Y_i = m + b_f F_i + b_g g_i + b_{fg} F_i g_i + e_i$ 

where *m* is intercept,  $F_i$  is the value of some "factor",  $g_i$  is the genotypic value, and  $e_i$  is random residual error

### WHAT COULD "F" BE?

- An environment (gene-environment interaction)
- Indicator of transmitting parent (imprinting models)
- Other genotype (gene-gene)
- ... etc.

#### A Genome-Wide Screen for Interactions Reveals a New Locus on 4p15 Modifying the Effect of Waist-to-Hip Ratio on Total Cholesterol

Study name	Main effect	Interaction term
FINRISK	-0-	<u> </u>
HBCS		<b>_</b>
NFBC1966	0	<b>O</b>
YFS		<u> </u>
KORAF3		
KORAF4	<u> </u>	<b>0</b>
RS-I	0	<b>(</b> )
RS-II	0	
EUROSPAN	<b>O</b>	<b>©</b>
TWINSUK		
KORCULA	0	
Stage 1 combined		
		0
NTR		
NTR2		- <b>o</b>
EGCUT		
LIFELINES		<u> </u>
SORBS		
Genmets	<b></b>	- <u>o</u>
Stage 2 combined	O	0
CoLaus		(O)
EPIC cohort		
EPIC cases		•
Stage 3 combined		<b>O</b>
All combined	(•)	
Effect size -0.	3 -0.1 0 0.1	0.3 -0.3 -0.1 0 0.1 0.3
) Number of individuals	Effect estimate	Standard error

- A meta-analysis of genomewide association (GWA) data from 18 population-based cohorts with European ancestry (maximum N = 32,225).
- Eight further cohorts (N = 17,102) for replication
- SNP rs6448771
  - demonstrated genome-wide significant interaction with waist-to-hip-ratio (WHR) on total cholesterol (TC) with a combined *P*-value of  $4.79 \times 10^{-9}$

## ANALYSIS OF IMPUTED GENOTYPIC DATA IN GWAS

- Short review of standard methods
- Methods for analysis of imputed data
- Open questions and problems

### IMPUTED DATA

We can not tell the exact genotype, but can estimate posterior probability distribution:  $P(g) = \{p_{AA}, p_{AB}, p_{BB}\}$ 

Directly typed SNPs: either AA, AB or BB. The probability distribution is *degenerate* (e.g. {0,1,0} that is to say AB) For imputed SNPs, the distribution is not degenerate

### IMPUTING: GUESS THE "?"

С

T C T T G G A C T G T C C T G T C C T G G A C T G T C C T G

Α

G

Т

G

Reference Haplotypes

Sample Genotypes



Zheng et al., 2011

### IMPUTING: GUESS THE "?"



### HOW CAN WE ANALYZE IMPUTED DATA?

- Instead of genotypes, we have probabilities that certain person at certain locus has certain genotype
- ???

### **BEST GUESS**

- Best guess: take the genotype with maximal posterior probability and treat it as if it was true, directly typed
- Drawback: biased estimates, reduced power

### REGRESSION ONTO PROBABILITIES

#### Use the model

 $E[Y_i] = m + b_1 P(g_i=1) + b_2 P(g_i=2)$ 

Note this is very similar to model for directly typed SNPs, with probabilities used instead of indicator variables. Different genetic models can be formulated in the same way.

### MAXIMUM LIKELIHOOD BASED ON PROBABILITIES

Define individual likelihood as  $L_i = SUM_{gi=\{0,1,2\}} P(g_i) P(Y_i | g_i)$ Where  $P(Y_i | g_i) = Normal(E[Y_i | g_i], s2)$ and  $E[Y_i | g_i] = m + b_1 I(g_i=1) + b_2 I(g_i=2)$ 

### MAXIMUM LIKELIHOOD BASED ON PROBABILITIES

Define joint likelihood as the product of individual likelihoods

Maximize the likelihood over the parameters involved

Maximum Likelihood Ratio test can be used to test nested models and draw statistical inferences

## POWER IN LARGE SAMPLES (SMALL EFFECTS)

Zheng et al., 2011



## POWER IN SMALL SAMPLES (LARGE EFFECTS)



## ANALYSIS OF IMPUTED GENOTYPIC DATA IN GWAS

- Short review of standard methods
- Methods for analysis of imputed data
- Open questions and problems

### **ADDITIONAL PROBLEMS**

- Multi-locus analysis: is problematic as no information about joint distribution of genotypes is retained after standard imputation procedures
- GxE analysis (and expected for other interaction analyses) represents a methodological challenge



#### λ's going all the way around 1

- Rotterdam study: population-based cohort used for genetic research for over 15 years
- In GWAS performed over many traits, always λ < 1.05
- GxE results for some traits:

	Environ			
	cov 1	cov 2	cov 3	cov 4
trait 1	1.13	1.13	1	1.14
trait 2	0.98	1.04	1.02	1.04
trait 3	1.12	1.22	1	1.09
trait 4	1.05	1.01	1.03	0.97
trait 5	1.1	1.09	1.07	1.01
trait 6	1.02	1.01	0.92	1.03
trait 7	0.94	0.95	0.89	1

Uppsala, 2011.05.05

Yuti Aukhenko



#### Solution: use robust (co)variances

Solution coming from Thomas Lumley Implemented in ProbABELv0.1-1 (Aulchenko *et al.*, BMC Bioinformatics, 2010)

	Environmental factor					
	cov 1	cov 2	cov 3	cov 4		
trait 1	1.03	1.04	1.03	1.02		
trait 2	1.03	1.01	1.03	1.02		
trait 3	1.02	1.04	1.03	1.02		
trait 4	1.04	1.03	1.03	1.01		
trait 5	1	1.02	1.03	1.01		
trait 6	1.03	1.01	1.02	1.01		
trait 7	1.02	1.03	1.03	1.01		

Still not 100% satisfactory! What about Mixed Models?

Uppsala, 2011.05.05

Yurii Aulchenko

### CONCLUSIONS

- Using regression onto genotype probabilities is a valid and powerful method for standard scenarios
- Use of ML/mixture method can give extra power in case of small samples and large effects
- Caution should be exercised in interaction analyses with imputed data