Advanced Gene Mapping Course

November 7-11, 2022 The Rockefeller University New York, NY

Lectures

Table of Contents

Genome-Wide Association Studies (GWAS) ¹	1
Data Quality Control – Next Generation Sequence and Genotype Array Data ²	15
Rare Variant Association Analysis ²	
Linkage Disequilibrium and its Application in Association Studies ³	40
Statistical Fine-Mapping in Genetic Association Studies ³	43
Fine-Mapping using Summary Statistics ³	50
Integrating GWAS with Functional Annotations ³	57
Transcriptome-Wide Association Studies (TWAS) ³	63
Multivariate Analysis in Genetic Association Studies ³	69
Ethics and Regulation of Human Subjects Research ⁴	70
Pleiotropy and Mediation Analysis ⁵	88
Mendelian Randomization ⁵	101
Generalized/Linear Mixed Models and Interaction ¹	109
Power/Sample Size Estimation ²	126
Special Lecture - Genotype Pattern Mining for Digenic Traits ⁶	133
Heritability ⁷	137
Population Genetics ⁷	153
Polygenic Risk Scores ⁷	161
Functional Annotation ⁷	167

Lectures given by: ¹Heather Cordell, ²Suzanne Leal, ³Gao Wang, ⁴Judy Matuk, ⁵Andrew DeWan, ⁶Jurg Ott, and ⁷Shamil Sunyaev



W

Genome-wide association studies (GWAS) - Part 1

Heather J. Cordell

Population Health Sciences Institute Faculty of Medical Sciences Newcastle University, UK heather.cordell@ncl.ac.uk

Genome-wide association studies (GWAS)

- ullet Popular (and highly successful) approach over past ~ 15 years
- Enabled by advances in high-throughput (microarray-based) genotyping technologies
- Idea is to measure the genotype at a set of single nucleotide polymorphisms (SNPs) across the genome, in a large set of unrelated individuals
 - Cases and controls
 - Or population cohort measured for relevant quantitative phenotypes (height, weight, blood pressure etc)
 - Or related individuals (family data) but need to analyse differently

GWAS (Part 1)

Genome-wide association studies (GWAS)

Association testing: case/control studies

- Collect sample of affected individuals (cases) and unaffected individuals (controls)
 - Or a else a sample of random "population" controls
 - Most of whom will not have the disease of interest
 - Examine the association (correlation) between alleles present at a genetic locus and presence/absence of disease
 - By comparing the distribution of genotypes in affected individuals with that seen in controls

Two individuals		1
Person 1	ACCTGTG <mark>T</mark> GCCCA <mark>A</mark> TGGCGTCCCATA <mark>C</mark> TATCGG ACCTGTG <mark>C</mark> GCCCA <mark>A</mark> TGGCGTCCCATA <mark>C</mark> TATCGG	
Person 2	ACCTGTG <mark>C</mark> GCCCA <mark>G</mark> TGGCGTCCCATA <mark>C</mark> TATCGG ACCTGTG <mark>C</mark> GCCCA <mark>G</mark> TGGCGTCCCATA <mark>G</mark> TATCGG	

• Test each SNP for association/correlation with disease or quantitative phenotype

GWAS (Part 1)

Heather Cordell (Newcastle)

Heather Cordell (Ne

Case/control studies

• Each person can have one of 3 possible genotypes at a diallelic genetic locus

Genotype	Cases	Controls
2 2	500 $(=a)$	200 (= b)
1 2	1100 (= c)	820 (= d)
1 1	$400 \ (=e)$	980 $(=f)$
Total	2000	2000

Case/control studies

• Each person can have one of 3 possible genotypes at a diallelic genetic locus

Genotype	Cases	Controls
2 2	500 $(= a)$	200 (= b)
1 2	1100 (= c)	$820 \ (= d)$
1 1	$400 \ (= e)$	980 $(= f)$
Total	2000	2000

• Test for association (correlation) between genotype and presence/ absence of disease using standard χ^2 test for independence on 2 df

Heather Cordell (Newcastle)	GWAS (Part 1)	5 / 40	Heather Cordell (Newcastle)	GWAS (Part 1)	5 / 40

Case/control studies

• Each person can have one of 3 possible genotypes at a diallelic genetic locus

Genotype	Cases	Controls
2 2	500 $(=a)$	200 (= b)
1 2	1100 (= c)	820 (= d)
1 1	$400 \ (=e)$	980 $(= f)$
Total	2000	2000

- Test for association (correlation) between genotype and presence/ absence of disease using standard χ^2 test for independence on 2 df
 - Defined as $\sum_{i=1,6} \frac{(O_i E_i)^2}{E_i}$ where O_i and E_i are observed and expected counts (calculated from the row and column totals) respectively
 - Generates a *p* value indicating how significant the association/ correlation appears to be

Heather Cordell (Newcastle)

\sim	/	
Case/	control	studies

• Each person can have one of 3 possible genotypes at a diallelic genetic locus

Genotype	Cases	Controls
2 2	500 $(=a)$	200 (= b)
1 2	1100 (= c)	820 (= d)
1 1	$400 \ (=e)$	980 $(= f)$
Total	2000	2000

- Test for association (correlation) between genotype and presence/ absence of disease using standard χ^2 test for independence on 2 df
 - Defined as $\sum_{i=1,6} \frac{(O_i E_i)^2}{E_i}$ where O_i and E_i are observed and expected counts (calculated from the row and column totals) respectively
 - Generates a *p* value indicating how significant the association/ correlation appears to be
- Two odds ratios can be estimated

 - OR $(2|2:1|1) = \frac{af}{be}$ OR $(1|2:1|1) = \frac{cf}{de}$

Odds ratios

- Odds of disease are defined as P(diseased)/P(not diseased)
 - Odds ratio OR (2|2:1|1) repesents the factor by which your odds of disease must be multiplied, if you have genotype 2|2 as opposed to 1|1
 i.e. the 'effect' of genotype 2|2

Odds ratios

- Odds of disease are defined as P(diseased)/P(not diseased)
 - Odds ratio OR (2|2:1|1) repesents the factor by which your odds of disease must be multiplied, if you have genotype 2|2 as opposed to 1|1
 i.e. the 'effect' of genotype 2|2
- Similarly, we can define the OR for 1|2 vs 1|1
 - As the factor by which your odds of disease must be multiplied, if you have genotype 1|2 as opposed to 1|1
 - ${\ensuremath{\, \bullet }}$ i.e. the 'effect' of genotype 1|2

Heather Cordell (Newcastle)	GWAS (Part 1)	6 / 40	Heather Cordell (Newcastle)	GWAS (Part 1)	6 / 40

Odds ratios

- Odds of disease are defined as P(diseased)/P(not diseased)
 - Odds ratio OR (2|2:1|1) repesents the factor by which your odds of disease must be multiplied, if you have genotype 2|2 as opposed to 1|1
 i.e. the 'effect' of genotype 2|2
- Similarly, we can define the OR for 1|2 vs 1|1
 - As the factor by which your odds of disease must be multiplied, if you have genotype 1|2 as opposed to 1|1
 - ullet i.e. the 'effect' of genotype 1|2

Heather Cordell (Newcastle)

- ORs are closely related (often pprox) genotype relative risks
 - The factor by which your probability of disease must be multiplied, if you have genotype 1|2 as opposed to 1|1 (say)
- If your genotype has no effect on your probability (and therefore on your odds) of disease, then the ORs=1.

VAS (Part 1)

• So the association test can be thought of as a test of the null hypothess that the ORs=1

Genotype relative risks

• If a disease is reasonably rare, the odds ratio approximates the genotype relative risk (GRR, RR)

Genotype	Penetrance	GRR	Odds	OR
1/1	0.01	1.0	0.01/0.99 = 0.0101	1.00
1/2	0.02	2.0	0.02/0.98 = 0.0204	2.02
2/2	0.05	5.0	0.05/0.95 = 0.0526	5.21

• If your genotype has no effect on your probability (and therefore your RR) of disease, then both the ORs and the GRRs=1.

GWAS (Part 1

Heather Cordell (Nev

Dominant/recessive effects

Dominant:

Genotype	Cases	Controls	Total
2 2 and 1 2	500 + 1100	200+820	700+1920
1 1	400	980	1380
Total	2000	2000	4000

Recessive:

Heather Cordell (Newcastle)

Genotype	Cases	Controls	Total
2 2	500	200	700
1 2 and $1 1$	1100 + 400	820+980	1920+1380
Total	2000	2000	4000

• Can also rearrange table to examine effects of alleles (1 df tests):

GWAS (Part 1)

	Counts in				
Allele	Cases	Controls			
2	2100 (=a)	1220 (=b)			
1	1900 ($=c$)	2780 (=d)			
Total	4000	4000			

Allelic OR = ad/bc

 χ^2 test statistic on 1 df = $\sum_i (O_i - E_i)^2 / E_i$ where O_i and E_i are the observed and expected values in cell *i*. ۲

Counting alleles

• Assumes HWE under null and multiplicative allelic effects under alternative: considers chromosomes as independent units

GWAS (Part 1)

- Better approach: use counts in previous genotype table to perform a Cochran-Armitage trend test
- Even better approach: use linear or logistic regression

Testing for association: quantitative traits

- Linear regression provides a natural test for quantitative traits
 - Testing the null hypothesis that the slope = 0



S (Part 1

Logistic regression

• Used in case/control studies

Heather Cordell (Newcastle)

- Outcome is affected or unaffected
- Model probability (and thus odds) of disease p as function of variable xcoding for genotype:

$$\ln \frac{p}{1-p} = \beta_0 + \beta_1 x \quad \equiv c + mx$$

11 / 40

• Use observed genotypes in cases and controls to estimate the values of regression coefficients β_0 and β_1

VAS (Part 1

• And to test whether $\beta_1 = 0$

Logistic regression

- Standard method used in standard epidemiological studies e.g. of risk factors such as smoking in lung cancer
- Main advantage is you can include more than one predictor in the regression equation e.g.

$$\ln \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

where x_1 , x_2 , x_3 code for

- genotypes at 3 loci
- measured environmental covariates (e.g. age, sex, smoking etc),
- genetic principal component scores (to adjust for population
- substructure),
- interactions between loci etc. etc.

Heather Cordell (Newcastle) GWAS (Part 1)

Testing for association

- All methods produce a test statistic and a p value at each SNP, indicating how significant the association/correlation observed appears to be
 - i.e. how likely it was to have occurred by chance
 - The threshold to declare 'genome-wide significance' is usually around $p=5\times 10^{-8}$
 - To account for multiple testing of many SNPs across the genome

GWAS (Part 1)

13 / 40

Testing for association

- All methods produce a test statistic and a p value at each SNP, indicating how significant the association/correlation observed appears to be
 - i.e. how likely it was to have occurred by chance
 - The threshold to declare 'genome-wide significance' is usually around $p=5\times 10^{-8}$
 - To account for multiple testing of many SNPs across the genome
- Alternative (Bayesian) methods produce a Bayes Factor

Heather Cordell (Newcastle)

- Indicates how likely the data is under the alternative hypothesis (of association between genotype and phenotype)
 - Compared to under the null hypothesis (of no association between genotype and phenotype)
- Requires you to make some prior assumptions regarding the likely strength of associations (i.e. the value of the β's)
- Choosing a sensible threshold (e.g. log₁₀ BF> 4) requires you to make some prior assumptions regarding what proportion of SNPs in the genome are likely to be associated with the phenotype

Manhattan Plots



- At any location showing 'significant' association, we expect to see several SNPs in the same region showing association/correlation with phenotype
 - Due to the correlation or linkage disequilibrium (LD) between neighbouring SNPs

Close-up of hit region



GWAS (Part 1)

Heather Cordell (Newcastle)

Historical Perspective: Complement Factor H in AMD

- First (?) GWAS was by Klein et al. (2005) Science 308:385-389
- Typed 116,204 SNPs in 96 cases (with age-related macular degeneration, AMD) and 50 controls
 - Very small sample size they were very lucky to find anything!

<u>GWAS (P</u>art 1)

- Luck was due to the fact the polymorphism has a very large effect (recessive OR=7.4)
- Klein et al. followed up on two SNPs passing threshold $(p < 4.8 \times 10^{-7})$
 - Plus a third SNP that just failed to pass significance threshold, but lay in same region as first SNP

16 / 40

18 / 40

Complement Factor H in AMD

• Of the 3 SNPs followed up:

Heather Cordell (Newcastle)

- One appeared to be due to genotyping errors: significance disappeared on filling in some missing genotypes
- First and third SNP lie in intron of Complement Factor H (CFH) gene
 - Lies in region previously implicated by family-based linkage studies
- Resequencing of the region identified a polymorphism of plausible functional effect
- Immunofluorescence experiments in the eyes of AMD patients supported the involvement of CFH in disease pathogenesis.

GWAS

GWAS really got going in around 2007

Heather Cordell (Newcastle

- See Visscher et al. (2012) AJHG 90:7-24 "Five Years of GWAS Discovery"
- And Visscher et al. (2017) AJHG 101:5-22 "10 Years of GWAS Discovery: Biology, Function and Translation"
- 2007/2008 saw a slew of high-profile GWAS publications
 - Breast cancer (Easton et al. 2007)
 - Rheumatoid Arthritis (Plenge et al. 2007)
 - Type 1 and Type 2 diabetes (Todd et al. 2007; Zeggini et al. 2008)
- Arguably the most influential was the Wellcome Trust Case Control Consortium (WTCCC) study of 7 different diseases
 - http://www.wtccc.org.uk/

WTCCC

Manhattan plots for 7 diseases

- Nature 447: 661-678 (2007)
- Considered 2000 cases for each of the following diseases:
 - Bipolar disorder, coronary artery disease, Crohn's disease, hypertension, rheumatoid arthritis, type 1 diabetes, type 2 diabetes
- Compared each disease cohort to common control panel
 - 3000 population-based controls
 - From 1958 birth cohort and National Blood Service
- Highly successful
 - WTCCC found 24 separate association signals
 - Including highly convincing signals in 5 out of the 7 diseases studied
 - All were replicated in subsequent independent follow-up studies



GWAS (Part 1)

Company Contracting the second secon

Lessons from WTCCC (and others)

- Typically used rather standard statistical/epidemiological methods (χ^2 tests, t tests, logistic regression etc.)
- Success largely due to:

Heather Cordell (Newcastle)

- An appreciation of the importance of large sample size (> 2000 cases, similar or greater number of controls)
- Stringent quality control procedures for discarding low-quality SNPs and/or samples
- Stringent significance thresholds $(p=5\times10^{-8})$ to account for multiple testing and/or low prior prob of true effect
- Importance of replication in an independent data set

GWAS (Part 1)

Short break

Quality Control

QC: call rates and heterozygosity

- Stringent QC checks are required for GWAS data
- Discard samples (people) deemed unreliable
 - Low genotype call rates, excess heterozygosity etc.
 - X chromosomal markers useful for checking gender
 - $\bullet\,$ Males should 'appear' homozygous at all X markers
 - Genome-wide SNP data useful for checking relationships and ethnicity
- Discard data from SNPs deemed unreliable
 - On basis of genotype call rates, Mendelian misinheritances, Hardy-Weinberg disequilibrium

GWAS (Part 1)

• Exclude SNPs with low minor allele frequency (MAF)



• 61 sample exclusions (low call-rate); 23 exclusions (heterozygosity)

GWAS (Part 1)

24 / 40

 SNP exclusions also made based on call-rates, MAF and Hardy-Weinburg equilibrium (HWE)

Heather Cordell (Newcastle)

Heather	Cordell	(Newcastle)		

QC: ethnicity tests



- Multidimensional scaling (with 210 HapMap individuals) identifies 33 samples with non-Caucasian ancestry
- Similar multivariate methods can be used to model more subtle population differences between samples...

Heather Cordell (Newcastle)

Multivariate Analysis

- Several related multivariate analysis techniques have been proposed for detecting population structure in genome-wide association studies
 - Principal components analysis (PCA)
 - Principal coordinates analysis (PCoA)
 - Multidimensional scaling (MDS)

Multivariate Analysis

- Several related multivariate analysis techniques have been proposed for detecting population structure in genome-wide association studies
 - Principal components analysis (PCA)
 - Principal coordinates analysis (PCoA)
 - Multidimensional scaling (MDS)

Heather Cordell (Newcastle)

Heat

- If population differences can be detected (and adjusted for) in association analysis, this offers a way to deal with the problem of population stratification
 - Population sampled actually consists of several 'sub-populations' that do not really intermix
 - Can lead to spurious false positives (type 1 errors) in case/control studies

GWAS (Part 1)

Multivariate Analysis

- Several related multivariate analysis techniques have been proposed for detecting population structure in genome-wide association studies
 - Principal components analysis (PCA)
 - Principal coordinates analysis (PCoA)
 - Multidimensional scaling (MDS)

Heather Cordell (Newcastle)

- If population differences can be detected (and adjusted for) in association analysis, this offers a way to deal with the problem of population stratification
 - Population sampled actually consists of several 'sub-populations' that do not really intermix
 - Can lead to spurious false positives (type 1 errors) in case/control studies
- These techniques can also be used in quality control (QC) procedures, to check for (and discard) population outliers

GWAS (Part 1)

Principal components analysis (PCA)		Principal Components Analysis		
<section-header></section-header>		 Price et al. (2006) Natu (2006) PLoS Genetics 2 Based on popn genet Idea is to form a large r to the genotype at a L M = 	ure Genetics 38:904-909; Patterso 2(12):e190 tics ideas from Cavalli-Sforza (1978) matrix M of SNP counts (0,1,2) c loci (=rows) for <i>n</i> individuals (=c $\begin{pmatrix}g_{11} & g_{12} & g_{1n} \\ g_{21} & g_{22} & g_{2n} \\ g_{31} & g_{32} & g_{3n} \\ & & & & \\ & & & & & \\ g_{L1} & g_{L2} & g_{Ln} \end{pmatrix}$	n et al. orresponding columns)
her Cordell (Newcastle) GWAS (Part 1)	27 / 40	Heather Cordell (Newcastle)	GWAS (Part 1)	28 / 4

Principal Components Analysis

• Subtract row means and normalise by function of row allele frequency $\sqrt{f_l(1-f_l)}$ to give matrix X

	(x_{11})	<i>x</i> ₁₂		x_{1n}
	x ₂₁	<i>x</i> ₂₂	•	x _{2n}
X =	<i>x</i> ₃₁	<i>x</i> ₃₂		x _{3n}
~			•	•
		•	•	•
	$\setminus x_{L1}$	<i>x</i> _{L2}	•	x _{Ln} /

• This matrix will be used as starting point for PCA

Heather Cordell (Newcastle)

Heather Cordell (Newcastle)

• In principal we could start with a different matrix – in particular not all PCA approaches would normalise by $\sqrt{f_i(1-f_i)}$

GWAS (Part 1)

Multivariate Analysis

 Estimate covariance matrix Ψ = X^TX between all pairs of individuals, with entries ψ_{ij} defined as the covariance (summing over SNPs) between column i and j of X

GWAS (Part 1)

ather Cordell (Newcastle

 Represents average genome-wide identity by descent (IBD) (estimated from identity by state, IBS)

Multivariate Analysis

- Estimate covariance matrix $\Psi = X^T X$ between all pairs of individuals, with entries ψ_{ij} defined as the covariance (summing over SNPs) between column *i* and *j* of *X*
 - Represents average genome-wide identity by descent (IBD) (estimated from identity by state, IBS)
 - Compute the eigenvectors \vec{v}_i and eigenvalues λ_i of matrix Ψ
 - Co-ordinate j of the kth eigenvector represents the ancestry of individual j along 'axis' k

GWAS (Part 1)

Multivariate Analysis

- Estimate covariance matrix Ψ = X^TX between all pairs of individuals, with entries ψ_{ij} defined as the covariance (summing over SNPs) between column i and j of X
 - Represents average genome-wide identity by descent (IBD) (estimated from identity by state, IBS)
 - Compute the eigenvectors v
 _j and eigenvalues λ_j of matrix Ψ
 Co-ordinate j of the kth eigenvector represents the ancestry of individual j along 'axis' k
- For technical details, see McVean (2009) PLoS Genetics 5;10:e1000686

Multivariate Analysis

- Estimate covariance matrix $\Psi = X^T X$ between all pairs of individuals, with entries ψ_{ij} defined as the covariance (summing over SNPs) between column *i* and *j* of *X*
 - Represents average genome-wide identity by descent (IBD) (estimated from identity by state, IBS)
 - Compute the eigenvectors $ec{v}_j$ and eigenvalues λ_j of matrix Ψ
 - $\bullet\,$ Co-ordinate j of the kth eigenvector represents the ancestry of individual j along 'axis' k
- For technical details, see McVean (2009) PLoS Genetics 5;10:e1000686
- Many genetics packages e.g. (PLINK) will allow you to calculate the top 10 (or more) PCs
 - Different geographic populations can often be well separated by just the first two or three PCs
 - Useful for outlier detection

Heather Cordell (Newcastle)

Heather Cordell (Newcastle

- For more subtle differences, you may need to calculate more PCs
 - And include them as covariates in the regression equation

GWAS (Part 1)

Post-GWAS QC can determine whether you have included 'enough'

Post GWAS QC: Q-Q Plots (good)

 Plot ordered test statistics (y axis) against their expected values under the null hypothesis (x axis)



S (Part 1

eather Cordell (N

Q-Q Plots (bad)



Population stratification

- A QQ plot showing constant inflation (straight line with slope > 1) can indicate population stratification/population substructure
- Simple solution: Genomic Control (Devlin and Roeder 1999)
 - Use your observed test statistics to estimate the slope (=inflation factor λ)
 - \bullet Divide each test statistic by λ to get an adjusted (deflated) test statistic

GWAS (Part 1)

- More complicated solution: use PCA/MDS or similar
- Even more complicated solution: use linear mixed models

Relatedness

Expected IBD sharing

- Assuming no inbreeding, the IBD state probabilities are:
- With genome-wide data, can also infer relationships based on average identity by descent (IBD) $\Psi = X^T X$ or identity by state (IBS)
 - Using 'thinned' subset of markers with high minor allele frequency (MAF) and in approximate linkage equilibrium
 - Simple relationships (PO, FS, MZ/duplicates) can identified with only a few hundred markers
 - More complicated relationships require 10,000-50,000 SNPs
- Various software packages, including PLINK, KING and TRUFFLE

GWAS (Part 1)

Heather Cordell (Newcastle)

Number of alleles shared IBD			
2	1	0	
1	0	0	
0	1	0	
1/4	1/2	1/4	
0	1/2	1/2	
0	1/2	1/2	
0	1/2	1/2	
0	1/4	3/4	
0	1/16	15/16	
1/16	6/16	9/16	
	Numb 2 1 0 1/4 0 0 0 0 0 0 1/16	$\begin{array}{c c} \text{Number of all}\\ \hline 2 & 1\\ \hline 1 & 0\\ 0 & 1\\ 1/4 & 1/2\\ 0 & 1/2\\ 0 & 1/2\\ 0 & 1/2\\ 0 & 1/2\\ 0 & 1/4\\ 0 & 1/16\\ 1/16 & 6/16\\ \end{array}$	

...

35 / 40

 A useful visualisation tool is to plot SE(IBD) vs mean(IBD) (as estimated across the genome)

GWAS (Part 1)

Heather Cordell (Newcastle)

34 / 40

• Or kinship coefficient $(\frac{1}{2}P(IBD=2)+\frac{1}{4}P(IBD=1))$ against P(IBD=0)

Full/half sibs and parent-offspring		CHD GWAS results (low QC)
n n n n n n n n n n n n n n		$\left[\begin{array}{c} 1\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\$
Heather Cordell (Newcastle) GWAS (Part 1)	36 / 40	Heather Cordell (Newcastle) GWAS (Part 1) 37 / 40



\sim			
I-ono	mawuda	moth hnh	VCIC
NICHO			רורעו

- Puts together data (or results) from a number of different studies
 Could analyse as one big study
 - But preferable to analyse using meta-analytic techniques

Heather Cordell (Newcastle)

• At each SNP construct an overall test based on the results

GWAS (Part 1)

(log ORs and standard errors) from the individual studies

Genome-wide meta-analysis

- Puts together data (or results) from a number of different studies
 Could analyse as one big study
 - But preferable to analyse using meta-analytic techniques
 At each SNP construct an overall test based on the results (log ORs and standard errors) from the individual studies
- Meta-analysis is often made easier by using *imputation*
 - Inferring (probabilistically) the genotypes at SNPs which have not actually been genotyped

GWAS (Part 1)

- On the basis of their known correlations with nearby SNPs that have been genotyped
- $\bullet\,$ Using a reference panel of people (e.g. 1000 Genomes) who have been genotyped at all SNPs

40 / 40

Heather Cordell (Newcastle)

Genome-wide meta-analysis

- Puts together data (or results) from a number of different studies
 Could analyse as one big study
 - But preferable to analyse using meta-analytic techniques
 - At each SNP construct an overall test based on the results
 - (log ORs and standard errors) from the individual studies
- Meta-analysis is often made easier by using *imputation*
 - Inferring (probabilistically) the genotypes at SNPs which have not actually been genotyped
 - On the basis of their known correlations with nearby SNPs that have been genotyped
 - Using a reference panel of people (e.g. 1000 Genomes) who have been genotyped at all SNPs
- Enables meta-analysis of studies that used different genotyping platforms
 - By imputing to generate data at a common set of SNPs
 - Ideally while accounting for the imputation uncertainty in the downstream statistical analysis
 - In practice often don't bother use post-imputation QC to remove poorly-imputed SNPS

Heather Cordell (Newcastle) GWAS (Part 1)

40 / 40







- · Genotype markers which can be used as DNA fingerprint
- Allows for Assessment of DNA quality
- Aids in determining the the genetic sex of study subjects
 To aid in identification of potential sample swaps
- Detects cryptic duplicates
- For family data
 - Aids in determining close familial relationships
 - Non-paternity
 - Sample swaps
 - Cryptic relationships



- Duplicate samples genotyped using arrays to detect inconsistencies
 - Can use duplicate samples that are inconsistent to adjust clusters to improve allele calls
 - Will not detect systematic errors
- Usually generated only for genotype array data
- Due to expense, duplicate samples are usually not generated for exome or whole genome sequencing studies









Variant Calling

- BAM files are large and take considerable resources
 - Storage is expensive
 - One 30x whole genome is ~80-90 gigabytes
 - A small study of 1,000 samples will consume 80 terabytes of disk space
- The cost of cloud computing to call variants
 - (Souilmi et al. 2015)
 - \$5 per exome
 - \$50 per genome
 - For 1,000 samples
 - \$5,000 exome
 \$50,000 genome



- Instead of obtaining VCF files
- Can obtain gVCF files to perform joint calling and complete the GATK pipeline
 - A whole genome gVCF
 - ~1 Gigabyte

 $-\ 1/100^{th}$ the size of a BAM file for one individual

12

Influences on Sequence Quality

- DNA guality
 - Age of sample
 - Extraction method
 - Source of sample
 - e.g., blood, skin punch, buccal
- Sequencing machines (read length)
- Median sequencing depth
- Alignment
- Variant calling method used
 - Single nucleotide variants and insertion/deletions
 - Structural variants
- 13

17

NGS Data Quality Control

- Extremely important to perform before data analysis
 - Poor data quality can increase type I and II errors - Due to inclusion of false positive variant sites or incorrect
 - genotype calls
- Protocols for data QC are still in their infancy - No set protocols for QC
- QC is data specific
 - Dependent on read depth
 - Batch effects
 - Availability of duplicate samples
- etc.

14

NGS Data Quality - Removal of Genotype Calls and Samples • Sequence depth of coverage - DP variant High DP could be an indication of copy number variants Which can introduce false positive variant calls
 » Due to down sampling in GATK maximum DP is 250 DP_genotype • Concerned if depth is too low or too high - Low insufficient reads to call a variant site - Remove genotypes with low read depth, e.g., ${\sf DP}\underline{{\leq}8}$ Genotype quality (GQ) score - Removal of sites with low genotype quality core, e.g., GQ<20 15 16 **VCF Example**



NGS Data Quality - Removal of Genotype Calls







NGS Data Quality – Removal of Genotype Calls and Samples

- Removal of sites with missing data
 e.g., missing > 10% of genotypes
- Removal of "novel" variant sites which only occur in one batch and the alternative allele is observed multiple times or the minor allele frequency (MAF) is high in overall sample
- Removal of sites that deviate from Hardy-Weinberg Equilibrium (HWE)
 - Must be performed by population, e.g., African American and European American
 - Related individuals should be removed from the sample before testing for deviations from HWE

19

NGS Data Quality Control

- GATK Variant Quality Score Recalibration (VQSR)
 - Used to determine variant sites of bad quality
 Variant site is a false positive call
- However even after this step
 - Concordance of duplicates (when available) and
 - and Ti/Tv ratios are often low
- Additional QC steps needs to be performed









Sequence Data QC Overview

- Variant level removal of variant sites
 - Low call rate

25

- i.e., missing call rate > 10%
- "Novel" variant sites observed >2 only in a single batch
- Deviation from Hardy-Weinberg-Equilibrium
 Population specific
 - Unrelated individuals
 - e.g., p<5 x 10⁻⁸, p<5x10⁻¹⁵

Data Clean – Assessing Sex Chromosomes

- When data is collected on study subjects they are asked about their gender/sex and not their genetic sex
 - Differences in gender/sex and genetic sex can be due to
 - Sample swaps
 - Study subjects who are not cisgender
- Some study subjects may have neither a XX nor XY karyotype
 - Turner syndrome X0
 - Klinefelter syndrome XXY

26



- their genetic sex are removed from the analysis
- This observation may be due to a sample swap
 - When samples are swapped
 - Phenotype data will be incorrect

 e.g., may be a case when labeled as a control









- Duplicate and related individuals can be detected
 - By examining <u>Identity-by-State (IBS)</u> adjusted for allele frequencies (p-hat) between all pairs of individuals within a sample
 - Identify-by-descent (IBD) sharing can be estimated

32



33

IBD Sharing Estimated Pairwise for all Individuals in a Samples

- PLINK (Purcell et al. 2007)
- Uses sequence (or genotype array) data to check IBD
 Prune markers to remove those in LD
- Prune markers to remove those in LD
 e.g., r²<0.1
- P-hat is calculated using the "population" allele frequency
- Used to approximates IBD sharing
- IBD is the number of alleles of alleles which are shared between a pair of individuals
 - Can either share 0, 1, and 2 alleles

34

Identifying Duplicate and Related Individuals Monozygote twins and duplicate samples will share

- 100% of their alleles IBD
 - IBD=2 is 1.0 (can be lower due to genotyping error)
- Siblings and child-parent pairs will share 50% of their alleles IBD
 - For parent-child IBD=1 is 1.0 (IBD=0 is 0 & IBD=2 is 0)
 - For sibs IBD=1 is ~0.50 (IBD=0 is ~0.25 & IBD=2 is ~0.25)
 - For more distantly related individuals the IBD measure will be lower

Identifying Duplicate and Related Individuals KING [Kinship-based INference for Gwas (Manichaikul et al. 2010)] can also be used to identify duplicate and related individuals KING is more robust to population substructure and admixture Prune markers for LD (e.g., r²<0.1) Provides kinship coefficients Duplicate samples Kinship coefficient equals 0.5 Siblings



Multiple Individuals observed that are distantly "Related"

- If individuals in sample come from different populations
- e.g., individuals from the same population within the sample will have inflated p-hat values due to incorrect allele frequencies Incorrectly appear to be related to each other
- "Relatedness" amongst many individuals can also be observed when batches are combined if they have different error rates Individuals from the same batch appear to be related
- . DNA contamination can cause "relatedness" between multiple individuals

38





40







• Individuals of different ancestry

- e.g., African American samples included with European Americans samples
- determine the ancestry for samples that are outliers Should not include HapMap/1000 genomes samples when calculating
- Batch effects







Detecting Genotyping Error – Examining HWE

- Testing for deviations from HWE not very powerful to detect genotyping errors
- The power to detect deviations from HWE dependent on: – Error rates
 - Underlying error model

Random

- Heterozygous genotypes -> homozygous genotypes
- Homozygous genotypes ->Heterozygous genotype
- Minor allele frequencies (MAF)

Detecting Genotyping Error – Examining HWE

- Controls and Cases are evaluated separately - Deviation found only in cases can be due to an association
- Test for deviation from HWE only in samples of the same ancestry
 - Population substructure can introduce deviations from HWE
- Do not include related individuals when testing for deviations from HWE
 - Can cause deviations from HWE

46

Detecting Genotyping Error – Examining HWE

- What criterion is used to remove variants due to a deviation from HWE
- $-\,$ GWAS studies have used 5.0 x 10^-7 to 5.0 x 10^{-15}
- Quantitative Traits

- Caution should be used removing markers which deviate from HWE may be due to an association

- Remove markers with extreme deviations from HWE and Flag markers with less extreme deviations from HWE
- When performing imputation need to be more stringent in removing variants which deviate from HWE

Sequence Data QC Overview

- Remove variant sites that fail VQSR
- Remove genotypes with low DP, GQ scores, etc.
- Remove variant sites with large percent of missing data
- Remove samples with missing large percent of missing data
- Evaluate genetic sex of individuals based upon X and Y chromosomal data
 - Sample mix-ups
 - Individuals with Turner or Klinefelter Syndrome



45

Sequence Data QC Overview

- Evaluate samples for cryptically related individuals and duplicates
 - Use variants which have been pruned for LD • e.g., r²<0.1
 - King or Plink algorithm
 - · Always remove duplicate individuals - Retaining only one in the sample
 - If sample includes related samples use linear mix models (LMM)/Generalized LMM (GLMM) to control for relatedness Best to perform even for data without related individuals
 - · If only a few related individuals can retain only one individual of a relative group if not using LMM or GLMM

49

51

Sequence Data QC Overview

• Detection of sample outliers

- Perform principal components analysis (PCA) or multidimensional scaling (MDS) to detect outliers
 - Use variants pruned for LD
 - e.g. r²<0.1
 - Use unrelated individuals and then project related individuals onto the PCs
- Due to population substructure/admixture and batch effects
- Remove effects by
- Additional QC
- Removal of outliers (can be determined by Mahalanobis distance) and\or
- Inclusion of MDS or PCA components in the association analysis

50

Sequence Data QC Overview

- Remove/flag variant sites that deviate from HWE in controls
 - HWE should be only be tested in unrelated individuals from the same population
- Post Analysis Quantile-Quantile (QQ) plots
 - To evaluate uncontrolled batch effects and population substructure/admixture

QQ Plots - Genome Wide Association Diagnosis

- Thousands of variants/genes are tested simultaneously
- The p-values of neutral markers follow the uniform distribution
- If there are systematic biases, e.g., population substructure, genotyping errors, there will be a deviation from the uniform distribution
- QQ plots offers an intuitive way to visually detect biases
- Observed p-values are ordered from largest to smallest and their $-\log_{10}(p)$ values are plotted on the y axis and the expected -log₁₀(p) values under the null (uniform distribution) on the x axis

52



Genomic Inflation Factor to Evaluate Inflation of the Test Statistic

- Genomic Inflation Factor (GIF): ratio of the median of the test statistics to expected median and is usually represented as λ
 - No inflation of the test statistic λ =1
 - Inflation λ>1
 - Deflation λ<1
 - · Can be observed when a study is underpowered
- Problematic to examine the mean of the test statistic
 - Can be large if many variants are associated
 - · Particularly if they have very small p-values
 - Should not be used



Phenotype	Covariate	Mean Chi-Square	GIF (λ)
BP		1.23829	1.16932
BP	Age	1.24119	1.18025
BP	Age-EV1	1.09471	1
BP	Age-EV2	1.0881	1
BP	Age-EV4	1.08385	1
BP	Age-EV10	1.09582	1.00402
BPI		1.14931	1.08921
BPI	Age	1.15139	1.08113
BPI	Age-EV1	1.05079	1.01148
BPI	Age-EV2	1.0428	1
BPI	Age-EV4	1.04204	1
BPI	Age-EV10	1.05421	1.01724
BPII		1.17283	1.25664
BPII	Age	1.17583	1.26996
BPII	Age-EV1	1.09874	1.15065
BPII	Age-EV2	1.09904	1.16425
BPII	Age-EV4	1.09502	1.14609
BPII	Age-EV10	1.10046	1.1418
BPII	Sex,Age-EV1	1.05958	1.06424
BPI	Sex.Age-EV4	1.05817	1.05323
BPII	Sex,Age-EV10	1.06338	1.05581























Ti/Ty Patios during OC Process						
in in ratios during QC Process						
	Known	Novel	All			
Before VQSR	2.95 ± 0.05	1.18 ± 0.29	2.86 ± 0.07			
Before additional QC	3.12 ± 0.03	2.01 ± 0.32	3.11 ± 0.03			
Genotype QC DP <u><</u> 8, GQ <u><</u> 20	3.18 ± 0.04	2.10 ±0.32	3.16 ± 0.03			
Remove sites missing >10% genotypes	3.39 ± 0.04	2.42 ± 0.52	3.39 ± 0.04			
Remove batch specific novel sites ≥2 N=17,835	3.39 ± 0.04	2.41 ± 0.53	3.39 ± 0.04			
Remove sites deviating from HWE p <u>≤</u> 5x10 ⁻⁸ N=4,414	3.41 ± 0.04	2.39 ± 0.54	3.40 ± 0.04			





Sequence Data QC

- Batch effects can sometimes be removed with additional QC
- Extreme outliers should be removed
- Additionally, MDS\PCA components can be included in the analysis to control for population substructure\admixture and batch effects
 - Unless correlated with the outcome (phenotype)
 - The MDS or PCA components should be recalculated after QC only including those samples included in the analysis
- Batch (dummy coding) may be included as a covariate in the analysis
 - Unless correlated with the outcome (phenotype)

69

Convenience Controls–Sequence Data

- Obtain BAM files and recall cases and control together
 - Can still have differential errors between cases and controls
 Check variant frequency by variant types in cases and control
 - Synonymous variants should have the same frequencies
 - Would not expect large differences in numbers of variants between cases and controls
- For single variants can compare difference in frequencies with gnomAD but is problematic
 - Differences in frequencies can be due to differences in ancestry and/or sequencing errors
 - Cannot adjust for confounders
 - e.g., sex, population substructure/admixture
- Don't perform an aggregate test using frequency information obtained from databases, e.g., gnomAD, TOPMed Bravo

Genotype Array Data

- Genotype Data QC Population Based Studies • Initially remove DNA samples from individuals who are missing
- >10% or their genotype dataFor variant sites with a minor allele frequency (MAF)>0.05
- Remove variants sites missing >5% of their genotype data
- For variant sites with a MAF<5%
 - Remove variant sites missing > 1% of their genotype data
- The genotypes for variant sites with missing data may have higher genotype error rates

71

72

26

• Can reduce the cost of a study

- Genotype data
- Type I error can be increased
 - Ascertainment from different population
 - Differential genotyping error
 - Even if performed at the same facility
- Proper QC can reduce or remove biases

Order of	Data	Cleaning-Genotype	Array [Data
----------	------	-------------------	---------	------

Remove samples missing >10% genotype data

Remove SNPs with missing genotype data

- If minor allele frequency >5%
 Remove markers with >5% missing genotypes
- If minor allele frequency <5%
- Remove markers with >1% missing genotypes
- Remove samples missing >3% genotype calls
- Check genetic sex of individuals based on X-chromosome markers & Y chromosome marker data (if available)
 - Remove individual whose reported gender/sex is inconsistent with genetic data
 - Could be due to a sample mix-up
- Check for cryptic duplicates and related individuals
 - Used "trimmed data set of markers which are not in LD
 - e.g. r2<0.1
 Remove duplicate samples

73

Order of Data Cleaning-Genotype Array

- Perform PCA or MDS to check for outliers
 - Use trimmed data set of markers which are not in LD
 e.g., r2<0.1
 - First with unrelated individuals and then project related individuals on the components
 - Remove outliers from data
 e.g., Mahalanobis distance
- Check for deviations from HWE
 - Separately in cases and controls
 - Only unrelated individuals
 - If more than one ancestry group
 Separately for each ancestry group
 As determined via PCA or MDS
- Examine QQ plots for potential problems with the data – e.g., not controlling adequately for population admixture



























- Combined multivariate & collapsing (CMC)
 Li & Leal, AJHG 2008
- Collapsing scheme which can be used in the regression framework
 - Can use various criteria to determine which variants to collapse into subgroups
 - Variant frequency
 - Predicted functionality

14



15



CMC





Methods to Detect Rare Variant Associations Weighted Approaches

- Group-wise association test for rare variants using the Weighted Sum Statistic (WSS)

 Variants are weighted inversely by their frequency in controls (rare
 - Variants are weighted inversely by their requery in controls (rare variants are up-weighted)
 Madsen & Browning, PLoS Genet 2009
- Kernel based adaptive cluster (KBAC)
 Adaptive weighting based on multilocus genotype
 Liu & Leal, PLoS Genet 2010

Methods to Detect Rare Variant Associations Maximization Approaches

- Variable Threshold (VT) method
 - Uses variable allele frequency thresholds and maximizes the test statistic
 Can also incorporate weighting based on functional information
 Price et al. AJHG 2010

RareCover

 Maximizes the test statistic over all variants with a region using a greedy heuristic algorithm

Bhatia et al. 2010 PLoS Computational Biology

20



Wu et al. 2011 AJHG

Optimal Test

• SKAT-O

22

 Maximizes power by adaptively using the data to combine a burden test and the sequence kernel association tests
 Lee et al. 2012 AJHG

21

19

Significance Level for Rare Variant Association Tests

- For exome data where individual genes are analyzed usually a Bonferroni correction for the number of genes tested is used
 There is very little to no linkage disequilibrium between genes
- Bonferroni correction used
 e.g., p≤2.5 x 10⁻⁶ (Correction for testing 20,000 genes)

Determine MAF Cut-offs for Aggregate Rare Variant Association Tests • MAF cut-offs are frequently used to determine which variants

- WAY CUL-OUS are frequently used to determine which variants to analyze in aggregate rare variant association tests
- MAF from controls should not be used
 - Increases in type I error rates
- Determine variant frequency cut-offs from databases - Using population frequencies for those understudy
 - gnomAD
 http://gnomad.broadinstitute.org/



- Same frequency of missing variant calls in cases and controls
 Decrease in power
- More variant calls missing for either cases or controls
- Increase in Type I error
- Decrease in power
- Remove variant sites which are missing genotypes, e.g., >10%
- Can impute missing genotypes using observed allele frequencies
 For the entire sample
 - Not based on case or control status
- Analyze imputed data using dosages



27


















Which uses the GRM





















Age-related Hearing Impairment (ARHI) (aka Presbycusis) • ARHI can impact quality of life and daily functioning • ARHI is one of the most common adult conditions – In the USA • ARHI affects 50% of individuals >75 years of age • It is estimated that 30-40 million will be affected with significant ARHI by 2030

Severe

44











Exclusion Criteria Obtained from ICD10, ICD9, & Self Report • Deafness • Early-onset hearing impairment • Otosclerosis • Meniere's • Labyrinthitis • Disorders of acoustic nerve • Bell's palsy • History of chronic suppurative and nonsuppurative otitis media

- Meningitis
- Encephalitis, myelitis, and encephalomyelitis
- Etc.

51















Analysis of Exome Data
Analysis limited to individuals of white European Ancestry
Sex, age, and two PCAs included as covariates

Age for cases first report of hearing difficulty & controls age at last visit
The PCAs where recalculated for only individuals included in the analysis
Using the pruned genotypes array data (r2<0.1)





- All variants with four or more alternative alleles observed in the sample analyzed
 - A very low minor allele frequency was used since it was hypothesized some of the variants may have large effect sizes

58























Results

- Replicated some previously reported ARHL genes – Some which had not been previously replicated
 - e.g., BAIAP2L2, CRIP3, KLHDC7B, MAST2, and SLC22A7
- Identified and replicated a new HL gene which has not been previously reported
 - Inner ear expression in humans and mice supports the involvement of gene in HL etiology
- Rare-variant aggregate analysis demonstrated the important contribution of Mendelian HL genes, i.e. *MYO6, TECTA*, and *EYA4* the genetics of ARHL

69

Results

- Rare variants for ARHL tend to have larger effect sizes than those for common variants
 - Rare variants should play an important role in risk prediction by increasing accuracy
- Although most of the studies findings were replicated in independent samples of white Europeans
 - Additional studies are necessary to elucidate whether these variants/genes play a role in the genetic etiology of ARHL in other populations

Linkage disequilibrium in genetic association studies

Gao Wang, Ph.D.

Advanced Gene Mapping Course, November 2022

The Gertrude H. Sergievsky Center and Department of Neurology Columbia University Vagelos College of Physicians and Surgeons

Genetic association studies

Identify genetic variants associated with complex traits

- Association does not imply causality
- Disease, quantitative traits, molecular phenotypes

2

in order to

1

- Understand biological mechanism
- Identify potential drug targets
- · Identify individuals with high disease risk

Sources of association signals

Causal association — meaningful

• Tested genetic variations influence traits directly

Linkage disequilibrium (LD) — useful

- Tested genetic variations associated with other nearby variations that influence traits
- Meaningful or misleading, in different contexts

Population stratification — misleading

- Tested genetic variations is unrelated to traits, but is associated due to sampling differences
- eg, minor allele frequency, disease prevalence

3

Sources of association signals: analysis tools

Causal association — meaningful

• Fine-mapping, colocalization, Mendelian randomization

Linkage disequilibrium (LD) — useful

- Meaningful: LD scores regression, polygenic risk scores (PRS), transcriptome-wide association studies (TWAS)
- Misleading: fine-mapping, LD pruning / clumping

Population stratification — misleading

• Principle component analysis, linear (mixed) models



Impact of LD on GWAS analysis

Polygenic: trait influenced by numerous genetic variants

- Misleading: stronger association due to more LD 'friends'
- Useful: whole-genome prediction with sparse models



A second thought on genomic inflation

Population stratification? Or, polygenic inheritance + LD?



Suggested reading: Yang et al (2011) EJHG

LD score regression (LDSC)

LD score regression model without population stratification



LD score regression (LDSC)

Separating h_g^2 and population stratification



A more powerful and accurate correction factor for GWAS summary statistics compared to genomic control approach.

- Bulik-Sullivan et al (2015) Nature Genetics the LDSC regression paper
- Zhu and Stephens (2017) AoAS a neat, alternative LDSC regression model derivation in supplemental material

10

12

LDSC application: heritability estimation

Narrow sense heritibility

• Proportion of phenotypic variation explained by additive genetic factors

Estimation strategy

- Pedigree design: genetic covariance and IBD sharing
- Population design: linear mixed models

Population design, summary statistics

- LDSC to estimate SNP-based heritability
- Stratified LDSC (S-LDSC) to partition heritability by functional annotations

11

9

Variance of height explained in GWAS





Fine-mapping with summary statistics: current methods and practical considerations

Gao Wang, Ph.D.

Advanced Gene Mapping Course, November 2022

The Gertrude H. Sergievsky Center and Department of Neurology Columbia University Vagelos College of Physicians and Surgeons



Figure: Benner et al. (2017) Am. J. Hum. Genet.

Association analysis summary statistics

z-scores from univariate association studies:

$$\hat{z}_j \coloneqq \hat{\beta}_j / s_j,$$

where

$$\hat{\beta}_j \coloneqq (\mathbf{x}_j^\mathsf{T} \mathbf{x})^{-1} \mathbf{x}^\mathsf{T} \mathbf{y} \quad s_j \coloneqq \sqrt{\hat{\sigma}_j^2 (\mathbf{x}_j^\mathsf{T} \mathbf{x})^{-1}}$$

- **Sufficient** statistics: $x^{\mathsf{T}}x, x^{\mathsf{T}}y, \hat{\sigma}_i^2$
- "Summary" statistics:
 - z-scores: \hat{z}
 - Genotypic correlation: \hat{R}

Reasons to work with summary statistics

Advantage over full data (genotypes and phenotypes):

- Easier to obtain and share with others
- Convenient to use: QC and data wrestling barely needed
- Computationally suitable for large-sample fine-mapping
 - $\mathcal{O}(p^2)$ (summary statistics) $\ll \mathcal{O}(np)$ (full data)
 - when sample size $n \gg$ variants in fine-mapped region p

Suggested reading: Pasaniuc and Price (2017) Nat. Rev. Genet.

3

$$\hat{z} \sim N(\hat{R}z, \hat{R})$$

Assumptions:

- 1. Heritability of any single SNP is small
- 2. \hat{R} is sample genotypic correlation for the same study
- 3. Genotypes used to computed z and \hat{R} are accurate

Properties of per SNP z scores

• z-score for a SNP depends on effects of both itself and other correlated SNPs:

$$\mathsf{E}(\hat{z}_j|\hat{\boldsymbol{R}}) = \sum_{i=1}^p r_{ij} z_j$$

GWAS marginal effects are biased due to LD!

• *z*-scores are correlated,

$$\operatorname{Cor}(\hat{z}_i, \hat{z}_k) = r_{ik}, \forall j, k$$

• Recall the previously discussed connection with LDSC

Fine-mapping via RSS model

"Single effect": z_l 's

$$\hat{z} \sim N(\hat{R}z, \hat{R})$$
 $z = \sum_{l=1}^{L} z_l$
 $z_l = \gamma_l z_l$
 $z_l \sim N(0, \omega_l^2)$
 $\gamma_l \sim \mathsf{Mult}(1, \pi)$

Z z_1 z_2 z_3 \equiv + + Suggested reading: Zou et al (2022) PLoS Genet.

Alternative models of GWAS summary statistics

The \hat{z} model:

5

The \hat{b}, \hat{s} model:

$$\hat{z} \sim N(Rz, R)$$

 $\hat{\boldsymbol{b}}|\hat{\mathbf{s}} \sim N(\hat{\mathbf{S}}\hat{\boldsymbol{R}}\hat{\mathbf{S}}^{-1}\boldsymbol{b},\hat{\mathbf{S}}\hat{\boldsymbol{R}}\hat{\mathbf{S}})$

- Both models can be easily written as SuSiE regression
 - \hat{z} model: lower MAF variants have larger effects
 - \hat{b}, \hat{s} model: effect sizes are the same regardless of MAF
- Fine-mapping using \hat{z} model: CAVIAR, FINEMAP
- Fine-mapping using \hat{b}, \hat{s} model: DAP-G

7



Impact of allele flips

What is allele flip?

- Different allele encoding between GWAS and LD reference
- e.g. AA=0, AC=1, CC=2 in GWAS; AA=2, AC=1, CC=0 in LD reference genotype
- A challenging problem coupled with strand flip, when merging sequence data from different platforms

10

Impact of allele flips



Addressing the allele flip challenge

- susieR::susie_rss() function implements a diagnosis
- bigsnpr::snp_match() function implements a basic allele matching for two sets of summary statistics
- Other resources
 - Allele flip illustration: https://statgen.us/ lab-wiki/compbio_tutorial/allele_qc
 - A powerful, multi-set data merger (by Yin Huang): https://cumc.github.io/xqtl-pipeline/ pipeline/misc/summary_stats_merger.html





Impact of mis-matched LD reference: credible sets







Benner et al. (2017) Am. J. Hum. Genet.

Fine-mapping in meta-analysis: overview



Kanai et al. (2022) Cell Genomics

16

15



Kanai et al. (2022) Cell Genomics

17

Fine-mapping in meta-analysis: diagnosis



Kanai et al. (2022) Cell Genomics

20

Fine-mapping in meta-analysis: diagnosis





Covariate adjustment in LD reference

Consider two GWAS regression analysis:

- 1. Evaluate SNP effect in Trait \sim SNP+Age+Sex+PCs
- Fit model Trait ~ Age+Sex+PCs, compute residual of Trait (remove covariates), and evaluate SNP effect in model Residual_Trait ~ SNP

Are these two analysis equivalent?

Covariate adjustment in LD reference

Consider two GWAS regression analysis:

- 1. Evaluate SNP effect in Trait \sim SNP+Age+Sex+PCs
- 2. Fit model Trait \sim Age+Sex+PCs, compute residual of Trait (remove covariates), and evaluate SNP effect in model Residual_Trait \sim SNP

They are not equivalent because covariates should also be removed from SNP data: Residual_Trait \sim Residual_SNP

Covariate adjustment in LD reference

Covariates should be removed from genotype before computing LD reference for fine-mapping





21

Fine-mapping with summary statistics: current methods and practical considerations

Gao Wang, Ph.D.

Advanced Gene Mapping Course, November 2022

The Gertrude H. Sergievsky Center and Department of Neurology Columbia University Vagelos College of Physicians and Surgeons



Figure: Benner et al. (2017) Am. J. Hum. Genet.

Association analysis summary statistics

z-scores from univariate association studies:

$$\hat{z}_j \coloneqq \hat{\beta}_j / s_j,$$

where

$$\hat{\beta}_j \coloneqq (\mathbf{x}_j^\mathsf{T} \mathbf{x})^{-1} \mathbf{x}^\mathsf{T} \mathbf{y} \quad s_j \coloneqq \sqrt{\hat{\sigma}_j^2 (\mathbf{x}_j^\mathsf{T} \mathbf{x})^{-1}}$$

- **Sufficient** statistics: $x^{\mathsf{T}}x, x^{\mathsf{T}}y, \hat{\sigma}_i^2$
- "Summary" statistics:
 - z-scores: \hat{z}
 - Genotypic correlation: \hat{R}

Reasons to work with summary statistics

Advantage over full data (genotypes and phenotypes):

- Easier to obtain and share with others
- Convenient to use: QC and data wrestling barely needed
- Computationally suitable for large-sample fine-mapping
 - $\mathcal{O}(p^2)$ (summary statistics) $\ll \mathcal{O}(np)$ (full data)
 - when sample size $n \gg$ variants in fine-mapped region p

Suggested reading: Pasaniuc and Price (2017) Nat. Rev. Genet.

3

$$\hat{z} \sim N(\hat{R}z, \hat{R})$$

Assumptions:

- $1. \ \mbox{Heritability of any single SNP is small}$
- 2. \hat{R} is sample genotypic correlation for the same study
- 3. Genotypes used to computed z and \hat{R} are accurate

Properties of per SNP z scores

 z-score for a SNP depends on effects of both itself and other correlated SNPs:

$$\mathsf{E}(\hat{z}_j|\hat{\boldsymbol{R}}) = \sum_{i=1}^p r_{ij} z_j$$

GWAS marginal effects are biased due to LD!

• *z*-scores are correlated,

$$\operatorname{Cor}(\hat{z}_i, \hat{z}_k) = r_{ik}, \forall j, k$$

• Recall the previously discussed connection with LDSC

Fine-mapping via RSS model

"Single effect": z_l 's

$$\hat{z} \sim N(\hat{R}z, \hat{R})$$

 $z = \sum_{l=1}^{L} z_l$
 $z_l = \gamma_l z_l$
 $z_l \sim N(0, \omega_l^2)$
 $\gamma_l \sim Mult(1, \pi)$

 Z
 Z1
 Z2
 Z3

 Image: Suggested reading:
 Image: PLoS Genet.
 Image: PLoS Genet.

Alternative models of GWAS summary statistics

The \hat{z} model:

5

The \hat{b}, \hat{s} model:

$$\hat{z} \sim N(Rz, R)$$

The *b*, s model.

 $\hat{\boldsymbol{R}}$) $\hat{\boldsymbol{b}}|\hat{\boldsymbol{s}} \sim N(\hat{\boldsymbol{S}}\hat{\boldsymbol{R}}\hat{\boldsymbol{S}}^{-1}\boldsymbol{b},\hat{\boldsymbol{S}}\hat{\boldsymbol{R}}\hat{\boldsymbol{S}})$

- Both models can be easily written as SuSiE regression
 - \hat{z} model: lower MAF variants have larger effects
 - \hat{b}, \hat{s} model: effect sizes are the same regardless of MAF
- Fine-mapping using \hat{z} model: CAVIAR, FINEMAP
- Fine-mapping using \hat{b}, \hat{s} model: DAP-G

7



Impact of allele flips

What is allele flip?

- Different allele encoding between GWAS and LD reference
- e.g. AA=0, AC=1, CC=2 in GWAS; AA=2, AC=1, CC=0 in LD reference genotype
- A challenging problem coupled with strand flip, when merging sequence data from different platforms

10

Impact of allele flips



Addressing the allele flip challenge

- susieR::susie_rss() function implements a diagnosis
- bigsnpr::snp_match() function implements a basic allele matching for two sets of summary statistics
- Other resources
 - Allele flip illustration: https://statgen.us/ lab-wiki/compbio_tutorial/allele_qc
 - A powerful, multi-set data merger (by Yin Huang): https://cumc.github.io/xqtl-pipeline/ pipeline/misc/summary_stats_merger.html





Impact of mis-matched LD reference: credible sets







Benner et al. (2017) Am. J. Hum. Genet.

Fine-mapping in meta-analysis: overview



Kanai et al. (2022) Cell Genomics

16

15



Kanai et al. (2022) Cell Genomics

17

Fine-mapping in meta-analysis: diagnosis



Kanai et al. (2022) Cell Genomics

Е 3 DENTIST-S outlier varian Ancestry AFR AMR EAS FIN NFE SAS Meta • 1,000 • 5,000 • 10,000 • 50,000 2 Missingness O Both ex

Fine-mapping in meta-analysis: diagnosis





Covariate adjustment in LD reference

Consider two GWAS regression analysis:

- 1. Evaluate SNP effect in Trait \sim SNP+Age+Sex+PCs
- 2. Fit model Trait \sim Age+Sex+PCs, compute residual of Trait (remove covariates), and evaluate SNP effect in model Residual_Trait \sim SNP

Are these two analysis equivalent?

Covariate adjustment in LD reference

Consider two GWAS regression analysis:

- 1. Evaluate SNP effect in Trait \sim SNP+Age+Sex+PCs
- 2. Fit model Trait \sim Age+Sex+PCs, compute residual of Trait (remove covariates), and evaluate SNP effect in model Residual_Trait \sim SNP

They are not equivalent because covariates should also be removed from SNP data: Residual_Trait \sim Residual_SNP

Covariate adjustment in LD reference

Covariates should be removed from genotype before computing LD reference for fine-mapping





21

Integrating GWAS with functional annotations

Gao Wang, Ph.D.

Advanced Gene Mapping Course, November 2022

The Gertrude H. Sergievsky Center and Department of Neurology Columbia University Vagelos College of Physicians and Surgeons



GWAS variants catelog by functional annotations



3

1

Functional enrichment in fine-mapped variants



Figure: Huang et al. (2017) Nature

2

Functional annotation in aggregated rare variant association analysis

Functional annotation filters in aggregated tests

Aggregated tests are sensitive to (mis-)classification of functional variants. Different sets can be evaluated in practice:

- Loss of function: start-loss, stop-gain, splice sites
- Damaging missense: start-loss, stop-gain, splice sites, nonsynonymous with REVEL score > 0.5
 - Ioannidis et al (2016) AJHG
- All: start-loss, stop-gain, splice sites, nonsynonymous

Annotations integrated to aggregated tests



Also see Li et al. (2022) Nature Methods

Annotations integrated to aggregated tests



Figure: Li et al. (2020) Nature Genetics

5

Chi-square GWAS statistic of variant j $\mathsf{E}[\chi_j^2] = 1 + \frac{Nh_g^2}{M} l_j \qquad \text{LD score of variant j}$ Total number of variants

A polygenic model: stratified LD score regression

$l_j = \sum_{k \neq j} r_{jk}^2 \quad \begin{array}{c} \text{LD score: sum of squared Pearson's} \\ \text{correlation coefficient between SNP j} \\ \text{and other (neighboring) SNPs} \end{array}$

7

A polygenic model: stratified LD score regression

Functional annotation in

analysis

common variant association



- Perform LDSC restricted to a functional category
- Enrichment: The proportion of SNP-heritability in the category divided by the proportion of SNPs

7

Cell-type enrichment in GWAS traits via S-LDSC



Integration approaches

- Integrate directly as range based binary annotations
 - Finucane et al (2015) Nature Genetics Stratified LDSC paper
- Possibility to work with variant specific continuous annotations
 - Gazal et al (2017) Nature Genetics
- Compute variant level annotations from epigenomic feature ranges
 - Deep Learning methods
 - Zhou et al (2015) Nature Genetics, Zhou et al (2018) Nature Genetics

9

A sparse model (a somewhat oligogenic view)

Generalized linear model for SNP effects given K annotations

$$\beta_j = (1 - \pi_j)\delta_0 + \pi_j g(\Theta)$$
$$\pi_j := \Pr(\gamma_j = 1 | \boldsymbol{\alpha}, \boldsymbol{d})$$
$$\log\left[\frac{\pi_j}{1 - \pi_j}\right] = \alpha_0 + \sum_{k=1}^K \alpha_k d_k$$

- α are log fold enrichment of functional genomic features
 - Suggested reading: Wen (2016) AoAS

Enrichment of DNase I in GTEx eQTLs



Figure: Wen et al. (2016) AJHG

Integrative fine-mapping with functional annotations



Integrating functional information prioritizes the left SNP.

12

Recall the toy example

Probability of association assuming one effect variable,

$$\frac{\mathsf{LR}_1}{\mathsf{LR}_1 + \mathsf{LR}_2} = 0.87 \quad \frac{\mathsf{LR}_2}{\mathsf{LR}_1 + \mathsf{LR}_2} = 0.13$$

Recall the toy example

Probability of association assuming one effect variable,

$$\frac{\mathsf{LR}_1}{\mathsf{LR}_1 + \mathsf{LR}_2} = 0.87 \quad \frac{\mathsf{LR}_2}{\mathsf{LR}_1 + \mathsf{LR}_2} = 0.13$$

What if we determine *a priori* that SNP 1 is **twice as important** as SNP 2?

$$\frac{2 \times LR_1}{2 \times LR_1 + LR_2} = 0.93 \quad \frac{LR_2}{2 \times LR_1 + LR_2} = 0.07$$

13

Fine-mapping with functional annotations

Recall the BVSR model

$$y = Xb + e$$

$$e \sim N(0, \sigma^2 I_n)$$

$$\gamma_j \sim \text{Bernoulli}(\pi)$$

$$b_{\gamma} | \gamma \sim g(\cdot)$$

$$b_{-\gamma} | \gamma \sim \delta_0$$

Key idea: π , prior inclusion probability, can be modelled by **enrichment** of functional annotations

13

Genome-wide approach with S-LDSC

- A single locus may not have enough power to leverage annotation enrichment
- Genome-wide evaluation of thousands of annotations can increase power of fine-mapping
 - Lead to new loci to discover
- Functional enrichment can be done under the same framework
 - Prioritize genomic features / tissues / cell-types
- Enrichment coefficient may be transferrable cross population
 - Weissbrod et al. (2021) medrxiv

15

Functionally informed fine-mapping in UK Biobank

In analyses of 49 UK Biobank traits, PolyFun + SuSiE identified >32% more fine-mapped variant-trait pairs compared to using SuSiE alone.



Figure: Weissbrod et al. (2020) Nat. Genet.



Complex phenotype prediction and transcriptome-wide association studies

Gao Wang, Ph.D.

Advanced Gene Mapping Course, November 2022

The Gertrude H. Sergievsky Center and Department of Neurology Columbia University Vagelos College of Physicians and Surgeons

- **1** Rationale and assumptions
- 2 Univariate TWAS methods (credits: Haky Im @ UChicago)

2

3

- **3** Multivariate TWAS methods
- Onnections between TWAS and fine-mapping, colocalization and Mendelian Randomization

Motivation: eQTLs are enriched in GWAS signals



Figure: Heinig (2018) Front. Cardiovasc. Med.

Transcriptome-wide association study (TWAS)	TWAS challenge: association vs causality
Contributions of <u>multiple</u> genetic variants to complex traits through their <u>impact</u> on molecular phenotypes	SNP GE TRAIT A. GE independent of trait Well-controlled: Supp. Table S9
Reference Panel Supplicit A = 0 = 0 0 A = 0 = 0 0 Individual TWAS 0 Summary-based TWAS Individual TWAS Inditin Transformation <th>SNP GE TRAIT B. Trait independent of GE Well-controlled: Supp. Table S9</th>	SNP GE TRAIT B. Trait independent of GE Well-controlled: Supp. Table S9
	SNP GE TRAIT C. All independent Subsumed by (A and B)
	SNP GE TRAIT D. Trait effects GE independently of SNPs With cis-GE component, equivalent to (C)
Figure: Gusev et al. (2016) Nat. Genet.	Figure: Gusev et al. (2016) Nat. Genet.
4	5

TWAS challenge: association vs causality



TWAS challenge: technical considerations

Ideal TWAS setup

- Homogenous population
- Tissue and cell-type specific
- Training data-set is large and complete (N > 200)

But in reality

- Cross population TWAS aplications
- Multiple tissue and cell-types
- Availability of individual level data vs summary statistics

7



Figure: Zhu and Zhou et al. (2020) Quantitative Biology

Univariate TWAS methods (credits: Haky Im @ UChicago)



These methods can also be used for Polygenic Risk Score (PRS) calculations

Simple regression method

LETTERS

Common polygenic variation contributes to risk of schizophrenia and bipolar disorder

The International Schizophrenia Consortium*

$$Y = \sum_{k=1}^{M} \hat{\beta}_{k}^{\text{GWAS}} X_{k}$$



10

9

Ridge regression / BLUP **REPORT** GCTA: A Tool for Genome-wide Complex Trait Analysis Jan Yang, 1* S. Hong Lee, ¹ Michael E. Goddard, ^{2,3} and Peter M. Visscher¹ AJHG 2011 $Y - \sum_{k=1}^{M} \hat{\beta}_{k}^{Ridge} X_{k}$ $|Y - \sum_{k=1}^{M} \hat{\beta}_{k}^{Ridge} X_{k}$ $|Y - \sum_{k} X_{k} \beta_{k}||_{2} + \lambda_{2} ||\beta_{2}||_{2}$

Other penalized regression

J. R. Statist. Soc. B (2005) 67, Part 2, pp. 301–320

Regularization and variable selection via the elastic net



Bayesian variable selection regression

OPEN ORCESS Freely available online

Polygenic Modeling with Bayesian Sparse Linear Mixed Models

Xiang Zhou¹*, Peter Carbonetto¹, Matthew Stephens^{1,2}*

$$Y = \sum_{k=1}^{M} \beta_k^L X_k + \sum_{k=1}^{M} \beta_k^S X_k + \epsilon$$
$$\beta_k^L \sim N(0, \sigma_L^2)$$
$$\beta_k^S \sim N(0, \sigma_S^2)$$

MultiBLUP: improved SNP-based prediction for complex traits Doug Speed and David J Balding Genome Res. published online June 24, 2014 Access the most recent version at doi:10.1101/gr.169375.113

13

Event RWAS / FUSION Eunctional Summary-based Imputation New! RWAS (Grishin et al.) models for TOGA ATAC-seq New! CONTENT (Thompson et al.) context-specific models for single-cell and bulk expression New! OBLY M models FUSION bia a suite of tools for performing transcriptome-wide and regulame-wide association studies (TWAS and RWAS). FUSION builds predictive models of the genetic component of a functional/holecular phenotype and predicts and tests that component for association with disease using GWAS summy statistics. The geal is to identify associations between a models from multiple studies to facilitate this analysis. Please cite the following manuscript for TWAS methods: Gusev et al. "Integrative approaches for large-scale transcriptome-wide association studies" 2016 Nature Genetics

Choice of methods: cross validation

12

Multivariate TWAS methods

Multivariate TWAS methods overview

Leverage similarity between molecular phenotypes



- UTMOST, Yu et al. (2019) Nature Genetics
- MR-JTI, Zhou et al. (2020) Nature Genetics

Multivariate TWAS hands-on exercise

statgen-setup launch --tutorial twas

Connections between TWAS and fine-mapping, colocalization and Mendelian Randomization

67





Figure: Zhu and Zhou (2022) Quantitative Biology
Multivariate analysis in genetic association studies

Gao Wang, Ph.D.

Advanced Gene Mapping Course, November 2022

The Gertrude H. Sergievsky Center and Department of Neurology Columbia University Vagelos College of Physicians and Surgeons 1 Motivation

1

- 2 Meta-analysis review
- **3** Meta-analysis: a multivariate regression prospective
- **4** Variant colocalization: variable selection in meta-analysis
- **5** Multivariate adaptive shrinkage and fine-mapping

Motivation

Beyond per trait per variant association studies

Statistical fine-mapping (multiple regressors)

• Identify non-zero effect variables by accounting for LD

Meta-analysis (multiple responses)

• Integrate information across multiple conditions / studies

"Causal" variants across multiple conditions?

• Cross-population fine-mapping; colocalization; pleiotropy; mediation; . . .



The problem

For a genetic variable analyzed in two conditions:

P("causal" in trait 1 & 2 | association data for 1 & 2)

5

6

The problem

For a genetic variable analyzed in two conditions:

P("causal" in trait 1 & 2 | association data for 1 & 2)

Denote data as D_1 and D_2 , and use indicator variables γ_1 , γ_2 for variable having effects in 1 and 2, and hyperparameters Θ :

$$P(\gamma_1 = 1, \gamma_2 = 1 | D_1, D_2, \Theta)$$

Multivariate relationships?



Meta-analysis review

Fixed effect and random effects models

Different assumptions on effects across studies

- Fixed effect model: all studies share a common effect size
- Random effects model: effect sizes are random variables *from an underlying distribution*

Fixed effect (FE) model

Let $\hat{\beta}_i$ be the observed effect size of study i, $1 \le i \le k$, and s_i^2 its variance. The true effect size is β . The observed effect is modelled as

$$\hat{\beta}_i \sim N(\beta, s_i^2),$$

with likelihood function

$$L(\beta) = P(\hat{\beta}|\beta) = \prod_{i}^{k} P(\hat{\beta}_{i}|\beta) \propto \prod_{i}^{k} \exp\left[-\sum_{i}^{k} \frac{(\hat{\beta}_{i}-\beta)^{2}}{2s_{i}^{2}}\right].$$

Fixed effect (FE) model

Let $\hat{\beta}_i$ be the observed effect size of study i, $1 \le i \le k$, and s_i^2 its variance. The true effect size is β . The observed effect is modelled as

$$\hat{\beta}_i \sim N(\beta, s_i^2)$$

with likelihood function

$$L(\beta) = P(\hat{\beta}|\beta) = \prod_{i}^{k} P(\hat{\beta}_{i}|\beta) \propto \prod_{i}^{k} \exp\left[-\sum_{i}^{k} \frac{(\hat{\beta}_{i}-\beta)^{2}}{2s_{i}^{2}}\right].$$

Let $w_i = 1/s_i^2$ be the weight of study i. The MLE of summary effect is

$$\hat{eta} = rac{\sum_{i}^{k} w_i eta_i}{\sum_{i}^{k} w_i}$$
 Inverse variance weighting

8

Random effects (RE) model

Let $\hat{\beta}_i$ be the observed effect size of study i, $1 \leq i \leq k$, and s_i^2 its variance. Let β_i be the true effect size of study i. The observed effect is modelled as

$$\hat{\beta}_i | \beta_i \sim N(\beta_i, s_i^2), \quad \beta_i \sim N(\beta, \sigma^2)$$

with likelihood function

$$P(\hat{\boldsymbol{\beta}}|\boldsymbol{\beta},\sigma^2) \propto \prod_{i}^{k} \frac{1}{s_i^2 + \sigma^2} \exp\left[-\sum_{i}^{k} \frac{(\hat{\boldsymbol{\beta}}_i - \boldsymbol{\beta})^2}{2(s_i^2 + \sigma^2)}\right].$$

Random effects (RE) model

Let $\hat{\beta}_i$ be the observed effect size of study i, $1 \leq i \leq k$, and s_i^2 its variance. Let β_i be the true effect size of study i. The observed effect is modelled as

$$\hat{\beta}_i | \beta_i \sim N(\beta_i, s_i^2), \quad \beta_i \sim N(\beta, \sigma^2)$$

with likelihood function

9

$$P(\hat{\boldsymbol{\beta}}|\boldsymbol{\beta},\sigma^2) \propto \prod_{i}^{k} \frac{1}{s_i^2 + \sigma^2} \exp\left[-\sum_{i}^{k} \frac{(\hat{\boldsymbol{\beta}}_i - \boldsymbol{\beta})^2}{2(s_i^2 + \sigma^2)}\right].$$

RE has weight $w_i^* = 1/(s_i^2 + \sigma^2)$; summary effect $\hat{\beta}$ can be similarly computed as FE, replacing w_i with w_i^* . σ^2 can be estimated (e.g. , MLE).

Multivariate model(s) for effect sizes

Meta-analysis: a multivariate regression prospective

Consider a parametric model on effect sizes across studies,

$$B_i | \gamma = 1 \sim MVN(0, U)$$

Consider 2 studies, *e.g.* height GWAS in Europeans and Africans.

Fixed-effect model multivariate analysis

Effect sizes are exactly the same between two studies,

$$U_{\mathsf{fixed}} = \sigma_0^2 \times \begin{bmatrix} 1 & \mathbf{1} \\ \mathbf{1} & \mathbf{1} \end{bmatrix}$$

Random effects model multivariate analysis

Effect sizes are different between two studies, but are from the same distribution,

$$U_{\rm random} = \sigma_0^2 \times \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

12

Other multivariate models

$$\mathcal{U}_{\mathsf{partially shared}} = \sigma_0^2 \times \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$$

where $|\rho| \leq 1.$ This contains the two meta-analysis models as special cases!

Other flexible multivariate models

More generally,

$$U = \begin{bmatrix} \sigma_1^2 & \sigma_{12}^2 \\ \sigma_{12}^2 & \sigma_2^2 \end{bmatrix}$$

- Pro: more generic than $U_{\rm fixed}$ and $U_{\rm random}$
- Con: 3 parameters to deal with, compared to one σ_0^2

73

Analogy to popular multivariate models (some necessary but, not sufficient)

• Colocalization correlation matrix:

$$\begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$$

• Condition specific correlation matrix:

$$\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$

Analogy to popular multivariate models (some necessary, but not sufficient)

• Mediation:

15

$$U_{\text{mediation}} = \sigma_0^2 \times \begin{bmatrix} 1 & \rho_{12} \\ \rho_{12} & \rho_2 \end{bmatrix}$$

- Genotype \rightarrow Trait 1 \rightarrow Trait 2.
- Effect on trait 2 should be smaller than that on trait 1.

The problem

Variant colocalization: variable selection in meta-analysis

For a genetic variable analyzed in GWAS and eQTL studies:

$$P(\gamma_g = 1, \gamma_e = 1 | D_g, D_e, \Theta)$$

17

Colocalization method: *coloc*

coloc [Giambartolomei *et al.* (2014) PLoS Genet.]

- On X: "one causal" assumption
- On Y: the null + 4 combinations given "one causal"
 - 1. In 1 but not 2
 - $2. \ \text{In} \ 2 \ \text{but not} \ 1 \\$
 - 3. In 1 and 2 but not the same variable
 - 4. In 1 and 2 and the same variable (colocalization)
 - 5. No association in both data 1 and 2 $\,$

Colocalization method: eCAVIAR

eCAVIAR [Hormozdiari et al. (2016) Am. J. Hum. Genet.]

- On X: multiple effect variables
- On Y: each effect variable can be
 - 1. In 1 but not 2
 - 2. In 2 but not 1
 - 3. In both 1 and 2
 - 4. No association in both data 1 and 2

eCAVIAR effects assumption	Colocalization method: <i>enloc</i>
Effect sizes are independent, $U = \begin{bmatrix} \sigma_g^2 & 0\\ 0 & \sigma_e^2 \end{bmatrix}$	<i>enloc</i> [Wen <i>et al.</i> (2017) PLoS Genet.] • Key difference: cross-condition effects not independent • eQTL signals are enriched in GWAS
20	21

Colocalization method: enloc

enloc [Wen et al. (2017) PLoS Genet.]

- Key difference: cross-condition effects not independent
- eQTL signals are enriched in GWAS

But how?

• Recall fine-mapping with functional annotation for *j*

$$\log\left[\frac{\pi}{1-\pi}\right] = \alpha_0 + \alpha \gamma_e$$

and in this context

$$\pi := P(\gamma_g = 1 | \gamma_e = 1$$

21

enloc two step procedure

- 1. Obtain $P(\gamma_g = 1)$ and $P(\gamma_e = 1)$ using fine-mapping
- 2. Fit the enrichment model via multiple imputation

22

Connections among colocalization methods	Connections among colocalization methods
 eCAVIAR is a special case of enloc with α = 0. coloc is a special case of "one causal" fine-mapping based enloc with fixed, high(!) α value by default. Recent coloc extension: coloc version 5, aka SuSiE-coloc removed the "one causal" assumption. Wallace (2021) PLoS Genetics https://chr1swallace.github.io/coloc/ 	 eCAVIAR is a special case of enloc with α = 0. coloc is a special case of "one causal" fine-mapping based enloc with fixed, high(!) α value by default. Recent coloc extension: coloc version 5, aka SuSiE-coloc removed the "one causal" assumption. Wallace (2021) PLoS Genetics https://chr1swallace.github.io/coloc/ Summary: pattern and scale of effect size correlations, represented as different prior models.
23	

Practical considerations

- Choice of prior
 - Best to estimate enrichment α from data
 - + $\alpha \in [0,5]$ suggested by $>4,000~{\rm GWAS}$ + GTEx data
- LD reference mismatch: underestimate α , thus power loss

Hukku et al. (2021) Am. J. Hum. Genet.

Multivariate adaptive shrinkage and fine-mapping



A naive mixture model

"FE and RE are equally likely for any variant":

$$U_{mixed} = 0.5 \times \begin{bmatrix} \sigma_0^2 & \sigma_0^2 \\ \sigma_0^2 & \sigma_0^2 \end{bmatrix} + 0.5 \times \begin{bmatrix} \sigma_0^2 & 0 \\ 0 & \sigma_0^2 \end{bmatrix}$$

Prior allows for possibility of both; data will determine where posterior lands.

A data-adaptive mixture model

Instead of making assumptions, can we learn from data:

• What are the latent structures for multivariate effects?

28

30

• How often does each structure appear?

and use these to construct the mixture model?

27

Patterns of sharing: factor analysis







Incorporating all possible patterns

Multivariate effects of a variant follows the *k*-th pattern with probability π_k :

$$U_{mixed} = \pi_1 \times \begin{bmatrix} 2.4 & 0.3 \\ 0.3 & 1.5 \end{bmatrix} + \pi_2 \times \begin{bmatrix} 1.6 & 0.001 \\ 0.001 & 0.02 \end{bmatrix} + \pi_3 \times \cdots$$

This is the Multivariate Adaptive Shrinkage Prior.

- Step 1: estimated π_k via EM algorithm using data across genome.
- Step 2: apply this prior to each variant in association mapping.



Application to multivariate fine-mapping



Figure: mvSuSiE fine-mapping with adaptive shrinkage model

32







- "The voluntary consent of the human subject is absolutely essential..."
- "The experiment should be conducted as to avoid all unnecessary physical and mental suffering and injury..."
- No experiment should be conducted where there is an a priori reason to believe that death or disabling injury will occur; except, perhaps, in those experiments where the experimental physicians also serve as subjects.
- "During the course of the experiment, the human subject should be at liberty to bring the experiment to an end if he has reached the physical or mental state where continuation of the experiment seems to him to be impossible."
- During the course of the experiment the scientist in charge must be prepared to terminate the experiment at any stage, if he has probable cause to believe...that a confinuation of the experiment is likely to result in injury, disability, or death to the experimental science.







3





Office for Human Research Protections

- 45 CFR 46 Subpart A ('Common Rule')
- Subpart B (Pregnant Women, Fetuses, and Nonviable/Questionable Viable Neonates),
- Subpart C (Prisoners),
- Subpart D (Minors)

Food & Drug Administration

(jurisdiction: clinical investigations of drugs, devices, biologics)

hirip

- 21 CFR 50: Protection of Human Subjects
- 21 CFR 56: Institutional Review Boards
- 21 CFR 312: Investigational Drugs
- 21 CFR 812: Investigational Devices

















15





Part II: What does 'designed to develop

or contribute to generalizable

knowledge' mean?

...designed to draw general conclusions:

✓ what we know about what is being tested is not

and ✓ the activity is not dependent on the unique

yet firmly established or accepted;

characteristics of the target population or system in

which it will be implemented







Exemption #4: Secondary research uses of identifiable private information or identifiable biospecimens can be exempt under this category, if at least one of the following criteria is met:

h r p

21











25



27







28

§ 46.111 Criteria for IRB approval of research.

(a) In order to approve research covered by this policy the IRB shall determine that all of the following requirements are satisfied:

h|r|p

hrp





- Who is included, who is excluded? Does it make scientific sense? Ethical sense?
- If applicable: Are children in a study involving a test article that hasn't first been tested in adults?
 Pregnant women before non-pregnant women?
- Costs or compensation that may impact 'fairness'
- Screening and recruitment?
- What about non-English speakers?



33

5. Informed consent will be appropriately documented or appropriately waived in accordance with §46.117

If not: Does the research meet one of the allowable criteria to waive documentation?





32

4. Informed consent will be sought from each prospective subject or the subject's legally authorized representative, in accordance with, and to the extent required by, §46.116

If not:

Are **ALL** the criteria for waiving informed consent or for altering/excluding specific elements of informed consent met?

hrp CONSULTING GROUP

hirid

- 6. When appropriate, the research plan makes adequate provision for monitoring the data collected to ensure the safety of subjects
- What data will be monitored for safety purposes? When? How?
- Who will be responsible for evaluating safety data? Is a DSMB needed?
- Stopping Rules?
- Communication plan of findings to investigators and IRBs (from the IRB of Record or Sponsor)



7. When appropriate, there are adequate provisions to protect the privacy of subjects...

Consider:

- Settings where recruitment, consent, and research procedures and interactions will occur
- Provisions to ensure privacy for each of the above
- · Provisions to ensure privacy when contacting or soliciting information from subjects



37

A closer look at data security: minimize the risk of disclosure or breach of data

- Obtaining the data What is the sensitivity of the data? Are all the data points that will be accessed or gathered for the research necessary to achieve the objectives of the research?
- Recording the data

What (if any) identifiers, including codes, will be recorded for the research?

- Storing the data
 - Where will paper research records, including signed consent forms, be stored? How will paper records be kept secure and restricted to authorized project personnel?

 - Where will the electronic research data be study be stored (Upiversity-provided database application like REDCap, IT hile server, etc.)? If there a key that links code numbers to identifiers, that list should be kept separate from the coded data, including copies of signed informed consent forms. Additionally, access to that list/key must be restricted authorized research personnel. h|r|p

39

...and to protect the confidentiality of subject data

General:

- How will the data/biospecimens be stored?
- If identifiers will be removed or replaced, is there a possibility that such information/biospecimens could be reidentified?
- Will the data/biospecimens be shared/transmitted/ transferred to a third party or otherwise disclosed or released? How?
- Is there a potential risk of harm to individuals if the data/biospecimens are lost, stolen, compromised, or otherwise used in a way contrary to the parameters of the study?

D

• Plans for data retention and destruction?

38



40



And (111.b) When some or all of the subjects are likely to be vulnerable to coercion or undue influence, such as children, prisoners, individuals with impaired decision-making capacity, or economically or educationally disadvantaged persons, additional safeguards have been included in the study to protect the rights and welfare of these subjects. (set aside issues with children, pregnant women/fetuses, prisoners,

regulations for which are codified in the Common Rule subparts---more on that in a moment)

· What are some considerations when determining if additional safeguards are necessary and sufficient?

• Examples:

- For economically disadvantaged...is there payment? What is the amount? schedule?
- For educationally disadvantaged...is the consent process particularly simplified? Should there be a witness to the consent process?





Yale From cross-phenotype associations to pleiotropy in human genetic studies Andrew DeWan, PhD, MPH Associate Professor of Epidemiology Director, Yale Center for Perinatal, Pediatric and Environmental Epidemiology Director, Yale Center for Perinatal, Pediatric and Environmental Epidemiology Work done in collaboration with Yasmmyn Salinas, PhD, MPH, Assistant Professor of Epidemiology Yale School of Public Health Work done in collaboration with Yasmmyn Salinas, PhD, MPH, Assistant Professor of Epidemiology Yale School of Public Health Yale School of Public Health

1

<section-header><section-header><section-header><section-header><section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item>

2



3

Examples in humans

- · Marfan syndrome
 - FBN1 (fibrillin-1)
 - thinness, joint hypermobility, limb elongation, lens dislocation, and increased susceptibility to heart disease.
- Holt-Oram syndrome,
 - TBX5 (transcription factor)
 - cardiac and limb defects
- Nijmegen breakage syndrome
 - NBS1 (DNA damage repair protein)
 - microcephaly, immunodeficiency, and cancer predisposition



Pleiotropy and complex disease comorbidity

- Examples of correlated (comorbid) disease
 - Obesity, hypertension, dyslipidemia, type 2 diabetes (metabolic disorder)
 - Depression, anxiety, personality disorders (psychiatric disorder)
 - Asthma, obesity (pro-inflammatory conditions)
- Why do certain disease occur together
- Causality
- Shared environmental risk factors
- Shared genetic risk factors





Pleiotropy and complex disease comorbidity

- Detecting shared genetics and/or molecular pathways between comorbid diseases can help us understand exactly how the etiology of the diseases overlap
- Etiologic overlaps:
 - provide opportunities for novel interventions that prevent or treat the comorbidity, rather than preventing/treating each disease separately
 - facilitate drug repurposing (that is, known drugs targeting a pleiotropic locus may be repurposed to treat other diseases controlled by that locus, precluding the need for the development and testing of a brand-new drug)

9

Pleiotropy in gene mapping

- Mapping a single genotype to multiple phenotypes has the potential to uncover novel links between traits or diseases
- It can also offer insights into the mechanistic underpinnings of known comorbidities
- It can increase power to detect novel associations with one or more phenotypes

Pleiotropy and complex disease comorbidity

- Pleiotropy-informed analyses consider multiple phenotypes together and take into account the correlation between the phenotypes
 - Analyzing multiple correlated phenotype (e.g. comorbid diseases) is equivalent to analyzing a single narrowly-defined phenotype with low heterogeneity

8

Abundant Pleiotropy in Human Complex Diseases and Traits

Shanya Sivakumaran,^{1,6} Felix Agakov,^{1,2,6} Evropi Theodoratou,^{1,6} James G. Prendergast,³ Lina Zgaga,^{1,4} Teri Manolio,⁵ Igor Rudan,¹ Paul McKeigue,¹ James F. Wilson,¹ and Harry Campbell^{1,*} The American Journal of Human Genetics *89*, 607–618, November 11, 2011

The American journal of Human Genetics 89, 607–618, November 11, 2011

	Genes		SNPs			
Disease Class	Pleiotropic (%)	Nonpleiotropic (%)	p Value*	Pleiotropic (%)	Nonpleiotropic (%)	p Value*
All (comparison group)	233 (16.9)	1147 (83.1)	-	77 (4.6)	1610 (95.4)	-
Immune-mediated phenotypes	106 (37.7)	175 (62.3)	< 0.0001	31 (8.3)	343 (91.7)	0.0066
Cancer	49 (34.8)	92 (65.2)	< 0.0001	8 (4.8)	158 (95.2)	0.8456
Metabolic syndrome	79 (28.5)	198 (71.5)	< 0.0001	30 (8.4)	327 (91.6)	0.0056
	_	CONN. , FADOL GONT, THATA	N SCOMA ST	CONES		
				Managed all stands and all stands		
	No.	CYERL ANDE		Andler Triggoenides Facts (ccoll HDL Cholesterol		
	Ret	Indiversal aneuryon Paranta aneuryon Paranta aneuryon Paranta aneuryon		HDL Cholesterol radio		

Souce ICAN 1 (CAN1) Patronin's deces (URRC) Serum magnesium (URRC) Menanthe lage at innet[(PLC) Heigh (SLC22AR)







15





- The overall approach is to:
- obtain univariate association p-values for each phenotype
- declare CP associations at genetic loci that are statistically significantly associated with each phenotype



















	Simulation s	scenari	OS
# traits associated	h _i ²	r _{Y1,Y2}	Pj
1	h ₁ ² =0.1%,h ₂ ² =0%	[-0.9,0.9]	P1 = P2 = 10%
			P1 = P2 = 20%
			P1 = 10%, P2 = 20%
			P1 = 20%, P2 = 10%
2	h ₁ ² = h ₂ ² = 0.1%	[-0.9,0.9]	P1 = P2 = 10%
			P1 = P2 = 20%
			P1 = 10%, P2 = 20%
			P1 = 20%, P2 = 10%
2	h ₁ ² = 0.1%,h ₂ ² = 0.05%	[-0.9,0.9]	P1 = P2 = 10%
			P1 = P2 = 20%
			P1 = 10%, P2 = 20%
			P1 = 20%, P2 = 10%



Problem: CP associations need not be indicative of pleiotropy































39



Mediation analysis: Assumptions Typically met in genetic epi studies! There must be no unmeasured: . Total Effect confounders of the total effect θ1 A confounders of the relationship Direct Effect between SNP A and the mediator M θ2 confounders of the relationship (M)between mediator M and phenotypic outcome Y Indirect Effect

















Guidelines for generating robust





Empirical searches for pleiotropic loci for asthma and obesity



		0	pu		
<pre>> med.fit<-glm(W1~rs1_2, data > out.fit<-glm(W2~W1+rs1_2, c > med.out<-mediate(med.fit,ou > summary(med.out)</pre>	=combined, fam lata=combined, t.fit, treat="rs1_2	ily=binomial(family=binom ", mediator='	"logit")) iial("logit")) 'W1", boot=	TRUE, boot.ci.type="bca", sim:	s=100
Causal Mediation Analysis					
Nonparametric Bootstrap Confi	dence Intervals	with the BCa	Method		
	Estimate	95% CI Lowe	r 95% CI U	pper p-value	
ACME (control)	0.02152	0.01823	0.03	<2e-16 ***	
ACME (treated)	0.02199	0.01868	0.03	<2e-16 ***	
ADE (control)	0.00723	0.00415	0.01	<2e-16 ***	
ADE (treated)	0.00771	0.00443	0.01	<2e-16 ***	
Total Effect	0.02022	0.02461	0.03		
Prop. Mediated (control)	0.73634	0.65429	0.84	<2e-16 ***	
Prop. Mediated (treated)	0.75247	0.67272	0.85	<2e-16 ***	
ACME (average)	0.02175	0.01847	0.03	-20-16 ***	
ADE (average)	0.00747	0.00426	0.01		
Dran Madiated (sugrama)	0.74441	0.66254	0.84	<2e-10 ***	













				Res	ults fo	r asthr	na		
Table 4. E	vidence for	Associati	ion of <i>Pl</i>	KCA with A	thma in Costa Rica and CAMP Number of Informative Families ¹ (number of offspring with 0/1 recoded genotype)				
Locati Marker (BP)*	Location (BP)*	Minor Allele	CR	САМР	CR	САМР	Costa Rica p Value ^{c.d}	CAMP Replication p Value ^{cal} (two-sided)	Joint p Value" (CR, CAMP two-sided)
rs732191	61779673	G	0.46	0.35	168 (117/51)	141 113/43	-0.0194	-0.0214 (-0.0428)	0.0036 (0.0067
rs9895580	61789701	С	0.47	0.35	168 (117/51)	141 114/43	-0.0171	-0.0160 (-0.0320)	0.0025 (0.0047
rs4411531	61793662	Α	0.29	0.12	88 (70/18)	25 (24/1)	-0.0058	-0.0058 (-0.0117)	0.0004 (0.0007
rs8080771	61824330	G	0.46	0.35	164 (116/48)	108 (90/29)	-0.0161	-0.0070 (-0.0140)	0.0011 (0.0021
rs11652956	61839798	G	0.29	0.12	83 (65/18)	23 (22/1)	-0.0101	-0.0111 (-0.0222)	0.0011 (0.0021
rs7221968	61848731	С	0.27	0.11	79 (63/16)	18 (17/1)	-0.0122	-0.0216 (-0.0432)	0.0024 (0.0045
rs7405806	61862056	А	0.49	0.31	164 (109/55)	90 (77/20)	-0.0309	-0.0009 (-0.0018)	0.0003 (0.0006
rs11079657	61862528	A	0.38	0.23	129 (94/35)	60 (56/8)	-0.0092	-0.0002 (-0.0004)	$2.6 \times 10^{-5_{44}}$ (5.0 × 10 ^{-5_{44}})
									+ +





Study design Two parts: Genome-wide search for cross-phenotype associations with asthma and body mass index Follow-up mediation analysis to dissect genome-wide significant CP associations

62



CC in PLINK CC in PLINK Estimation of genetic correlation using BOLT-REML Univariate association analyses using linear mixed effects models in BOLT-LMM Search for overlapping signals between asthma and BMI Search for overlapping signals between asthma and BMI Assessment of asthma-BMI relationship in the UK Biobank GWA sample Assessment of potential confounders of the asthma-BMI relationship Follow-up mediation analysis in 'mediation' R Package

65

Phenotype definitions

- BMI at baseline (kg/m²):
 - calculated based on height and weight measurements collected by trained UK Biobank staff at the recruitment sites
- Asthma diagnosed prior to baseline (yes/no):
 - ascertained via the question "Has a doctor ever told you that you had asthma?"
 - Note: In mediation analyses, two subgroups were created based on age-at-diagnosis































<figure>











- Mimic randomized trial using genetic data as instruments for exposures
- Leverages information on genetic variants that segregate randomly at conception
- If an association between the instrument and outcome is detected, a causal relationship for this association is strengthened



8

MR Assumptions • The genetic instrument (G) is associated with the exposure (X)

- The genetic instrument is not accepted with any confounder (1)
- The genetic instrument is not associated with any confounder (U) of the exposure-outcome association
- The genetic instrument is conditionally independent of the outcome (Y) given the exposure and confounders



9











One-sample vs. two-sample designs

Assumption/Issue	One-sample	Two-sample
Instrument variable related to risk factor	Weak instrument biases towards the confounded regression result	Weak instrument biases towards the null
Confounders	Can (and should) check this for measured confounders	Not often possible when using summary statistics
Pleiotropy	Multiple methods to explore this issue (including MR-Egger)	Multiple methods to explore this issue (including MR-Egger) and may be more powerful with large consortium datasets since methods tend to be statistically inefficient
Subgroup analyses	Possible if large sample sizes and data on relevant risk factors are available	Only possible if individual level data are available
Bias from adjustments made in GWAS	N/A as all adjustments made in the same set of subjects	Summary data may or may not have been adjusted

14

Selecting genetic variants for an instrument Single or multiple variants Current recommendation is to select variant(s) that are significantly associated with the exposure at the genome-wide level Want a strong genetic instrument to avoid weak instrument bias A single variant or variants with modest effects in small samples are likely to have low power and can suffer from bias

• If selecting multiple variants these should not be in LD and assumes negligible gene-gene interaction among variants

15

Instrument strength

Measured using the F statistic in the regression of the IV on the exposure

$$F = \frac{N-K-1}{K} * \frac{R^2}{1-R^2}$$

 R^2 proportion of the variance of the exposure explained by IV N: sample size

K: number of genetic variants

General Rule: F < 10 is an indication of a weak instrument







17

Testing MR: 2 stage least squares (2SLS)

- Single continuous instrument (GRS)
- Only for one sample method
- Assumes a linear relationship between exposure and outcome
- Regress X on G
- Calculate genetically predicted values of X
- Regress Y on genetically predicted values of X
- Fix the standard errors (e.g. sandwich estimator)

19



21



For each variant calculate the Wald ratio:

 $\widehat{\beta}_j = \frac{\widehat{\Gamma}_j}{\widehat{\gamma}_j}$

Combine into an overall estimate using a

 $\widehat{\beta}_{IVW} = \frac{\Sigma_j \widehat{\gamma}_j^2 \sigma_{Yj}^{-2} \widehat{\beta}_j}{\Sigma_j \widehat{\gamma}_j^2 \sigma_{Yj}^{-2}}$

formula from meta-analysis literature:

- One or two sample designs
- Tends to give more reliable results in the presence of heterogeneity and when using large number of instruments
- Fixed (assumes no heterogeneity across SNP) or random effects meta-analysis

20



22






Body mass index and risk of dying from a bloodstream infection: A Mendelian randomization study

Tormod Rogne^{1,2,3}*, Erik Solligård^{1,3}, Stephen Burgess^{1,5}, Ben M. Brumptong^{6,7,8}, Julie Paulseno¹, Hallie C. Prescotti^{0,11}, Randi M. Mohus^{1,3}, Lise T. Gustadg^{1,12}, Arme Mehl², Bjorn O. Asvold^{6,51}, Andrew T. DeWan^{1,24}, Jan K. Damäso^{1,1,154} PLOS Medicine | https://doi.org/10.1371/journal.pmed.1003413 November 16, 2020

Assess the causal association between BMI and risk of and mortality from BSI by overcoming the limitations of previous observational studies by conducting an MR study in a general population of approximately 56,000 participants in Norway with 23 years of follow-up

25

Study Population

- The Trondelag Health Study (HUNT) is a series of cross-sectional surveys carried out in Nord-Trondelag County, Norway
- 130,000 inhabitants who are representative of the general Norwegian population in terms of morbidity, mortality, sources of income and age distribution
- Based on HUNT2 survey conducted in 1995-1997 with 65,236 participants, 55,908 of whom had complete data for the analysis

26

Characteristic	Total population (n = 55,908)	BSI incidence (n = 2,547)	BSI death (n = 451)
Age (years) ⁵	48.3 (36.5-62.3)	63.6 (52.9-71.4)	67.3 (57.1-74.5)
Male sex*	26,324 (47.1)	1,345 (52.8)	263 (58.3)
BMI (kg/m ²) ^A	26.3 (4.1)	27.7 (4.5)	27.9 (4.8)
Median follow-up time (years) ⁶	21.1 (17.1-21.8)	13.8 (8.4-18.3)	13.3 (7.7-17.9)
Self-reported cancer"	1,955 (3.7)	144 (6.2)	24 (5.9)
Smoking*			
Never	23,594 (43.0)	876 (35.2)	156 (35.6)
Previous	15,133 (27.6)	893 (35.8)	164 (37.4)
Current	16,117 (29.4)	723 (29.0)	118 (26.9)
Physical activity"			
None	3,821 (7.6)	243 (11.9)	54 (15.4)
Slight	15,662 (31.0)	714 (34.9)	117 (33.3)
Moderate	17,167 (34.0)	693 (33.9)	116 (33.1)
High	13,810 (27.4)	397 (19.4)	64 (18.2)
Education*			
≤9 years	19,033 (35.7)	1,305 (55.8)	240 (58.8)
10-12 years	23,468 (44.0)	762 (32.6)	125 (30.6)
≥13 years	10,832 (20.3)	274 (11.7)	43 (10.5)
BMI, body mass index; BSI, bloodstream inf ^mean (standard deviation) ⁹ median (25th-75th percentiles), or ¹ (1%) BSI incidence is based on first occurs	ection. Data are presented as	ed Education defined as follows: <9 vs	ears ("arimum school 7-10 years, continuation
n (s), for inductive is based on mist occar	rence, outerwate, and occurrence is us	tu. Duucation utimtu as ionows. So ji	cars (primary scroot 7-10 years, communitation

collegs. A levels'), and \geq 15 years ('university or other past-secondary education, less than 4 years' and 'university(college 4 years or more'). Activity defines none(''no light or signoron activity'), which coll light activity/work and no vigorous activity'), moderate (" \geq 31 hight activity/work or <1 h vigorous activit or hight (' \geq 1 h vigorous activity)/work').

27

Outcome

- Linked to all prospectively recorded blood cultures at the two community hospitals in the catchment area (Levanger and Namsos Hospitals) as well as St. Olav's Hospital in Trondheim (tertiary referral center)
- Data on blood cultures were available from January 1, 1995 through the end of 2017
- Date of death and emigration out of Nord-Trondelag County were obtained from the Norwegian population registry
- BSI was defined as a positive blood culture of pathogenic bacteria
- BSI mortality was defined as death within 30 days of BSI diagnosis

28

Genetic Instrument

- Based on a BMI meta-analysis of ~700,000 individuals (Yesgo Lat AL (2018) Meta-analysis)
- 939 of 941 SNPs identified as associated with BMI (p<5x10⁻⁸, two SNPs did not pass imputation quality control)
- Genetic risk score (GRS) was calculated for BMI using the --score command in PLINK (version 1.9) and weighted based on the effect estimates from the meta-analysis
- GRS (939 variants) explained 4.2% of the variation in BMI in the population (F-statistic = 2,461)

Analysis Methods

- Fractional polynomial model (suggestion of a nonlinear relationship between BMI and BSI)
- 2-stage least squares (with sandwich estimator) for analyses assuming
- a linear relationship between exposure and outcome
- Sensitivity analyses
- MR Egger (random effects)
 INW
- Weighted median
- 2-sample (using Yengo et al. for SNP-exposure associations)























- 5. Assumptions

INTRODUCTION

Explicitly state assumptions for the main analysis (e.g. relevance, exclusion, independence, homogeneity) as well assumptions for any additional or sensitivity analysis.

37

6. Statistical methods: main analysis Describe statistical methods and statistics used. a) Describe how quantitative variables were handled in the analyses (i.e., scale, units, model). b) Describe the process for identifying genetic variants and weights to be included in the analyses (i.e, independence and model). Consider a flow diagram. c) Describe the MR estimator, e.g. two-stage least squares, Wald ratio, and related statistics Detail the included covariates and, in case of two-sample MR, whether the same covariate set was used for adjustment in the two samples. d) Explain how missing data were addressed.e) If applicable, say how multiple testing was dealt with. 7. Assessment of assumptions Describe any methods used to assess the assumptions or justify their validity. 8. Sensitivity analyses ribe any sensitivity analyses or additional analyses performed. RESULTS 10. Descriptive data r two-sample Mendelian randomization: Provide information on the similarity of the genetic variant-exposure associat between the exposure and outcomes samples. Provide information on extent of sample overlap between the exposure and outcome data sources. d) For two
 i. Pro Pro

38

ii.









- Exposure randomization may not be truly random
- Unmeasured or residual confounding (population stratification, parental genotype associated with outcome)
- · Weak instrument bias resulting from measurement error for the exposure of interest
- Adaptation to the exposure
- · Inconsistent results and selective publication

ent variable analysis: time to call Mendelian randomizatio





Wewcastle

W

Genome-wide association studies (GWAS) - Part 2

More advanced topics:

Linear Mixed Models and G×G or G×E interactions

Heather J. Cordell

Population Health Sciences Institute Faculty of Medical Sciences Newcastle University, UK heather.cordell@ncl.ac.uk

Linear Mixed Models (LMMs)

 Linear Mixed Models have been used for many years in the plant and animal breeding communities

 In the mid 1990s they became popular in the human genetics field, mostly for performing linkage analysis and estimating heritability

GWAS (Part 2

• Using family (pedigree) data i.e. related individuals

Linear Mixed Models (LMMs)

- Linear Mixed Models have been used for many years in the plant and animal breeding communities
- In the mid 1990s they became popular in the human genetics field, mostly for performing linkage analysis and estimating heritability
 - Using family (pedigree) data i.e. related individuals
- In recent years they have become popular in the genetic association studies field for:
 - Testing for association while accounting for varying degrees of relatedness
 - Close family relationships

Heather Cordell (Newcastle)

• Distant relationships and population stratification/substructure

Linear Mixed Models (LMMs)

- Linear Mixed Models have been used for many years in the plant and animal breeding communities
- In the mid 1990s they became popular in the human genetics field, mostly for performing linkage analysis and estimating heritability
 - Using family (pedigree) data i.e. related individuals
- In recent years they have become popular in the genetic association studies field for:
 - Testing for association while accounting for varying degrees of relatedness
 - Close family relationships
 - Distant relationships and population stratification/substructure
 - Estimating the heritability accounted for various partitions of SNPs: • All SNPs typed on a GWAS panel
 - All typed SNPs and others in LD with them

 - · Partitions of SNPs in various functional categories

GWAS (Part 2)

Linear Mixed Models (LMMs)

- Linear Mixed Models have been used for many years in the plant and animal breeding communities
- In the mid 1990s they became popular in the human genetics field, mostly for performing **linkage analysis** and estimating **heritability**
 - Using family (pedigree) data i.e. related individuals
- In recent years they have become popular in the **genetic association** studies field for:
 - Testing for association while accounting for varying degrees of relatedness
 - Close family relationships
 - \bullet Distant relationships and population stratification/substructure
 - Estimating the heritability accounted for various partitions of SNPs:
 All SNPs typed on a GWAS panel
 - All typed SNPs and others in LD with them
 - Partitions of SNPs in various functional categories
 - Investigating genetic correlations between different traits

Heather Cordell (Newcastle) GWAS (Part 2)

Linear Mixed Models (LMMs)

- Linear Mixed Models have been used for many years in the plant and animal breeding communities
- In the mid 1990s they became popular in the human genetics field, mostly for performing linkage analysis and estimating heritability
 Using family (pedigree) data i.e. related individuals
- In recent years they have become popular in the **genetic association** studies field for:
 - Testing for association while accounting for varying degrees of relatedness
 - Close family relationships
 - Distant relationships and population stratification/substructure
 - Estimating the heritability accounted for various partitions of SNPs:
 All SNPs typed on a GWAS panel
 - All typed SNPs and others in LD with them
 - Partitions of SNPs in various functional categories
 - Investigating genetic correlations between different traits

GWAS (Part 2)

Predicting trait values in a new individual

Heather Cordell (Newcastle)

Heather Cordell (Newcastle)

Linear Mixed Models (LMMs)

- A linear mixed model is a statistical model in which the dependent variable is a linear function of both fixed and random independent variables
 - Known respectively as fixed and random effects
 - $\bullet\,$ Fixed effects are considered 'fixed' at their measured values
 - Random effects are considered to be sampled from a distribution
- Recall the usual linear regression model

$$y = mx + c$$
 or $y = \beta_0 + \beta_1 x$

• This model may also be written

Heather Cordell (Newcastle)

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

- y_i refers to the trait value of person i
- x_i refers to the measured value of person *i*'s predictor variable

GWAS (Part 2)

- ϵ_i refers to the displacement from the regression line
- \bullet i.e. the discrepency between the observed and the predicted y value

Linear Regression



Linear Mixed Models (LMMs)

- In linear regression we have $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$
 - Here β_0 and β_1 are fixed effects while ϵ_i is a random error
 - x_i is the 'loading' of the fixed effect that someone has (based on their genotype)
- In matrix notation we can write this model:

• or
$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

Heather Cordell (Newcastle)

Linear Mixed Models (LMMs)

- In linear regression we have y_i = β₀ + β₁x_i + ε_i
 Here β₀ and β₁ are fixed effects while ε_i is a random error
 - x_i is the 'loading' of the fixed effect that someone has (based on their genotype)
- In matrix notation we can write this model:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

• or $\mathbf{y} = \mathbf{X} \boldsymbol{eta} + \boldsymbol{\epsilon}$

er Cordell (N

- A LMM takes the form $\mathbf{y} = \mathbf{X} oldsymbol{eta} + \mathbf{Z} \mathbf{u} + oldsymbol{\epsilon}$
- where u corresponds to a vector of random effects
 with loadings specified in Z

Linear Mixed Models (LMMs)

- E.g. suppose 2 fixed effects β_1 and β_2 , and 3 random effects (plus *n* random errors)
- Then $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}$ corresponds to:

Heather Cordell (Newcastle)

Heather Cordell (Newcastle)

y ₁]	x ₁₁	x ₁₂		<i>z</i> ₁₁	z_{12}	z ₁₃ -]		ϵ_1
<i>y</i> 2		<i>x</i> ₂₁	<i>x</i> ₂₂	[<i>B</i> .]	<i>z</i> ₂₁	<i>z</i> ₂₂	Z23	$\int u_1$	1	ϵ_2
	=	•		$\begin{vmatrix} \rho_1 \\ \beta_2 \end{vmatrix} +$	•			<i>u</i> ₂	+	
yn _		<i>x</i> _{n1}	x _{n2}		_ <i>Z</i> _{n1}	z _{n2}	<i>z</i> _{n3}			ϵ_n

LMMs in genetics

• or $y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + u_1 z_{i1} + u_2 z_{i2} + u_3 z_{i3} + \epsilon_i$

LMMs in genetics

- In genetics we generally work with two equivalent forms of LMM
- One is: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}$
 - The random effect $u_{\rm l}$ corresponds to a scaled additive effect of causal variant (locus) l
 - Assuming many (m) such causal variants all across the genome
 - Z is a standardized genotype matrix i.e. z_{il} takes value

$$\left(\frac{-2f_l}{\sqrt{2f_l(1-f_l)}}, \frac{(1-2f_l)}{\sqrt{2f_l(1-f_l)}}, \frac{2(1-f_l)}{\sqrt{2f_l(1-f_l)}}\right)$$

if individual *i* has genotype (qq, Qq, QQ)

• where f_l is the frequency of allele Q at locus l

GWAS (Part 2)

LMMs in genetics

 The other form is: y = Xβ + g + ε
 Where g_i = Σ^m_{l=1} z_{il}u_l is the total genetic effect in individual i, summed over all the causal loci The other form is: y = Xβ + g + ε
 Where g_i = Σ^m_{l=1} z_{il}u_l is the total genetic effect in individual i, summed over all the causal loci • In this form, g_i can be considered as a random effect operating In this form, g_i can be considered as a random effect operating in individual i in individual i • The vector of random effects **g** takes distribution $\mathbf{g} \sim N(0, \mathbf{G}\sigma_a^2)$ • The vector of random effects ${f g}$ takes distribution ${f g}\sim {\it N}(0,\,{f G}\sigma_a^2)$ $\bullet~$ Where ${\bf G}$ is the genetic relationship matrix (GRM) • Where **G** is the genetic relationship matrix (GRM) between individuals - i.e. their IBD sharing at the causal loci between individuals - i.e. their IBD sharing at the causal loci • $\sigma_a^2 = m\sigma_u^2$ is the total additive genetic variance • $\sigma_a^2 = m \sigma_u^2$ is the total additive genetic variance • G=ZZ'/m • **G**=**ZZ'**/*m* • For family data (close relatives), the expected values of the elements of **G** equal the expected IBD sharing • i.e. twice the kinship coefficients • Thus G is just equal to twice the kinship matrix

• Models their expected relatedness at the causal loci (and elsewhere)

Use of LMMs in genetics

- The formulation $\mathbf{y} = \mathbf{X}\beta + \mathbf{g} + \boldsymbol{\epsilon}$ is known as the Animal Model and has been used extensively in plant and animal breeding
 - Mostly to predict the *breeding values* g_i in order to inform breeding strategies
 - E.g. to increase milk yield, meat production etc. etc.
 - Similar approaches could be used for *prediction* of trait values given genotype data
- In the mid 1990s it became popular in human genetics as the backbone of variance components linkage analysis
- Now commonly used in association analysis (GWAS)
 - To correct for relatedness, when testing for association

Testing for association using LMMs

- Idea is to test a fixed SNP effect β₁
 While including a random effect γ_i that models relatedness
- Fit regression model: $y_i = \beta_0 + \beta_1 x_i + \gamma_i$
 - y is the trait value
 - x is a variable coding for genotype at the test SNP
 - (e.g. an allele count, coded 0, 1, 2 for genotypes 1/1, 1/2, 2/2)

GWAS (Part 2

• $\gamma_i = \mathbf{g}_i + \epsilon_i$

Heather Cordell (Newcastle)

Heather Cordell (Newcastle) GWAS (Part 2)

× /

Testing for association using LMMs

- \bullet Idea is to test a fixed SNP effect β_1
 - While including a random effect γ_i that models relatedness
- Fit regression model: $y_i = \beta_0 + \beta_1 x_i + \gamma_i$
- y is the trait value

Heather Cordell (Newcastle)

- x is a variable coding for genotype at the test SNP
- (e.g. an allele count, coded 0, 1, 2 for genotypes 1/1, 1/2, 2/2) • $\gamma_i = g_i + \epsilon_i$
- We assume $\gamma \sim MVN(0,\,{\rm V})$ where variance/covariance matrix ${\rm V}$ follows standard variance components model
 - Variance/covariance matrix structured as:

$$V_{ij} = \sigma_a^2 + \sigma_e^2 \quad (i = j)$$

$$V_{ij} = 2\Phi_{ij}\sigma_a^2 \quad (i \neq j)$$

• σ_a^2 , σ_e^2 represent the additive polygenic variance (due to all loci) and the environmental (=error) variance, respectively

Testing for association using LMMs

- LMMs were first (?) applied in human genetics by Boerwinkle et al. (1986) and Abney et al. (2002)
- Chen and Abecasis (2007) implemented them via the "FAmily based Score Test Approximation" (FASTA) in the MERLIN software package
 - Closely related to earlier QTDT method (Abecasis et al. 2000a;b) which implements a slightly more general/complex model
 - FASTA was also implemented in GenABEL, along with a similar test called GRAMMAR (Aulchenko et al. 2007)

12 / 38

ather Cordell (Ne

Estimating the genetic relationship matrix

- These early implementations calculated the kinship matrix Φ on the basis of known (theoretical) kinships constructed from known pedigree relationships
- Amin et al. (2007) proposed instead *estimating* the kinships based on genome-wide SNP data
 - Ideally we want to use G=ZZ'/m, the genetic relationship matrix (GRM) between individuals at the causal loci
 - Since we don't know the causal loci, we approximate **G** by **A**, the overall GRM between individuals
 - Various different ways to estimate this, usually based on scaled (by allele frequency) matrix of *identity-by-state* (IBS) sharing

Estimating the genetic relationship matrix

- Once you move to estimating the GRM, you are no longer limited to using family data
- Kang et al. (2010) and Zhang et al. (2010) suggested applying the approach to apparently unrelated individuals
 - As a way of accounting for population substructure/stratification
 - Also proposed applying to binary traits (case/control coded 1/0)
 - Implemented in EMMAX and TASSEL software, respectively

WAS (Part 2)

14 / 38

Estimating the genetic relationship matrix

- Once you move to estimating the GRM, you are no longer limited to using family data
- Kang et al. (2010) and Zhang et al. (2010) suggested applying the approach to apparently unrelated individuals
 - As a way of accounting for population substructure/stratification
 - Also proposed applying to binary traits (case/control coded 1/0)
 - Implemented in EMMAX and TASSEL software, respectively
- Subsequently a number of other publications/software packages have implemented essentially the same model
 - FaST-LMM (Lippert et al. 2011)
 - GEMMA (Zhou and Stephens 2012)
 - GenABEL (GRAMMAR-Gamma) (Svishcheva et al. 2012)
 - MMM (Pirinen et al. 2013)
 - MENDEL (Zhou et al. 2014)
 - RAREMETALWORKER
 - GCTA

Heather Cordell (Newcastle)

DISSECT
Heather Cordell (Newcastle)

Software implementations

- Main difference between them is the precise computational tricks used to speed up the calculations
 - And the convenience/ease of use

Heather Cordell (Newcastle

• See comparison in Eu-Ahsunthornwattana et al. (2014) PLoS Genetics 10(7):e1004445

Software implementations

- Main difference between them is the precise computational tricks used to speed up the calculations
 - And the convenience/ease of use

Heather Cordell (Newcastle)

- See comparison in Eu-Ahsunthornwattana et al. (2014) PLoS Genetics 10(7):e1004445
- BOLT-LMM (Loh et al. 2016) uses a slightly different approach, based on a Bayesian implementation of LMM formulation 1:

$$\mathsf{y} = \mathsf{X}eta + \mathsf{Z}\mathsf{u} + \epsilon$$

- One of the first mixed model packages that worked for really large-scale (e.g. UK Biobank) datasets
- Now potentially (?) superseded by fastGWA module in GCTA
- And by REGENIE (https://doi.org/10.1038/s41588-021-00870-7), which uses a slightly different formulation based on analysing the residuals following a whole-genome blockwise ridge regression
 - Again based on LMM formulation 1: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}$ GWAS (Part 2)

Binary traits

- For binary traits, coding cases and controls as a 1/0 quantitative trait is not optimal
 - Though in practice it seems to work reasonably well
- LTMLM (Hayeck et al. 2015) and LEAP (Weissbrod et al. 2015) instead use an underlying liability model to improve power

GWAS (Part 2

Assuming known disease prevalence

Heather Cordell (Newcastle)

Binary traits

- For binary traits, coding cases and controls as a 1/0 quantitative trait is not optimal
 - Though in practice it seems to work reasonably well
- LTMLM (Hayeck et al. 2015) and LEAP (Weissbrod et al. 2015) instead use an underlying liability model to improve power
 - Assuming known disease prevalence

Heather Cordell (Newcastle)

- Chen et al. (2016) showed that high levels of population stratification can invalidate the analysis, when applied to a case/control sample
 - Resulting in a mixture of inflated and deflated test statistics
 - Developed GMMAT software to address this problem
 - See also CARAT software (Jiang et al. 2016, AJHG 98:243-55)

GWAS (Part 2

Binary traits

- SAIGE software (Zhou et al. 2018, AJHG 50(9):1335-1341) implements a mixed model test that deals with large case-control imbalance, as you might see (for example) in UK Biobank
- REGENIE also implements this same saddle point approximation (SPA) test
 - Along with an approximate Firth penalized likelihood-ratio test

Heather Cordell (Newcastle)

Elucidating genetic architecture

- Seminal paper by Yang et al. (2010) [Nat Genet 42(7):565-9]
- Showed that by framing the relationship between height and genetic factors as an LMM, 45% of variance could be explained by considering 294,831 SNPs simultaneously
 - So-called 'SNP heritability' or 'chip heritability'
 - Demonstrated that modelling effects at all genotyped SNPs explained the 'known' heritability (\approx 80%) much better than just the top SNPs from GWAS
- Moreover, if you estimate effects of additional SNPs in LD with the genotyped SNPS, the variance explained goes up to 84% (s.e. 16%), consistent with 'known' value
- Subsequently many papers have shown similar results for a variety of complex traits

Heather Cordell (Newcastle) GW

Heather Cordell (Newcastle)

Elucidating genetic architecture

Basic idea is to use formulation

eather Cordell (Newcastle)

$$\mathbf{y} = \mathbf{X}\boldsymbol{eta} + \mathbf{g} + \mathbf{e}$$

with $\mathbf{g} \sim N(0, \mathbf{A}\sigma_a^2)$ and $\epsilon \sim N(0, \mathbf{I}\sigma_e^2)$ so $\mathbf{V} = \mathbf{A}\sigma_a^2 + \mathbf{I}\sigma_e^2$

- A is the GRM between individuals, estimated using all genotyped SNPs
- σ_a^2 and σ_e^2 estimated using REML (or MLE)
- Thus we can estimate heritability accounted for by the genotyped SNPs as $\sigma_a^2/(\sigma_a^2+\sigma_e^2)$
- Implemented in several software packages including GCTA and DISSECT
 - ALBI software (Schweiger et al. 2016, AJHG 98:1181-1192) can then be used to construct accurate confidence intervals for the heritability

GWAS (Part 2)

Partitioning variance

- The same formulation can be used to partition the variance explained by different subsets of SNPs
 - Yang et al. (2010) partitioned variance onto each of the 22 autosomes using formulation

 $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \sum_{c=1}^{22} \mathbf{g}_{c} + \boldsymbol{\epsilon}$ with $\mathbf{V} = \sum_{c=1}^{22} \mathbf{A}_{c}\sigma_{c}^{2} + \mathbf{I}\sigma_{e}^{2}$,

where \mathbf{g}_{c} is a vector of effects attributed to the *c*th chromosome, and \mathbf{A}_{c} is the GRM estimated from SNPs on the *c*th chromosome

- $\bullet\,$ Slight adjustment is needed for estimating variance explained by SNPs on chromosome X
- Similar partitioning can be used to examine subsets of SNPs defined in other ways e.g. according to MAF or functional annotation

Other approaches

- Some recent work has focussed on achieving similar ends
 i.e. estimating
 - heritability explained by sets of SNPs
 - genetic correlations across traits

using summary statistics only

- Bulik-Sullivan et al. (2015) [Nat Genet 47:291-295]
- Bulik-Sullivan et al. (2015) [Nat Genet 47:1236-1241]
 - Clever idea that allows the variance component parameters to be estimated via a simple regression on 'LD Scores'

Short break	Gene-gene (and gene-environment) interactions
	 GWAS have been extraordinarily successful at detecting genetic locations harboring genes associated with complex disease But the SNPs identified do not account for the known (estimated) heritability for most disorders Could G×G and G×E effects account for part of the 'missing heritability'? Zuk et al. (2012) PNAS 109:1193-1198

Heather Cordell (Newcastle)

\sim	/		
(pene_gene)	and	gene-environment	Interactions
UCITC-gene	anu	gene-environment	

GWAS (Part 2)

- GWAS have been extraordinarily successful at detecting genetic locations harboring genes associated with complex disease
 - But the SNPs identified do not account for the known (estimated) heritability for most disorders
 - Could G×G and G×E effects account for part of the 'missing heritability'?
 - Zuk et al. (2012) PNAS 109:1193-1198

Heather Cordell (Newcastle)

Heather Cordell (Newcastle)

- Effects operating through interactions may not be visible unless you stratify by or take account of the interacting genetic (or environmental) factors
 - By modelling interactions, we hope to increase our power to detect loci with weak marginal effects

Gene-gene (and gene-environment) interactions

GWAS (Part 2)

- GWAS have been extraordinarily successful at detecting genetic locations harboring genes associated with complex disease
 - But the SNPs identified do not account for the known (estimated) heritability for most disorders
 - \bullet Could G×G and G×E effects account for part of the 'missing heritability'?
 - Zuk et al. (2012) PNAS 109:1193-1198
- Effects operating through interactions may not be visible unless you stratify by or take account of the interacting genetic (or environmental) factors
 - By modelling interactions, we hope to increase our power to detect loci with weak marginal effects

23 / 38

Phenomenon of biological interest?

Heather Cordell (Newcastle)

• Identifying genes that interact to cause disease could help us understand the mechanisms and pathways in disease progression

Definition of (pairwise) interaction

- Statistical interaction most easily described in terms a of (logistic) regression framework
 - Supppose x₁ and x₂ are binary factors whose presence/absence (coded 1/0) may be associated with a disease outcome
 - Logistic regression models their effect on the log odds of disease as:

$$\log \frac{p}{1-p} = \beta_0 + \beta_1 x_1$$

$$\log \frac{p}{1-p} = \beta_0 + \beta_2 x_2$$
Marginal effect of factor 1
Marginal effect of factor 2

 $\log \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \qquad \log \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2$ Main effects of factors 1 and 2 Main effects and interaction term

• For quantitative traits, use linear regression (replace $\log \frac{p}{1-p}$ with y)

GWAS (Part 2)

ullet For modelling as an LMM, add in a random effect γ

Heather Cordell (Newcastle)

Heather Cordell (Newcastle)

Interaction

• Expected trait values (log odds of disease) take the form:

	Factor 2		
Factor 1	1	0	
1	$\beta_0 + \beta_1 + \beta_2 + \beta_{12}$	$\beta_0 + \beta_1$	
0	$\beta_0 + \beta_2$	β_0	

• β_0 , β_1 , β_2 , β_{12} are regression coefficients (numbers) that can be estimated from real data

GWAS (Part 2

Interaction			Interaction							
Expected trait values (log odds of disease) take the form:		• Expected trait values (log odds of disease) take the form:					form:			
		Factor 2]				Factor 2		
	Factor 1	1	0			Facto	r 1	1	0	
	1	$\beta_0 + \beta_1 + \beta_2 + \beta_{12}$	$\beta_0 + \beta_1$			1		$\beta_0 + \beta_1 + \beta_2 + \frac{\beta_{12}}{\beta_{12}}$	$\beta_0 + \beta_1$	
	0	$\beta_0 + \beta_2$	β_0			0		$\beta_0 + \beta_2$	β_0	i
• $\beta_0, \beta_1, \beta_2, \beta_1$	• β_0 , β_1 , β_2 , β_{12} are regression coefficients (numbers) that can be									

 $\beta_0, \beta_1, \beta_2, \beta_{12}$ are regression coefficients (numbers) that cleatimated from real data

 $\bullet~$ Having factor 1 adds β_1 to your trait value

- β_0 , β_1 , β_2 , β_{12} are regression coefficients (numbers) that can be estimated from real data
 - Having factor 1 adds β_1 to your trait value
 - $\bullet~$ Having factor 2 adds β_2 to your trait value

ather Cordell (Ne

Interaction

• Expected trait values (log odds of disease) take the form:

	Factor 2	
Factor 1	1	0
1	$\beta_0 + \beta_1 + \beta_2 + \frac{\beta_{12}}{\beta_{12}}$	$\beta_0 + \beta_1$
0	$\beta_0 + \beta_2$	β_0

- $\beta_0, \beta_1, \beta_2, \beta_{12}$ are regression coefficients (numbers) that can be estimated from real data
 - Having factor 1 adds β_1 to your trait value
 - Having factor 2 adds β_2 to your trait value
 - Having both factors adds an additional β_{12} to your trait value \Rightarrow Implies that the overall effect of two variables is greater (or less) than the 'sum of the parts'
 - The 'effect' of factor 2 is different in the presence/absence of factor 1

Interaction

• Expected trait values (log odds of disease) take the form:

	Factor 2				
Factor 1	1	0			
1	$\beta_0 + \beta_1 + \beta_2 + \beta_{12}$	$\beta_0 + \beta_1$			
0	$\beta_0 + \beta_2$	β_0			

- β_0 , β_1 , β_2 , β_{12} are regression coefficients (numbers) that can be estimated from real data
 - Having factor 1 adds β_1 to your trait value
 - $\bullet~$ Having factor 2 adds $\beta_{\rm 2}$ to your trait value
 - Having both factors adds an additional β_{12} to your trait value \Rightarrow Implies that the overall effect of two variables is greater (or less) than the 'sum of the parts'
 - The 'effect' of factor 2 is different in the presence/absence of factor 1
- Suppose no main effects ($\beta_1 = \beta_2 = 0$)

ather Cordell (Newcastle)

	Factor 2					
Factor 1	1	0				
1	$\beta_0 + \beta_{12}$	β_0				
0	Bo	βn				

Trait value only differs from baseline if both factors present

GWAS (Part 2)

Heather Cordell (Newcastle)

GWAS (Part 2)

Gene-gene interaction (epistasis)

- However SNPs are not binary, but rather take 3 levels according to the number of copies (0,1,2) of the susceptibility allele possessed
- Most general 'saturated' (9 parameter) genotype model allows all 9 penetrances to take different values
 - Via modelling log odds in terms of:
 - A baseline effect (β_0)
 - Main effects of locus G (β_{G_1} , β_{G_2})
 - Main effects of locus $H(\beta_{H_1}, \beta_{H_2})$
 - 4 interaction terms

Heather Cordell (Newcastle)

		Locus H	
Locus G	2	1	0
2	$\beta_0 + \beta_{G_2} + \beta_{H_2} + \beta_{22}$	$\beta_0 + \beta_{G_2} + \beta_{H_1} + \beta_{21}$	$\beta_0 + \beta_{G_2}$
1	$\beta_0 + \beta_{G_1} + \beta_{H_2} + \beta_{12}$	$\beta_0 + \beta_{G_1} + \beta_{H_1} + \beta_{H_1}$	$\beta_0 + \beta_{G_1}$
0	$\beta_0 + \beta_{H_2}$	$\beta_0 + \beta_{H_1}$	β_0

• Corresponds in statistical analysis packages to coding x_1 , x_2 (0,1,2) as a "factor"

Gene-gene interaction

- Alternatively we can assume additive effects of each allele at each locus:
 - Corresponds to fitting

$$\log \frac{p}{1-p} = \beta_0 + \beta_G x_1 + \beta_H x_2 + \frac{\beta_{GH} x_1 x_2}{1-p}$$

with x_1 , x_2 coded (0,1,2)

		Locus H	
Locus G	2	1	0
2	$\beta_0 + 2\beta_G + 2\beta_H + 4\beta_{GH}$	$\beta_0 + 2\beta_G + \beta_H + 2\beta_{GH}$	$\beta_0 + 2\beta_G$
1	$\beta_0 + \beta_G + 2\beta_H + 2\beta_{GH}$	$\beta_0 + \beta_G + \beta_H + \beta_{GH}$	$\beta_0 + \beta_G$
0	$\beta_0 + 2\beta_H$	$\beta_0 + \beta_H$	β_0

Change of scale

- Transformations of outcome variable *y* can change whether or not the predictor variables interact
 - Due to definition of interaction as departure from a linear model for the effects of x₁ and x₂, for predicting y
 - Two SNPs that interact on the log odds scale may not interact on the penetrance scale (and vice versa)
 - Makes biological interpretation of resulting interaction model difficult

Change of scale

- Transformations of outcome variable *y* can change whether or not the predictor variables interact
 - Due to definition of interaction as departure from a linear model for the effects of x₁ and x₂, for predicting y
 - Two SNPs that interact on the log odds scale may not interact on the penetrance scale (and vice versa)
 - Makes biological interpretation of resulting interaction model difficult
- Much discussion in the literature
 - Siemiatycki and Thomas (1981) Int J Epidemiol 10:383-387; Thompson (1991) J Clin Epidemiol 44:221-232
 - Phillips (1998) Genetics 149:1167-1171; Cordell (2002) Hum Molec Genet 11:2463-2468
 - McClay and van den Oord (2006) J Theor Biol 240:149-159; Phillips (2008) Nat Rev Genet 9:855-867
 - Clayton DG (2009) PLoS Genet 5(7): e1000540; Wang, Elston and Zhu (2010) Hum Hered 70:269-277

Heather Cordell (Newcastle)	GWAS (Part 2)	28 / 38	Heather Cordell (Newcastle)	GWAS (Part 2)	28 / 38

Change of scale

- Transformations of outcome variable *y* can change whether or not the predictor variables interact
 - Due to definition of interaction as departure from a linear model for the effects of x₁ and x₂, for predicting y
 - Two SNPs that interact on the log odds scale may not interact on the penetrance scale (and vice versa)
 - Makes biological interpretation of resulting interaction model difficult

Much discussion in the literature

Heather Cordell (Newcastle)

- Siemiatycki and Thomas (1981) Int J Epidemiol 10:383-387; Thompson (1991) J Clin Epidemiol 44:221-232
- Phillips (1998) Genetics 149:1167-1171; Cordell (2002) Hum Molec Genet 11:2463-2468
- McClay and van den Oord (2006) J Theor Biol 240:149-159; Phillips (2008) Nat Rev Genet 9:855-867
- Clayton DG (2009) PLoS Genet 5(7): e1000540; Wang, Elston and Zhu (2010) Hum Hered 70:269-277
- Bottom line is, little direct correspondence between statistical interaction and biological interaction
 - In terms of whether, for example, gene products physically interact

Change of scale

- Transformations of outcome variable *y* can change whether or not the predictor variables interact
 - Due to definition of interaction as departure from a linear model for the effects of x₁ and x₂, for predicting y
 - Two SNPs that interact on the log odds scale may not interact on the penetrance scale (and vice versa)
 - Makes biological interpretation of resulting interaction model difficult
- Much discussion in the literature
 - Siemiatycki and Thomas (1981) Int J Epidemiol 10:383-387; Thompson (1991) J Clin Epidemiol 44:221-232
 - Phillips (1998) Genetics 149:1167-1171; Cordell (2002) Hum Molec Genet 11:2463-2468
 - McClay and van den Oord (2006) J Theor Biol 240:149-159; Phillips (2008) Nat Rev Genet 9:855-867
 - Clayton DG (2009) PLoS Genet 5(7): e1000540; Wang, Elston and Zhu (2010) Hum Hered 70:269-277
- Bottom line is, little direct correspondence between statistical interaction and biological interaction
 - In terms of whether, for example, gene products physically interact
- However, existence of statistical interaction does imply both loci are "involved" in disease in some way

ather Cordell (Ne

Change of scale

- Transformations of outcome variable y can change whether or not the predictor variables interact
 - Due to definition of interaction as departure from a linear model for the effects of x_1 and x_2 , for predicting y
 - Two SNPs that interact on the log odds scale may not interact on the penetrance scale (and vice versa)
 - Makes biological interpretation of resulting interaction model difficult

Much discussion in the literature

Heather Cordell (Newcastle)

- Siemiatycki and Thomas (1981) Int J Epidemiol 10:383-387; Thompson (1991) J Clin Epidemiol 44:221-232
- Phillips (1998) Genetics 149:1167-1171; Cordell (2002) Hum Molec Genet 11:2463-2468
- McClay and van den Oord (2006) J Theor Biol 240:149-159; Phillips (2008) Nat Rev Genet 9:855-867
- Clayton DG (2009) PLoS Genet 5(7): e1000540; Wang, Elston and Zhu (2010) Hum Hered 70:269-277
- Bottom line is, little direct correspondence between statistical interaction and biological interaction
- In terms of whether, for example, gene products physically interact
- However, existence of statistical interaction does imply both loci are "involved" in disease in some way
 - Good starting point for further investigation of their (joint) action GWAS (Part 2)

Gene-environment $(G \times E)$ interactions

The same regression model

 $\log \frac{p}{1-p} = \beta_0 + \beta_G x_1 + \beta_H x_2 + \beta_{GH} x_1 x_2$

can be used to model interaction between a genetic factor G and an environmental factor H

• With the environmental variable x_2 coded in binary fashion (e.g. smoking) or quantitatively (e.g. age)

Gene-environment $(G \times E)$ interactions

The same regression model

Heather Cordell (Newcastle)

$$\log \frac{p}{1-p} = \beta_0 + \beta_G x_1 + \beta_H x_2 + \frac{\beta_{GH} x_1 x_2}{1-p}$$

can be used to model interaction between a genetic factor G and an environmental factor H

- With the environmental variable x_2 coded in binary fashion (e.g. smoking) or quantitatively (e.g. age)
- Focus of analysis is often risk estimation
 - Estimating genetic risks in particular environments
 - Estimating effect of environmental factor on particular genetic background
 - Important for treatment/screening strategies and public health interventions
- For G×G, focus of interest is more related to
 - Increasing power to detect an effect (by taking into account the effects of other genetic loci)
 - Modelling the biology, especially related to the joint action of the loci

Testing association and/or interaction

Go back to binary coding of genetic (and/or environmental) factors

$$\log \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2$$

• 3df test of $\beta_1 = \beta_2 = \beta_{12} = 0$ tests for association at both loci (or both variables), allowing for their possible interaction

Testing association and/or interaction

• Go back to binary coding of genetic (and/or environmental) factors

$$\log \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2$$

- 3df test of β₁ = β₂ = β₁₂ = 0 tests for association at both loci (or both variables), allowing for their possible interaction
- 2df test of $\beta_2 = \beta_{12} = 0$ tests for association at locus 2, while allowing for possible interaction with locus (or variable) 1

Testing association and/or interaction

• Go back to binary coding of genetic (and/or environmental) factors

$$\log \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2$$

- 3df test of β₁ = β₂ = β₁₂ = 0 tests for association at both loci (or both variables), allowing for their possible interaction
- 2df test of $\beta_2 = \beta_{12} = 0$ tests for association at locus 2,
- while allowing for possible interaction with locus (or variable) 1
- 1df test of $\beta_{12}=$ 0 tests the interaction term alone

Heather Cordell (Newcastle)	GWAS (Part 2)	30 / 38	Heather Cordell (Newcastle)	GWAS (Part 2)	30 / 38

Testing association and/or interaction

• Go back to binary coding of genetic (and/or environmental) factors

$$\log \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2$$

- 3df test of β₁ = β₂ = β₁₂ = 0 tests for association at both loci (or both variables), allowing for their possible interaction
- 2df test of $\beta_2 = \beta_{12} = 0$ tests for association at locus 2,

Heather Cordell (Newcastle)

- while allowing for possible interaction with locus (or variable) 1 • 1df test of $\beta_{12} = 0$ tests the interaction term alone
- Depending on circumstances, any of these tests may be a sensible option

Testing association and/or interaction

• Go back to binary coding of genetic (and/or environmental) factors

$$\log \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2$$

- 3df test of β₁ = β₂ = β₁₂ = 0 tests for association at both loci (or both variables), allowing for their possible interaction
- 2df test of $\beta_2 = \beta_{12} = 0$ tests for association at locus 2,
- while allowing for possible interaction with locus (or variable) 1
- 1df test of $\beta_{12} = 0$ tests the interaction term alone
- Depending on circumstances, any of these tests may be a sensible option
- Most tests of interaction/joint action can be thought of as a version of one or other of these tests
 - Although different tests vary in their precise details
 - And their relationship to the logistic regression formulation not always clearly described
 - See Howey and Cordell (2017)
 - https://pubmed.ncbi.nlm.nih.gov/28852712/

$G \times G$ versus $G \times E$ in the context of GWAS

- Typically GWAS measure thousands if not millions of genetic variants
 But only a few (tens or at most 100s) of environmental factors
- Feasible to consider all G×E combinations
- All pairwise G×G combinations possible, but much more time consuming
 - And leads to greater multiplicity of tests
 - Also, why stop at 2-way interactions?
 - Could look at all 3 way, 4 way etc. combinations
 - Scale of problem quickly gets out of hand
 - \bullet Less obvious reason to do this for $\mathsf{G}{\times}\mathsf{E}{\ldots}$

$G \times G$ in the context of GWAS

- Many recent publications have focussed on finding clever computational tricks to speed up exhaustive search procedure
 BOOST (Wan et al. (2010) AJHG 87:325-340)
 - SIXPAC (Prabhu and Pe'er (2012) Genome Res 22:2230-2240)
 - SIXPAC (Prablu and Pe er (2012) Genome Res 22.2250-2240)
 Kam-Thong et al. (2012) Hum Hered 73:220-236 (GPUs)

GWAS (Part 2

- Fraanberg et al. (2012) Fluin Fleted 75:220-230 (GF0)
 Fraanberg et al. (2015) PLOS Genetics 11(9):e1005502
- "Discovering genetic interactions in large-scale association studies by stage-wise likelihood ratio tests"

Heather Cordell (Newcastle)

Heather Cordell (Newcastle)

GWAS (Part 2)

31 / 38 Heather Cordell (Newcastle)

$G \times G$ in the context of GWAS

- Many recent publications have focussed on finding clever computational tricks to speed up exhaustive search procedure
 - BOOST (Wan et al. (2010) AJHG 87:325-340)
 - SIXPAC (Prabhu and Pe'er (2012) Genome Res 22:2230-2240)
 - Kam-Thong et al. (2012) Hum Hered 73:220-236 (GPUs)
 - Fråanberg et al. (2015) PLOS Genetics 11(9):e1005502
 "Discovering genetic interactions in large-scale association studies by stage-wise likelihood ratio tests"
- Or have proposed filtering based on single-locus significance or other (biological or statistical) considerations
 - Reduces multiple testing burden, improves interpretability

$G\!\times\!G$ in the context of GWAS

- Many recent publications have focussed on finding clever computational tricks to speed up exhaustive search procedure
 - BOOST (Wan et al. (2010) AJHG 87:325-340)
 - SIXPAC (Prabhu and Pe'er (2012) Genome Res 22:2230-2240)
 - Kam-Thong et al. (2012) Hum Hered 73:220-236 (GPUs)
 - Fråanberg et al. (2015) PLOS Genetics 11(9):e1005502
 "Discovering genetic interactions in large-scale association studies by stage-wise likelihood ratio tests"
- Or have proposed filtering based on single-locus significance or other (biological or statistical) considerations
 - Reduces multiple testing burden, improves interpretability
- Or have proposed testing at the gene level rather than the SNP level
 Ma et al. (2013) PLoS Genet 9(2): e1003321

S (Part 2

- Compared 4 different tests that combine P values from pairwise (SNP × SNP) interaction tests
- Showed that the truncated tests did best
- Presented an application only considering gene pairs known to exhibit protein-protein interactions

Case-only analysis

- Piergorsh et al. 1994; Yang et al. 1999; Weinberg and Umbach 2000
- Several authors have shown that, for binary predictor variables, a test of the interaction term β_{12} in the logistic regresssion model

$$\log \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2$$

can be obtained by testing for correlation (association) between the genotypes at two separate loci, within the sample of cases

- Gains power from making assumption that genotypes (alleles) at the two loci are uncorrelated in the population
 - So only really suitable for unlinked or loosely linked loci (since closely linked loci are likely to be in LD)
- Alternatively contrast the genotype correlations in cases with those seen in controls (--fast-epistasis in PLINK)

GWAS (Part 2)

Heather Cordell (Newcastle)

Heather Cordell (Newcastle)

Testing correlation between loci

- A similar idea is implemented in EPIBLASTER (Kam-Thong et al. 2011; EJHG 19:465-571)
- Wu et al. (2010) (PLoS Genet 6:e1001131) also proposed a similar approach – though needs adjustment to give correct type I error rates
- See also Joint Effects (JE) statistics (Ueki and Cordell 2012; PLoS Genetics 8(4):e1002625)
- All these methods test whether correlation exists (case-only) or is different in cases and controls (case/control)
 - Via testing a log OR for association between two loci
 - However, the log OR for association (λ) encapsulates a slightly different quantity between the different methods
- All implemented (along with standard logistic and linear regression) in CASSI

GWAS (Part 2)

http://www.staff.ncl.ac.uk/richard.howey/cassi/

Empirical evidence for $\mathsf{G}{\times}\mathsf{G}$ interactions

- Epistasis among HLA-DRB1, HLA-DQA1, and HLA-DQB1 in multiple sclerosis (Lincoln et al. 2009 PNAS 106:7542-7547)
- HLA-C and ERAP1 in psoriasis (Strange et al. 2010)
- HLA-B27 and ERAP1 in ankylosing spondylitis (Evans et al. 2011)
- BANK1 and BLK in SLE (Castillejo-Lopez et al. 2012)
- Gusareva et al. (2014) found a reasonably convincing (partially replicating) interaction between SNPs on chromosome 6 (*KHDRBS2*) and 13 (*CRYL1*) in Alzheimer's disease
- Dai et al. (2016) [AJHG 99:352-365] identified 3 loci simultaneously interacting with established risk factors gastresophageal reflux, obesity and tobacco smoking, with respect to risk for Barrett's esophagus

Empirical evidence for $G \times G$ interactions

- Hemani et al. 2014 (Nature 508:249-253) found 501 instances of epistatic effects on gene expression, of which 30 could be replicated in two independent samples
 - Many SNPs are close together, could represent haplotype effects?
 - Or the effect of a single untyped variant?
 - See caveats in

Heather Cordell (Newcastle

Heather Cordell (Newcastle)

- Wood et al. (2014) Nature 514(7520):E3-5. PMID:25279928
- Fish et al. (2016) Am J Hum Genet 99(4):817830. PMID:27640306
- The Hemani et al. paper was subsequently retracted (https://www.nature.com/articles/s41586-021-03766-y)

Empirical evidence for $G \times E$ interactions

- Myers et al. (2014) Hum Mol Genet 23(19): 5251-9 "Genome-wide Interaction Studies Reveal Sex-Specific Asthma Risk Alleles"
- Small et al. (2018) Nat Genet 50(4): 572-580 "Regulatory Variants at KLF14 Influence Type 2 Diabetes Risk via a Female-Specific Effect on Adipocyte Size and Body Composition"
- Sung et al. (2019) Hum Molec Genet 28(15): 2615-2633 "A multi-ancestry genome-wide study incorporating gene-smoking interactions identifies multiple new loci for pulse pressure and mean arterial pressure."



Empirical evidence for $G \times E$ interactions

• Favé et al. (2018) Nat Commun 9(1): 827 "Gene-by-environment Interactions in Urban Populations Modulate Risk Phenotypes"





Why Estimate Sample Sizes and/or Power?

· To avoid wasting time and money

 Does not make sense to perform an inadequately powered study for which it is unlikely to to correctly reject the null hypothesis due to inadequate sample size

· Collaborations can aid in increasing sample sizes

- Caveats
 » Disease definition may not be the same between studies
- » Study subjects may be drawn for different populations
- Processing of genetic material maybe not be consistent
 Almost always necessary for grant proposals
- Can be denied funding if unable to demonstrate planned study has adequate power
 - Realistic disease models are necessary when performing power calculations
 - Correctly adjust alpha for multiple testing which will be performed

 e.g., use genome-wide significant level of 5 x 10⁸ for GWAS studies
 - e.g., use genome-wide significant level of 5 x 10° for GWAS studies

2

4

Power and Sample Size Estimation for Case-Control Data

- The correct α must be use for sample size estimation/power analysis
- Type I (α) the probability of rejecting the null hypothesis of no association when it is true
- Due to multiple testing a more stringent value than α =0.05 is used in order to control the Family Wise Error Rate

Power and Sample Size Estimation for Case-Control Data

- GWAS of common variants where each variant is test separately

 α=5 X10.⁸ (Bonferroni Correction for testing 1,000,000 variant sites)
 - Shown to be a good approximation for the effective number of tests
 Valid even when more than 1,000,000 variant sites tested
 - Effective number of tests is dependent of the linkage disequilibrium (LD) structure
- Single variant tests using whole genome sequence data

 Many more rare variants than common variants
 Lower levels of LD between rare variants than between common variants
 - The number of effective tests for rare variants than between common variants
 The number of effective tests for rare variants is higher than for analysis limited to common variants
 - α is yet to be determined for association analysis of whole genome sequence data

An Example of Determining Genome-wide Significance Levels for Common Variants

- Using genotypes from the Wellcome Trust Case-Control Consortium
- Dudbridge and Gusnato, Genet Epidemiol 2008
- Estimated a genome-wide significance threshold for the UK European population
- By sub-sampling genotypes at increasing densities and using permutation to estimate the nominal p-value for a 5% familywise error
- Then extrapolating to infinite density
- The genome wide significance threshold estimate ~7.2X10⁻⁸
- Estimate is based on LD structure for Europeans
 - Not sufficiently stringent for populations of African Ancestry

Power and Sample Size Estimation for Aggregate Rare Variant Tests

- For gene-based rare variant aggregate methods a Bonferroni correction for the number of genes/regions tested is used

 e.g., 20,000 genes significance level α=2.5 x 10⁻⁶
 - Can use a less stringent criteria
 - Not all genes have two or more variants
 » Divide 0.05 by number of genes tested
 - If units other than genes are used
 - A more stringent criteria may be necessary
- For rare variants very low levels of LD between variants in separate genes
- Therefore, a Bonferroni correction is not overly stringent
 - The number of tests ≅ effective number tests
 This would not be the case for variants in LD

3

Power and Sample Size Estimation for Replication **Studies**

- For replication studies can base the significance level (α)
- On the number of genes/variants being brought from the discovery (stage I) study
- To replication (stage II)
- For example, if it is hypothesized that 20 genes and 80 independent variants will be brought to stage II (replication) - A Bonferroni correct can be made for performing 100 tests - An α = 5.0 x $10^{-3}\,can$ be used for a family wise error rate of 0.05

Estimating Power/Sample Sizes For Single Variant Tests

- Can be obtained analytically
- Information necessary
 - Prevalence
 - Risk allele frequency
 - Effect size (odds ratio-for case control data)
 - Genetic model for the susceptibility variant
 - Recessive (y1=1)
 - Dominant (y2=y1) Additive (v2=2v1-1)
 - Multiplicative (y2=y1²)

8

Estimating Power/Sample Sizes For Individual Variants

- Usually, information on disease prevalence is known from epidemiological data
- A range of risk allele allele frequencies and effect sizes are used
- A variety of genetic models can also used
 - Dominant Additive
 - Multiplicative

Armitage Trend Test

• Power and Sample size

- Calculated under different models
- Where y is the relative risk
 - Multiplicative » $\gamma_2 = \gamma_1^2$ – Additive

» γ2=2γ1-1 Dominant

- » $\gamma_2 = \gamma_1$
- Recessive » γ1=1

9

7



Gamma is the Relative Risk not the Odd Ratio

- Most software for power calculations/sample size estimation use the relative risk (\boldsymbol{y}) and not the odds ratio
- The relative risk only approximates the odds ratio when disease is rare (Prevalence ~< 0.1%)

- The relative risk is not appropriate for common traits when a case-control design is used

Correspo	ndence Between	the Odds	Ratio and	Relative Risk							
	Domir	nant Mode	el								
	Disease Prevalence	1/2* RR=1.5	2/2 ^{**} RR=1.5								
	0.01	1.51	1.51								
	0.10	1.59	1.59								
	0.20	1.71	1.71								
	Multiplicative Model										
	Disease Prevalence	1/2 RR=1.5	2/2 RR=2.25								
	0.01	1.51	2.28								
	0.10	1.59	2.61								
	0.20	1.71	3.25								
Marker minor a D' and r ² =1 *1/2 genotype **2/2 genotype	allele and disease allele – heterozygous (one co e - homozygous for the	frequency 0.0 py of the alterr alternative alle	1 native allele) le								



<section-header>









Power Association With Errors (PAWE)

- http://compgen.rutgers.edu/pawe/
- · Implements the linear trend test
- · Four different error models can be used - See online documentation for complete explanation
- Can either perform:

with the tag SNP

N/r²=N'

same level as when r²=1

• Can estimate sample size when $r^2 \neq 1$

- Valid only for multiplicative model

- (Pritchard and Przeworski, 2001)

• Power calculation almost always assume that r²=1

- Power calculations for a fixed sample size
- Sample size calculations for a fixed power
- The genotype frequencies can be generated either using a:

Linkage Disequilibrium (LD)

• Power will be reduced if causal variant is not in perfect LD (r²=1)

For whole genome sequence data this should be the case since

usually the causal variant would be included in the data

- Can adjust sample size when $r^2 < 1$ to increase power to the

- Genetic model free method or
- Genetic model-based method

20

Quanto

- · Provides sample size and power calculations for
- · Genetic and environmental main effects
- Interactions
 - Gene x gene Gene x environment
- Sample & power calculations can be carried for: Case-control
 - Unmatched
 - Matched
 - Case-sibling
 - Case-parent (trios)
 - QuantitativeQualitative

 - Independent sample of individuals Quantitative traits
 - Assumption sampled from a random population
- Can only be run under windows
- https://pphs.usc.edu/download-quanto/

21

Power Analysis for Rare Variant Aggregate **Association Tests**

- · Many unknown parameters must be modeled - Allelic architecture within a genetic region
 - Varied across genes and populations
 - Effects of variants within a region
 - · Fixed or varied effect sizes of causal variants
 - Bidirectional effect of variants · Proportion of non-causal variants
- Power estimated empirically
- · Simplified assumptions can be made to obtain analytical estimates
 - All variants have the same effect size
 - No non-causal variants within a region that is analyzed in aggregate

Simplistic Analytical Power Calculation for Rarevariant Aggregate Association Analysis

- Assumption
 - All rare variants are causal and have the same effect size
- Although usual not be correct
- Provides a gestalt of the power for a given samples or sample size for a given power
- Use aggregate of allele frequencies
 - For example, assume a cumulative allele frequency of 0.025 - Use an exome-wide significant level e.g., 2.5x10⁻⁶
- · Provide disease prevalence and penetrance model
- Perform calculations in the same manner as was described for single variants

22



Empirical Power Calculations

- Examples
 - 5,000 replicates are generated each with 20,000 cases and 20,000 controls
 - The power is the proportion of replicates with p-value less than the specified threshold, e.g., 5x10⁻⁸
 - For rare-variant aggregate tests all autosomal genes are generated and those genes with more than two rare variants (e.g., predicted loss of function) are analyzed
 - The power is the proportion of genes that were tested with p-value which is below a specified threshold, e.g., 2.5x10⁻⁶

26





28





- Power and sample size calculations for binary and quantitative traits
- User specify proportion of variants that increase or lower risk









33



- It is unknown which genes are important in disease etiology
 Correct allelic architecture is unknown
- Can get a better understanding of power to detect
- associations by generating variants for the entire exome
- Use a variety of disease models
- Odds ratios
- Proportion of pathogenic variants
- Analyze of all genes
- e.g., those with 2 or more variant sites
- Determine power as the proportion of genes that meet exome-wide significance (e.g., alpha=2.5x10⁻⁶)



Power Analysis

- For tests of individual variants
 - Power depended on sample size, disease prevalence, minor allele frequency, genetic model and variant effect size
- For rare variants (aggregate association tests)
 - Also dependent on the allelic architecture
 - Cumulative variant frequency within analyzed region
 Proportion of causal variants
 - How much contamination from non-causal variants
 - Effect sizes the same the same or different across gene regions – Effects of variants in the same or different directions
 - » Protective and detrimental for binary traits
 - » Increase and decrease quantitative trait values





- Power will not only vary between traits greatly
- The power to detect an association will also vary
- drastically between genes for the same complex trait - For some causal genes even with hundreds of thousands of
 - samples power will be low
- While for other causal genes a few thousand samples may be sufficient



37

38



Genotype Pattern Mining For Digenic Traits

Advanced Gene Mapping Course, November 2022

Jurg Ott, Ph.D., Professor Emeritus Rockefeller University, New York <u>https://lab.rockefeller.edu/ott/</u> <u>ott@rockefeller.edu</u> PH +1 646 321 1013





https://lab.rockefeller.edu/ott/ EM: ott@rockefeller.edu PII: +1 646 321 1013

Jurg Ott, PhD

Research Interests

Development of analysis methods for genetic data, genetic linkage and association analysis

Implementation in computer programs, dissemination on website

Collaboration with researchers world-wide on their data

Recent publications: [1-5] (#1 now freely available from https://www.jurgott.org/linkage/LinkageHandbook.pdf)

- 1. Terwilliger DJ, Ott J. Handbook of human genetic linkage. Johns Hopkins University Press. 1994.
- Imai-Okazaki A, Li Y, Horpaopan S, Riazalhosseini Y, Garshasbi M, Mosse YP, et al. Heterozygosity mapping for human dominant trait variants. Hum Mutat. 2019 Apr 24;40(7):5996-1004.
- Professor Emeritus and Director Laboratory of Statistical Genetics Rockefeller University New York: NY 10065 24;40(7):996-1004. 3. Horpaopan S, Fann CSJ, Lathrop M, Ott J. Shared genomic segment analysis with equivalence rescher Leidemind. 2020 Oct;44(7):741-47.
 - Okazaki A, Yamazaki S, Inoue I, Ott J. Population genetics: past, present, and future. Human genetics. 2020:1-10.
 - Okazaki A, Horpaopan S, Zhang Q, Randesi M, Ott J. Genotype pattern mining for pairs of interacting variants underlying digenic traits. Genes. 2021;12(8):1160.

Ott "Genotype Patterns"

2

Topics

	Frequent Pattern Mining	ARREN CA
	ase-control association analysis	0000 8:8:8
	Main effects in genetic association studies	
	Interaction effects in case-control data	
M	ining consumer databases	
	The Apriori algorithm	
	Newer algorithms: eclat, fpgrowth	
-	Analysis of AMD dataset	

Main association effects

- □ Consider two DNA variants with minor alleles *A* and *T*. Even when on the same chromosome, the frequency of *A*-*T* chromosomes (haplotypes) is the product of allele frequencies, $P(A-T) = P(A) P(T) \rightarrow$ linkage equilibrium. Variants very close together: $P(A-T) \neq P(A) P(T) \rightarrow$ LD, linkage disequilibrium.
- □ Disease variant vs marker variant: Different genotype frequencies in cases and controls → genetic association.
- Recessive traits: Variants close to disease tend to be homozygous (homozygosity mapping; Lander & Botstein, Science 1987;236:1567-70).
- Dominant traits: Variants close to disease tend to be heterozygous (Imai-Okazaki et al, Hum Mutat 2019;40:996-1004): P(het) > 1 - f, P(het, popul) = 2f(1 - f), f = MAF





Multiple Hits ... Digenic Diseases Ming & Muenke (2002) Am | Hum Genet 71, 1017 (review)

Schaffer A (2013) J Med Genet 50, 641-52 (review)

EFFECT AND	GE	ne 1	Gene 2			
Phenotype	Mutation	Phenotype	Mutation	Phenotype		
Synergistic:						
RP	ROM1 ^{+/G80insG}	Normal	$RDS^{+/L18SP}$	Normal		
RP	ROM1 ^{+/L114insG}	Normal	$RDS^{+/L185P}$	Normal		
Bardet-Biedl	BBS2 ^{Y24X/Q59X}	Normal	BBS6 ^{+/Q147X}	Normal		
Deafness	GJB2 ^{+/35delG}	Normal	$GJB6^{+/-}$	Normal		
Deafness	GJB2 ^{+/167delT}	Normal	$GJB6^{+/-}$	Normal		
Hirschsprung	$RET^{+/164711}$	Normal	EDNRB ^{+/S305N}	Normal		
Severe insulin resistance	PPARG ^{+/A553delAAAiT}	Normal	PPP1R3A ^{+/C1984delAG}	Normal		
Modifier:						
Juvenile-onset glaucoma	MYOC ^{+/G399V}	Adult-onset glaucoma	CYP1B1 ^{+/R368H}	Normal		
Usher 1	USH3 ^{mut/mut}	Usher 3	MYO7A ^{+/delG (exon 25)}	Normal		
Congenital nonlethal JEB	COL17A1 ^{R1226X/L855X}	Juvenile JEB	LAMB3 ^{+/R635X}	Normal		
More severe ADPKD	$PKD1^{+/mut}$	Less severe ADPKD	PKD2 ^{+/2152delA}	Less severe ADPKE		
More severe hearing loss	DFNA1	Mild hearing loss	DFNA2	Mild hearing loss		
WS2/OA	$MITF^{+/944delA}$?WS2	$TYR^{+/R402Q}$	Normal		
More severe WS2/OA	MITF ^{+/944delA}	?WS2	TYR ^{R402Q/R402Q}	Normal		

Ott "Genotype Patterns"

How to analyze interaction effects?

Hyperlipidemia data: 5 relevant genes, ~200 variants in each ge interactions in each pair of variants. Work with LR chi-square!	ene, look for
--	---------------

CASES	V	ariant	1		CONTROLS	V	/arian	t 1	Data	chi-sq	df
Var 2	GG	GT	TT		Var 2	GG	GT	TT	cases	3.3591	4
AA					AA				controls	3.6658	4
AC					AC				both	1.4255	4
CC					CC				heterogeneity	5.5994	4
$\chi^2_{\text{Heterogen}}$	eity =	χ^2_{Cas}	$ses + \gamma$	$\ell^2 c$	$c_{ontrols} - \chi^2_{bot}$	h					

Var 1 ->		GG			GT			TT		Source	chi-sq	df
Var 2 ->	AA	AC	CC	AA	AC	CC	AA	AC	CC	Var 1 main	0.4196	2
cases										Var 2 main	48.1979	2
controls										Interaction	5.5994	4
2	2		2		2					Total table	54.2169	8
Interaction	-λ.	Total -	·λv	arl —	L Var	2				·		

More sophisticated analysis by logistic regression (Cordell, Nat Rev Genet 2009;10:392-404).

Ott "Genotype Patterns"

Different Levels of Genetic Interactions

Okazaki & Ott (2022) Trends in Genetics 38 (10):1013-1018

1. Traditionally, disease association has been carried out on the level of alleles or **genotypes**. The total number of pairs can be prohibitively large. While this level of analysis generally requires the most effort, it also entails the highest level of precision in the sense that disease-causing elements can be directly traced down to nucleotides.

2. Working with pairs of **variants** provides some economy of computational effort but may 'dilute' a signal from a single genotype pair when all nine genotype pairs in a pair of variants are analyzed jointly.

3. Finally, focusing on pairs of **genes** represents the most economical approach but is also the most imprecise among the three strategies. Also, focusing on genes disregards susceptibility elements outside of genes. Distant-acting transcriptional enhancers have been known for over 10 years to affect susceptibility to human disease and noncoding RNAs have been shown to be associated with many diseases, for example, cardiac hypertrophy.

Ott "Genotype Patterns"

Finding disease-associated pairs of genotypes

- Multifactor Dimensionality Reduction (MDR) Ritchie MD, Hahn LW, Moore JH. Power of multifactor dimensionality reduction for detecting gene-gene interactions ... Genet Epidemiol 2003;24:150–157
- Exhaustive evaluation of all pairs of genotypes at all pairs of variants
- Applying off-the-shelf pattern search algorithms Chee C-H, Jaafar J, Aziz IA, Hasan MH, Yeoh W. Algorithms for frequent itemset mining: a literature review. Artificial Intelligence Review. 2019;52(4):2603-21
- Construction of Bayesian network Guo Y, Zhong Z, et al. Epi-GTBN: An approach of epistasis mining based on genetic Tabu algorithm and Bayesian network. BMC Bioinform 2019;20:444
- Sophisticated computational approaches Titarenko SS, Titarenko VN, Aivaliotis G, Palczewski J. Fast implementation of pattern mining algorithms ... Journal of Big Data. 2019; 13(6):37

Ott "Genotype Patterns'

8

6

7

1. MDR



2. Exhaustive search for interacting SNPs

- Discovering Genetic Factors for psoriasis through exhaustively searching for significant second order SNP-SNP interactions"
- Kwan-Yeung Lee, Kwong-Sak Leung, Nelson L. S. Tang & Man-Hon Wong. Sci Rep 2018;8:15186

Abstract: To deal with the enormous search space, our search algorithm is accelerated with eight biological plausible interaction patterns and a precomputed look-up table. After our search, we have discovered several SNPs having a stronger association to psoriasis when they are in combination with another SNP...

Ott "Genotype Patterns"

All pairs of SNPs

Ueki & Cordell (2012) *PLoS Genet* 8(4): e1002625

- □ Interaction statistic, χ^2 (1 df). Implemented in *plink* with option --fast-epistatisis joint-effects
- □ Applied to schizophrenia data: 2,164 males, 853,934 SNPs
- □ Trend genotype test (*plink*), permutation testing with 10,000 replicates: 5 SNPs with p < 0.05 by permutation and Bonferroni.
- Interaction tests for all pairs of SNPs, disregarding the 5 SNPs significant in trend test.

even though n	$p_{Bon} <$	#SNPs	same chr.
obtained as $853.929 \times n$.	0.01	259	41
\square General result?	0.05	452	63

Ott "Genotype Patterns"

3. Frequent Pattern Mining

- Thirty years ago, supermarkets started collecting huge amounts of consumer data at their cashiers. Consumer habits – if someone buys bread, how likely will they also buy milk and wine?
- Apriori algorithm (Agrawal et al, ACM SIGMOD Conference on Management of Data 1993; 207-216): Efficient search for frequent sets of items ("itemsets") purchased by one consumer ("transaction"). Development of association rules, that is, conditional probabilities P(Y | X), with Y and X being items or itemsets.
- Research published in conference proceedings, rarely in traditional journals.
- In the absence of strong main effects, we need to directly search for genotype patterns (at two or more variants) with different frequencies in cases and controls, without consulting main effects.
- Zhang Q, Long Q, Ott J. AprioriGWAS, a new pattern mining strategy for detecting genetic variants associated with disease through interaction effects. PLoS Comput Biol. 2014 Jun;10(6):e1003627
- Other implementations of search algorithms, e.g. *fpgrowth* (<u>https://borgelt.net/software.html</u>). Huge memory demands: Using Linux desktop with 512 GB of memory.

Ott "Genotype Patterns"

12

10

4. Bayesian networks

- □ Guo Y et al. Epi-GTBN: An approach of epistasis mining based on ... Bayesian network. BMC Bioinform. 2019;20:444
- Like many other approaches, Epi-GTBN employs a Bayesian network, that is, a probabilistic model to represent actions and interactions among variants and phenotypes.
- □ Authors analyzed a well-known dataset on age-related macular degeneration (AMD), which has been investigated by various other researchers. For analysis by Epi-GTBN, to reduce the computational burden, only the 1,039 SNPs with smallest p-values (p < 0.01) out of the original 103,611 SNPs were retained.</p>
- □ Results were comparable to those obtained elsewhere.
- Focusing on variants with strong main effects is fallacious! Frequencies of genotype patterns depend on main and interaction effects: Strong main effects are likely to lead to strong (significant) genotype patterns.

Ott "Genotype Patterns"

13

5. Newer Algorithms

 <u>http://www.philippe-fournier-viger.com/spmf/</u>
 Titarenko SS et al. Fast implementation of pattern mining algorithms Journal of Big Data. 2019; 13(6):37

- □ 1) Superb documentation, freely downloadable. Large memory requirements although not as large as for *fpgrowth* (Borgelt).
- 2) Apriori principle vs. evaluating all genotype pairs.
- □ Schizophrenia data: 853,934 SNPs vs 344,831,940,768 pairs (diff. chrom)
- □ For a given pattern (pair) of genotypes, X, its relation with phenotype Y is specified by a 2 x 2 table.

-	C .		Phenotype	x present	X absent
ш	Support,	s = a + c	Y = 2, cases	a	b
	Confidence,	c = a/(a + c)	Y = 1, controls	С	d

 Current computer approaches: Work with bitwise operations (1 word = 4 bytes = 32 bits) and with multiple threads in a single machine.

Ott "Genotype Patterns"

14

AMD data: Genotype pattern analysis

Klein et al (2005) Science 308 (5720:385-389

	Search for patterns (genotype pairs)	supp	conf	chisq	pPerm	OR	ch1	ch2	
	with minimum support of 40.	40	90.0	47.2604	0.01	18.5	3	4	
	Perform 1000 random permutations	40	90.0	47.2604	0.01	18.5	3	4	
	for <i>p</i> -value estimation (corrected for	53	81.1	42.5124	0.098	3 9.5	1	5	
	multiple testing).	56	78.6	39.481	0.315	5 8.1	6	7	
	Find <i>m</i> = 18,044,794 patterns.	56	78.6	39.481	0.315	5 8.1	6	7	
	Two patterns are significant, $p = 0.015$,								
	compared with the best $p = 0.60$ in sing	le-		alpha	Perm	Bonf	FDF	R-BY	
	variant analysis by trend test.				0	2		2	
	Expect many more significant genotyp	e		0.01	0	9		3	
_	patterns than single-variant results.			0.02	2	11		11	
	Different ways of establishing significa Bonferroni correction: EDR depends or	nce:		0.03	2	16		13	
	large number <i>m</i> of "null" results			0.04	2	18		61	
	Compare confidence of 90% with 43% of	0.05	2	19	11	1905			
	cases ("null" confidence) in data.								
	Ott "Capatura Dattarra" 15								






















































27





We start from the effect of selection on

quantitative traits

Historically, observations on selection gave first clue about polygenicity of quantitative traits

• We will look at the same questions from a different perspective a little later





































Additive variance $Var(Y) = \sum 2\beta_j^2 p_j q_j + \overline{\epsilon^2}$ Additive genetic variance (variance due to independent effects of alleles) $V_A = \sum 2\beta_j^2 p_j q_j$ 46























aa

A:











62

1. Common variants of weak effect

Likely reasons for missing heritability

- 2. Rare variants of larger effect
 - 3. Epistatic interactions













		AA	Aa	аа	
	Genotypes	0	1	2	
	Normalized genotypes	$\frac{0 - E(X)}{\sqrt{Var(X)}}$	$\frac{1 - E(X)}{\sqrt{Var(X)}}$	$\frac{2 - E(X)}{\sqrt{Var(X)}}$	
	Normalized genotypes	$\frac{-2q}{\sqrt{2pq}}$	$\frac{p-q}{\sqrt{2pq}}$	$\frac{2p}{\sqrt{2pq}}$	
		In these not	ations, $V_A = \sum \beta^2$		
68					

$$\begin{split} Y_i &= \sum_j \beta_j X_{ij} + \varepsilon \\ \text{Xs} - \text{Normalized genotype of individual } i \text{ at SNP} j \end{split}$$
In the matrix form: $\overline{y} &= X \overline{\beta} + \varepsilon \\ GRM &= \frac{1}{N} X X^T \end{cases}$ Now, assume that allelic effects are random $\beta \sim N(0, \sigma^2) \qquad V_A = M \cdot \sigma^2 \end{split}$



















$$\beta_j \sim N\left(0, \left[2p_j(1-p_j)\right]^S \sigma_\beta^2\right)\pi + \phi(1-\pi)$$





Evidence in favor of the highly polygenic model $\int_{0}^{0} \int_{0}^{0} \int_{0}^$





































Mutation rate is variable along the genomeImage: Colspan="2">Image: Colspan="2">Image: Colspan="2">Image: Colspan="2">Image: Colspan="2">Image: Colspan="2">Image: Colspan="2">Image: Colspan="2">Image: Colspan="2"Image: Colspan="2"</t















Demographic history





Basic facts about human genetic variation

- Nucleotide diversity (density of nucleotide differences between two randomly chosen chromosomes) is about 0.001
- Most common SNPs are very old (~300-400K years old)
- Protein coding regions are showing clear signs of selection (reduced diversity and excess of rare alleles)

15







- Selection coefficient (s) is the expected relative loss of fitness due to the sequence variant
- Variants with selection coefficients less than ~1/Ne are insensitive to selection. This is the drift barrier

14

Methods of mathematical population genetics

16

Diffusion approximation

Random walk that does not jump long distances can be approximated by a diffusion process

$$\frac{\partial \phi(x, p, t)}{\partial t} = -\frac{\partial M \phi(x, p, t)}{\partial x} + \frac{1}{2} \frac{\partial^2 V \phi(x, p, t)}{\partial x^2}$$





Site Frequency Spectrum

Non-synonymous

20

0.6

0.5

0.4 Lobortion 0.2























35

What happens if we incorporate drift?

- 1) The approach fails if selection is weak
- 2) The approach fails if mutational target is small
- 3) These considerations are important for regional constraint scores
- 4) Overall, the approach is non-informative in case of recessivity

34











 $|s_{\rm dn}|$



Stabilizing selection is the most common type of selection on a quantitative trait









Koch & Sunv





Effect sizes of individual variants are very small

- Genotype at a single locus carries very little information about phenotype.
- It does not mean that one cannot predict phenotype from genotype.
- Accuracy (r²) of an ideal genetic predictor equals heritability.











Why estimate genetic risk? • An estimate of the long-term risk at birth • Genetic risk can be combined with biomarkers and clinical features • Genetics explains about 50% of risk. One cannot predict risk any better than that but 50% is a non-trivial proportion of risk

<image><image><image><image><image><image><image><image>







12

















































- NPS accounts for the correlation of sampling errors in GWAS summary statistics.
- NPS provides an extensible framework to estimate the shrinkage curve from training data.
- NPS is best-suited to take advantage of the high density of markers and imputation accuracy in latest GWAS datasets.

The preprint is available in BioRxiv: Chun et al. AIHG 2020 "Non-parametric polygenic risk prediction via partitioned GWAS summary statistics." Software is available at: https://github.com/sgchun/nps

Identifying functionally significant variants



.





One of most significant types of variants usually leading to the complete loss of function.

Nonsense variants are enriched in sequencing artifacts

Important considerations: i) location along the gene, ii) does the variant cause NMD? iii) is the variant in a commonly skipped exon?

Tool: LOFTEE

4































































• Rare variant studies in search of drug targets

















































 Mostly coding • Are in "putatively causative" genes









- Variants involved in common forms alter regulatory sequence of these genes.
- This in turn induces changes in gene expression; regulatory variants are *eQTLs*.


















