

Advanced Gene Mapping Course


April 22-26, 2024
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	10
Rare Variant Association Analysis ²	24
Generalized/Linear Mixed Models and Interaction ¹	37
Power/Sample Size Estimation ²	49
Imputation ²	56
Linkage Disequilibrium and its Application in Association Studies ³	61
Statistical Fine-Mapping in Genetic Association Studies ³	64
Fine-Mapping using Summary Statistics ³	92
Integrating GWAS with Functional Annotations ³	97
Multivariate Analysis in Genetic Association Studies ³	102
Transcriptome-Wide Association Studies (TWAS) ³	110
Special Lecture - Genotype Pattern Mining for Digenic Traits ⁴	115
Pleiotropy and Mediation Analysis ⁵	117
Mendelian Randomization ⁵	131
Pleiotropy – review of exercise results.....	138
Mendelian Randomization – review of exercise results.....	140
Ethics and Regulation of Human Subjects Research ⁶	142
Polygenic Risk Scores ⁷	151
Population Genetics ⁷	158
Functional Annotation ⁷	167

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

Newcastle University 

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

1

Genome-wide association studies (GWAS)

- Popular (and highly successful) approach over past ~ 18 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

2

Genome-wide association studies (GWAS)

Two individuals

Person 1 ACCTGTGTGCCCAATGGGGTCCCATACTATCGG
 ACCTGTGGGCCAATGGGGTCCCATACTATCGG

Person 2 ACCTGTGGGCCAATGGGGTCCCATACTATCGG
 ACCTGTGGGCCAATGGGGTCCCATAGTATCGG

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

3

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

4

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

5

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

6

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 $\frac{(O-E)^2}{E}$ where O and E 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 Cordill (Newcastle) GWAS (Part 1) 5 / 40

7

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 $\frac{(O-E)^2}{E}$ where O and E 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) = $\frac{a/f}{c/d}$
 - OR(1|2: 1|1) = $\frac{e/d}{c/f}$

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

8

Odds ratios

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

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

9

Odds ratios

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

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

10

Odds ratios

- Odds of disease are defined as P(diseased)/P(not diseased)
 - Odds ratio OR(2|2: 1|1) represents the factor by which your odds of disease must be multiplied, if you have genotype 2|2as opposed to 1|1
 - i.e. the 'effect' of genotype 2|2
 - Similarly, we can define the OR for 1|2vs 1|1
 - As the factor by which your odds of disease must be multiplied, if you have genotype 1|2as opposed to 1|1
 - i.e. the 'effect' of genotype 1|2
 - ORs are closely related (often =) genotype relative risks
 - The factor by which your probability of disease must be multiplied, if you have genotype 1|2as 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.
 - So the association test can be thought of as a test of the null hypothesis that the ORs=1

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

11

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.

Heather Cordill (Newcastle) GWAS (Part 1) 7 / 40

12

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:

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):

Heather Cordill (Newcastle) GWAS (Part 1) 8 / 40

13

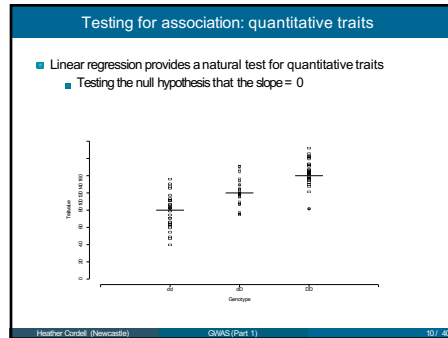
Counting alleles

Allele	Counts in		Allelic OR = ad/bc
	Cases	Controls	
2	2100 (=a)	1220 (=b)	
1	1900 (=c)	2780 (=d)	
Total	4000	4000	

- χ^2 test statistic on 1 df = $\sum (O_i - E_i)^2 / E_i$, where O_i and E_i are the observed and expected values in cell i .
 - Assumes HWE under null and multiplicative allelic effects under alternative: considers chromosomes as independent units
 - Better approach:** use counts in previous genotype table to perform a Cochran-Armitage trend test
 - Even better approach:** use linear or logistic regression

Heather Cordill (Newcastle) GWAS (Part 1) 9 / 40

14



15

Logistic regression

- Used in case/control studies
 - Outcome is affected or unaffected
 - Model probability (and thus odds) of disease p as function of variable x coding for genotype:

$$\ln \frac{p}{1-p} = \beta_0 + \beta_1 x \equiv c + mx$$
 - Use observed genotypes in cases and controls to estimate the values of regression coefficients β_0 and β_1
 - And to test whether $\beta_1 = 0$

Heather Cordill (Newcastle) GWAS (Part 1) 11 / 40

16

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 Cordill (Newcastle) GWAS (Part 1) 12 / 40

17

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

Heather Cordill (Newcastle) GWAS (Part 1) 13 / 40

18

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** 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. $\log_{10} BF > 4$) requires you to make some prior assumptions regarding what proportion of SNPs in the genome are likely to be associated with the phenotype

Heather Cordell (Newcastle) GWAS (Part 1) 13 / 40

19

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

Heather Cordell (Newcastle) GWAS (Part 1) 14 / 40

20

Close-up of hit region

Heather Cordell (Newcastle) GWAS (Part 1) 15 / 40

21

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! 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

Heather Cordell (Newcastle) GWAS (Part 1) 16 / 40

22

Complement Factor H in AMD

- Of the 3 SNPs followed up:
 - 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.

Heather Cordell (Newcastle) GWAS (Part 1) 17 / 40

23

GWAS

- GWAS really got going in around 2007
 - Visscher et al. (2012) AJHG 90:7-24 "Five Years of GWAS Discovery"
 - Visscher et al. (2017) AJHG 101:5-22 "10 Years of GWAS Discovery: Biology, Function and Translation"
 - Abdelaloui et al. (2023) AJHG 110:179-194 "15 Years of GWAS Discovery: Realizing the promise"
- 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/>

Heather Cordell (Newcastle) GWAS (Part 1) 18 / 40

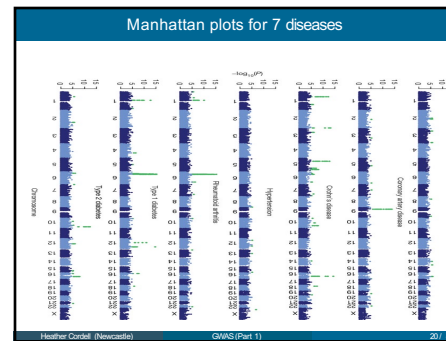
24

WTCCC

- 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

Heather Cordell (Newcastle) GWAS (Part 1) 19 / 40

25



26

Lessons from WTCCC (and others)

- Typically used rather standard statistical/epidemiological methods (χ^2 tests, t tests, logistic regression etc.)
- Success largely due to:
 - 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** ($p = 5 \times 10^{-8}$) to account for multiple testing and/or low prior prob of true effect
 - Importance of **replication** in an independent data set

Heather Cordell (Newcastle) GWAS (Part 1) 21 / 40

27

Short break

Heather Cordell (Newcastle) GWAS (Part 1) 22 / 40

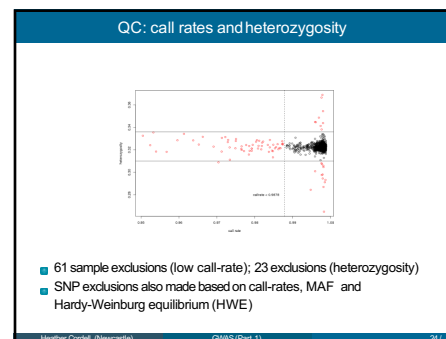
28

Quality Control

- 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
 - 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
 - Exclude SNPs with low minor allele frequency (MAF)
- See tutorials at:
 - <https://pubmed.ncbi.nlm.nih.gov/21085122/>
 - <https://pubmed.ncbi.nlm.nih.gov/29484742/>

Heather Cordell (Newcastle) GWAS (Part 1) 23 / 40

29



30

QC: ethnicity tests

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

Heather Cordell (Newcastle) GWAS (Part 1) 25 / 40

31

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) GWAS (Part 1) 26 / 40

32

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)
- 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

Heather Cordell (Newcastle) GWAS (Part 1) 27 / 40

33

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)
- 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) gross population outliers

Heather Cordell (Newcastle) GWAS (Part 1) 28 / 40

34

Principal components analysis (PCA)

Genes mirror geography within Europe

J Novembre et al. (2008) Nature 456(7218):98-101, doi:10.1038/nature07331

Heather Cordell (Newcastle) GWAS (Part 1) 27 / 40

35

Principal Components Analysis

- Price et al. (2006) Nature Genetics 38:904-909; Patterson et al. (2006) PLoS Genetics 2(12):e190
 - Based on popn genetics ideas from Cavalli-Sforza (1978)
- Idea is to form a large matrix M of SNP counts (0,1,2) corresponding to the genotype at L loci (=rows) for n individuals (=columns)

$$M = \begin{matrix} & \begin{matrix} g_{11} & g_{12} & \dots & g_{1n} \end{matrix} \\ \begin{matrix} g_{21} \\ g_{31} \\ \vdots \\ g_{L1} \end{matrix} & \begin{matrix} g_{22} & g_{32} & \dots & g_{2n} \\ g_{32} & g_{32} & \dots & g_{3n} \\ \vdots & \vdots & \ddots & \vdots \\ g_{L2} & g_{L2} & \dots & g_{Ln} \end{matrix} \end{matrix}$$

Heather Cordell (Newcastle) GWAS (Part 1) 28 / 40

36

Principal Components Analysis

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

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ x_{31} & x_{32} & \dots & x_{3n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{L1} & x_{L2} & \dots & x_{Ln} \end{bmatrix}$$
- This matrix will be used as starting point for PCA
 - In principal we could start with a different matrix—in particular not all PCA approaches would normalise by $\sqrt{f_i(1-f_i)}$

Heather Cordell (Newcastle) GWAS (Part 1) 29 / 40

37

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)

Heather Cordell (Newcastle) GWAS (Part 1) 30 / 40

38

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 and eigenvalues λ_j of matrix Ψ
 - Co-ordinate j of the k th eigenvector represents the ancestry of individual j along 'axis' k

Heather Cordell (Newcastle) GWAS (Part 1) 31 / 40

39

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 and eigenvalues λ_j of matrix Ψ
 - Co-ordinate j of the k th eigenvector represents the ancestry of individual j along 'axis' k
- For technical details, see McVean (2009) PLoS Genetics 5;10:e1000686

Heather Cordell (Newcastle) GWAS (Part 1) 31 / 40

40

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 and eigenvalues λ_j of matrix Ψ
 - Co-ordinate j of the k th 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
 - For more subtle differences, you may need to calculate more PCs
 - And include them as covariates in the regression equation
 - Post-GWAS QC can determine whether you have included 'enough'

Heather Cordell (Newcastle) GWAS (Part 1) 31 / 40

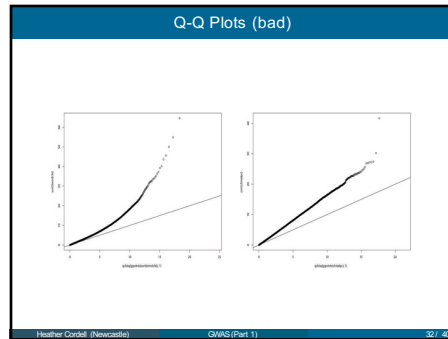
41

Post GWAS QC: Q-Q Plots (good)

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

Heather Cordell (Newcastle) GWAS (Part 1) 31 / 40

42



43

Population stratification

- A Q-Q 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)
 - Divide each test statistic by λ to get an adjusted (deflated) test statistic
- More complicated solution: use PCA/MDS or similar Even
- more complicated solution: use linear mixed models

Heather Cordell (Newcastle) GWAS (Part 1) 33 / 40

44

Relatedness

- 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 be identified with only a few hundred markers
 - More complicated relationships require 10,000-50,000 SNPs
- Various software packages, including PLINK, KING and TRUFFLE

Heather Cordell (Newcastle) GWAS (Part 1) 34 / 40

45

Expected IBD sharing

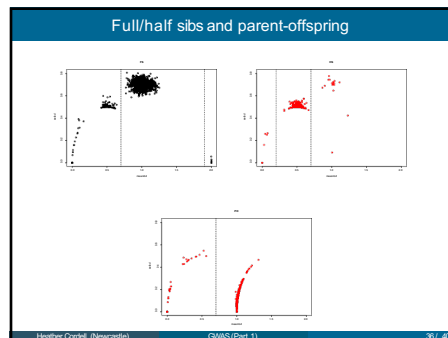
- Assuming no inbreeding, the IBD state probabilities are:

Relationship	Number of alleles shared IBD		
	2	1	0
MZ twins	1	0	0
Parent-Offspring	0	1	0
Full siblings	1/4	1/2	1/4
Half siblings	0	1/2	1/2
Grandchild-grandparent	0	1/2	1/2
Uncle/aunt-nephew/niece	0	1/2	1/2
First cousins	0	1/4	3/4
Second cousins	0	1/16	15/16
Double 1st cousins	1/16	6/16	9/16

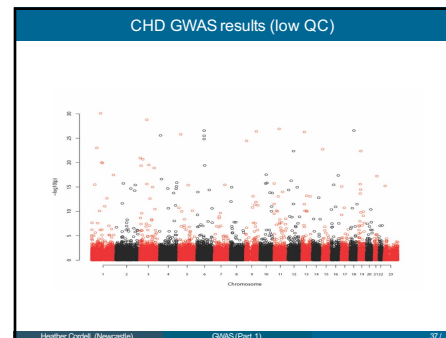
- A useful visualisation tool is to plot SE(IGD) vs mean(IGD) (as estimated across the genome)
 - Or kinship coefficient $\{ \frac{E(IGD=2)}{E(IGD=1)} \}$ against $P(IGD=0)$

Heather Cordell (Newcastle) GWAS (Part 1) 35 / 40

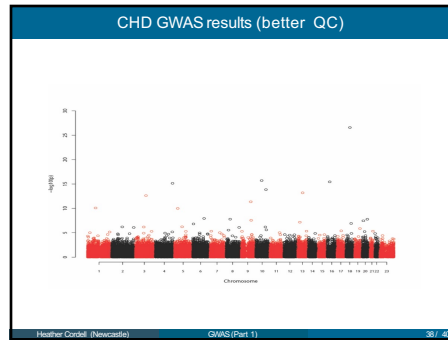
46



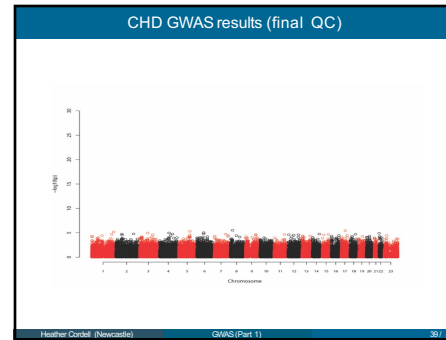
47



48



49



50

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

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

51

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

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

52

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) 41 / 41

53

Data Quality Control NGS and Genotype Array Data

Suzanne M. Leal, Ph.D.
sm13@cumc.columbia.edu

© 2024 Suzanne M. Leal

1

DNA Collection

- Blood samples
 - For unlimited supply of DNA
 - Transformed cell lines
 - is expensive
 - Whole genome amplification
 - Allows for the creation of large amounts of DNA from initial small DNA sample
 - Perform WGA on each sample three or more times and use pooled samples
 - Can experience lower call rates and higher genotyping error rates
 - Not recommend for whole genome sequencing or copy number variant (CNV) analysis
- Buccal Swabs
 - Small amounts of DNA
 - DNA not stable
- Saliva (Origene collection kit)

Measurement of DNA Concentrations

- Nanodrop
- Picogreen

2

Effect of Genotyping Error – Same Error Rates for Cases and Controls

- For family-based association studies - Trios
 - Can increase both type I and II error
- Population based studies
 - Increases type II error only

Quantitative Traits

If genotyping error is not correlated with trait values type II errors will be increased

3

Effects of Genotyping Error – Different Error Rates for Cases and Controls

- Cases and controls are sequenced/genotyped
 - At different times
 - Different institutions
 - Or one group, e.g., case or control, is predominately sequenced/genotyped in the same batch
- Can lead to different genotyping error rates in cases and controls
 - In this situation both type I and II error can be increased
- If sequencing/genotyping cases and controls
 - Randomize cases and controls so they are spread evenly across batches

Quantitative Traits

If genotyping error is correlated with trait values, it will also increase type I and II errors, e.g., individuals with elevated systolic blood pressure are genotyped in one batch and those with systolic blood pressure within the normotensive range in another batch

4

Genotype SNPs (~20-96) before Exome or Whole Genome Sequencing

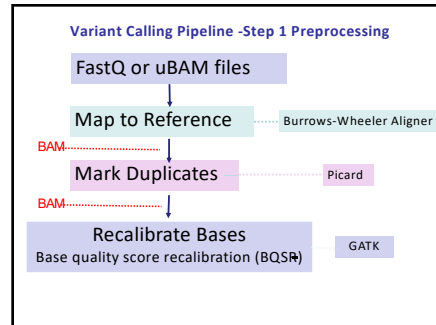
- Genotype markers which can be used as DNA fingerprint
- Allows for Assessment of DNA quality
- Aids in determining the genetic sex of study subjects
 - To aid in identification of potential sample swaps
- Detects cryptic duplicates
- For family data
 - Aids in determining close familial relationships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

5

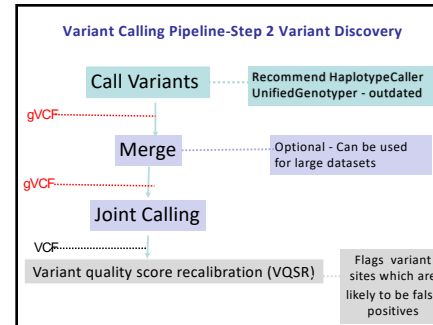
Detecting Genotyping Errors

- 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

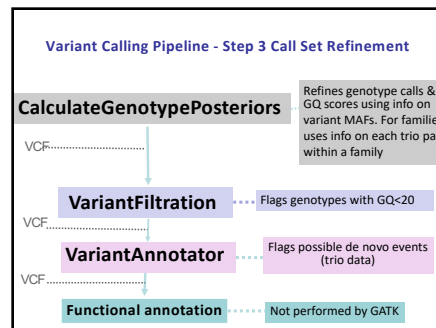
6



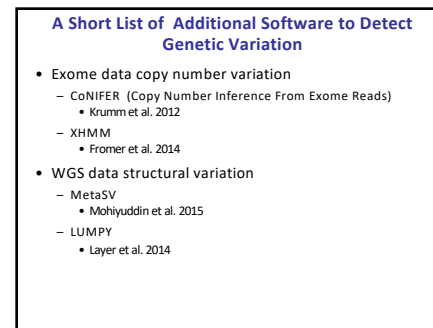
7



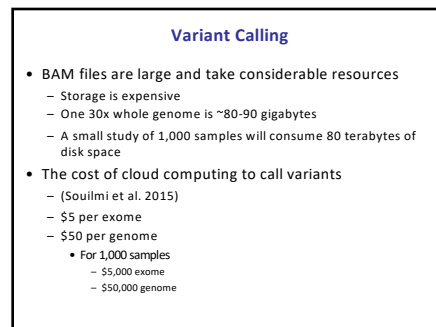
8



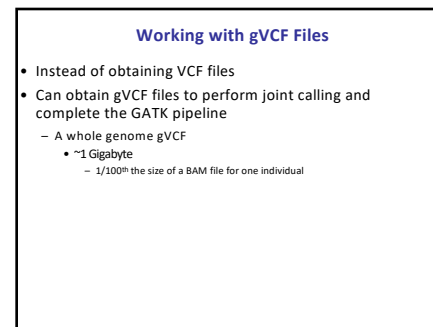
9



10



11



12

Influences on Sequence Quality

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

13

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., $DP < 8$
 - Genotype quality (GQ) score
 - Removal genotypes with a low genotype quality core, e.g., $GQ < 20$
 - Bcftools
 - Can be used to remove variants sites and genotypes which do not meet quality control criteria

15

VCF Example

```

##fileformat=VCF v.4.2
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##INFO=
##
```

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

19

NGS Data Quality Control

- Values which are used for DP (genotype), GQ, and missing data cut offs are based upon
 - Concordance rates
 - If there are duplicate samples are available
 - Ti/Tv ratios
 - By individual
 - By batch
 - Entire data set
 - Amount of data removed
 - QC can remove substantial amounts of data which should be avoided
 - e.g., >15% of variant sites

20

Transition/Transversion (Ti/Tv) Ratios

- Transition
 - Purine → Purine
 - Pyrimidine → Pyrimidine
- Transversion
 - Purine → Pyrimidine
 - Pyrimidine → Purine

AKA Ts/Tv ratios

21

Transition/Transversion (Ti/Tv) Ratios

- Ti/Tv Ratios
 - Whole genome ~2.0
 - Exome novel ~2.7
 - Exome known ~3.5
- Ti/Tv ratios can be calculated by
 - Sample or
 - Dataset
- Ti/Tv ratios can be evaluated for subsets of data
 - e.g., by batch

22

Sequence Data QC Overview

- Variant and genotype call level
 - Evaluation of batch effects
- Genotype call level – Removal of genotype calls
 - Low or high depth of coverage DP < 8
 - Low genotype quality score GQ < 20
- Removal of individual samples
 - >20% missing data
 - After taking the intersect of capture arrays
 - Samples without phenotype information

23

Sequence Data QC Overview

- Variant level – removal of variant sites
 - Low call rate
 - 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., $p < 5 \times 10^{-8}$, $p < 5 \times 10^{-15}$

24

QC – 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 XO
 - Klinefelter syndrome XXY

25

QC – Assessing Sex Chromosomes

- Study subjects labeled as females with an excess of homozygous genotypes on the X chromosome can denote
 - That their genetic sex is male
 - Turner Syndrome

26

QC – Assessing Sex Chromosomes

- Study subjects labeled as males with an excess of heterozygous SNPs* on the X chromosome can denote
 - That their genetic sex is female
 - Klinefelter syndrome
- Note: Individuals who are XY will also be heterozygous for markers in the pseudoautosomal regions
- Availability of Y chromosome data
 - Can greatly aid in determining genetic sex and if an individual has Turner or Klinefelter syndrome

*Both genetic males and females have two alleles for each locus on the X chromosome in the dataset, although genetic males are hemizygous

27

Data Clean – Assessing Sex Chromosomes

- Individuals whose labeled gender/sex does not match 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

28

Checking for Duplicate and Related Individuals

- Duplicate samples are sometimes included in a study as part of quality control to detect inconsistencies
 - Will not detect systematic errors
 - Usually not included in exome and whole genome sequencing studies
 - Intentional duplicates can easily be removed before data quality control
- Cryptic duplicates (unintentional)
 - DNA sample aliquoted more than once
 - Individual ascertained more than once for a study
 - e.g. The same individual undergoes the same operation more than once and is ascertained each time
- Individuals who are related to each other may participate in the same study
 - Unknown to the investigator
 - Or be part of the study design

29

Duplicate and Related Individuals Need to be Identified

- For duplicate samples
 - Only one can be retained
- For related individuals
 - PCA is performed first with unrelated individuals and related individuals are then projected onto the PCs of unrelated individuals
 - Mixed-models need to be used to analyze the data if related individuals are included*
 - Case-Control
 - Generalized linear mixed models (GLMM)
 - Quantitative traits
 - Linear mixed models (LMM)
- If related individuals are ignored in the analysis type I error rates can be inflated

*If only a few related individuals in sample, may wish to remove them or use LMM/GLMM to control type I errors. Must use LMM/GLMM if related individuals are included in the dataset. If possible, opt for LMM/GLMM since it can help to control type I error due to other types of structure in the data, even when no closely related individuals are included in the analysis.

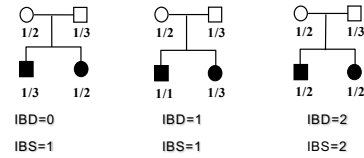
30

Identifying Duplicate and Related Individuals

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

31

Identity by Descent (IBD)/Identity-by-State (IBS)



32

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 linkage disequilibrium (LD)
 - e.g., $r^2 < 0.1$
- \hat{p} -hat is calculated using the “population” allele frequency
- Used to approximate 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

33

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

34

Identifying Duplicate and Related Individuals

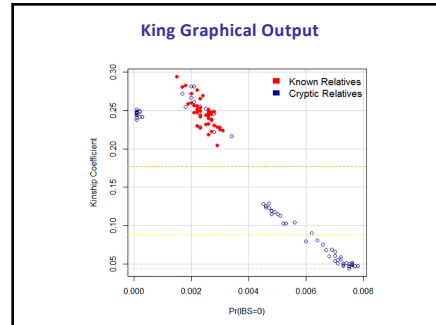
- KING [Kinship-based Inference for Gwas (*Manichaikul et al. 2010*)] can also be used to identify duplicate and related individuals
 - KING is more robust to population substructure and admixture
 - Prune markers for LD (e.g., $r^2 < 0.1$)
 - Provides kinship coefficients
 - Duplicate samples
 - Kinship coefficient equals 0.5
 - Siblings
 - Kinship coefficient equals 0.25

35

UK Biobank Related Individuals > Kinship Coefficient 0.0625

White European		African		Asian	
# of Relatives	# of relatives	# of relatives	# of individuals	# of relatives	# of individuals
1	86089	1	715	1	743
2	18491	2	153	2	115
3	3691	3	26	3	33
4	707	4	10	4	4
5	165	5	3	5	4
6	40	6	5		
7	9	7	5		
8	5	8	4		
9	1	9	1		
10	11	10	4		
11	2	11	2		
12	2	12	3		
16	1	17	2		
19	1	19	3		
25	1	19	2		
30	1	21	1		
3085	1	23	1		
		-	-		
		-	-		
		-	-		
		-	-		
		390	1		
		391	1		
		393	1		
		396	1		

36



37

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

Principal Components Analysis (PCA) / Multidimensional Scaling (MDS)

- Can be used to identify outliers
- Population substructure
 - Individuals from different ancestry
 - e.g., African American samples included in samples of European Americans
- Batch effects
- Use a subset of markers which have been LD pruned
 - Only very low levels of LD between marker loci
 - e.g., $r^2 < 0.1$
 - MAF cutoff dependent on sample size
 - e.g. MAF > 0.01
 - Can use lower MAF for large sample sizes

39

Principal Components Analysis (PCA) / Multidimensional Scaling (MDS)

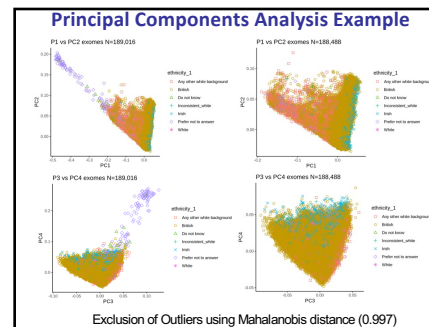
- Unrelated individuals are used to generate PC plots
 - Related individuals are projected onto the PC plots
- Plot 1st component vs. 2nd component
 - Additional PCs should also be plotted
 - e.g., PCs 1-10
- Mahalanobis distance can be used to determine outliers
 - e.g., < 1

40

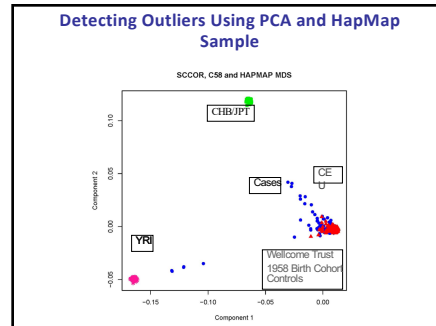
PCA/MDS Can be Used to Identify Outliers

- Individuals of different ancestry
 - e.g., African American samples included with European Americans samples
 - Can use samples from HapMap/1000 genomes to help to determine the ancestry for samples that are outliers
 - Should not include HapMap/1000 genomes samples when calculating components to control for population substructure/admixture
- Batch effects

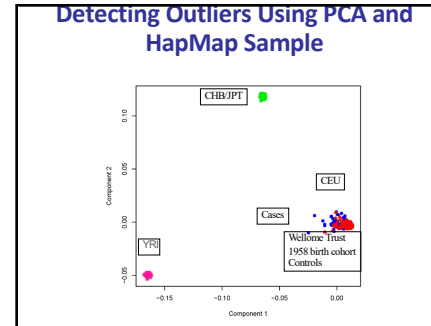
41



42



43



44

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)

45

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 5.0×10^{-7} to 5.0×10^{-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

47

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

48

Sequence Data QC Overview

- Evaluate samples for cryptically related individuals and duplicates
 - Use variants which have been pruned for LD
 - e.g., $r^2 < 0.1$
 - 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

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., $r^2 < 0.1$
 - 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

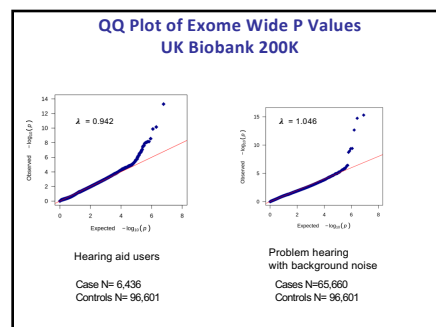
- Remove/flag variant sites that deviate from HWE in controls
 - HWE should be only be tested in unrelated individuals from the same population
- Post Analysis - Quantile-Quantile (QQ) plots
 - To evaluate uncontrolled batch effects and population substructure/admixture

51

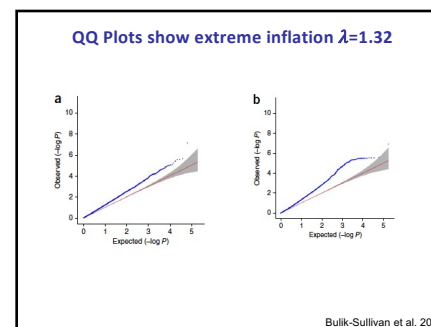
QQ Plots - Genome Wide Association Diagnosis

- 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 $-\log_{10}(p)$ values under the null (uniform distribution) on the x axis

52



53



Bulik-Sullivan et al. 2015

54

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 $\lambda=1$
 - Inflation $\lambda>1$
 - Deflation $\lambda<1$
 - 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

55

Phenotype	Covariate	Mean Chi-Square	GIF (λ)
BP		1.23829	1.16932
BP	Age	1.24119	1.18025
BP	Age-EV1	1.05471	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-EV4	1.0428	1
BPI	Age-EV10	1.04204	1
BPI	Sex, Age-EV1	1.05421	1.01724
BPII		1.17283	1.25664
BPII	Age	1.17583	1.26996
BPII	Age-EV1	1.09874	1.15085
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	1.09558	1.06424
BPII	Sex, Age-EV4	1.05817	1.03327
BPII	Sex, Age-EV10	1.06338	1.05581

56

Evaluating Reason for Inflated λ

- LD score regression (LDSC) can be used to determine if the observed λ is inflated due to
 - Problems in the data
 - Population substructure/admixture
 - Batch effects/genotyping errors
 - Polygenicity
 - Many associated loci each with a very small effect size
- LDSC is performed and the intercept is examined
 - If intercept is >1 than inflation is due to population substructure, etc.
 - If intercept is ~ 1 than $\lambda < 1$ is due to polygenicity

57

- Bulik-Sullivan et al. (2015) performed simulation studies using LDSC regression to evaluate polygenicity

Panels a & c data were simulated with population substructure. The $\lambda=1.32$ (a) & LDSC intercept = 1.30 (c)

Panels b & d data were simulated with polygenicity with 0.1% of variants having a causal effect. The $\lambda=1.32$ and LDSC intercept = 1.006

Panel a: LDSC regression line for population substructure, showing a positive slope and intercept of 1.30.

Panel b: LDSC regression line for polygenicity, showing a positive slope and intercept of 1.006.

Panel c: Scatter plot of observed λ^2 vs LD Score (M²) for population substructure, showing a positive linear relationship.

Panel d: Scatter plot of observed λ^2 vs LD Score (M²) for polygenicity, showing a positive linear relationship.

58

Post Analysis QC

- Observe in Manhattan plots individual associated variants with no surrounding associated variants

59

Post Analysis QC

- Most variants are in LD with neighboring SNPs
- Genotyping error can cause a variant site not to be in LD with any of its neighbors
- Genotyping error can also cause a spurious associations
- A lone associated variant site can be due to genotyping error

60

Post Analysis QC

- Imputation can be used to determine if for the variant site there is genotype error
- Variant site is imputed
 - Check how accurate variant is imputed
 - R² or INFO score
 - If imputation accuracy is high e.g. R²>0.8
 - Check the correlation between the imputed variant and sequence or genotype array data
 - R²
 - If r² is low there is genotyping error
- The variant site should be removed

61

Example Project Description

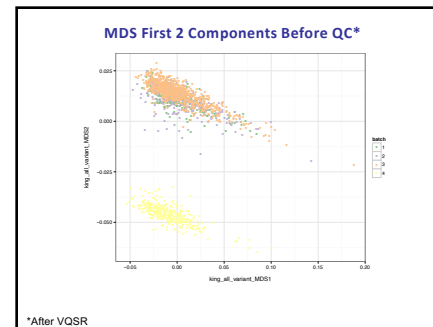
- 1,667 Samples
- Seven cohorts
- Two sequencing centers
 - Center 1
 - Two capture arrays
 - NimbleGen V2Refseq 2010 (CA1): 1082
 - Batch 1 and 3
 - NimbleGen bigexome 2011 (CA2): 234
 - Batch 2
 - Center 2
 - One capture array
 - Agilent SureSelect
 - Batch 4
- Four batches
- No intentional duplicate samples

62

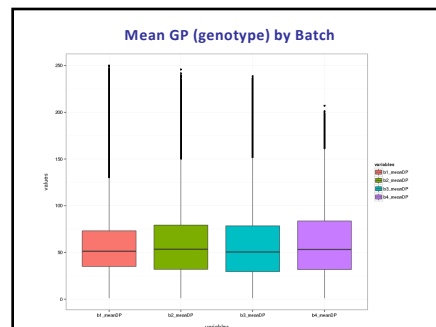
Example Project Description

- Intersection of the three capture arrays used
 - NimbleGen V2Refseq 2010
 - Batch 1 and 3
 - NimbleGen bigexome 2011
 - Batch 2
 - Agilent Sure Select
 - Batch 4
- Sequencing machine
 - Illumina HiSeq
- Sequence alignment
 - BWA
- Multi-sample variant calling
 - GATK

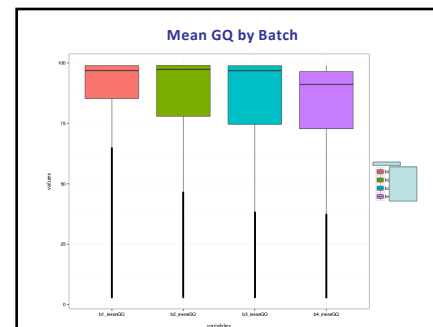
63



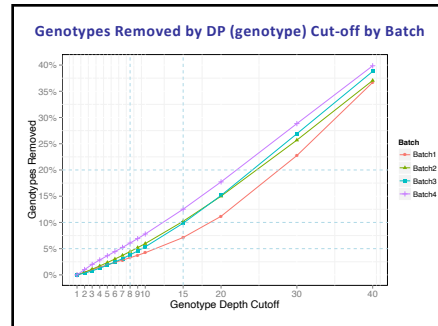
64



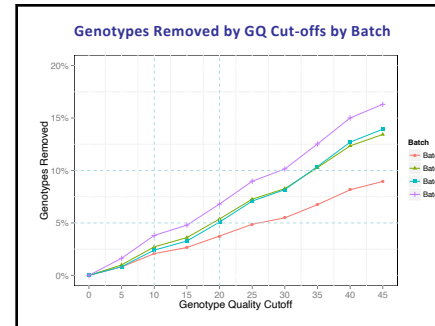
65



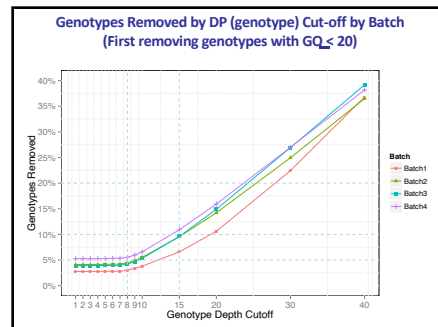
66



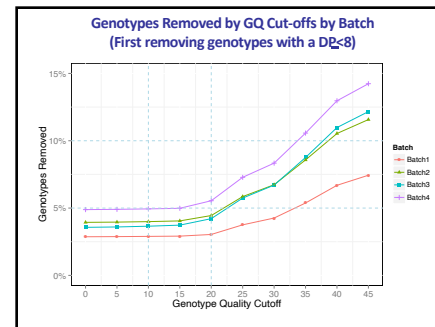
67



68



69



70

Missing Rate Criteria & Sites Removed

	Variant sites removed if missing >10% of their genotypes	Variant sites removed if missing >5% of their genotypes
Percent of genotype data removed		
Before QC*	2.5%	3.9%
After QC	12.9%	18.3%

Variant sites missing >10% of their data were removed

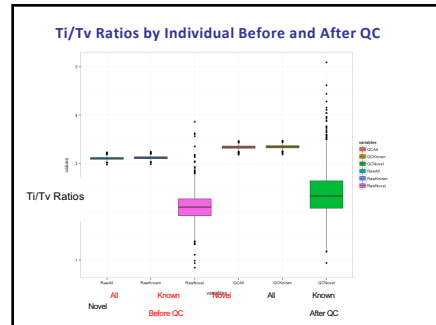
*After VQSR

71

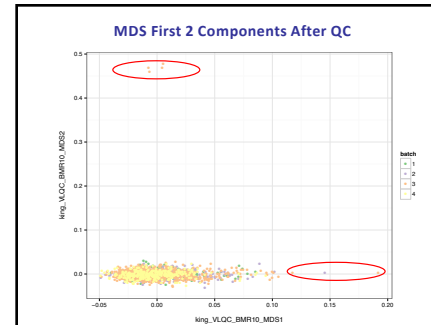
Ti/Tv Ratios during QC Process

	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 $DP \geq 8, GQ \geq 20$	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 $p < 2 \times 10^{-8}$ N=17,835	3.39 ± 0.04	2.41 ± 0.53	3.39 ± 0.04
Remove sites deviating from HWE $q < 5 \times 10^{-8}$ N=4,414	3.41 ± 0.04	2.39 ± 0.54	3.40 ± 0.04

72



73



74

- ### 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)

75

- ### Convenience Controls
- Can reduce the cost of a study
 - Genotype data
 - Ascertainment from different population
 - Differential genotyping error
 - Even if performed at the same facility
 - Proper QC can reduce or remove biases

76

- ### 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

77

- ### Genotype Array Data
- #### Genotype Data QC – Population Based Studies
- Initially remove DNA samples from individuals who are missing >10% or their genotype data
 - For variant sites with a minor allele frequency (MAF) > 0.05
 - Remove variant 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

78

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. $r^2 < 0.1$
 - Remove duplicate samples

79

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., $r^2 < 0.1$
 - 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

80

Order of Data Cleaning-Genotype Array

- Examine QQ plots
 - e.g., not controlling adequately for population admixture
 - Inflated test statistics Deflated p-values
- Examine Manhattan to detect associated variants which are not in LD with other variants
 - Genotyping errors causing spurious associations

81

Complex Trait Association Analysis of Rare Variants Obtained from Sequence Data

Suzanne M. Leal, Ph.D.
Sergievsky Family Professor of Neurological Sciences
Director of the Center for Statistical Genetics
Columbia University
smi3@columbia.edu

© 2024 Suzanne M. Leal

1

Complex Diseases (Traits)

Top 10 leading causes of death in the United States

Cause of Death	2015 (per 100,000)	2016 (per 100,000)
Heart disease	208.4	202.2
Cancer	159.2	155.2
Stroke	97.3	95.2
Chronic lower respiratory diseases	84.7	82.4
Diabetes	75.1	73.3
Alzheimer's disease	55.1	53.3
Kidney disease	45.4	43.3
Intentional self-harm	35.4	33.3
Suicide	25.4	23.3
Unintentional injuries	15.4	13.3

Genetic and environmental contribution to complex disorders

T.A. Manolio, et al. JAMA Intern Med. 2012

2

Heritability

- **Broad-sense heritability**
 - Considers all genetic factors
 - Phenotype = Genetics + Environmental Noise
 - $Y = G + E$
 - $\text{Var}(Y) = \text{Var}(G) + \text{Var}(E)$
 - $H^2 = \text{Var}(G)/\text{Var}(Y)$
- **Narrow-sense heritability**
 - Considers only additive contributions
 - Phenotype = Additive Genetics + Environmental Noise
 - $Y = A + E$
 - $\text{Var}(Y) = \text{Var}(A) + \text{Var}(E)$
 - $h^2 = \text{Var}(A)/\text{Var}(Y)$

3

Height Heritability

- The variance of human height is about $\sim 25 \text{ cm}^2$
 - Adjusted for sex
- Total Variation
 - $\sim 20 \text{ cm}^2$ due to genetics
 - $\sim 5 \text{ cm}^2$ due to other factors (noise)
- The heritability of height is $\sim 20/\sim 25 \sim 80\%$
- The heritability of height has been estimated using a variety of study types, e.g. twin, sibpairs
- Karolinska

4

Heritability for Common Traits

Human height heritability is $\sim 80\%$

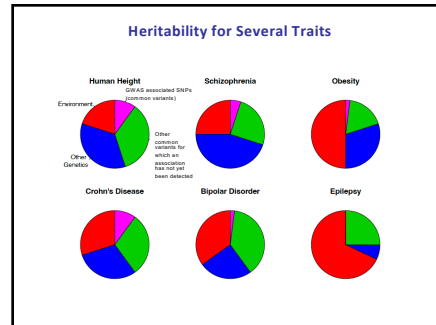
- Strongly associated common variation explain 21—29%
 - Those that statistically significant
- All common variation explains 60% of height heritability (h^2)

5

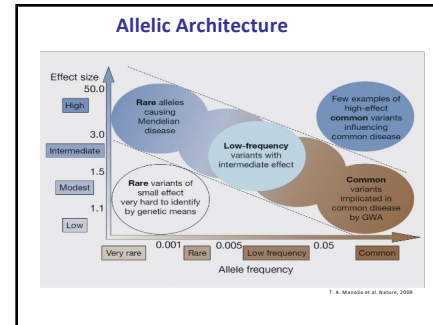
Heritability for Several Traits

Area in blue is the so called missing heritability

6



7



8

Complex Disease – Common Variant Associations

- Disease susceptibility is conferred by variants which are common within populations
 - Variants are old and widespread
- These variants have modest phenotypic effect
- This model is supported by many replicated examples
 - Age Related Macular Degeneration (Klein et al. 2005)
 - Complement factor H (CFH) gene

9

Studying Complex Traits – Common Variant Associations

- Hundreds of thousands of Single nucleotide polymorphism (SNPs) genotyped and analyzed
 - Indirect mapping
 - Markers usually had a minor allele frequency (MAF) > 0.05
 - Usually not pathogenic – tag SNPs
 - In linkage disequilibrium with disease susceptibility variant

10

Complex Trait – Common Variant Associations

- Although highly successful in identifying thousands of complex trait loci
- Usually pathogenic susceptibility variant(s) not identified

NHGRI GWA catalogue
www.genome.gov/GWAstudies
www.ebi.ac.uk/ftpp/gwas/

11

Complex Disease – Rare Variant Associations

- Complex traits are the result of multiple rare variants
 - Although first thought to large effects, there effect sizes are usually small
- Although these variants are rare, e.g., MAF<0.005
 - Collectively they may be quite common
- Direct tests of this hypothesis where first reported >15 years ago
 - Dallas Heart Study
 - Small sample ~1,200 individuals
 - Multi-ancestry
 - Used “extreme” sampling
 - Plasma low density lipoprotein levels (Cohen et al. 2004)
 - NPC1L1

12

Rationale for Rare Variant Aggregate Association Tests

- Testing individual variants with low effect sizes and minor allele frequencies (MAFs)
 - Underpowered to detect associations
- Testing variants in aggregate increases MAFs
 - Improving the power to detect associations

13

Caveats - Aggregate Rare Variant Association Tests

- Misclassification of variants can reduce power
 - Inclusion of non-causal variants
 - Exclusion of causal variants
- Analysis can be performed using region boundaries for
 - Genes
 - Genes within pathways
 - Regulatory regions
 - As determined for example by
 - FANTOM5 CAGE profiles to identify promoter regions (Noguchi et al. 2017)
 - STAAR pipeline that combines multiple *in silico* annotations (Li et al. 2020)
- Unlikely a sliding window approach will work
 - Size of window unknown and will differ across the genome

14

Analysis of Rare Variants

- For biobank sized datasets higher frequency rare variants, e.g., 0.5% can be analyzed individually
 - Using same same methods implemented for common variants

Example
 $\alpha = 5 \times 10^{-8}$
 Disease prevalence 5%
 $1 - \beta = 0.80$

*Note: a more stringent significance criterion may be necessary for genome-wide sequence data. Due to a larger number of effective tests compared to analysis of common variant GWAS panels

15

A Few Rare Variant Association Tests

- Combined Multivariate Collapsing (CMC)
 - Li and Leal AJHG 2008
- Burden of Rare Variants (BRV)
 - Auer, Wang, Leal Genet Epidemiol 2013
- Weighted Sum Statistic (WSS)
 - Madsen and Browning PLoS Genet 2009
- Kernel based adaptive cluster (KBAC)
 - Liu and Leal PLoS Genet 2010
- Variable Threshold (VT)
 - Price et al. AJHG 2010
- Sequence Kernel Association Test (SKAT)
 - Wu et al. AJHG 2011
- SKAT-O
 - Lee et al. AJHG 2012

Fixed Effect Tests: CMC, BRV, WSS, KBAC, VT
 Random Effect Test: SKAT
 Optimal test: SKAT-O

16

Types of Aggregate Analyses

- Frequency cut offs used to determine which variants to include in the analysis
 - Rare Variants (e.g., MAF<0.05% frequency)
 - Rare and low (MAF=0.05-5%) frequency variants
- Maximization approaches
- Tests developed to detection associations when variants effects are bidirectional
 - e.g., protective and detrimental
- Incorporate weights based upon annotation
 - Frequency
 - e.g., gnomAD
 - Functionality
 - CADD c-scores

17

Methods to Detect Rare Variant Associations Using Variant Frequency Cut-offs

- 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

18

CMC

- Define covariate X_j for individual j as

$$X_j = \begin{cases} 1 & \text{if rare variants present} \\ 0 & \text{otherwise} \end{cases}$$
- Compute Fisher exact test for 2x2 table

Number of cases for which one or more rare variants are observed e.g., nonsynonymous variants freq. $\leq 1\%$

	X=1	X=0
cases		
controls		

Number of cases without a rare variants

Number of controls for which one or more rare variants are observed

Number of controls without a rare variants

Can also use same coding in a regression framework

19

CMC

- Example of coding used in regression framework:
 - Binary coding $X_j = \begin{cases} 1 & \text{if rare variants present} \\ 0 & \text{otherwise} \end{cases}$
 - Gene region with 5 variant sites

	Individual	Coding
1	1	1
2	1	1
3	0	0

Rare Variant Sites
 Green bars: Major allele is observed in the study subject
 Red bars: Minor allele has been observed

20

Methods to Detect Rare Variant Associations Using Variant Frequency Cut-offs

- Gene- or Region-based Analysis of Variants of Intermediate and Low frequency (GRANVIL)
 - Aggregate number of rare variants used as regressors in a linear regression model
 - Can be extended to case-control studies
 - Morris & Zeggini 2010 Genet. Epidemiol
 - Test also referred to as MZ

21

GRANVIL

- Example of coding used in regression framework
 - Gene region with 5 variant sites – data available on all sites

Individual 1: Coded 2/5 (0.4)

Individual 2: Coded 2/5 (0.4) Note same coding for heterozygous and homozygous genotypes

Individual 3: Coded 1/2 (0.5)

Burden Rare Variant (BRV) extension (Auer et al. 2013 Genet Epidemiol)
 Individual 1: Coded 2
 Individual 2: Coded 3
 Individual 3: Coded 1

22

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 up-weighted)
 - Madsen & Browning, PLoS Genet 2009
- Kernel based adaptive cluster (KBAC)
 - Adaptive weighting based on multilocus genotype
 - Liu & Leal, PLoS Genet 2010

23

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

24

Methods to Detect Associations with Protective & Detrimental Variants within a Region

- C-alpha
 - Detects variant counts in cases and controls that deviate from the expected binomial distribution
 - For qualitative traits only
 - Neale et al. 2011 PLoS Genet
- Sequence Kernel Association Test (SKAT)
 - Variance components score test performed in a regression framework
 - Can also incorporate weighting
 - Wu et al. 2011 AJHG

25

Optimal Test

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

26

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., $p < 2.5 \times 10^{-6}$ (Correction for testing 20,000 genes)

27

Determine MAF Cut-offs for Aggregate Rare Variant Association Tests

- MAF cut-offs 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 under study
 - gnomAD
 - <http://gnomad.broadinstitute.org/>

28

Problem of Missing Genotypes for Aggregate Rare Variant Association Tests

- 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

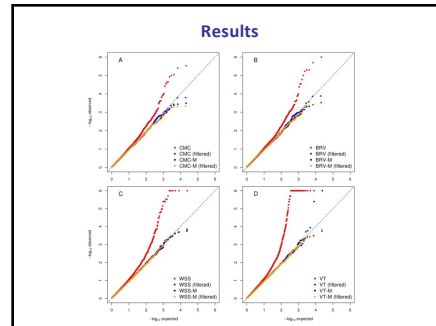
29

Dosages

- Genotypes are no longer assigned 0 (1/1), 1 (1/2) or 2 (2/2)
 - Due to uncertainty
- Each genotype is assigned a probability
 - Probabilities sum to 1
- For example
 - Probability of 0 (1/1) genotype is 0.98 and 1 (1/2) genotype is 0.015
- The dosage can be estimated for this example as follows

$$\begin{aligned}
 0 \times 0.98 &= 0 \\
 1 \times 0.015 &= 0.015 \\
 2 \times 0.005 &= 0.01 \\
 \text{Dosage} &= 0.025
 \end{aligned}$$
- Instead of using the most likely genotype the dosage is used

30



31

Rare Variant Aggregate Methods

- Ideally should be performed in a regression framework to adjust for covariates
 - Logistic
 - Linear regression

- Almost all rare variant aggregate methods have been extended to be implemented within a regression framework
- Some have also been implemented in a linear mixed model (LMM)/generalized LMM (GLMM)

32

Analyzing Quantitative Variants

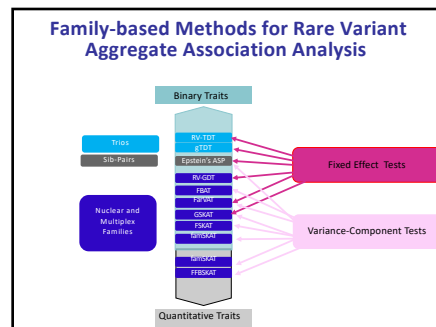
- Most rare variant aggregate analysis methods can be performed on quantitative traits
- If phenotype data includes outliers or deviates from normality
 - Can increase type I errors

33

Analyzing Quantitative Variants

- For data that deviates from normality
 - Quantile-quantile normalization
- For data that includes outliers
 - Winsorize
- Don't winsorize and then normalize
- Instead of analyzing quantitative trait values
 - Residual can be generated
 - Adjusting for confounders

34



35

Linear Mixed Model (LMM) & generalized LMM (GLMM) Analysis of Related & Unrelated Individuals

- LMM is an extension of the linear model to allow for both fixed & random effects and also allows for non-independence of samples
 - Early implementations calculated the kinship matrix Φ on the basis of known relationships
 - Amin et al. (2007) proposed to estimate kinships based on genome-wide variant data
 - The generalized relationship matrix (GRM) can be estimated for all individuals using for example identical-by-descent (IBD) sharing
- Extended to binary (case-control) traits - GLMM

36

**LMM and GLMM:
Analysis of Related & Unrelated Individuals**

- Can be applied to analyze families, cryptically related, & unrelated individuals
 - e.g., UK Biobank
 - 500K study subjects of which 30.3% are 3rd degree relatives & 4.5% sib-pairs
- More recent implementation for large scale data using a variety of methods
 - BOLT-LMM (Loh et al. 2015)
 - FastGWA (Jiang et al. 2019)
 - SAIGE (Zhao et al. 2015)*
 - REGIE (Mbatchou et al. 2020) *
 - SMMAT (Chen et al. 2019)**
- *Can be used to analyze data where case to control ratio is very unbalanced
 - e.g., 20 cases for every control
- **Cannot be used for UK Biobank Scale data

37

**LMM and GLMM:
Analysis of Related & Unrelated Individuals**

- To allow for use with biobank sized datasets
- REGIE does not use the GRM
 - It uses whole genome regression, i.e., the ridge regression
 - In essence, it includes all the SNVs as covariates in the null model
 - Performed by blocks to avoid having to load the entire genome in memory
 - » Using different effect size differences per block
- This large-scale approximation may not control type I error for individuals that are closely related
 - e.g., when only families are being analyzed
 - Can use for example SMMAT
 - Which uses the GRM

38


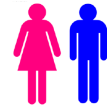
**LMM and GLMM:
Analysis of Related & Unrelated Individuals**

- A few programs which can perform rare variant aggregate analysis
 - REGIE - Burden test, SKAT, & SKAT-O
 - SMMAT - Burden, SKAT, & SKAT-O
 - rvtests (Zhan 2020) implements BOLT-LMM to perform burden association analysis

39

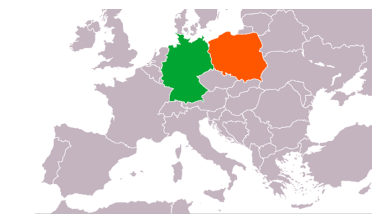
Rare Variant Association Analysis - Confounders

- Control for covariates in the analysis which are potential confounders
 - Age
 - Sex
 - Batch
 - Body Mass Index (BMI)
 - Smoking pack years
 - Population substructure

40

Confounder - Population Substructure and Admixture



41

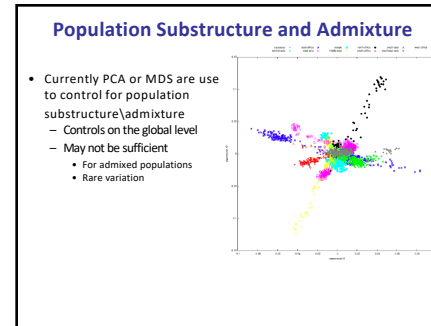
Population Substructure and Admixture

- If proportion of cases and controls sampled from each population is different
 - Can occur due to
 - Disease frequency is different between populations
 - Sloppy sampling
- Population substructure\admixture can cause detection of differences in variant frequencies within a gene which is due to sampling and not disease status
 - False positive findings can be increased

42



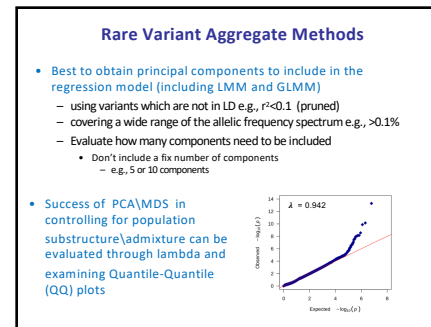
43



44

- ### Rare Variant Aggregate Association Analysis
- When analyzing different populations, e.g.,
 - Africans
 - Europeans
 - When analyzing data from different source
 - Analyze each group separately
 - Meta-analysis can be used to combine the results from each group

45



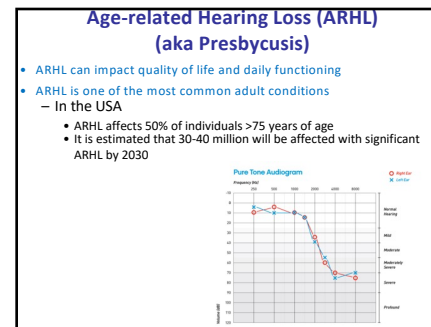
46

Part II

Example of a Rare Variant Association Study

Analysis of UK Biobank Exome Data to Study the Etiology of Late-onset Hearing Loss

47



48

Goals of the Study

- Using data from the UK Biobank to detect associations between self-reported measures of ARHL and genetic variants
 - H-aid** self-reported hearing aid use (f.3393: "Do you use a hearing aid most of the time?")
 - H-diff** self-reported hearing difficulty (f.2247: "Do you have any difficulty with your hearing?")
 - H-noise** self-reported hearing difficulty with background noise (f.2257: "Do you find it difficult to follow a conversation if there is background noise e.g., TV, radio, children playing?")
 - H-both** individuals with both H-diff and H-noise
- With an emphasis of understanding the role that rare variation plays in ARHL
 - Current analysis - exome sequence data

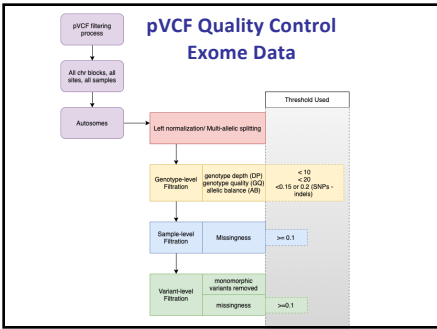
49

UK Biobank

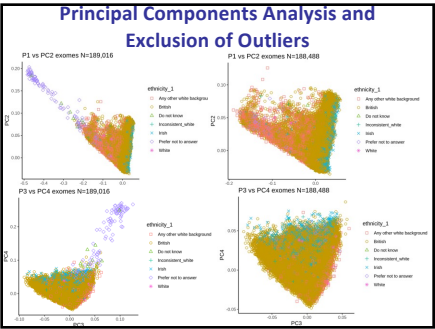
- 500,000 individuals randomly sampled
 - Aged 40-69 at time of enrollment
 - To be followed for at least 20 years
 - Predominantly white Europeans
 - Also includes South Asians and individuals of African Ancestry and smaller number of individuals of a few other ancestries
- Extensive phenotype data
 - Qualitative and quantitative traits
 - ICD-10 and ICD-9 codes
 - Self reports
 - Cognitive test
 - Brain MRIs
 - NMR-metabolomics data
- Genetic Data
 - Genotype and imputed data
 - Exome sequence data
 - Whole genome sequence data
 - Telomere length data

*Data showcase can be used to examine phenotypes and sample sizes available

50



51



52

Exclusion Criteria Obtained from ICD10, ICD9, & Self Report

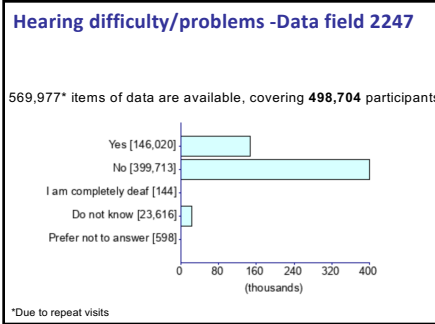
- Deafness
- Early-onset hearing impairment
- Otosclerosis
- Meniere's
- Labyrinthitis
- Disorders of acoustic nerve
- Bell's palsy
- History of chronic suppurative and nonsuppurative otitis media
- Meningitis
- Encephalitis, myelitis, and encephalomyelitis
- Etc.

53

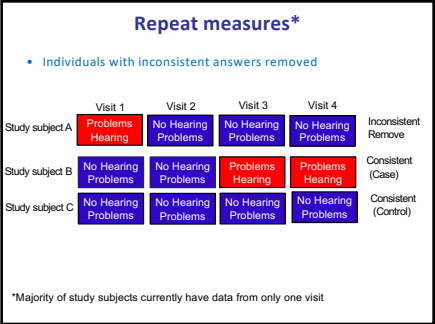
Defining Cases and Controls

- Based on answers obtained from a touch screen
- Cases - self-reported hearing difficulty
 - f.2247: "Do you have any difficulty with your hearing?"
- Controls - did **not** have any self-reported HL or ID10/9 HL codes

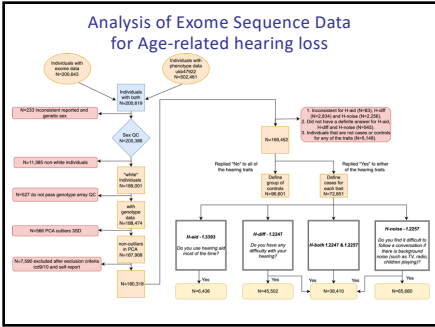
54



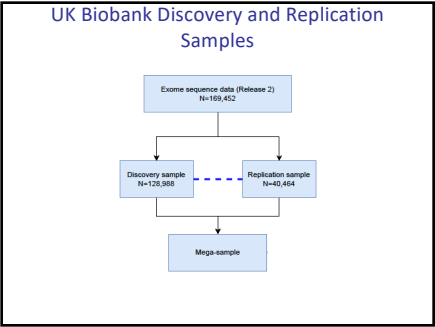
55



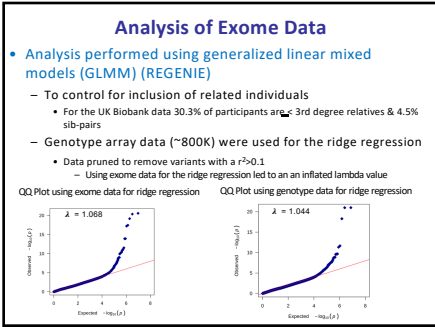
56



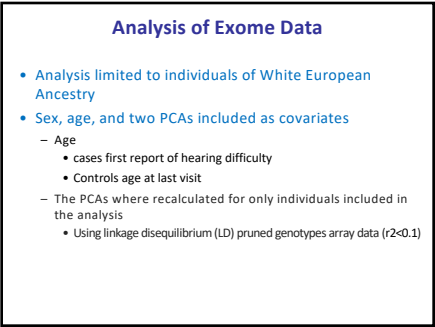
57



58



59



60

Selection of Variants to Include in Rare Variant Aggregate Association Tests

Annotation File	Mask File	AAF file
1:55039839:T:C PCSK9 LoF	Mask1 LoF	1:55039839:T:C 1.53e-05
1:55039842:G:A PCSK9 missense	Mask2 LoF,missense	1:55039842:G:A 2.19e-06
1:55039839:T:C PCSK9 CADD30	Mask1 CADD score > 30	1:55039839:T:C 1.53e-05
1:55039842:G:A PCSK9 CADD20	Mask2 CADD score > 20	1:55039842:G:A 2.19e-06

REGENIE will use information from the annotation and alternative allele frequency (AAF) files to build the Masks (variants to be included in the association testing)

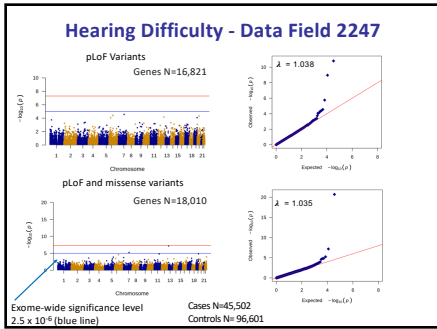
67

Rare Variant Aggregate Analysis

- Exome sample was split
 - Second release of 150K exome were used as the discovery sample.
 - First release of 50K exome were used as the replication sample
- Entire exome sample (200K) was also analyzed*
- Discovery sample significance level
 - $p < 2.5 \times 10^{-6}$
 - 0.05/20,000 Bonferroni correction for testing 20,000 genes
- Replication sample significant level
 - $p < 0.05$
 - Empirical p-values generated
 - Permutation used to adjust for the number of phenotypes and genes brought to replication (pLoF and pLoF & missense)

*No replication sample available for these findings

68



69

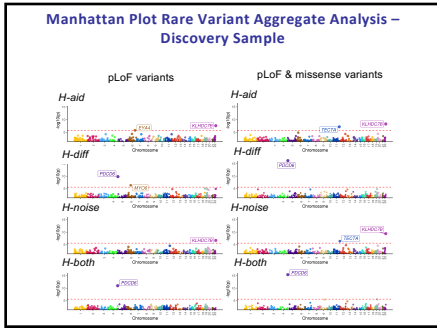
Rare Variant Aggregate Analysis – Discovery and Replication Samples

Discovery Sample Rare-variant aggregate association analysis with age-related hearing traits

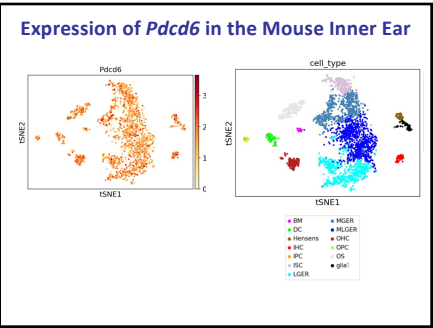
Type of variants	Gene	H-aid		H-noise		H-diff		H-both	
		SE	P	SE	P	SE	P	SE	P
pLoF	SH3BP1	0.20	2.86×10^{-6}	0.46	1.03×10^{-5}	0.40	1.10×10^{-5}	0.21	2.03×10^{-5}
	SLC6A1	0.40	3.14×10^{-6}	0.34	7.08×10^{-6}	0.41	1.01×10^{-5}	0.41	1.01×10^{-5}
pLoF & missense	PCSK9	0.46	1.76×10^{-6}	0.29	1.73×10^{-5}	0.47	1.03×10^{-5}	0.47	1.03×10^{-5}
	PCSK9	0.46	1.76×10^{-6}	0.29	1.73×10^{-5}	0.47	1.03×10^{-5}	0.47	1.03×10^{-5}
pLoF & missense	SH3BP1	0.40	3.14×10^{-6}	0.34	7.08×10^{-6}	0.41	1.01×10^{-5}	0.41	1.01×10^{-5}
	SLC6A1	0.40	3.14×10^{-6}	0.34	7.08×10^{-6}	0.41	1.01×10^{-5}	0.41	1.01×10^{-5}

Genes associated to an exome-wide significance level ($p < 2.5 \times 10^{-6}$) with hearing aid (H-aid), hearing difficulty (H-diff), hearing difficulty with background noise (H-noise), and the combined trait (H-both). Using rare-variant aggregate association tests (pLoF or missense + pLoF variants with a MAJ < 0.1) in geneAD v2.1.1, were analyzed in the discovery and/or replication samples of white European individuals from the UK Biobank. The p-values for replicated associations [empirical p-values < 0.05 adjusting for genes (pLoF and missense & pLoF) and traits brought to replication] are shown in red.

70



71



72

Overview



- 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, *PDCD6* which has not been previously reported
 - Inner ear expression in humans and mice supports the involvement of gene in HL etiology
 - *PDCD6* is a cytoplasmic Ca²⁺ binding protein with an important role in apoptotic cell death
- Rare-variant aggregate analysis demonstrated the important contribution of Mendelian HL genes, i.e. *MYO6*, *TECTA*, and *EYA4* the genetics of ARHL
- 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
- For additional information see
 - Cornejo-Sanchez et al. (2023) Eur J Hum Genet PMID: 36788145

73

Overview/Future Direction

- The entire exome sequence data set of White Europeans has been analyzed
 - Revealing many additional known Mendelian nonsyndromic HL genes
- Mendelian genes (although not necessarily the same variants) play an important role in ARHL
- Performing Mendelian Randomization and testing for pleiotropy (vertical & horizontal) to evaluate associations between ARHL and comorbidities
 - e.g., dementia, depression
- Analysis of UK Biobank and All of Us WGS data including structural variants and performing rare variant aggregate tests outside of the coding regions

74

Newcastle University  

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

1

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

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

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

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

3

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

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

4

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

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

5

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
 - Predicting trait values in a new individual

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

6

Population stratification and relatedness

Genes mirror geography within Europe

J Novembre et al. (2008) *Nature* 456(7218):98-101. doi:10.1038/nature07331

Heather Cordell (Newcastle) GWAS (Part 2) 3 / 38

7

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
 - Fixed effects are considered 'fixed' at their measured values
 - Random effects are considered to be sampled from a distribution

Heather Cordell (Newcastle) GWAS (Part 2) 4 / 38

8

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
 - 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

$$y = mx + c \quad \text{or} \quad y = \beta_0 + \beta_1x$$
- This model may also be written

$$y_i = \beta_0 + \beta_1x_i + E_i$$
 - y_i refers to the trait value of person i
 - x_i refers to the measured value of person i 's predictor variable
 - E_i refers to the displacement from the regression line
 - i.e. the discrepancy between the observed and the predicted y value

Heather Cordell (Newcastle) GWAS (Part 2) 4 / 38

9

Linear Regression

Heather Cordell (Newcastle) GWAS (Part 2) 5 / 38

10

Linear Mixed Models (LMMs)

- In linear regression we have $y_i = \beta_0 + \beta_1x_i + E_i$
 - Here β_0 and β_1 are fixed effects while E_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:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} E_1 \\ E_2 \\ \vdots \\ E_n \end{bmatrix}$$
- or $y = X\beta + \epsilon$

Heather Cordell (Newcastle) GWAS (Part 2) 6 / 38

11

Linear Mixed Models (LMMs)

- In linear regression we have $y_i = \beta_0 + \beta_1x_i + E_i$
 - Here β_0 and β_1 are fixed effects while E_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:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} E_1 \\ E_2 \\ \vdots \\ E_n \end{bmatrix}$$
- or $y = X\beta + \epsilon$
- A LMM takes the form $y = X\beta + Z\alpha + \epsilon$
 - where α corresponds to a vector of random effects
 - with loadings specified in Z

Heather Cordell (Newcastle) GWAS (Part 2) 6 / 38

12

Linear Mixed Models (LMMs)

- E.g. suppose 2 fixed effects β_1 and β_2 , and 3 random effects (plus n random errors)
- Then $y = X\beta + Zu + \epsilon$ corresponds to:

y_1	x_{11}	x_{12}		z_{11}	z_{12}	z_{13}		u_1		ϵ_1
y_2	x_{21}	x_{22}		z_{21}	z_{22}	z_{23}		u_2		ϵ_2
\vdots	\vdots	\vdots		\vdots	\vdots	\vdots		\vdots		\vdots
y_n	x_{n1}	x_{n2}		z_{n1}	z_{n2}	z_{n3}		u_3		ϵ_n
- or $y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + u_1 z_{i1} + u_2 z_{i2} + u_3 z_{i3} + \epsilon_i$

13

LMMs in genetics

- In genetics we generally work with two equivalent forms of LMM
- One is: $y = X\beta + Zu + \epsilon$
 - The random effect u corresponds to a scaled additive effect of **causal variant (locus) l**
 - We assume there are many (m) such causal variants all across the genome
 - Considering it to be a random effect (within a population of interest) could be thought of as taking a Bayesian perspective

14

LMMs in genetics

- In genetics we generally work with two equivalent forms of LMM
- One is: $y = X\beta + Zu + \epsilon$
 - The random effect u corresponds to a scaled additive effect of **causal variant (locus) l**
 - We assume there are many (m) such causal variants all across the genome
 - Considering it to be a random effect (within a population of interest) could be thought of as taking a Bayesian perspective
 - Z is a standardized genotype matrix i.e. z_i takes value

$$\sqrt{\frac{2f_l}{2f_l(1-f_l)}}, \sqrt{\frac{(1-2f_l)}{2f_l(1-f_l)}}, \sqrt{\frac{2(1-f_l)}{2f_l(1-f_l)}}$$
 if individual i has genotype (qq, Qq, QQ)
 - where f_l is the frequency of allele Q at locus l

15

LMMs in genetics

- The other form is: $y = X\beta + g + \epsilon$
 - Where $g_i = \sum_{l=1}^m z_{il} u_l$ 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 individual i
 - The vector of random effects g takes distribution $g \sim N(0, G\sigma_g^2)$
 - Where G is the genetic relationship matrix (GRM) between individuals – i.e. their IBD sharing **at the causal loci**
 - $\sigma_g^2 = m\sigma_u^2$ is the total additive genetic variance
 - $G = ZZ^T/m$

16

LMMs in genetics

- The other form is: $y = X\beta + g + \epsilon$
 - Where $g_i = \sum_{l=1}^m z_{il} u_l$ 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 individual i
 - The vector of random effects g takes distribution $g \sim N(0, G\sigma_g^2)$
 - Where G is the genetic relationship matrix (GRM) between individuals – i.e. their IBD sharing **at the causal loci**
 - $\sigma_g^2 = m\sigma_u^2$ is the total additive genetic variance
 - $G = ZZ^T/m$
- 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)

17

Use of LMMs in genetics

- The formulation $y = X\beta + g + \epsilon$ 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

18

Testing for association using LMMs

- 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_i is the trait value
 - x_i 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)
 - $\gamma_i = g + E$

Heather Cordell (Newcastle) GWAS (Part 2) 11 / 38

19

Testing for association using LMMs

- 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_i is the trait value
 - x_i 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)
 - $\gamma_i = g + E$
 - We assume $\gamma \sim MVN(0, \mathbf{V})$ where variance/covariance matrix \mathbf{V} follows standard variance components model
 - Variance/covariance matrix structured as:
$$V_{ij} = \sigma_g^2 + \sigma_e^2 \delta_{ij}$$

$$V_{ij} = 2\Phi_{ij}\sigma_g^2 \delta_{ij}$$
 - σ_g^2, σ_e^2 represent the additive polygenic variance (due to all loci) and the environmental (=error) variance, respectively

Heather Cordell (Newcastle) GWAS (Part 2) 11 / 38

20

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 QTD 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)

Heather Cordell (Newcastle) GWAS (Part 2) 12 / 38

21

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 $\mathbf{G} = \mathbf{Z}\mathbf{Z}'/m$, the genetic relationship matrix (GRM) between individuals **at the causal loci**
 - Since we don't know the causal loci, we approximate \mathbf{G} by \mathbf{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

Heather Cordell (Newcastle) GWAS (Part 2) 13 / 38

22

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

Heather Cordell (Newcastle) GWAS (Part 2) 14 / 38

23

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
 - DISSECT

Heather Cordell (Newcastle) GWAS (Part 2) 14 / 38

24

Software implementations

- Main difference between them is the precise computational tricks used to speed up the calculations
 - And the convenience/ease of use
 - See comparison in Eu-Ahsunthornwattana et al. (2014) PLoS Genetics 10(7):e1004445

Heather Cordell (Newcastle) GWAS (Part 2) 15 / 38

25

Software implementations

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

$$y = X\beta + Zu + E$$
 - 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**, which uses a slightly different formulation based on analysing the residuals following a whole-genome blockwise ridge regression
 - Again based on LMM formulation 1: $y = X\beta + Zu + E$
 - See also **LDAK-KVIK**

Heather Cordell (Newcastle) GWAS (Part 2) 15 / 38

26

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) GWAS (Part 2) 16 / 38

27

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
- 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)

Heather Cordell (Newcastle) GWAS (Part 2) 16 / 38

28

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
- See also **LDAK-KVIK**

Heather Cordell (Newcastle) GWAS (Part 2) 17 / 38

29

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 (= 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) GWAS (Part 2) 18 / 38

30

Elucidating genetic architecture

- Basic idea is to use formulation

$$y = X\beta + g + \epsilon$$
 with $g \sim N(0, A\sigma_g^2)$ and $\epsilon \sim N(0, I\sigma_e^2)$ so $V = A\sigma_g^2 + I\sigma_e^2$
 - A is the GRM between individuals, estimated using all genotyped SNPs
 - σ_g^2 and σ_e^2 estimated using REML (or MLE)
 - Thus we can estimate heritability accounted for by the genotyped SNPs as $\sigma_g^2 / (\sigma_g^2 + \sigma_e^2)$
- 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

Heather Cordell (Newcastle) GWAS (Part 2) 19 / 38

31

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

$$y = X\beta + \sum_{c=1}^{22} g_c + \epsilon \quad \text{with } V = \sum_{c=1}^{22} A_c\sigma_c^2 + I\sigma_e^2$$
 where g_c is a vector of effects attributed to the c th chromosome, and A_c is the GRM estimated from SNPs on the c th chromosome
 - 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

Heather Cordell (Newcastle) GWAS (Part 2) 20 / 38

32

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'
 - See LDSC software (<https://github.com/bulik/ldsc>)

Heather Cordell (Newcastle) GWAS (Part 2) 21 / 38

33

Short break

Heather Cordell (Newcastle) GWAS (Part 2) 22 / 38

34

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) GWAS (Part 2) 23 / 38

35

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
 - 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

Heather Cordell (Newcastle) GWAS (Part 2) 23 / 38

36

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
- 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
- Phenomenon of biological interest?
 - Identifying genes that interact to cause disease could help us understand the mechanisms and pathways in disease progression

Heather Cordell (Newcastle) GWAS (Part 2) 23 / 38

37

Definition of (pairwise) interaction

- Statistical interaction most easily described in terms of a (logistic) regression framework
 - Suppose x_1 and x_2 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:

$$\log \frac{p}{1-p} = \beta_0 + \beta_1 x_1$$
 Marginal effect of factor 1

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

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

$$\log \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$
 Main effects and interaction term
 - For quantitative traits, use linear regression (replace $\log \frac{p}{1-p}$ with y)
 - For modelling as an LMM, add in a random effect γ

Heather Cordell (Newcastle) GWAS (Part 2) 24 / 38

38

Interaction

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

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

Heather Cordell (Newcastle) GWAS (Part 2) 25 / 38

39

Interaction

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

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

Heather Cordell (Newcastle) GWAS (Part 2) 25 / 38

40

Interaction

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

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

Heather Cordell (Newcastle) GWAS (Part 2) 25 / 38

41

Interaction

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

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

Heather Cordell (Newcastle) GWAS (Part 2) 25 / 38

42

Interaction

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

		Factor 2	
Factor 1	1	$\beta_0 + \beta_1 + \beta_2 + \beta_{12}$	$\beta_0 + \beta_1$
	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
 - 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 ($\beta_1 = \beta_2 = 0$)

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

 - Trait value only differs from baseline if both factors present

Heather Cordell (Newcastle) GWAS (Part 2) 25 / 38

43

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, β_G^2)
 - Main effects of locus H (β_H, β_H^2)
 - 4 interaction terms

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

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

Heather Cordell (Newcastle) GWAS (Part 2) 26 / 38

44

Gene-gene interaction

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

$$\log \frac{p}{1-p} = \beta_0 + \beta_G x_1 + \beta_H x_2 + \beta_{GH} x_1 x_2$$
 with x_1, x_2 coded (0,1,2)

		Locus H		
Locus G	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

Heather Cordell (Newcastle) GWAS (Part 2) 27 / 38

45

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

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

46

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
 - 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) 29 / 38

47

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
 - 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

Heather Cordell (Newcastle) GWAS (Part 2) 29 / 38

48

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
 - Siemiatycki and Thomas (1981) Int J Epidemiol 10:383-387; Thompson (1991) J Clin Epidemiol 44:221-232
 - Phillips (1998) Genetics 148:1167-1171; 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

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

49

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
 - Siemiatycki and Thomas (1981) Int J Epidemiol 10:383-387; Thompson (1991) J Clin Epidemiol 44:221-232
 - Phillips (1998) Genetics 148:1167-1171; 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

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

50

Gene-environment (G×E) interactions

- The same regression model

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

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)

Heather Cordell (Newcastle) GWAS (Part 2) 29 / 38

51

Gene-environment (G×E) interactions

- The same regression model

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

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)
- 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 $G \times G$, 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

Heather Cordell (Newcastle) GWAS (Part 2) 29 / 38

52

Testing association and/or interaction

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

$$\log \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$
 - 3df test of $\beta_1 = \beta_2 = \beta_3 = 0$ tests for association **at both loci** (or both variables), allowing for their possible interaction

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

53

Testing association and/or interaction

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

$$\log \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$
 - 3df test of $\beta_1 = \beta_2 = \beta_3 = 0$ tests for association **at both loci** (or both variables), allowing for their possible interaction
 - 2df test of $\beta_2 = \beta_3 = 0$ tests for association at locus 2, **while allowing for possible interaction with locus (or variable) 1**

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

54

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
 - 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

55

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
 - 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

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

56

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
 - 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/>

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

57

G×G versus G×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
 - Less obvious reason to do this for G×E...

Heather Cordell (Newcastle) GWAS (Part 2) 31 / 38

58

G×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)
 - Fr'aaenberg 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) GWAS (Part 2) 32 / 38

59

G×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)
 - Fr'aaenberg 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

Heather Cordell (Newcastle) GWAS (Part 2) 32 / 38

60

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

Heather Cordell (Newcastle) GWAS (Part 2) 32 / 38

61

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 regression 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)

Heather Cordell (Newcastle) GWAS (Part 2) 33 / 38

62

Testing correlation between loci

- A similar idea is implemented in EPIBLASTER (Kam-Thong et al. 2011; EUHG 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 (*A*) encapsulates a slightly different quantity between the different methods
- All implemented (along with standard logistic and linear regression) in CASSI
 - <http://www.staff.ncl.ac.uk/richard.howe/cassi/>

Heather Cordell (Newcastle) GWAS (Part 2) 34 / 38

63

Empirical evidence for G×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 gastroesophageal reflux, obesity and tobacco smoking, with respect to risk for Barrett's esophagus

Heather Cordell (Newcastle) GWAS (Part 2) 35 / 38

64

Empirical evidence for G×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
 - Wood et al. (2014) Nature 514(7520):E3-5. PMID:25279928
 - Fish et al. (2016) Am J Hum Genet 99(4):817-830. PMID:27640306
- The Hemani et al. paper was **subsequently retracted** (<https://www.nature.com/articles/s41586-021-03766-v>)

Heather Cordell (Newcastle) GWAS (Part 2) 36 / 38

65

Empirical evidence for G×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."

Heather Cordell (Newcastle) GWAS (Part 2) 37 / 38

66

Empirical evidence for G×E interactions

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

ARTICLE

The figure consists of two panels, (a) and (b). Panel (a) is a box plot showing gene expression levels for three genes: AK, AG, and GS. The y-axis is labeled 'Expression level' and ranges from -4 to 4. For each gene, there are two box plots: a blue one for 'Low-NE exposure' and a red one for 'High-NE exposure'. Panel (b) is a network diagram with nodes representing genes. The nodes are labeled with gene IDs: SARC2, SARC22, SARC21, SARC23, SARC24, SARC25, SARC26, SARC27, SARC28, SARC29, SARC30, SARC31, SARC32, SARC33, SARC34, SARC35, SARC36, SARC37, SARC38, SARC39, SARC40, SARC41, SARC42, SARC43, SARC44, SARC45, SARC46, SARC47, SARC48, SARC49, SARC50, SARC51, SARC52, SARC53, SARC54, SARC55, SARC56, SARC57, SARC58, SARC59, SARC60, SARC61, SARC62, SARC63, SARC64, SARC65, SARC66, SARC67, SARC68, SARC69, SARC70, SARC71, SARC72, SARC73, SARC74, SARC75, SARC76, SARC77, SARC78, SARC79, SARC80, SARC81, SARC82, SARC83, SARC84, SARC85, SARC86, SARC87, SARC88, SARC89, SARC90, SARC91, SARC92, SARC93, SARC94, SARC95, SARC96, SARC97, SARC98, SARC99, SARC100. The nodes are interconnected by lines representing interactions. One node, SARC21, is highlighted in yellow.

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

67

Power Analysis for Single and Rare Variant Aggregate Association Analyses

Suzanne M. Leal, Ph.D.
Sergievsky Family Professor of Neurological Sciences
Director of the Center for Statistical Genetics
Columbia University
smi3@columbia.edu

© 2024 Suzanne M. Leal

1

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 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 from different populations
 - Processing of genetic material may 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 5×10^{-8} for GWAS studies

2

Power and Sample Size Estimation for Case-Control Data

- The correct α must be used 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 $\alpha=0.05$ is used in order to control the Family Wise Error Rate

3

Power and Sample Size Estimation for Case-Control Data

- GWAS of common variants where each variant is test separately
 - $\alpha=5 \times 10^{-8}$ (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 is higher than for analysis limited to common variants
 - α is yet to be determined for association analysis of whole genome sequence data

4

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% family-wise error
- Then extrapolating to infinite density
- The genome wide significance threshold estimate $\sim 7.2 \times 10^{-8}$
- Estimate is based on LD structure for Europeans
 - Not sufficiently stringent for populations of African Ancestry

5

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 $\alpha=2.5 \times 10^{-6}$
 - 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 is effective number tests
 - This would not be the case for variants in LD

6

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 $\alpha = 5.0 \times 10^{-3}$ can be used for a family wise error rate of 0.05

7

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 ($y=1$)
 - Dominant ($y=y_1$)
 - Additive ($y=2y_1-1$)
 - Multiplicative ($y=y_1^2$)

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

9

Armitage Trend Test

- Power and Sample size
 - Calculated under different models
 - Where y is the relative risk
 - Multiplicative
 - $y = y_1^2$
 - Additive
 - $y = 2y_1 - 1$
 - Dominant
 - $y = y_1$
 - Recessive
 - $y = 1$

10

Gamma is the Relative Risk not the Odds Ratio

- Most software for power calculations/sample size estimation use the relative risk (γ) and not the odds ratio
- The relative risk only approximates the odds ratio when disease is rare (Prevalence $\sim < 0.1\%$)
 - The relative risk is not appropriate for common traits when a case-control design is used

11

Correspondence Between the Odds Ratio and Relative Risk

Dominant Model

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 allele and disease allele frequency 0.01
 D' and $r^2=1$
 *1/2 genotype - heterozygous (one copy of the alternative allele)
 **2/2 genotype - homozygous for the alternative allele

12

Genetic Association Study (GAS) Power Calculator

- http://csg.snh.umich.edu/abecasis/cats/gas_power_calculator/index.html
- A one-stage study power calculator
 - Which was derived from CaTs
 - Which is to perform two-stage genome wide association studies
 - Skol et al. 2006
- Cochran Armitage Trend Test
- Displays graphs of the results

13

GAS Power Calculator

14

Genetic Power Calculator

- <http://zzz.bwh.harvard.edu/gnc/>
- S Purcell & P Sham
- Uses the methods described in Sham PC et al. (2000) Am J Hum Genet 66:1616-1630
 - VC QTL linkage for sibships
 - VC QTL association for sibships
 - VC QTL linkage for sibships conditional on the trait
 - TDT for discrete traits
 - Case-Control for discrete traits
 - TDT for quantitative traits
 - Case-Control quantitative traits
- Although input is the relative risk
 - Displays odds ratios

15

Genetic Power Calculator

Case - control for discrete traits

High risk allele frequency (A)	: 0.01 (0 - 1)
Prevalence	: 0.2 (0.0001 - 0.9999)
Genotype relative risk Aa	: 1.5 (> 1)
Genotype relative risk AA	: 1.5 (> 1)
D-prime	: 1 (0 - 1)
Marker allele frequency (B)	: 0.01 (0 - 1)
Number of cases	: 10000 (0 - 1000000)
Controls : case ratio	: 1 (> 0)

Unselected controls? (* see below)

User-defined type I error rate : 0.0000000 (0.00000001 - 0.5)
 User-defined power: detection B : 0.80 (0 - 1)
 (1 = type II error rate)

Process Reset

Created by *Shawn Purcell* 24 Oct 2008

16

17

Power Association With Errors (PAWE)

- <http://compeng.rutgers.edu/pawe/>
- Implements the linear trend test
- Four different error models can be used
 - See online documentation for complete explanation
- Can either perform:
 - Power calculations for a fixed sample size
 - Sample size calculations for a fixed power
- The genotype frequencies can be generated either using a:
 - Genetic model free method or
 - Genetic model-based method

18

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

19

Linkage Disequilibrium (LD)

- Power will be reduced if causal variant is not in perfect LD ($r^2=1$) with the tag SNP
- Can adjust sample size when $r^2 < 1$ to increase power to the same level as when $r^2=1$
- Can estimate sample size when $r^2 \neq 1$
 - $N/r=r^2N$
 - Valid only for multiplicative model
 - (Pritchard and Przeworski, 2001)
- Power calculation almost always assume that $r^2=1$
- For whole genome sequence data this should be the case since usually the causal variant would be included in the data

20

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

21

Simplistic Analytical Power Calculation for Rare-variant 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., 2.5×10^{-6}
- Provide disease prevalence and penetrance model
- Perform calculations in the same manner as was described for single variants

22

Empirical Power Calculations

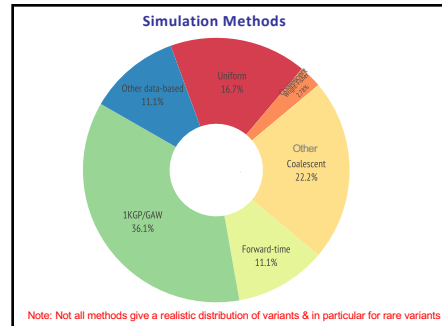
- A variety of methods can be used to generate variant data to empirically estimate power
- Variant data is generated
 - Based upon a penetrance model samples of cases and controls are generated
 - Or a quantitative trait is generated based upon the genetic variance
- Multiple replicates are generated and analyzed
 - To determine the power

23

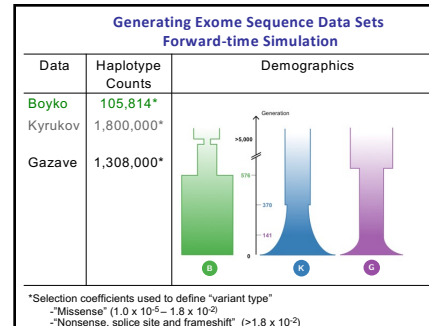
Empirical Power Calculations

- 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., 5×10^{-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., 2.5×10^{-6}

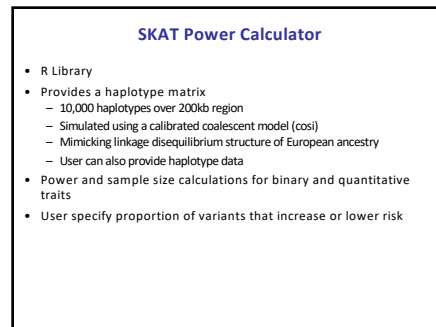
24



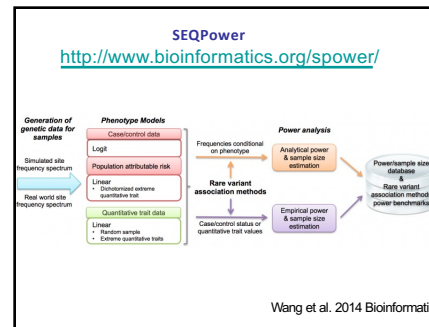
25



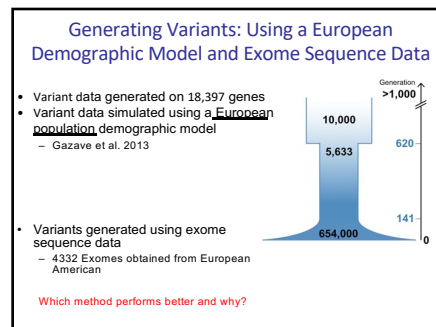
26



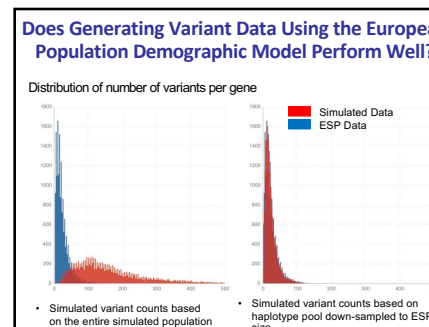
27



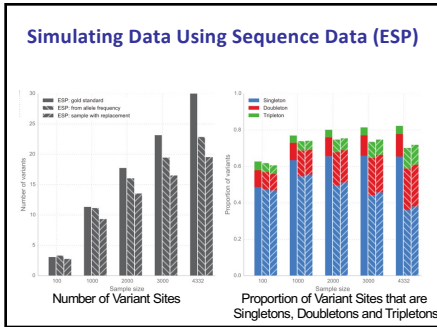
28



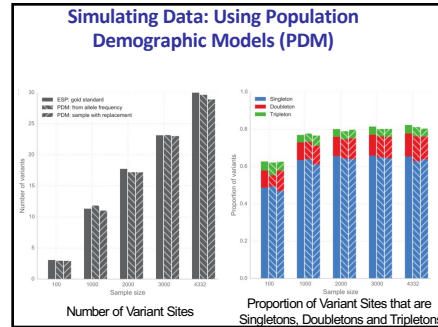
29



30



31



32

- ### Simulation Studies to Evaluate Power for Rare Variant Association Studies
- 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., $\alpha=2.5 \times 10^{-6}$)
 - If addition regions besides genes are analyzed
 - A more stringent α value should be used

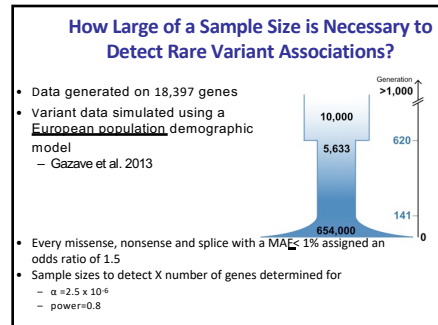
33

- ### 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

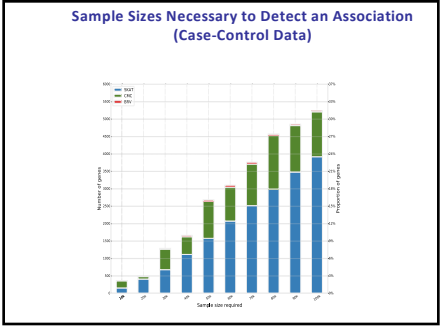
34

- ### Power Analysis Rare Variants (Aggregate Association Tests)
- 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

35



36



37

Imputing and Analyzing Imputed Genotype data

Suzanne M. Leal, Ph.D.

© 2024 Suzanne Leal

1

Motivation for Imputation of Genotype Data

- Obtain genotypes for variant sites that are not genotyped
 - Additional variants can be tested for associations
 - Providing additional power to tag causal variant sites
 - Potential inclusion of causal variants that are unavailable on genotyping arrays
 - Aids in fine-mapping
- Considerably less expensive than generating whole genome sequence data
 - Does come at a cost of accuracy
 - In particular for very rare variants
 - Imputed data will be available for very rare variants if
 - For a variant site the alternative allele has been observed ~8X in the reference panel in order for it to be imputed

2

Imputation of Variants

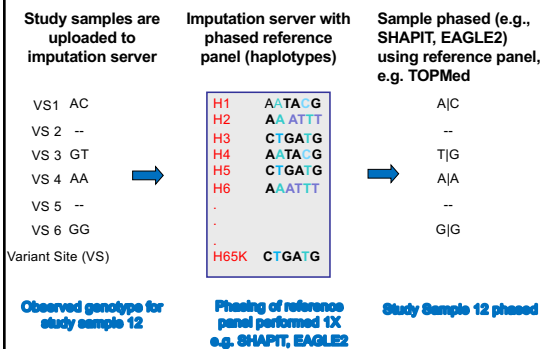
- Can be performed locally or on an imputation server
- Imputation locally has its limitation due to availability of a reference panel
 - Internal data
 - 1000 genomes
 - Haplotype reference consortium (HRC)
 - Only part of this dataset is made publicly available
- Smaller imputation panels will impact the ability to impute lower frequency and rare variants
 - Additionally, regardless of variant MAF a decrease in the size and diversity of imputation panel will lead to a decrease in the imputation accuracy

3

Phasing and performing imputation using an Imputation Server

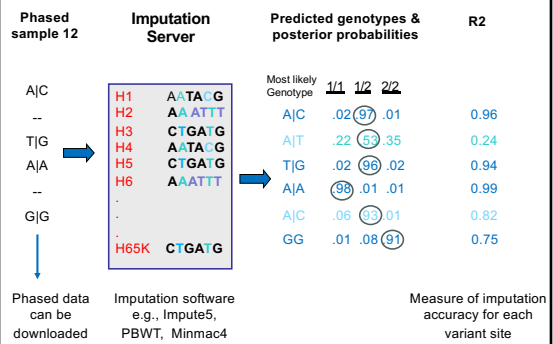
4

Imputation Step 1 Phasing



5

Step 2 Imputation



6

Measures of Imputation Accuracy

- R^2 /INFO
 - Measures of imputation accuracy
 - Most programs report R^2
 - Impute provides INFO scores
- r^2 is the correlation between the dosage and genotype obtained from sequence or genotype array data
 - Must have imputed data and sequence or genotype array data for the same person to estimate r^2 .

7

Step 3 Analysis of Imputed Data

- Variants are filtered according to R^2 values
 - e.g., analyze variants with an $R^2 > 0.8$
- Most likely genotypes are not analyzed instead dosages are analyzed
- The dosage can be estimated as follows for variant site 1 sample 12: A|C with prior probabilities $1/1 = 0.02$, $1/2 = 0.97$, & $2/2 = 0.01$ ($R^2 = 0.96$)

Genotype 1/1	$0 \times 0.02 = 0.0$
Genotype 1/2	$1 \times 0.97 = 0.97$
<u>Genotype 2/2</u>	<u>$2 \times 0.01 = 0.02$</u>
Dosage	0.99
- The dosage for variant site 2 sample 12: A|T with prior probabilities $1/1 = 0.22$, $1/2 = 0.53$, & $2/2 = 0.35$ ($R^2 = 0.23$)

Genotype 1/1	$0 \times 0.22 = 0.0$
Genotype 1/2	$1 \times 0.53 = 0.53$
<u>Genotype 2/2</u>	<u>$2 \times 0.35 = 0.70$</u>
Dosage	1.23

8

Imputation Panels

- 1000 Genomes Phase 3*
 - 2,504 reference samples
 - 26 populations from Africa, the Americas, Europe, East Asia, & South Asia
- African Genome Resource
- Asthma among African-ancestry Populations in the Americas (CAAPA)
- Genome Asia Pilot (GAsP)
- HAPMAP2
- Haplotype Reference Consortium (HRC) *
 - 32,470 reference samples (39,635,008 variants)
 - Predominately European Ancestry

*Commonly used imputation panels

9

Imputation Reference Panels

- Multi-ethnic HLA
- Southeast Asian Reference Database (SEAD)
- The Trans-Omics for Precision Medicine (TOPMed)*
 - Version R3 133,597 reference samples (445,600,184 variants)
 - ~40% European, ~29% African/African American, ~19% Hispanic/Latino, ~8% Asian, & ~4% other/unknown)
- UK10K
- Westlake Biobank for Chinese (WBBC)

*Commonly used imputation panels

10

Imputation Servers

- Michigan (US)
 - Reference panels include, HRC, 1,000 Genomes, etc.
 - Phasing EAGLE2
 - Imputation Minmax4
 - <https://imputationserver.sph.umich.edu/index.html#!>
- NHLBI (US)
 - Reference panel TOPMed
 - Phasing EAGLE2
 - Imputation Minmax4
 - <https://imputation.biocatalyst.nih.gov/#/>

11

Imputation Servers

- Sanger (UK)
 - Reference panels include HRC, 1,000 Genomes, etc.
 - Phasing SHAPEIT or EAGLE2
 - Imputation PBWT
 - <https://www.sanger.ac.uk/tool/sanger-imputation-service/>
- Westlake (People's Republic of China)
 - Reference panels include 1000 Genomes, GAsP, SEAD, & WBBC
 - Phasing SHAPEIT2
 - Imputation Minmax4
 - <https://imputationserver.westlake.edu.cn/index.html>

12

What Impacts Imputation Quality?

- Reference (imputation) panel
 - Sample size
 - Larger samples
 - Increase imputation accuracy
 - Ability to impute rare variants
 - Ancestry diversity
- Target sample
 - Density of markers
 - Genotype quality
 - Ancestry and representation on the imputation panel
 - The population's linkage disequilibrium structure

Note: Since each target sample is phased and imputed separately using the pre-phased imputation panel on the imputation server, sample size of the target sample does not impact imputation accuracy

13

How Well do 1000 Genomes, HRC, and TOPMed Imputation Reference Panels Perform?

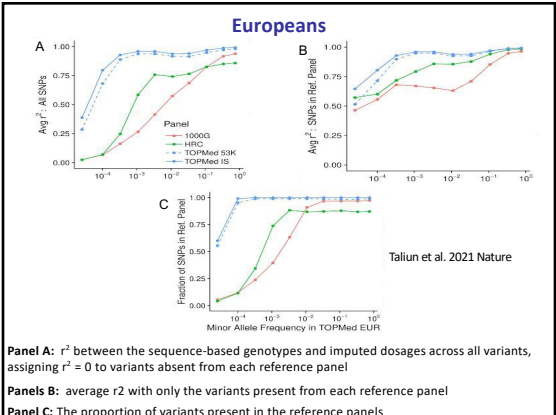
- Reference Panels
 - 1000 Genomes Phase 3
 - 2,504 reference samples
 - 26 populations from Africa, the Americas, Europe, East Asia, & South Asia
 - HRC v1.1 2016
 - 32,470 reference samples (39,635,008 variants)
 - Predominately European Ancestry
 - TOPMed (Version r2)
 - 97,256 reference samples (308,107,085 variants)
 - Diverse population from the USA 48.49% European, 25.95% African/African American, 17.57% Hispanic/Latino/Admixed Americans, 1.22% East Asian, 0.66 South Asians, 6.11% other/unknown)
 - TOPMed (53K)
 - 53,831 reference samples

14

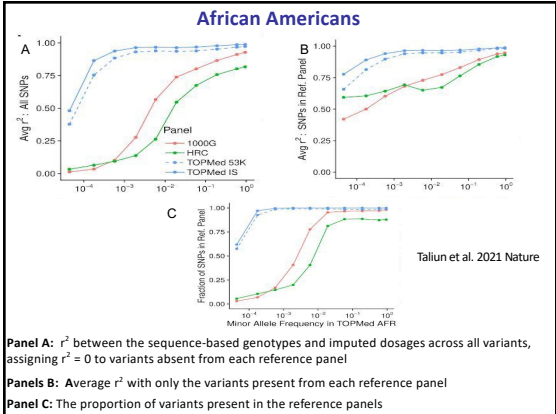
How Well do 1000 Genomes, HRC, and TOPMed Imputation Reference Panels Perform?

- Target Sample
 - 100 ancestry specific samples,
 - e.g. Europeans, African-Americans, & South Asians
 - Obtained from BioMe
 - Samples are not included in any of the reference panels

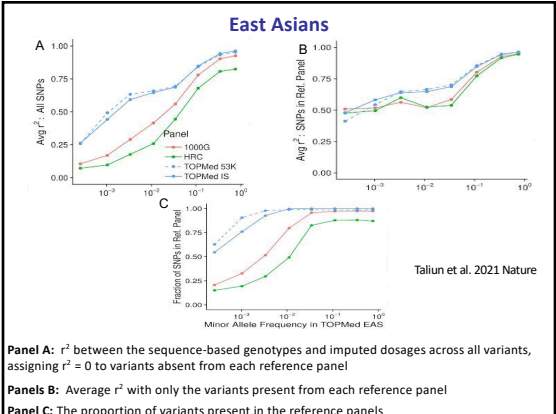
15



16



17



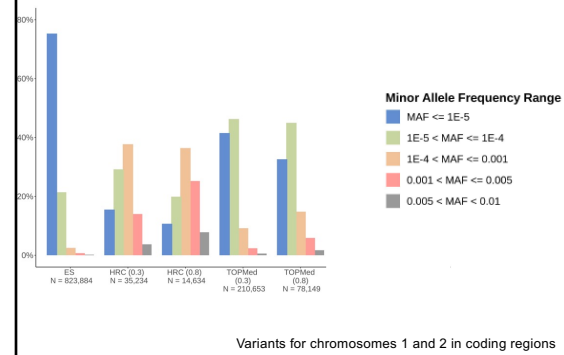
18

Comparison of Rare Variant Distributions

- Unrelated white European UK Biobank study participants (N=168,206) with
 - Release 2 exome sequence
 - Genotype array data available
- Imputed variants using both HRC and TOPMed (v2)
- Comparison of variant distributions
 - Exome sequence (ES) data
 - HRC imputed data $r^2 > 0.3$ and $r^2 > 0.8$
 - TOPMed imputed data $r^2 > 0.3$ and $r^2 > 0.8$

19

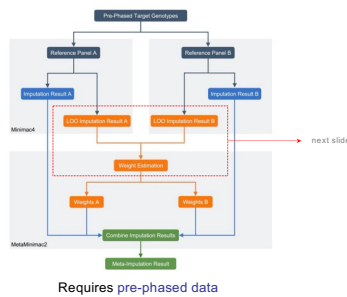
Distribution of Rare Variants



20

Meta-imputation (I)

Use, in turn, two or more reference panels*, then combine the results



*The reference panels must use the same genome build

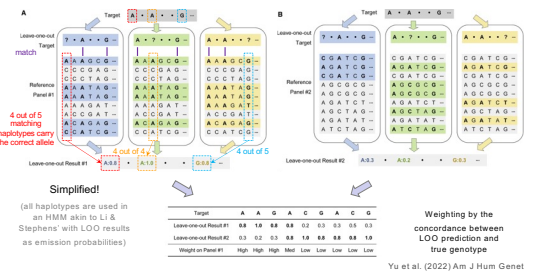
Yu et al. (2022) Am J Hum Genet

21

Meta-imputation (II)

Obtain region-specific weights via leave-one-out (LOO) in an HMM

Imputation using 1st reference panel Imputation using 2nd reference panel



Yu et al. (2022) Am J Hum Genet

22

Imputation of Variants without using an Imputation Server

- Imputation locally has its limitations due to availability of reference panels
 - Internal data
 - 1000 genomes
 - HRC
 - Only part of this dataset is made publicly available to download to use locally
- Can be computationally intensive to phase and impute genotypes locally
- All haplotype phasing and imputation software used on imputation servers are publicly available
- Due to data sharing limitations in particular within the European Union
 - It may not be possible to use imputation servers which are located in the US, UK or China

23

Using Imputation to Detect Genotyping Errors

- Can provide information on genotyping error by comparing the genotype of the imputed variant with genotypes obtained from array or sequence data
 - Would suggest there is genotype error if for the imputed data the R^2 (measure of imputation accuracy) is high
 - But the r^2 (correlation) between the imputed variant and the genotypes obtained from sequence or array data is low.
 - Association analysis results obtained for the imputed variant and the same variant obtained from genotyping array or sequencing vary greatly even though the R^2 value is high for the imputed variant
 - Suggest that there is probably genotyping error for the variant obtained from genotyping array or sequence data
- The variant obtained from array or sequence data can be replaced with the imputed variant

24

Combining data obtained from different genotyping arrays

- Variants that don't overlap between arrays can be imputed
 - As well as variants not available on any of the arrays
- Caution should be used because the imputation quality can vary between datasets
 - Influenced by different error rates between datasets
 - Principal components analysis (PCA) can be used to determine if the potential problems
 - If additional quality control is necessary
- If there are more cases or controls for a particular dataset
 - Type I errors can be increased

25

Linkage disequilibrium in genetic association studies

Gao Wang, Ph.D.
Advanced Gene Mapping Course, May 2024
The Gertrude H. Sergievsky Center and Department of Neurology
Columbia University Vagelos College of Physicians and Surgeons

1

1

Genetic association studies (recap)

Identify genetic variants **associated** with **complex traits**

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

in order to

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

2

2

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

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

4

4

Linkage disequilibrium (LD)

LD: the sharing of certain combinations of variants

- Formally, equivalent to *Haplotype structure*
- There are several measures of LD but largely irrelevant to our learning objectives
- In gene-mapping, let's simply understand LD as Pearson's correlation between variants

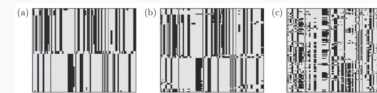
5

5

Linkage disequilibrium (LD)

Levels of LD is a result of chromosomal "shuffling"

- Segregation and Recombination



Each row is a variant site

- Shuffle within rows does not change *marginal MAF*.
- Multi-loci MAF, i.e., *haplotype frequency*, will change.

6

6

Why do we care about LD?

When obviously LD is an issue

- Many variants will look "similar" by genotype but have different biological function — mapping "causal" variants is challenging

When LD is useful

- Can leverage non-causal genetic variables to **predict phenotypes** when causal variant is not observed in data
- Can leverage variants that are LD to **infer each other's genotype** to complete missing genotype data
 - also, association study summary statistics

7

Impact of LD on GWAS analysis

Oligogenic: trait influenced by a few genetic variants

- Misleading: difficult to identify causal variants
- Useful: tag SNPs in array based GWAS design

8

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

9

A second thought on genomic inflation

Population stratification? Or, polygenic inheritance + LD?

Suggested reading: Yang et al (2011) EHG

10

LD score regression (LDSC)

LD score regression model without population stratification

$$E[\chi^2_j] = 1 + \frac{N h^2}{M} I_j$$

N : Sample size
 h^2 : Narrow sense heritability
 M : Total number of variants
 I_j : LD score of variant
 $I_j = \sum_{k \neq j} r_{jk}^2$: LD score: sum of squared Pearson's correlation coefficient between SNP j and other (neighboring) SNPs

11

LD score regression (LDSC)

Separating h^2_g and population stratification

$$E[\chi^2_j] = N\alpha + 1 + \frac{N h^2_g}{M} I_j$$

$N\alpha + 1$: Regression intercept
 $\frac{N h^2_g}{M} I_j$: LD score of variant (Regression slope)

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

12

LDSC application: heritability estimation

Narrow sense heritability

- 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

13

13

Variance of height explained in GWAS

Yengo et al. (2022) Nature

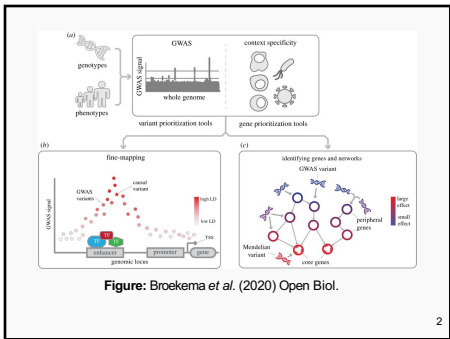
14

14

Statistical fine-mapping in genetic association studies

Gao Wang, Ph.D.
 Advanced Gene Mapping Course, May 2024
 The Gertrude H. Sergievsky Center and Department of Neurology
 Columbia University Vagelos College of Physicians and Surgeons

1



1

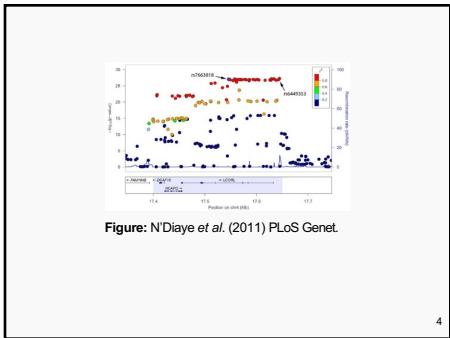
2

Correlated variables in association studies

Due to a phenomenon called **linkage disequilibrium (LD)**

$\text{cor}(x_1, x_2) = 0.9$

3



3

4

Objectives

Statistical fine-mapping **aids in** the identification of causal variants, in order to

- interpret association signals (pinpoint to genes)
- understand biological function of a variant
- elucidate genetic architecture of complex and molecular phenotypes

5

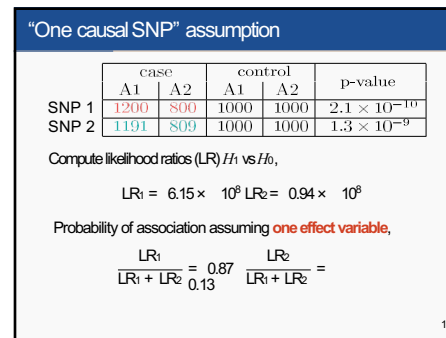
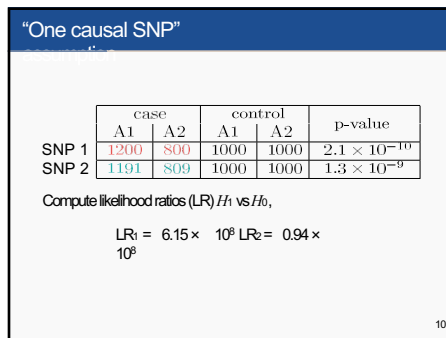
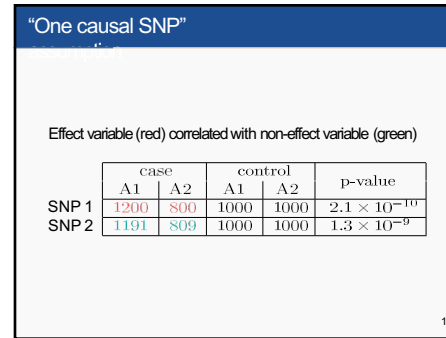
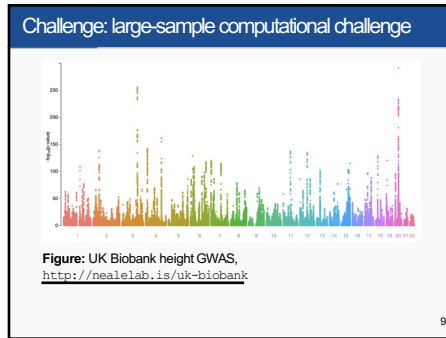
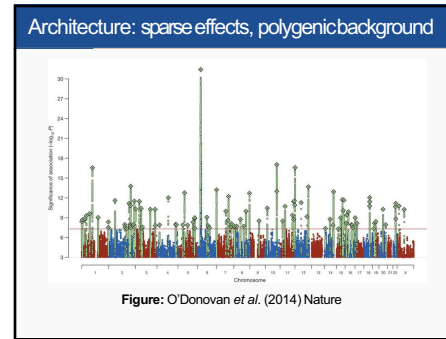
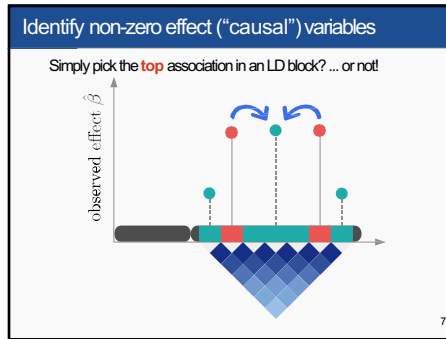
Identify non-zero effect ("causal")

Simply pick the **top** association in an LD block? Maybe?

7

5

6



11

12

Per variable contingency table analysis, R code

```
# returns likelihood ratio of H1 vs H0
get_2x2_lr = function(tbl) {
  tbl = as.table(matrix(tbl, 2, 2,
    dimnames=list(status=c('case', 'control'),
      genotype=c('minor_allele', 'major_allele'))))
  test = MASS::loglm(~status+genotype, data=tbl)
  return(exp(test$lr1 / 2))
}
lr1 = get_2x2_lr(c(1200, 800, 1000, 1000))
lr2 = get_2x2_lr(c(1191, 809, 1000, 1000))
```

13

A "single effect" Bayesian variable

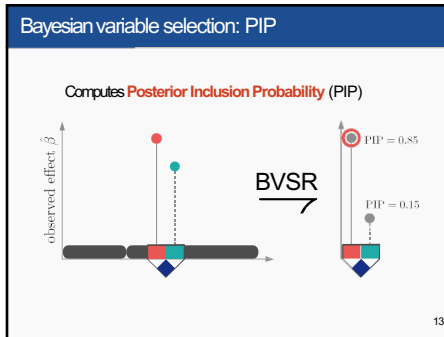
Use Bayes Factor, and compute **posterior inclusion probability**

case		control		p-value
A1	A2	A1	A2	
1200	800	1000	1000	2.1×10^{-10}
1191	809	1000	1000	1.3×10^{-9}

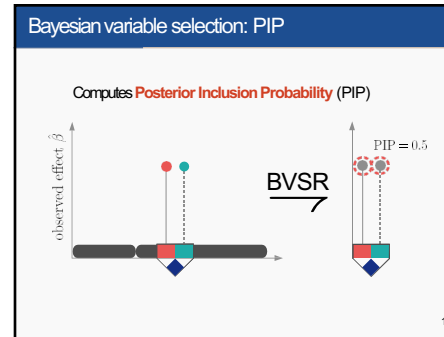
$$PIP_1 = \frac{BF_1}{BF_1 + BF_2} = 0.85$$

$$PIP_2 = \frac{BF_2}{BF_1 + BF_2} = 0.15$$

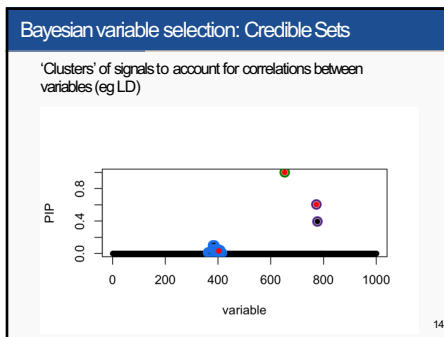
14



15



16



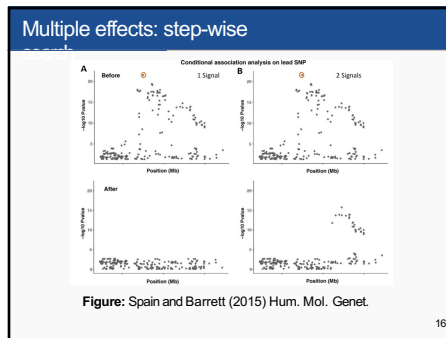
17

Bayesian variable selection: Credible Sets

- **95% credible set S**: $\Pr(\text{effect variable in } S) \geq 95\%$
- e.g., "Single effect" model:

$$\sum_{j \in S} PIP_{(j)} \geq 95\%$$
 where $PIP_{(j)}$'s are in descending order.
- Formal definition: Wang et al. (2020) J. R. Stat. Soc. B

18



19

A simple frequentist conditional analysis

Forward selection algorithm

1. For each SNP fit a simple linear regression model
2. Select the SNP_j that has the largest model likelihood
3. Form residuals $\hat{y} := y - X_j \hat{b}_j$, and repeat

20

A simple frequentist conditional analysis

Forward selection algorithm

1. For each SNP fit a simple linear regression model
2. Select the SNP_j that has the largest model likelihood
3. Form residuals $\hat{y} := y - X_j \hat{b}_j$, and repeat

A greedy algorithm to choose the "best" SNPs, but is incapable of capturing multiple SNPs in LD

21

To quantify uncertainty

Bayesian forward selection algorithm

1. For each SNP_j fit a simple Bayesian linear regression model to get Bayes Factor BF_j
2. Form weight for each SNP, $w_j \propto BF_j$
3. Form residuals $\hat{y} := y - \sum_j w_j X_j \hat{b}_j$, and repeat

22

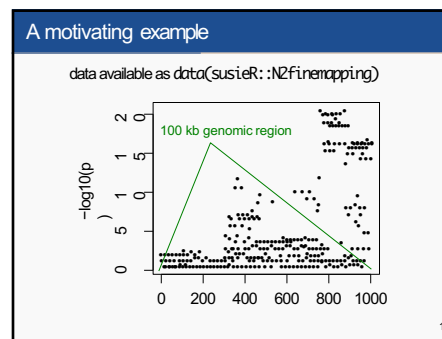
To quantify uncertainty

Bayesian forward selection algorithm

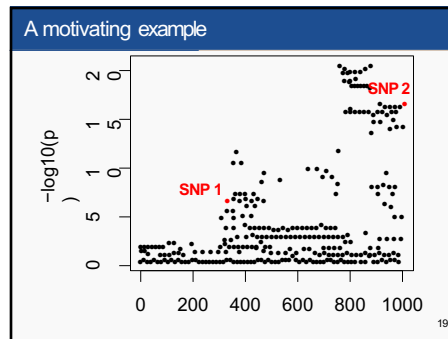
1. For each SNP_j fit a simple Bayesian linear regression model to get Bayes Factor BF_j
2. Form weight for each SNP, $w_j \propto BF_j$
3. Form residuals $\hat{y} := y - \sum_j w_j X_j \hat{b}_j$, and repeat

What if a "bad decision" is made early on?

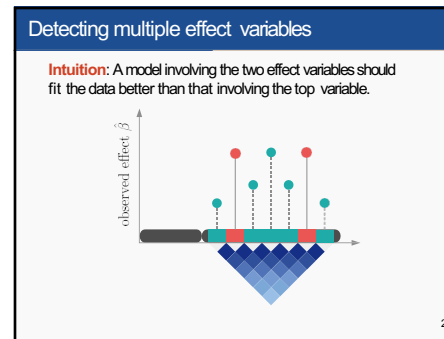
23



24



25



26

Bayesian Variable Selection Regression (BVSr)

Fine-mapping is a particular multiple regression problem:

$$y_{n \times 1} = X_{n \times p} b_{p \times 1} + e_{n \times 1}$$

- b is sparse: most of its elements are 0's
- Columns of X are very correlated

21

27

Why BVSr?

- Other sparse variable selection regression may not work
 - designed to minimize prediction errors, e.g. LASSO

22

28

Why BVSr?

- Other sparse variable selection regression may not work
 - designed to minimize prediction errors, e.g. LASSO
- Bayesian variable selection regression (BVSr)
 - can evaluate significance of effect variables
 - can quantify **uncertainty** in variables selected

22

29

Software	Test type	Hard constraint	Use auxiliary methods	Maximum number of variables	Fast estimation?	Good search	Multioutput
BAYESPIE	q and binary	No	No	Fixed	No	Education	Bayes factor
BayesPIE	q and binary	No	No	Fixed	No	Education	Bayes factor
SNPTEST	q and binary	No	No	Fixed	No	Education	Bayes factor
BAYESPIE	q and binary	No	No	Compared	No	MC/MC	Bayes factor and PP
SNPTEST	Binary	Yes	No	Compared	Yes	MC/MC	Bayes factor and PP
SNPTEST	q	No	No	Compared	No	MC/MC	Bayes factor and PP
SNPTEST	q	Yes	Yes	1, fixed and compared	Yes	Education	Bayes factor and PP
Finemapping	Multioutput	Yes	No	Compared	No	Greedy	PP
Finemapping	Multioutput	Yes	No	Compared	No	Greedy	Bayes factor and PP
Finemapping	Binary	Yes	No	20	No	MC/MC	PP
SNPTEST	q and binary	Yes	No	1	No	Education	Bayes factor
SNPTEST	q and binary	No	Yes	1	Yes	Education	Bayes factor and PP
SNPTEST	q and binary	No	Yes	Fixed	No	Education	probability, confidence set and PP
FINEMAPPING	q, binary and any	No	Yes	Fixed and compared	Yes	Education and MC/MC	Bayes factor and PP
FINEMAPPING	q and binary	No	No	Fixed	Yes	Education	Bayes factor and PP
FINEMAPPING	q and binary	No	Yes	Fixed	No	Shogun	Bayes factor and PP, confidence set
FINEMAPPING	q and binary	No	Yes	Fixed and compared	No	Education	Bayes factor and PP

Figure: Schaid et al. (2018) Nat. Rev. Genet.

23

30

BVSR model

$$y = Xb + e$$

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

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

$$b_j | y \sim g(\cdot)$$

$$b_{-j} | y \sim \tilde{\alpha}$$

y : model configurations; π : prior inclusion probability.

24

31

BVSR results

Assess combinations of variables

SNPs	1	2	3	4	5	...	Probability
model configurations	1	0	1	0	0	...	0.25
	1	0	0	1	0	...	0.25
	0	1	1	0	0	...	0.25
	0	1	0	1	0	...	0.25

- $PIP_j := Pr(z_j \text{ is non-zero})$

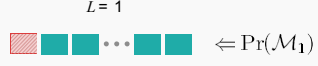
PIP = (0.5, 0.5, 0.5, 0.5, 0, ...)

25

32

Assessing multi-effects configurations


$L = 1$



$\Leftarrow Pr(\mathcal{M}_1)$

26

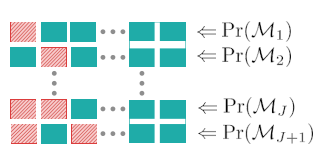
33



$\Leftarrow Pr(\mathcal{M}_1)$

$\Leftarrow Pr(\mathcal{M}_2)$

34



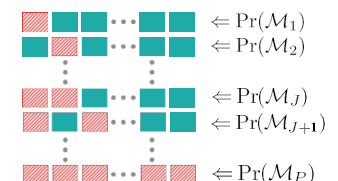
$\Leftarrow Pr(\mathcal{M}_1)$

$\Leftarrow Pr(\mathcal{M}_2)$

$\Leftarrow Pr(\mathcal{M}_J)$

$\Leftarrow Pr(\mathcal{M}_{J+1})$

35



$\Leftarrow Pr(\mathcal{M}_1)$

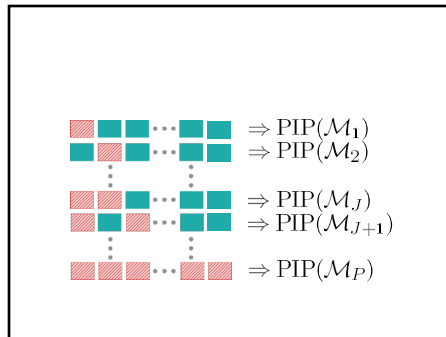
$\Leftarrow Pr(\mathcal{M}_2)$

$\Leftarrow Pr(\mathcal{M}_J)$

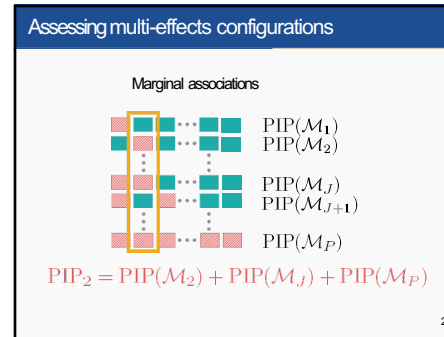
$\Leftarrow Pr(\mathcal{M}_{J+1})$

$\Leftarrow Pr(\mathcal{M}_P)$

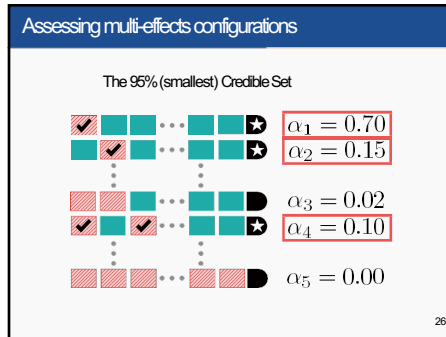
36



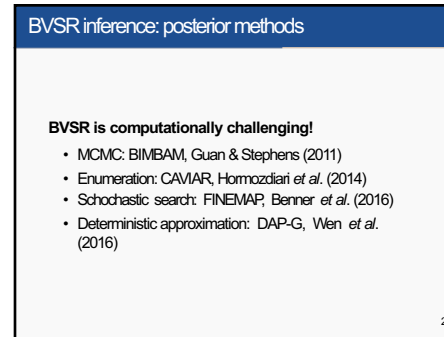
37



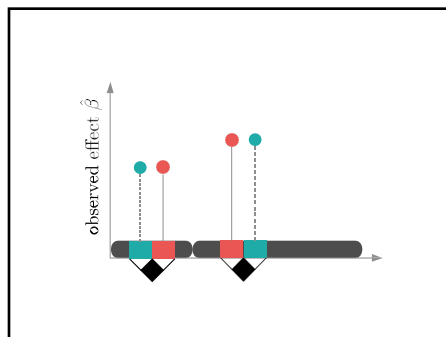
38



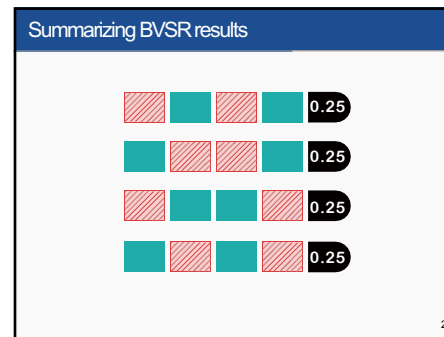
39



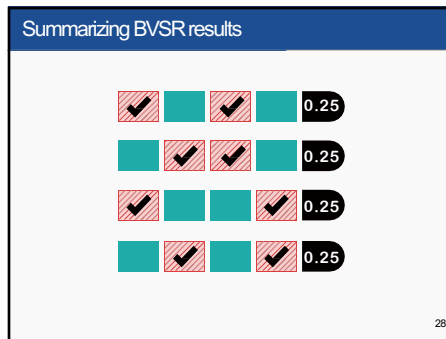
40



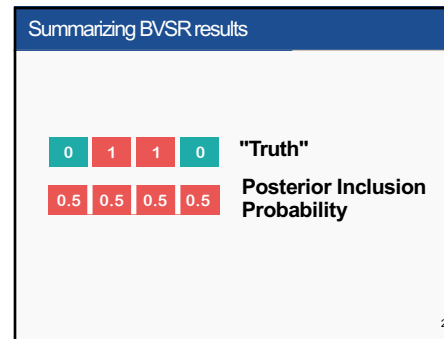
41



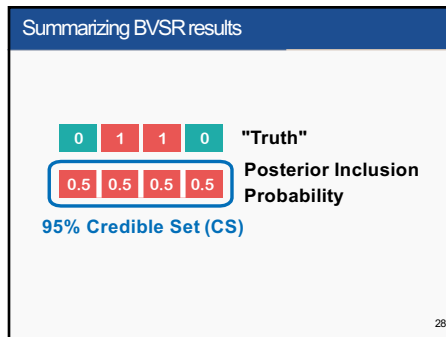
42



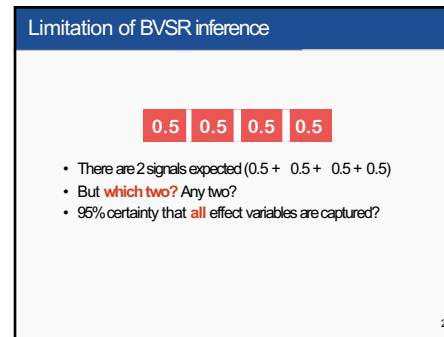
43



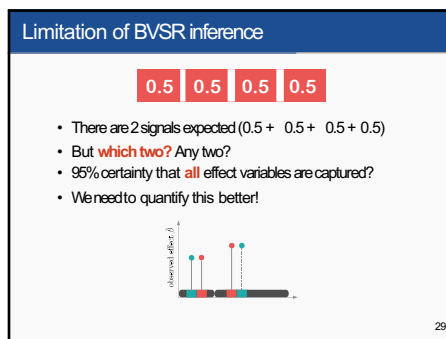
44



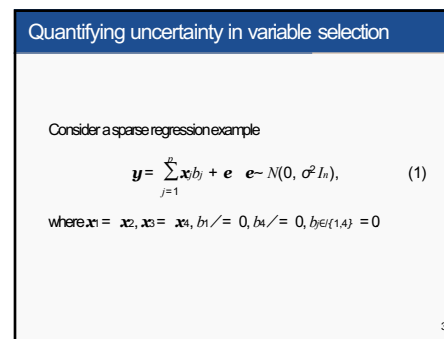
45



46



47



48

Quantifying uncertainty in variable selection

Consider a sparse regression example

$$y = \sum_{j=1}^n x_j b_j + e \quad e \sim N(0, \sigma^2 I_n), \quad (1)$$

where $x_1 = x_2, x_3 = x_4, b_1 \neq 0, b_2 \neq 0, b_3 \neq 0, b_4 \neq 0, b_{j \in \{1,4\}} = 0$

We are interested in making the following statement,

$$(b_1 \neq 0 \text{ or } b_2 \neq 0) \text{ AND } (b_3 \neq 0 \text{ or } b_4 \neq 0).$$

30

49

Quantifying uncertainty in variable selection

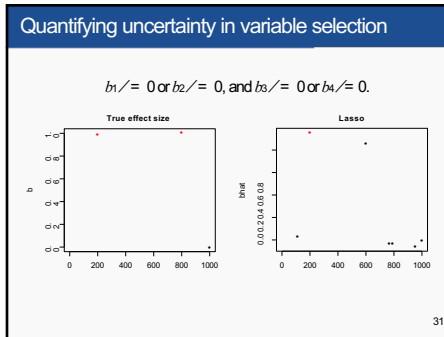
We are interested in making the following statement,

$$(b_1 \neq 0 \text{ or } b_2 \neq 0) \text{ AND } (b_3 \neq 0 \text{ or } b_4 \neq 0).$$

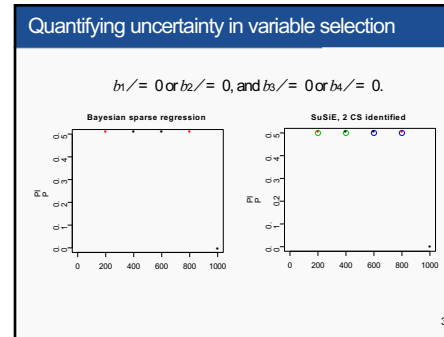
1. There are two independent variables with non-zero effect
2. x_1 and x_2 (and x_3 and x_4) are too similar to distinguish
3. yet they can be prioritized relative to each other

30

50



51



52

The Sum of Single Effects model (SuSIE)

$$y = Xb + e$$

$$b = \sum_{l=1}^L b_l$$

Wang et al. (2020) J. R. Stat. Soc. B

32

53

The Sum of Single Effects model (SuSIE)

$$y = Xb + e$$

$$b = \sum_{l=1}^L b_l$$

A variational approximation to posterior under SuSIE

$$q(b_1, \dots, b_L) = \prod_l q_l(b_l)$$

- b_1, \dots, b_L are treated as independent a posteriori.
- Do not assume q_l factorizes over the elements of b_l .

32

54

A fast Bayesian variable selection

Iterative Bayesian forward selection algorithm (IBSS)

- For each iteration l
 - For each SNP j fit $y = X_j^{(l)} + e$ get $BF_j^{(l)}$
 - Form weight for each SNP $w_j^{(l)} \propto BF_j^{(l)}$
 - Form residuals $y' := y - \sum_j w_j^{(l)} X_j^{(l)}$ and repeat
- Until converge

Coordinate ascent algorithm; convergence based on evidence lower bound (ELBO)

33

55

SuSiE model, formal notation

"single effect": b 's

$$y = Xb + e$$

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

$$b = \sum_{l=1}^L b_l$$

$$b_l = \gamma_l \beta_l$$

$$\gamma_l \sim \text{Mult}(1, \beta_l)$$

$$\beta_l \sim N(0, \sigma_l^2)$$

$$\sigma_l^2 \geq 0$$

A mean-field approximation $q(b_1, \dots, b_L) = \prod_l q_l(b_l)$

- b_1, \dots, b_L are treated as independent a posteriori.
- Do not assume q_l factorizes over the elements of b_l .

34

56

IBSS algorithm, formal notation

Algorithm Iterative Bayesian forward selection

Require: single effect regressor $SER(y, X) \rightarrow (\alpha, \mu, \sigma^2)$

- Initialize α, μ, σ^2 for $l = 1, \dots, L$.
- repeat
- for l in $1, \dots, L$ do
- $n \leftarrow y - \sum_{j \neq l} \alpha_j$
- $(\hat{\alpha}_l, \hat{\mu}_l, \hat{\sigma}_l^2) \leftarrow SER(n, X_l)$
- $\alpha_l \leftarrow \hat{\alpha}_l$
- until converged
- return $\alpha, \mu, \dots, \alpha, \mu$.

35

57

SuSiE model yields single-effect CS

36

58

SuSiE model yields single-effect CS

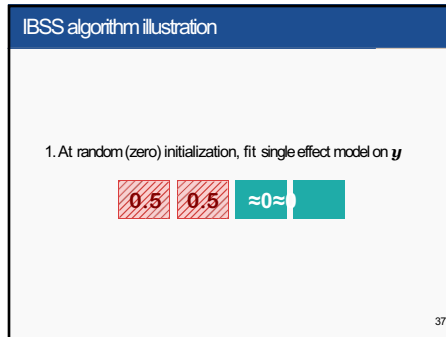
36

59

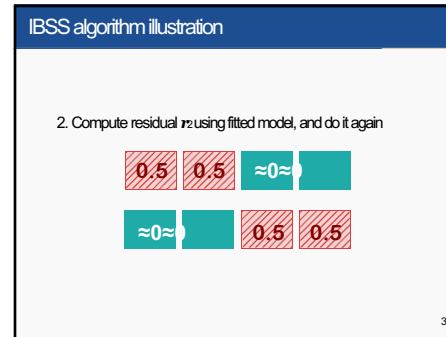
SuSiE model yields single-effect CS

36

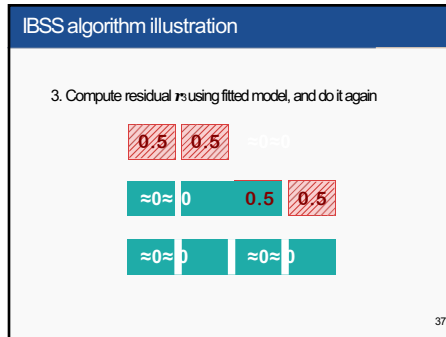
60



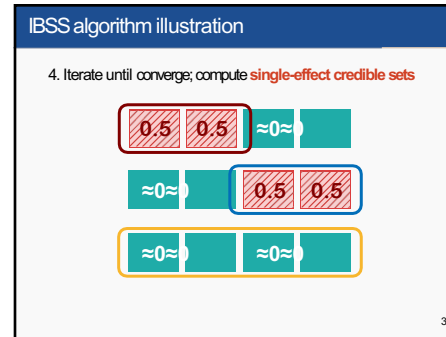
61



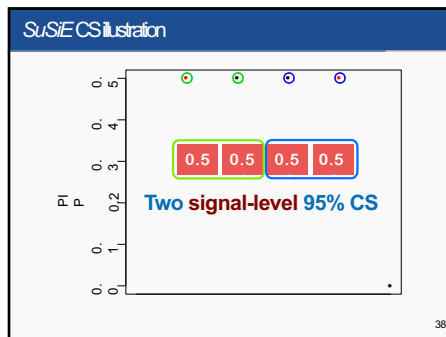
62



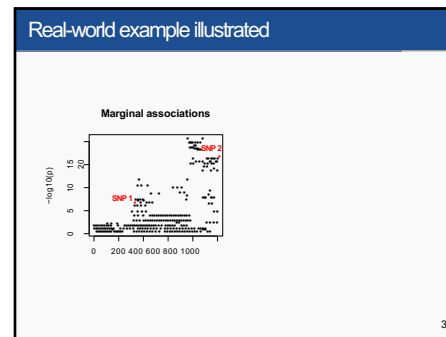
63



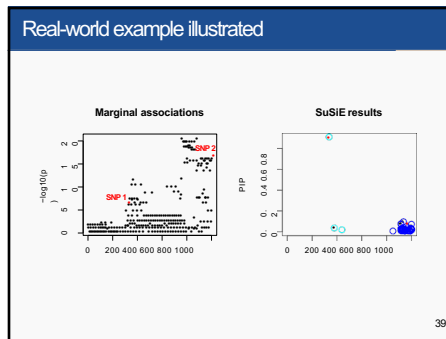
64



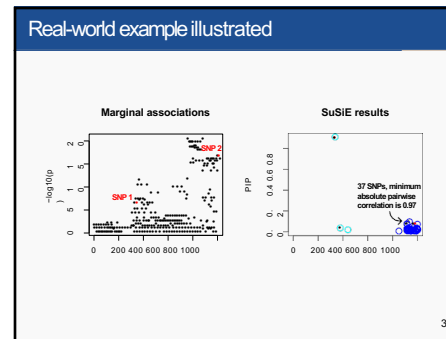
65



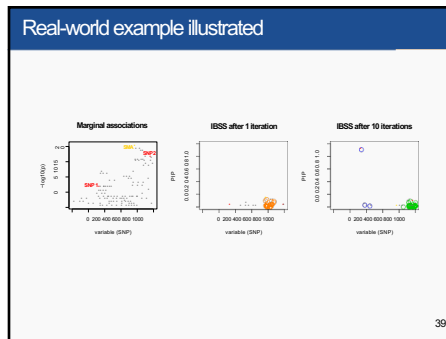
66



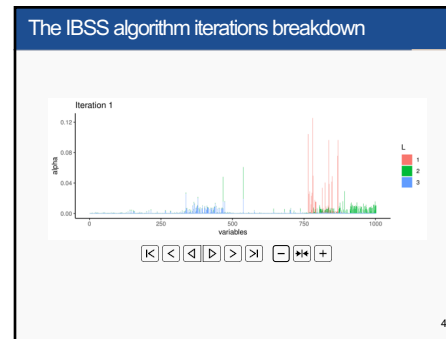
67



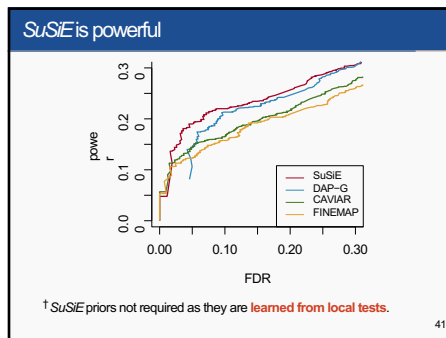
68



69



70



71

SuSIE is fast

Speed comparison (3 causal variables; unit: sec.)

Method	Avg.	Min.	Max.
SuSIE [†]	0.64	0.34	2.28
DAP-G	2.87	2.23	8.87
FINEMAP	23.01	10.99	48.16
CAVIAR	2907.51	2637.34	3018.52

† An R implementation of SuSIE. Others are implemented in C++.

42

72

Similar model, different problems

- X is gene expression, y is tissue / cell type
- X is pathway, y is gene-set
- X is functional annotation, y is GWAS effect size
- X is "step matrix", y is spatially-structured variable

43

73

The "changepoint" problem

Data is piecewise constant, e.g. copy number variation

44

74

The "changepoint" problem

Can be modelled as linear combination of step functions

44

75

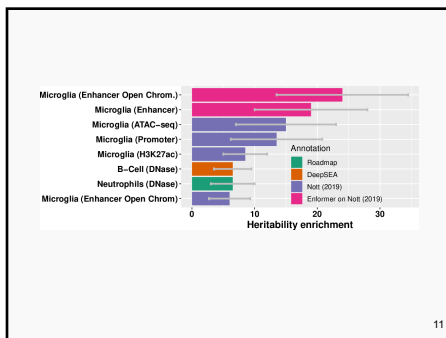
Example: simulated DNA copy number variation

SuSE vs Circular Binary Segmentation Cohen et al. (2004) Biostatistics

Notice the benefit of quantifying uncertainty in this example

45

76



77

A sparse model (a somewhat oligogenic)

Generalized linear model for SNP effects given K annotations

$$\beta_j = (1 - \pi_j)\alpha_0 + \pi_j g(\Theta)$$

$$\pi_j := \Pr(y_j = 1 | \alpha \mathbf{d})$$

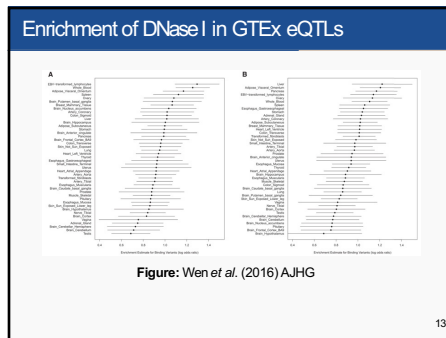
$$\log \frac{\pi_j}{1 - \pi_j} = \alpha + \sum_{k=1}^K \alpha_k d_j$$

α are log fold enrichment of functional genomic features

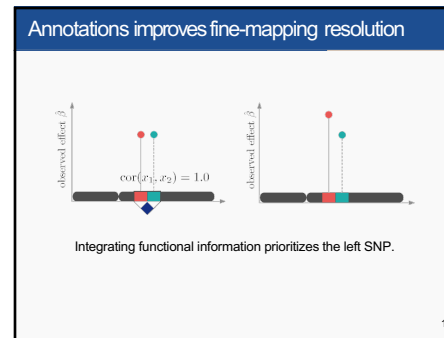
- Suggested reading: Wen (2016) AoAS

12

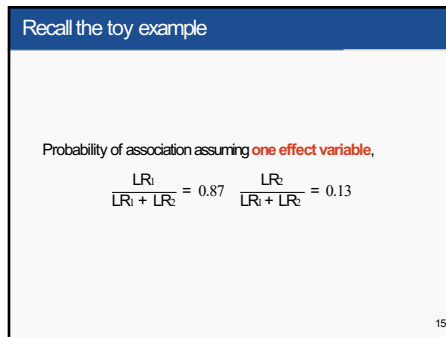
78



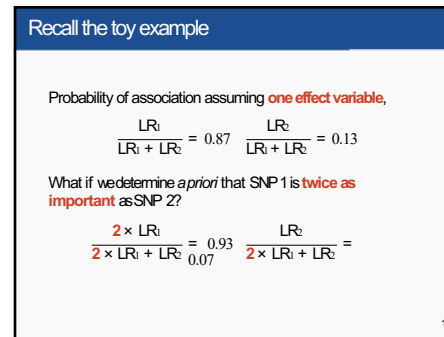
79



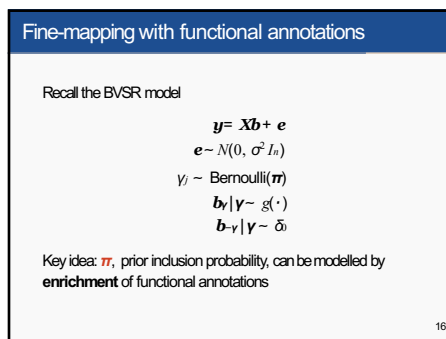
80



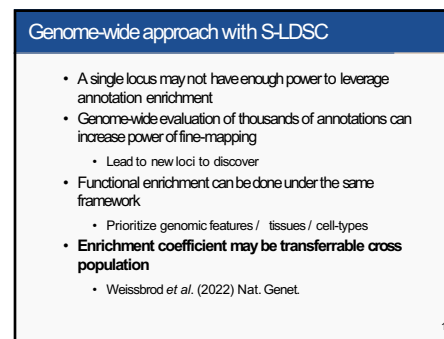
81



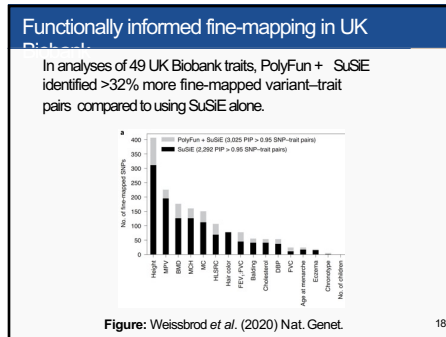
82



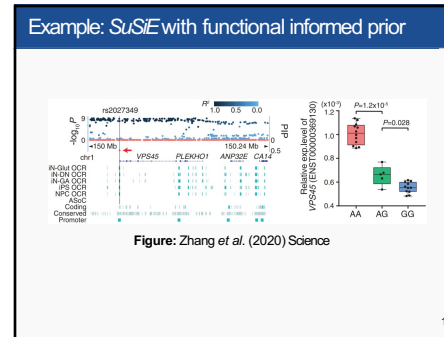
83



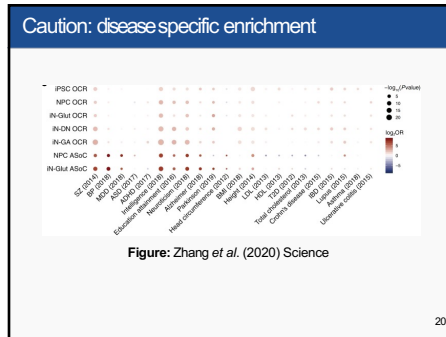
84



85



86



87

Multivariate analysis in genetic association studies

Geo Wang, Ph.D.
Advanced Gene Mapping Course, May 2024
The Gertrude H. Sergievsky Center and Department of Neurology
Columbia University Vagelos College of Physicians and Surgeons

88

- 1 Motivation
- 2 Meta-analysis review
- 3 Meta-analysis: a multivariate regression perspective
- 4 Multivariate adaptive shrinkage and fine-mapping

89

Motivation

90

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)

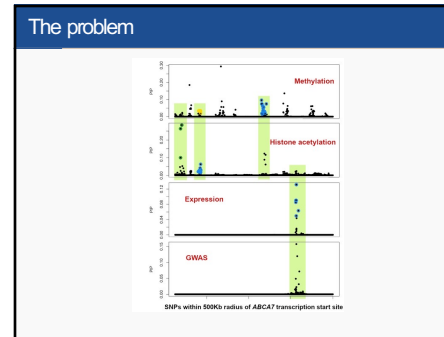
- Integrate information across multiple conditions / studies

“Causal” variants across multiple conditions?

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

3

91



92

The problem

For a genetic variable analyzed in two conditions:

$$P(\text{“causal” in trait 1 \& 2} \mid \text{association data for 1 \& 2})$$

5

93

The problem

For a genetic variable analyzed in two conditions:

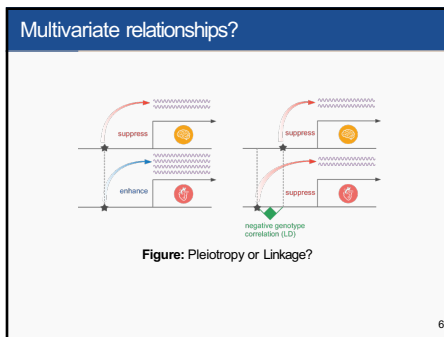
$$P(\text{“causal” in trait 1 \& 2} \mid \text{association data for 1 \& 2})$$

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 :

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

5

94



95

Meta-analysis review

96

Meta-analysis: a multivariate regression prospective

97

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

7

98

Fixed effect (FE) model

Let $\hat{\beta}_i$ be the observed effect size of study i , $1 \leq i \leq k$, 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

$$L(\beta) = P(\hat{\beta} | \beta) = \prod_i P(\hat{\beta}_i | \beta) \propto \prod_i \exp - \frac{(\hat{\beta}_i - \beta)^2}{2s_i^2}$$

8

99

Fixed effect (FE) model

Let $\hat{\beta}_i$ be the observed effect size of study i , $1 \leq i \leq k$, 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

$$L(\beta) = P(\hat{\beta} | \beta) = \prod_i P(\hat{\beta}_i | \beta) \propto \prod_i \exp - \frac{(\hat{\beta}_i - \beta)^2}{2s_i^2}$$

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

$$\hat{\beta} = \frac{\sum_i w_i \hat{\beta}_i}{\sum_i w_i} \quad \text{Inverse variance weighting}$$

8

100

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

$$P(\hat{\beta} | \beta, \sigma^2) \propto \prod_i \frac{1}{s_i^2 + \sigma^2} \exp - \frac{(\hat{\beta}_i - \beta)^2}{2(s_i^2 + \sigma^2)}$$

9

101

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

$$P(\hat{\beta} | \beta, \sigma^2) \propto \prod_i \frac{1}{s_i^2 + \sigma^2} \exp - \frac{(\hat{\beta}_i - \beta)^2}{2(s_i^2 + \sigma^2)}$$

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).

9

102

Multivariate model(s) for effect

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.

10

103

Fixed-effect model multivariate analysis

Effect sizes are exactly the same between two studies,

$$U_{\text{fixed}} = \sigma^2 \times \begin{matrix} i & 1 \\ 1 & 1 \\ 1 & 1 \end{matrix}$$

11

104

Random effects model multivariate analysis

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

$$U_{\text{random}} = \sigma^2 \times \begin{matrix} i & 1 \\ 1 & 0 \\ 0 & 1 \end{matrix}$$

12

105

Other multivariate models

$$U_{\text{partially shared}} = \sigma^2 \times \begin{matrix} i & 1 \\ 1 & \rho \\ \rho & 1 \end{matrix}$$

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

13

106

Other flexible multivariate models

More generally,

$$U = \begin{matrix} i & 1 \\ \sigma_{12}^2 & \sigma_2^2 \\ \sigma_{12}^2 & \sigma_2^2 \end{matrix}$$

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

14

107

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

- Colocalization correlation matrix:

$$\begin{matrix} i & 1 \\ 1 & \rho \\ \rho & 1 \end{matrix}$$
- Condition specific correlation matrix:

$$\begin{matrix} i & 1 & i & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{matrix}$$

15

108

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

- Mediation:

$$U_{\text{mediation}} = \begin{matrix} & i & 1 & 1 \\ \sigma^2 & \times & & \\ \rho_{12} & & \rho_{12} & \\ \rho & & & \rho \end{matrix}$$

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

16

109

The problem

For a genetic variable analyzed in GWAS and eQTL studies:

$$P(y_g = 1, y_e = 1 | D_g, D_e, \Theta)$$

17

110

Colocalization method: *coloc*

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

- On X: "one causal" assumption
- On Y: the null + 4 combinations given "one causal"
 - In 1 but not 2
 - In 2 but not 1
 - In 1 and 2 but not the same variable
 - In 1 and 2 and the same variable (colocalization)
 - No association in both data 1 and 2

18

111

Colocalization method: *eCAVIAR*

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

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

19

112

eCAVIAR effects assumption

Effect sizes are independent,

$$U = \begin{matrix} & i & 1 \\ \sigma_g^2 & & 0 \\ \sigma_e^2 & & \sigma_e^2 \end{matrix}$$

20

113

Colocalization method: *enloc*

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

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

21

114

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?

- Through a simple logistic link **using eQTL as an annotation** for j

$$\log \frac{\pi}{1 - \pi} = \alpha + \alpha \gamma_e$$

and in this context

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

21

115

enloc two step procedure

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

22

116

Connections between colocalization methods

- eCAVIAR is a special case of *enloc* with $\alpha = 0$.
- 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://chrishwallace.github.io/coloc/>

23

117

Connections between colocalization methods

- eCAVIAR is a special case of *enloc* with $\alpha = 0$.
- 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://chrishwallace.github.io/coloc/>

Summary: **pattern** and **scale** of effect size correlations, represented as different **prior** models.

23

118

Practical considerations

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

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

24

119

Multi-trait colocalization

The screenshot shows the HyPrColoc interface with several hypothesis selection options (e.g., 'No association with any of the traits', 'One trait has a CV in the region') and enrichment model configurations (e.g., 'HyPrColoc', 'SuSIE-Enrich').

Figure: HyPrColoc, Foley *et al.* (2021) Nat. Comm.

Assuming a single causal variant in the loci.

25

120

Incorporating all possible patterns

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

$$U_{mixed} = \pi_1 \times \begin{pmatrix} i & 2.4 & 0.3 & 1 \\ 0.3 & 1.5 & & \end{pmatrix} + \pi_2 \times \begin{pmatrix} i & 1.6 & 0.001 & 1 \\ 0.001 & 0.02 & & \end{pmatrix} + \pi_3 \times \dots$$

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.

31

127

Multivariate effect size sharing in eQTLs

Figure: Quantitative characterization of eQTL effects heterogeneity in GTEx

32

128

Application to multivariate fine-mapping

Figure: mvSuSiE fine-mapping with adaptive shrinkage model

Zou et al. (2023) bioRxiv

33

129

Multi-trait fine-mapping methods & challenges

	mvSuSiE	CAFEI	PAINTOR	WTFESS	BayesSUI	fmshin	mvCovar	HyP-Covar	mtm
>5 traits integrated									
>10 traits integrated									
Multiple causal alleles									
Individual level data									
Summary statistics									
Missing data									
Trait specific LD									
Correlated effects									
Trait specific effects									
Arbitrary heterogeneous effects									
Arbitrary multi-trait correlation									
Correlated traits									
Partial sample overlap									
Functional annotation									
Trait specific functional annotation									
Genetic architecture									
Heritability									

Reference: CAFEI: Arvanitis et al (2022); PAINTOR: Kichaev et al (2017); WTFESS: Levin et al (2016); BayesSUI: Zhou et al (2017); fmshin: Hernandez et al (2021); mvCovar: Luffano et al (2020); HyP-Covar: Froy et al (2021); mtm: Gumbastromer et al (2016).

34

130

Comparison to other methods

35

131

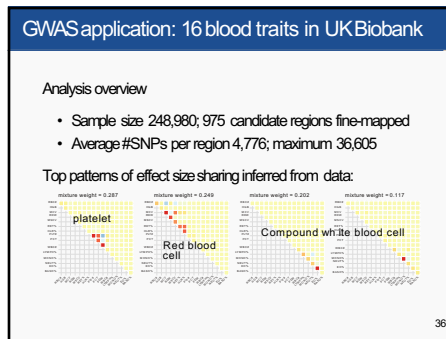
GWAS application: 16 blood traits in UK Biobank

Analysis overview

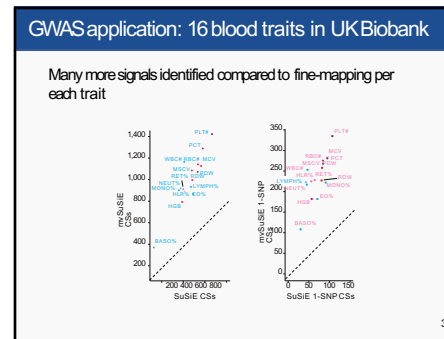
- Sample size 248,980; 975 candidate regions fine-mapped
- Average #SNPs per region 4,776; maximum 36,605

36

132



133



134

Complex phenotype prediction and transcriptome-wide association studies

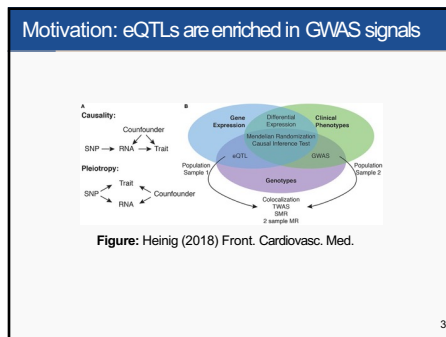
Gao Wang, Ph.D.
 Advanced Gene Mapping Course, May 2024
 The Gertrude H. Sergievsky Center and Department of Neurology
 Columbia University Vagelos College of Physicians and Surgeons

1

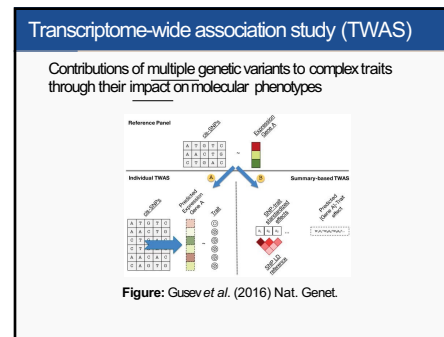
135

2

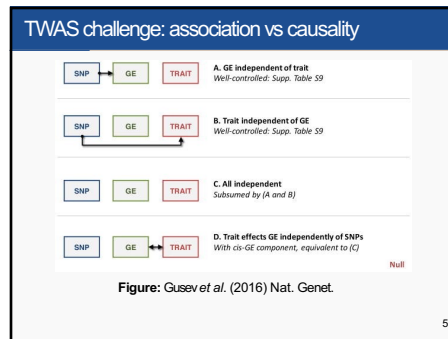
136



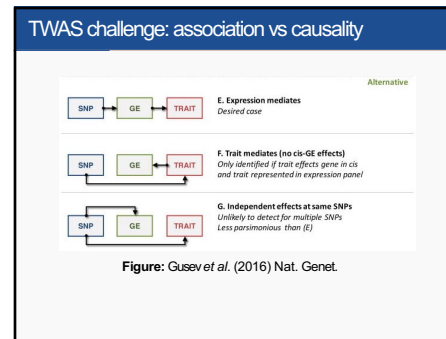
137



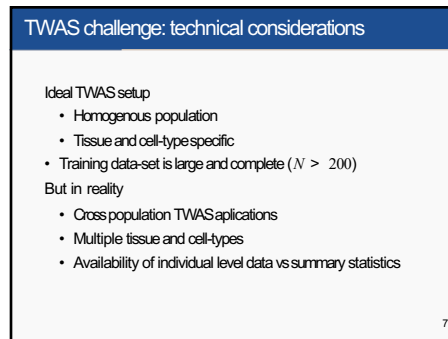
138



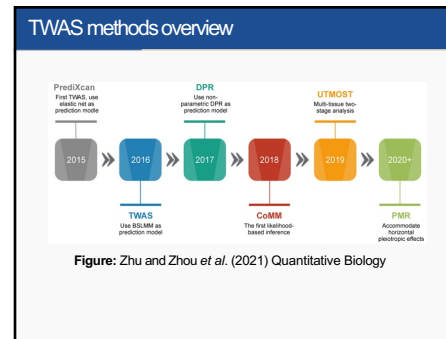
139



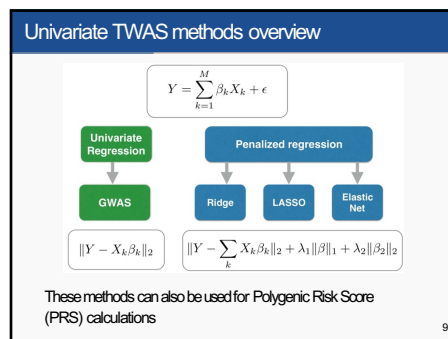
140



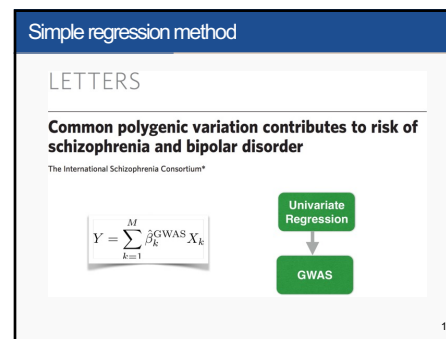
141



142



143



144

Ridge regression / BLUP

REPORT
GCTA: A Tool for Genome-wide Complex Trait Analysis
 Jian Yang,^{1*} S. Hong Lee,¹ Michael E. Goddard,^{2,3} and Peter M. Visscher¹
 AJHG 2011

Penalized regression
 ↓
 Ridge

$$Y = \sum_{k=1}^M \beta_k^{\text{Ridge}} X_k$$

$$\|Y - \sum_k X_k \beta_k\|_2 + \lambda_2 \|\beta_2\|_2$$

11

145

Other penalized regression

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

Regularization and variable selection via the elastic net
 Hui Zou and Trevor Hastie
 Stanford University, USA

Penalized regression
 ↓
 LASSO Elastic Net

$$Y = \sum_{k=1}^M \beta_k^{\text{E-N}} X_k$$

$$\|Y - \sum_k X_k \beta_k\|_2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta_2\|_2$$

12

146

Bayesian variable selection regression

OPEN ACCESS Freely available online **PLOS GENETICS**

Polygenic Modeling with Bayesian Sparse Linear Mixed Models
 Xiang Zhou^{1*}, 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 Biol. published online June 24, 2014
 Access the most recent version at doi:10.1101/gb.169075.113

13

147

Choice of methods: cross validation

TWAS / FUSION
 Functional Summary-based Imputation

- New TWAS (Stratton et al.) models for cis-eQTLs
- New FUSION (Thompson et al.) context-specific models for single-cell and bulk expression
- New eQTLs v8 models

FUSION is a suite of tools for performing transcriptome-wide and region-wide association studies (TWAS and RWAS). FUSION builds predictive models of the genetic component of a functional/molecular phenotype and predicts and tests that component for association with disease using GWAS summary statistics. The goal is to identify associations between a GWAS phenotype and a functional phenotype that was only measured in reference data, to provide uncollected predictive results from multiple studies in full-scale GWAS analysis.

Please cite the following manuscript for TWAS methods:
 Zhou et al. "Integrative approaches for large-scale transcriptome-wide association studies" 2016 Nature Genetics

14

148

Likelihood based approach

Figure: CoMM, Yeung et al. (2019)

Also see Yuan et al. (2022) likelihood based Mendelian Randomization

15

149

Multivariate TWAS methods overview

Leverage similarity between molecular phenotypes

44 molecular phenotypes
 44 tissues
 11 million SNPs

- UTMOST, Yu et al. (2019) Nature Genetics
- MR-JTI, Zhou et al. (2020) Nature Genetics
- mr.mash, Morgante et al. (2023) PLoS Genetic (to appear)

16

150

Multivariate TWAS hands-on exercise

statgen-setup launch --tutorial twas

21

157

Missing regulation in eQTL and GWAS

The missing link between genetic association and regulatory function

... by applying a gene-based approach we found limited evidence that the baseline expression of trait-related genes explains GWAS associations, whether using colocalization methods (8% of genes implicated), transcription-wide association (2% of genes implicated), or a combination of regulatory annotations and distance (4% of genes implicated). These results contradict the hypothesis that most complex trait-associated variants coincide with homeostatic expression QTLs, suggesting that better models are needed. The field must confront this deficit and pursue this 'missing regulation.'

22

158

TWAS and fine-mapping: variable selection

23

159

TWAS and fine-mapping: variable selection

Figure: Zhao et al. (2022) bioRxiv

24

160

TWAS and colocalization: pleiotropy

Figure: Jordan et al. (2019) Genome Biology

25

161

TWAS + colocalization: pleiotropy

Image credit: Haky Im @UChicago
 • "Locus level" colocalization: Hukku et al. (2022) AJHG;
 Okamoto et al. (2023) AJHG.

26

162

TWAS and colocalization: statistical framework

$$M = \mu_M \mathbf{1} + G\beta_E + e_M, e_M \sim N(0, \sigma_M^2 \mathbf{I})$$

$$Y = \mu_Y \mathbf{1} + \gamma M + G\beta_Y + e_Y, e_Y \sim N(0, \sigma_Y^2 \mathbf{I})$$

- "locus level", $Pr(\gamma \neq 0 | \text{Data}) \propto Pr(\gamma \neq 0) P(\text{Data})$
- $Pr(\gamma \neq 0) = Pr(\text{coloc}) \times Pr(\text{twas})$
- Data: z-score from TWAS.
- Key idea: Test $\gamma = 0$, not to estimate γ which is Mendelian Randomization.

27

163

TWAS and Mendelian randomization

Figure: Zhu and Zhou (2021) Quantitative Biology

TWAS can be viewed as two-sample MR — using various IV selection methods.

28

164

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

Gao Wang, Ph.D.
 Advanced Gene Mapping Course, May 2024
 The Gertrude H. Sergievsky Center and Department of Neurology
 Columbia University Vagelos College of Physicians and Surgeons

1

Association analysis summary statistics

z-scores from univariate association studies:

$$\hat{z}_j := \hat{\beta}_j / s_{\beta_j}$$

where

$$\hat{\beta}_j := (\mathbf{x}_j^T \mathbf{x}_j)^{-1} \mathbf{x}_j^T \mathbf{y} \quad s_{\beta_j} := \sqrt{\hat{\sigma}_j^2 (\mathbf{x}_j^T \mathbf{x}_j)^{-1}}$$

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

2

1

2

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 problems
 - $O(p^2)$ (summary statistics) \ll $O(np)$ (full data)
 - when sample size $n \gg$ variants in fine-mapped region p

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

3

3

Regression with Summary Statistics (RSS)

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

Assumptions:

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

4

4

Properties of per SNP z scores

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

$$E(\hat{z}_j | \hat{\mathbf{R}}) = \sum_{i=1}^n r_{ij} z_i$$

GWAS marginal effects are biased due to LD!

- z-scores are correlated,

$$\text{Cor}(\hat{z}_j, \hat{z}_k) = r_{jk}, \forall j, k$$

- Recall the previously discussed connection with LDSC

5

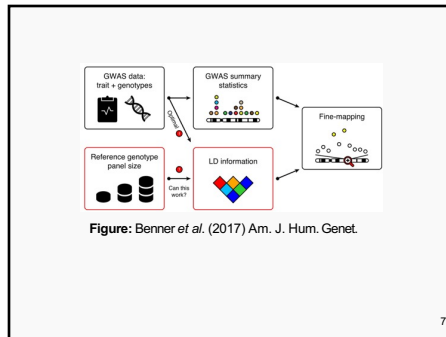
5

Summary of summary statistics

- \mathbf{X} , genotype matrix
- \mathbf{Y} , phenotype matrix, can be multiple traits
- $\mathbf{X}^T \mathbf{Y}$, association results — effect size estimate
- \mathbf{X}^X , LD matrix
- $\mathbf{X} \mathbf{X}^T$, genomic relatedness matrix, reflects kinship
- \mathbf{Y}^Y , trait correlation, relevant in multi-trait analysis and integration

6

6



7

Fine-mapping via RSS model

"Single effect": z_i 's

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

$$\mathbf{z} = \sum_{l=1}^L \mathbf{z}_l$$

$$\mathbf{z}_l = \mathbf{Y}_l \gamma_l$$

$$\gamma_l \sim N(0, \omega_l^2)$$

$$\mathbf{Y}_l \sim \text{Mult}(1, \boldsymbol{\pi})$$

Suggested reading:
Zou et al (2022) PLoS Genet.

8

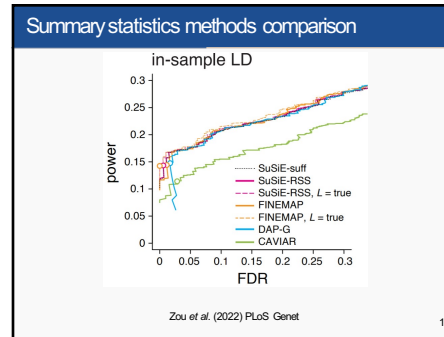
$\hat{\beta}$ and $SE(\hat{\beta})$ based models

The \hat{z} model: $\hat{z} \sim N(\mathbf{R}^T \mathbf{z}, \mathbf{R})$

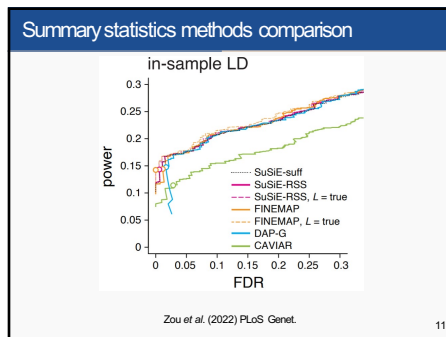
The \hat{b}, \hat{s} model: $\hat{b}, \hat{s} \sim N(\mathbf{S}^T \mathbf{R} \mathbf{S}^{-1} \mathbf{b}, \mathbf{S} \mathbf{R} \mathbf{S}^T)$

- 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
- \hat{b}, \hat{s} model takes sample size into consideration
 - No longer have to assume small effect per SNP
- \hat{z} model: CAVIAR, FINEMAP (2016)
- \hat{b}, \hat{s} model: Zhu and Stephens (2017) AoAS; FINEMAP (2018 10.1101/318618), SuSIE-RSS (Zou et al. 2022)

9



10



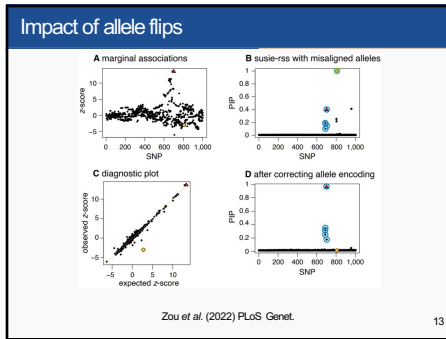
11

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

12



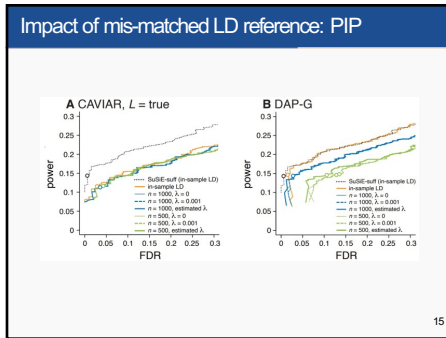
13

Addressing the allele flip challenge

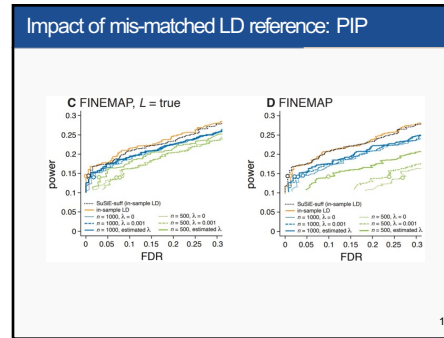
- `susier::susie_rss()` function implements a diagnosis
- `bigsnpr::snp_match()` function implements a basic allele matching for two sets of summary statistics
- Other resources
 - Allele flip illustration: https://statgen.us/lab-wiki/combio_tutorial/allele_fc
 - A powerful, multi-set data merger (by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

14

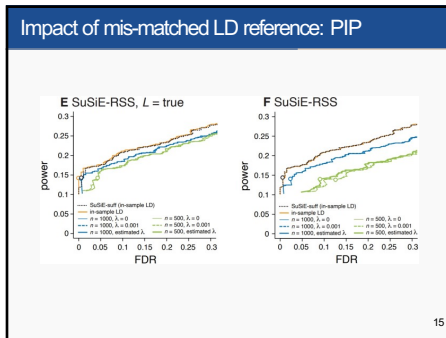
14



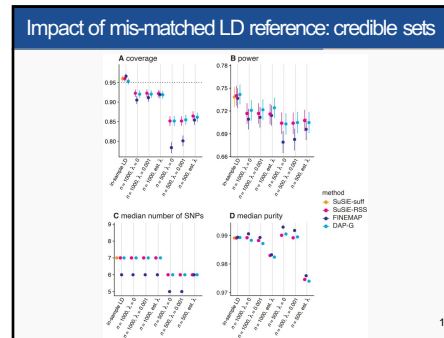
15



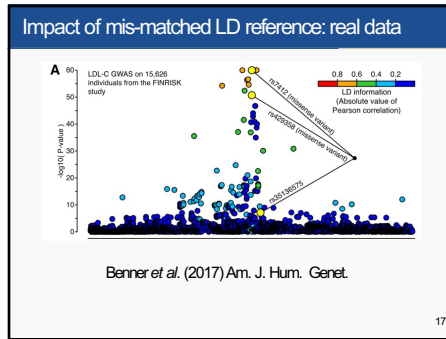
16



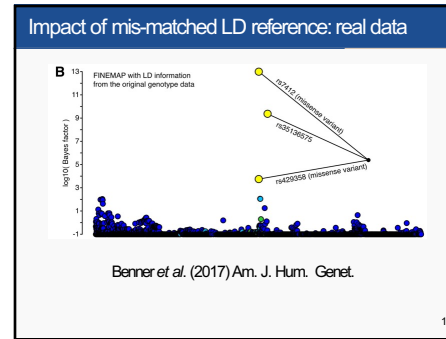
17



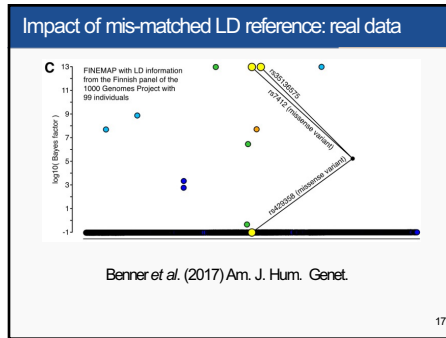
18



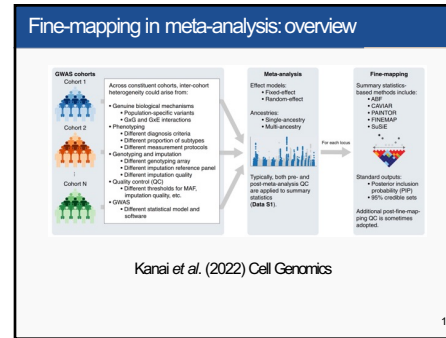
19



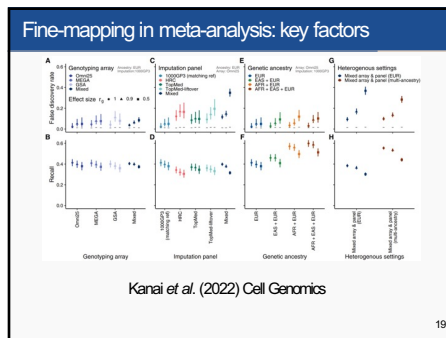
20



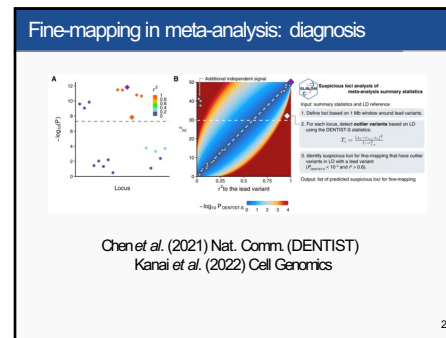
21



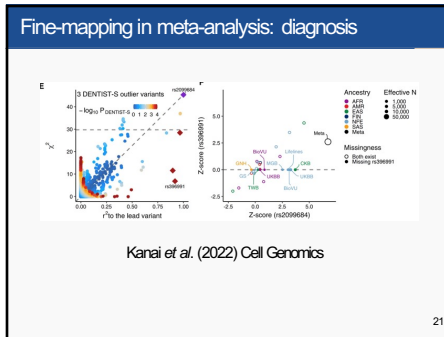
22



23



24



25

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?

More technical details see McCaw *et al.* (2020) Biometrics

26

Covariate adjustment in LD reference

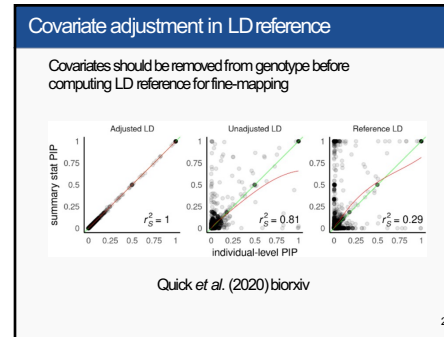
Consider two GWAS regression analysis:

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

More technical details see McCaw *et al.* (2020) Biometrics

27



28

Integrating GWAS with functional annotations

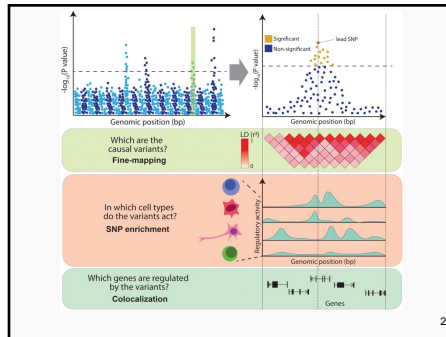
Gao Wang, Ph.D.
Advanced Gene Mapping Course, May 2024
The Gertrude H. Sergievsky Center and Department of Neurology
Columbia University Vagelos College of Physicians and Surgeons

1

1

Non-coding functional annotation in GWAS

2

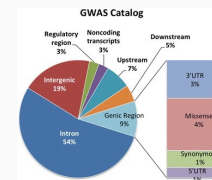


2

3

GWAS variants catalog by functional annotations

Most GWAS variants are non-coding



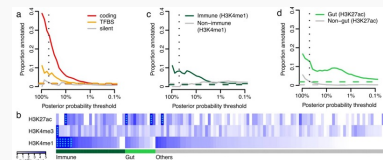
Lee et al. (2018) *Human Genetics*

3

4

Functional enrichment in fine-mapped variants

Signals concentrated in tissue / cell specific functional area



4

5

Functional annotation in aggregated rare variant association analysis

6

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

5

7

Annotations integrated to aggregated tests

Figure: Li et al. (2020) Nature Genetics

Also see Li et al. (2019) AJHG; Li et al. (2022) Nature Methods

6

8

Annotations integrated to aggregated tests

Figure: Li et al. (2020) Nature Genetics

6

9

Rare xQTL can improve PRS for complex traits

Figure: Small et al. (2022) AJHG

Also see Li et al. (2017) Nature; Ferraro et al. (2020) Science

7

10

Functional annotation in common variant association analysis

11

A polygenic model: stratified LD score regression

$$E[\chi^2_j] = 1 + \frac{Nl_j^2}{M} \left(\frac{1}{L_j} + \frac{1}{M} \right)$$

Labels for the equation:

- Chi-square GWAS statistic of variant j
- Sample size
- Narrow sense heritability
- LD score of variant j
- Total number of variants

$$l_j = \sum_{k \neq j} r_{jk}^2$$

LD score: sum of squared Pearson's correlation coefficient between SNP j and other (neighboring) SNPs.

8

12

A polygenic model: stratified LD score regression

Chi-square GWAS statistic of variant

$$E[\chi^2] = 1 + \frac{N h^2}{M} l_j$$

Sample size

Narrow sense heritability

LD score of variant

Total number of variants

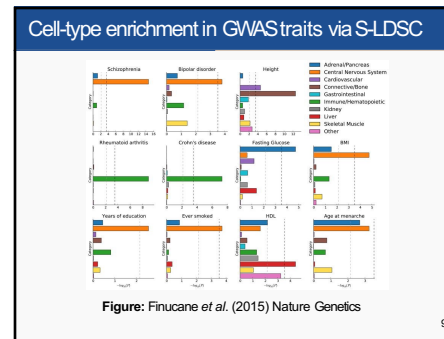
$$l_j = \sum_{k \neq j} r_{jk}^2$$

LD score: sum of squared Pearson's correlation coefficient between SNP and other (neighboring) SNPs

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

8

13



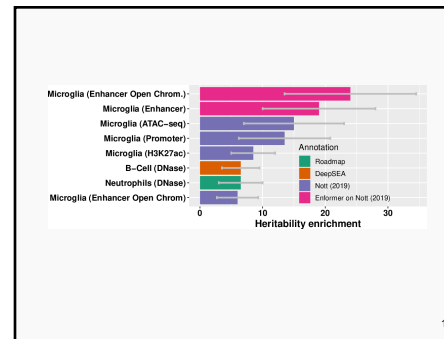
14

Integration approaches

- Integrate directly as range based binary annotations
 - Finucane *et al.* (2015) Nature Genetics — Stratified LDSC paper
- Extension: variant specific continuous annotations
 - Gazal *et al.* (2017) Nature Genetics
- Tissue specific variant level annotations independent of GWAS results
 - Deep Learning methods
 - Zhou *et al.* (2015) Nature Genetics, Zhou *et al.* (2018) Nature Genetics, Lai *et al.* (2022) PLoS Comp Bio
 - Avsec *et al.* (2021) Nature Methods

10

15



16

A sparse model (a somewhat oligogenic view)

Generalized linear model for SNP effects given K annotations

$$\beta_j = (1 - \pi_j)\delta_0 + \pi_j g(\Theta)$$

$$\pi_j := \Pr(y_j = 1 | \alpha \mathbf{d})$$

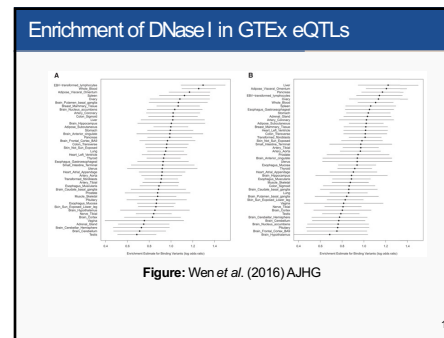
$$\log \frac{\pi_j}{1 - \pi_j} = \alpha_0 + \sum_{k=1}^K \alpha_k d_k$$

α are log fold enrichment of functional genomic features

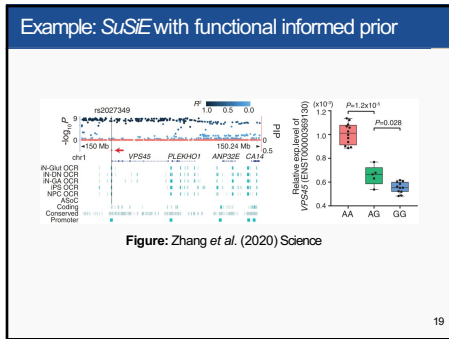
- Suggested reading: Wen (2016) AoAS

12

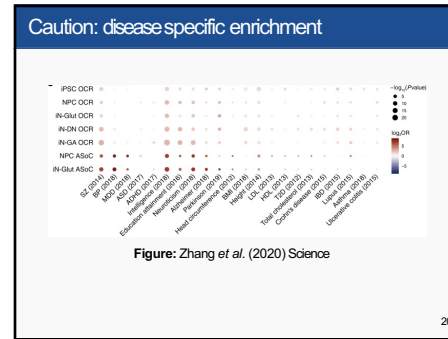
17



18



25



26

Multivariate analysis in genetic association studies

Gao Wang, Ph.D.
Advanced Gene Mapping Course, May 2024
The Gertrude H. Sergievsky Center and Department of Neurology
Columbia University Vagelos College of Physicians and Surgeons

1

1

- 1 Motivation
- 2 Meta-analysis review
- 3 Meta-analysis: a multivariate regression perspective
- 4 Multivariate adaptive shrinkage and fine-mapping

2

2

Motivation

3

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

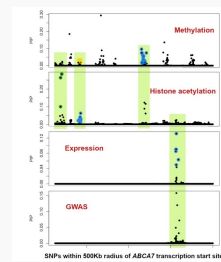
“Causal” variants across multiple conditions?

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

3

4

The problem



4

5

The problem

For a genetic variable analyzed in two conditions:

$$P(\text{“causal” in trait 1 \& 2} \mid \text{association data for 1 \& 2})$$

5

6

The problem

For a genetic variable analyzed in two conditions:

$P(\text{"causal" in trait 1 \& 2} \mid \text{association data for 1 \& 2})$

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 :

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

5

7

Meta-analysis review

9

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

7

11

Multivariate

Figure: Pleiotropy or Linkage?

6

8

Meta-analysis: a multivariate regression perspective

10

Fixed effect (FE) model

Let $\hat{\beta}_i$ be the observed effect size of study i , $1 \leq i \leq k$, 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

$$L(\beta) = P(\hat{\beta} \mid \beta) = \prod_i P(\hat{\beta}_i \mid \beta) \propto \prod_i \exp - \frac{(\hat{\beta}_i - \beta)^2}{2s_i^2}$$

8

12

Fixed effect (FE) model

Let $\hat{\beta}_i$ be the observed effect size of study i , $1 \leq i \leq k$, 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

$$L(\beta) = P(\hat{\beta} | \beta) = \prod_i^k P(\hat{\beta}_i | \beta) \propto \prod_i^k \exp - \frac{(\hat{\beta}_i - \beta)^2}{2s_i^2}$$

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

$$\hat{\beta} = \frac{\sum_i^k w_i \hat{\beta}_i}{\sum_i^k w_i} \quad \text{Inverse variance weighting}$$

8

13

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

$$P(\hat{\beta} | \beta, \sigma^2) \propto \prod_i^k \frac{1}{s_i^2 + \sigma^2} \exp - \frac{(\hat{\beta}_i - \beta)^2}{2(s_i^2 + \sigma^2)}$$

9

14

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

$$P(\hat{\beta} | \beta, \sigma^2) \propto \prod_i^k \frac{1}{s_i^2 + \sigma^2} \exp - \frac{(\hat{\beta}_i - \beta)^2}{2(s_i^2 + \sigma^2)}$$

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).

9

15

Multivariate model(s) for effect

Consider a parametric model on effect sizes across studies,

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

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

10

16

Fixed-effect model multivariate analysis

Effect sizes are exactly the same between two studies,

$$U_{\text{fixed}} = \sigma^2 \times \begin{matrix} & i & 1 \\ \begin{matrix} i & 1 \\ 1 & 1 \end{matrix} & & \end{matrix}$$

11

17

Random effects model multivariate analysis

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

$$U_{\text{random}} = \sigma^2 \times \begin{matrix} & i & 1 \\ \begin{matrix} i & 0 \\ 0 & 1 \end{matrix} & & \end{matrix}$$

12

18

Other multivariate models

$$U_{\text{partially shared}} = \sigma^2 \times \begin{matrix} i & 1 \\ 1 & \rho \\ \rho & 1 \end{matrix}$$

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

13

19

Other flexible multivariate models

More generally,

$$U = \begin{matrix} i & 1 \\ \sigma_1^2 & \sigma_{12}^2 \\ \sigma_{12}^2 & \sigma_2^2 \end{matrix}$$

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

14

20

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

- Colocalization correlation matrix:

$$\begin{matrix} i & 1 \\ 1 & \rho \\ \rho & 1 \end{matrix}$$
- Condition specific correlation matrix:

$$\begin{matrix} i & 1 & i & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{matrix}$$

15

21

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

- Mediation:

$$U_{\text{mediation}} = \sigma^2 \times \begin{matrix} i & 1 & \rho_{12} \\ 1 & \rho_{12} & \rho_2 \\ \rho_{12} & \rho_2 & \rho_2 \end{matrix}$$
- Genotype \rightarrow Trait 1 \rightarrow Trait 2.
- Effect on trait 2 should be smaller than that on trait 1.

16

22

The problem

For a genetic variable analyzed in GWAS and eQTL studies:

$$P(\gamma_S = 1, \gamma_e = 1 | D_S, D_e, \Theta)$$

17

23

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. In 2 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

18

24

Colocalization method: *eCAVIAR*

eCAVIAR [Homozidiari *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

19

25

eCAVIAR effects assumption

Effect sizes are independent,

$$U = \begin{matrix} & i & & 1 \\ & \sigma_g^2 & 0 & \\ & 0 & \sigma_e^2 & \end{matrix}$$

20

26

Colocalization method: *enloc*

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

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

21

27

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?

- Through a simple logistic link **using eQTL as an annotation** for j

$$\log \frac{\pi_j}{1 - \pi_j} = \alpha + \alpha \gamma_j$$

and in this context

$$\pi_j := P(Y_g = 1 | Y_e = 1)$$

21

28

enloc two step procedure

1. Obtain $P(Y_g = 1)$ and $P(Y_e = 1)$ using fine-mapping
2. Fit the enrichment model via **multiple imputation**

22

29

Connections between colocalization methods

- *eCAVIAR* is a special case of *enloc* with $\alpha = 0$.
- *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 *SuSE-coloc* removed the "one causal" assumption.
 - Wallace (2021) *PLoS Genetics*
 - <https://chrswallace.github.io/coloc/>

23

30

Connections between colocalization methods

- eCAVIAR is a special case of *enloc* with $\alpha = 0$.
- 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://chriswallace.github.io/coloc/>

Summary: **pattern** and **scale** of effect size correlations, represented as different **prior** models.

23

31

Practical considerations

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

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

24

32

Multi-trait

The screenshot shows a web interface for HyPrColoc. It features a table with columns for 'Hypothesis', 'Number of SNPs', 'Enrichment model', and 'Number of SNPs'. The 'Hypothesis' column lists various models such as A_1 , A_2 , A_3 , A_{1+2} , A_{1+2+3} , $A_{1+2+3+4}$, $A_{1+2+3+4+5}$, and $A_{1+2+3+4+5+6}$. The 'Enrichment model' column shows mathematical expressions for each hypothesis, such as $\frac{1}{\alpha} \frac{\beta_1}{\beta_2}$ for A_1 . The 'Number of SNPs' column shows the number of SNPs for each hypothesis, ranging from 1 to 6. Below the table, there is a caption: 'Figure: HyPrColoc, Foley *et al.* (2021) Nat. Comm. Assuming a single causal variant in the loci.'

25

33

Multivariate adaptive shrinkage and fine-mapping

34

More phenotypes, more complications

The diagram shows three different patterns of sharing between phenotypes. The first pattern is a full matrix where all elements are red, indicating that all phenotypes share all SNPs. The second pattern is a lower triangular matrix where only the elements on and below the diagonal are red, indicating that each phenotype shares SNPs with itself and all other phenotypes. The third pattern is a sparse matrix where only a few elements are red, indicating that only a few SNPs are shared between a few phenotypes. Below the diagram is a caption: 'Figure: Plausible patterns of sharing'.

26

35

Major

- For a given variant: the less assumption made on multivariate effects, the more parameters to estimate.
 - FE and RE models are restrictive but easy to fit.
- Different variants: may fit in different multivariate effect models

27

36

A naive mixture model

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

$$U_{mixed} = 0.5 \times \begin{matrix} i & 1 \\ \sigma_0^2 & 0 \\ \sigma_0^2 & 0 \\ \sigma_0^2 & 0 \end{matrix} + 0.5 \times \begin{matrix} i & 1 \\ \sigma_0^2 & 0 \\ 0 & \sigma_0^2 \\ 0 & \sigma_0^2 \end{matrix}$$

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

28

37

A data-adaptive mixture model

Instead of making assumptions, can we **learn from data**:

- What are the latent structures for multivariate effects?
- How often does each structure appear?

and use these to construct the mixture model?

29

38

Patterns of sharing: factor analysis

Decomposing effect estimates, $\hat{B} = LF + E$

Figure: Sparse factor analysis of GTEx data

30

39

Incorporating all possible patterns

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

$$U_{mixe} = m_1 \times \begin{matrix} i & 1 \\ 2.4 & 0.3 \\ 0.3 & 1.5 \end{matrix} + m_2 \times \begin{matrix} i & 1 \\ 1.6 & 0.001 \\ 0.001 & 0.02 \end{matrix} + m_3 \times \dots$$

This is the Multivariate Adaptive Shrinkage Prior.

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

31

40

Multivariate effect size sharing in eQTLs

Figure: Quantitative characterization of eQTL effects heterogeneity in GTEx

32

41

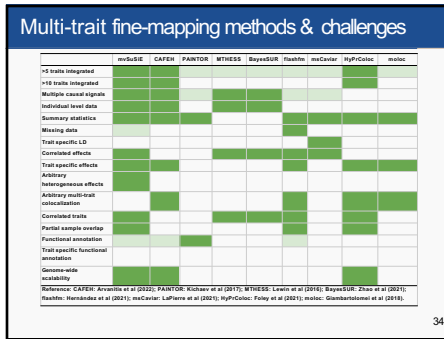
Application to multivariate fine-mapping

Figure: mvSuSIE fine-mapping with adaptive shrinkage model

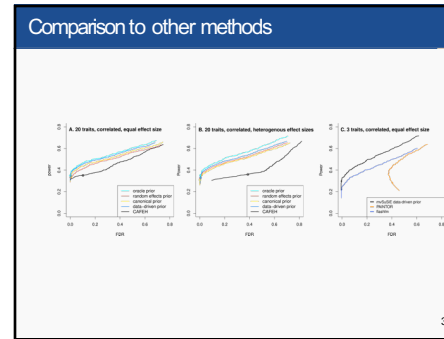
Zou et al. (2023) bioRxiv

33

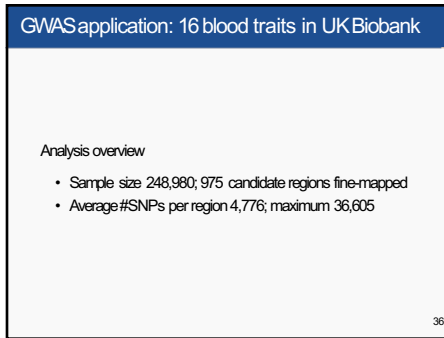
42



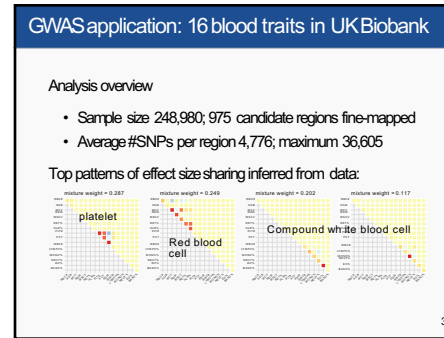
43



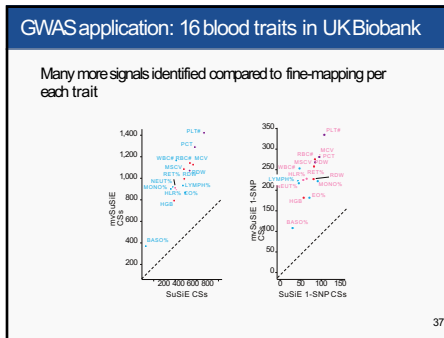
44



45



46



47

Complex phenotype prediction and transcriptome-wide association studies

Gao Wang, Ph.D.
 Advanced Gene Mapping Course, May 2024
 The Gertrude H. Sergievsky Center and Department of Neurology
 Columbia University Vagelos College of Physicians and Surgeons

1

Motivation: eQTLs are enriched in GWAS signals

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

2

Transcriptome-wide association study (TWAS)

Contributions of multiple genetic variants to complex traits through their impact on molecular phenotypes

Figure: Gusev et al. (2016) Nat. Genet.

3

TWAS challenge: association vs causality

Figure: Gusev et al. (2016) Nat. Genet.

4

TWAS challenge: association vs causality

Figure: Gusev et al. (2016) Nat. Genet.

5

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 applications
- Multiple tissue and cell-types
- Availability of individual level data vs summary statistics

6

TWAS methods overview

PrediXcan (2015): First TWAS, use logistic risk as prediction model.
TWAS (2016): Use SLDMM as prediction model.
DPR (2017): Use non-parametric DPR as prediction model.
CoMM (2018): The first likelihood-based inference.
UTMOST (2019): Multi-tissue two-stage analysis.
PMR (2020+): Accommodate horizontal pleiotropic effects.

Figure: Zhu and Zhou *et al.* (2021) Quantitative Biology

7

Univariate TWAS methods overview

$$Y = \sum_{k=1}^M \beta_k X_k + \epsilon$$

Univariate Regression leads to **GWAS** with loss function $\|Y - X_k \beta_k\|_2$.
Penalized regression includes **Ridge**, **LASSO**, and **Elastic Net** with loss function $\|Y - \sum_k X_k \beta_k\|_2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta_2\|_2$.

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

8

Simple regression method

LETTERS

Common polygenic variation contributes to risk of schizophrenia and bipolar disorder
The International Schizophrenia Consortium*

$$Y = \sum_{k=1}^M \beta_k^{GWAS} X_k$$

9

Ridge regression / BLUP

REPORT

GCTA: A Tool for Genome-wide Complex Trait Analysis
Jian Yang,^{1*} S. Hong Lee,¹ Michael E. Goddard,^{2,3} and Peter M. Visscher¹
AJHG 2011

$$Y = \sum_{k=1}^M \beta_k^{Ridge} X_k$$

Loss function: $\|Y - \sum_k X_k \beta_k\|_2 + \lambda_2 \|\beta_2\|_2$

10

Other penalized regression

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

Regularization and variable selection via the elastic net
Hui Zou and Trevor Hastie
Stanford University, USA

Loss function: $\|Y - \sum_k X_k \beta_k\|_2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta_2\|_2$

11

Bayesian variable selection regression

OPEN ACCESS Freely available online PLOS GENETICS

Polygenic Modeling with Bayesian Sparse Linear Mixed Models
Xiang Zhou^{1*}, Peter Carbonetto¹, Matthew Stephens^{1,2*}

$$Y = \sum_{k=1}^M \beta_k^L X_k + \sum_{k=1}^M \beta_k^S X_k + \epsilon$$

$$\beta_k^L \sim N(0, \sigma_L^2)$$

$$\beta_k^S \sim N(0, \sigma_S^2)$$

MultBLUP: 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

12

Choice of methods: cross validation

TWAS / FUSION

Functional Summary-based Imputation

- Need RIN6 (Grubic et al.) models for TSS ATAC-seq
- Need CONTEXT (Thompson et al.) context-specific models for single-cell and bulk expression
- Need cTEU v8 models

FUSION is a suite of tools for performing transcriptome-wide and region-wide association studies (TWAS and RWAS). FUSION builds predictive models of the genetic component of a functional/molecular phenotype and tests that component for association with disease using GWAS summary statistics. The goal is to identify associations between a GWAS phenotype and a functional phenotype that was only measured in reference data. We provide precomputed predictive 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

13

13

Likelihood based approach

Figure: CoMM, Yeung et al. (2019)

Also see Yuan et al. (2022) likelihood based Mendelian Randomization

14

14

Multivariate TWAS methods overview

Leverage similarity between molecular phenotypes

- UTMOST, Yu et al. (2019) Nature Genetics
- MR-JTI, Zhou et al. (2020) Nature Genetics
- mr.mash, Morgante et al. (2023) PLoS Genetic (to appear)

16

15

Multivariate TWAS method:

16

16

An omnigenic view of genetic regulations

Figure: Liu et al. (2019) Cell

17

17

Multi-omic Strategies for TWAS

Mediator-enriched TWAS

A. MeTWAS scheme

1. Model mediators M_1, \dots, M_m SNPs local to mediators X_{M_1}, \dots, X_{M_m} Intensities of mediators M_1, \dots, M_m Incidence or prevalence of trait
2. Find imputed $\hat{M} = (M_1, \dots, M_m)$ mtRNA expression of gene Y_G
3. Model Y_G with fixed \hat{M} and random X_G
4. TWAS test of association

Figure: Bhattacharya et al. (2020) PLoS Genet.

19

18

Multi-omic Strategies for TWAS

Distal-eQTL Prioritization via Mediation Analysis

B. DePMA scheme

1. If $H_0: TME = \alpha \beta w = 0$ is rejected

2. Estimate w_j with full X_d

3. TWAS test of association

Append X_d with X_g

SNPS local to gene G X_g

Distal eQTL X_d

Mediators local to G M_1, \dots, M_n

mRNA expression of gene G X_g

Incidence or prevalence of trait

Figure: Bhattacharya *et al.* (2020) PLoS Genet.

19

Multi-omic Strategies for TWAS

C. Example biological mechanism leveraged by MOSTWAS

Regulation via TF binding, chromatin or methylation state

Mediation of distal eQTL

Regulatory element

Gene X local to distal eSNP

Distal eSNP

mRNA of gene G

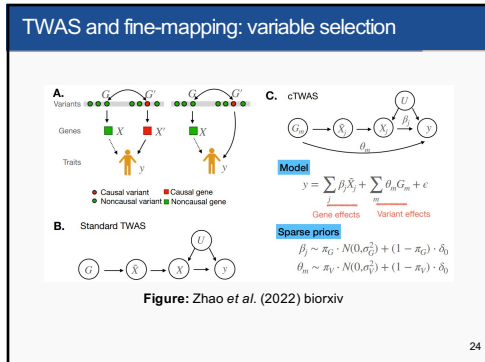
Regulatory element

eGene G under study

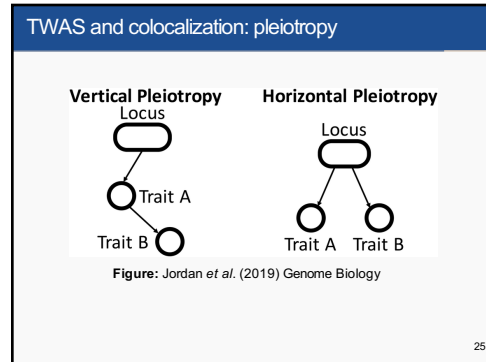
Figure: Bhattacharya *et al.* (2020) PLoS Genet.

20

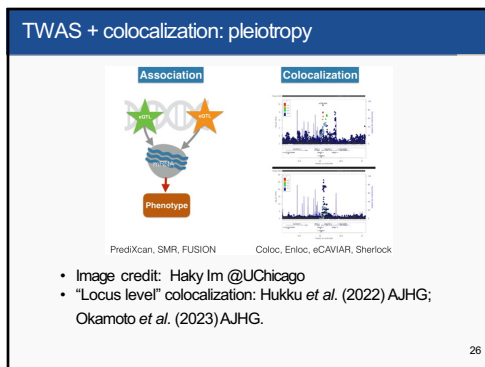
Deep learning to predict molecular traits



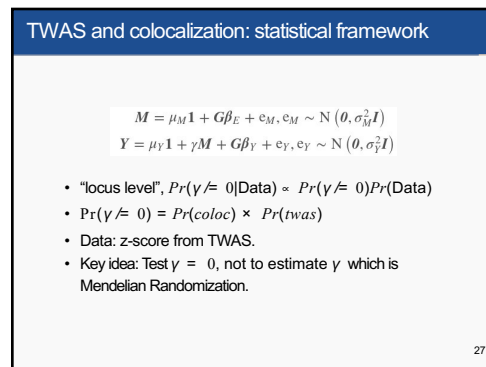
25



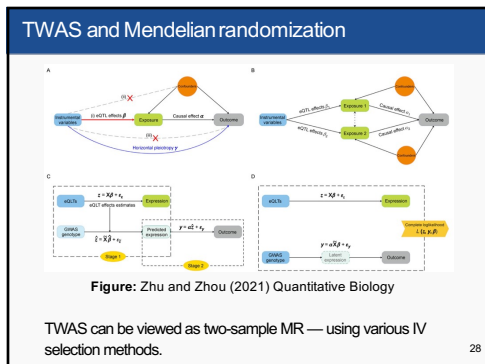
26



27



28




29

Genotype Pattern Mining For Digenic Traits

Advanced Gene Mapping Course, April 2024

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



1

Research Interests



Development of analysis methods for genetic data, genetic linkage and association analysis.
 Current topics: Digenic disease mapping, disease prediction based on genotype patterns.
 Implementation in computer programs, dissemination on website
 Collaboration with researchers world-wide on their data
 Recent publications: [1-7] #1 now freely available from <https://github.com/jurgott/handbook>

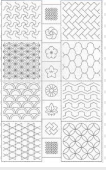
1. Terwilliger, J.D. and Ott, J. (1994) *Handbook of human genetic linkage* Johns Hopkins University Press
2. Horppanen, S. et al. (2020) Slated genomic segment analysis with equivalence testing. *Gene Epidemiol* 44, 741-747. DOI:10.1002/gepi.12235
3. Okazaki, A. et al. (2020) Population genetics: past, present, and future. *Human genetics*, 1-10. DOI:10.1007/s00439-020-02208-5
4. Okazaki, A. et al. (2021) Genotype pattern mining for pairs of interacting variants underlying digenic traits. *Genet* 12, 1161. DOI:10.3390/genet12081160
5. Okazaki, A. and Ott, J. (2022) Machine learning approaches to explore digenic inheritance. *Trends Genet*. DOI:10.1016/j.tig.2022.04.009
6. Ott, J. and Park, T. (2022) Overview of frequent pattern mining. *Genomics Inform* 20, e39. DOI:10.5898/gi.22074
7. Zhang, Q. et al. (2023) A multi-threaded approach to genotype pattern mining for detecting digenic disease genes. *Front Genet* 14, 1222517. DOI:10.3389/fgene.2023.1222517

Ott "Genotype Patterns" 2

2

Topics

- Science develops independently in different fields
 - Frequent Pattern Mining
 - Human gene mapping
- Mining consumer databases
 - The *Apriori* algorithm (30 years ago)
 - Newer algorithms: *eclat*, *fpgrowth*
- Case-control association analysis
 - GWAS: Main effects in genetic association studies
 - Digenic traits (20 years ago)
 - MDR, Multifactor Dimensionality Reduction (20 years ago)
 - Differences in interaction between cases and controls
 - *AprioriGWAS* (10 years ago)
 - Newest approach, *Gpairs* program
 - Analysis of AMD dataset



Ott "Genotype Patterns" 3

3

Frequent Pattern Mining

<https://www.philippe-fournier-viger.com/spmf/>

- Thirty years ago, supermarkets started collecting huge amounts of consumer data at their cashiers. Consumer habits – if someone buys bread and milk, how likely will they also buy wine?
- **Apriori algorithm** (Agrawal et al, *ACM SIGMOD Conference on Management of Data* 1993, 207-216). Efficient search for frequent sets of items ("itemsets", patterns) purchased by a consumer ("transaction"). (1) Development of **association rules**, that is, conditional probabilities $P(Y|X)$, with Y and X being items or itemsets. (2) **Apriori property**: "If an itemset is infrequent, all its supersets will be infrequent". Recursive search for longer patterns.
- Research published in conference proceedings, less so in traditional journals.
- Other implementations of search algorithms, e.g. *fpgrowth* (written in C) (<https://borgelt.net/software.html>), SPMF (in java). Huge memory demands.

Ott "Genotype Patterns" 4

4

Digenic Traits

Ming & Muenke (2002) *Am J Hum Genet* 71, 1017 (review)
 Schaffer A (2013) *J Med Genet* 50, 641-52 (review)

EFFECT AND PHENOTYPE	GENE 1		GENE 2	
	Mutation	Phenotype	Mutation	Phenotype
Synergistic:				
RP	<i>ROM1</i> ^{C104delC}	Normal	<i>RDS</i> ^{A183P}	Normal
RP	<i>ROM1</i> ^{V114delG}	Normal	<i>RDS</i> ^{A183P}	Normal
Bardet-Biedl	<i>BBS2</i> ^{288G>S}	Normal	<i>BBS6</i> ^{951A>T}	Normal
Deafness	<i>GJB2</i> ^{D36AG}	Normal	<i>GJB6</i> ^{47C}	Normal
Deafness	<i>GJB2</i> ^{D36AGT}	Normal	<i>GJB6</i> ^{47C}	Normal
Hirschsprung	<i>RET</i> ^{946T}	Normal	<i>EDNRB</i> ^{R319S}	Normal
Severe insulin resistance	<i>PPARG</i> ^{V253M>A>A>T}	Normal	<i>PP1R3A</i> ^{V191M>A>M>G}	Normal
Modifier:				
Juvenile-onset glaucoma	<i>MYOC</i> ^{S139V}	Adult-onset glaucoma	<i>CYP11B1</i> ^{R36H1}	Normal
Usher 1	<i>USH1B</i> ^{909delG}	Usher 3	<i>MYO7A</i> ^{G46E>G46E>12D}	Normal
Congenital nonlethal JEB	<i>COL17A1</i> ^{G1128A>S15X}	Juvenile JEB	<i>LAMB3</i> ^{R1665S}	Normal
More severe ADPKD	<i>PKD1</i> ¹⁰⁴⁴	Less severe ADPKD	<i>PKD2</i> ^{G1526A>A}	Less severe ADPKD
More severe hearing loss	<i>DFNA1</i>	Mild hearing loss	<i>DFNA2</i>	Mild hearing loss
WS2/OA	<i>MITF</i> ^{R646A}	WS2	<i>TYR</i> ^{R362Q}	Normal
More severe WS2/OA	<i>MITF</i> ^{R646A}	WS2	<i>TYR</i> ^{R362Q>R402Q}	Normal

Ott "Genotype Patterns" 5

5

Genetic Interactions between Variants

Okazaki & Ott (2022) *Trends in Genetics* 38 (10):1013-1018; DOI:10.1016/j.tig.2022.04.009

1. Traditionally, disease association has been carried out at 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 degree 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" 6

6

Finding disease-associated pairs of variants or 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
- Zhang Q, Long Q, Ott J. **AprioriGWAS**, *PLoS Comput Biol*. 2014;10(6):e1003627
Apriori applied to GWAS: 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.
- 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

Ott "Genotype Patterns"

7

7

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 pre-computed 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"

8

8

Gpairs program: All pairs of genotypes, schizophrenia data

<https://lab.rockefeller.edu/ott/programs/GPM>
https://github.com/jureott/gpm_prog

- Schizophrenia case-control data: 1,044 cases and 2,052 controls genotyped for 892,850 SNPs. Pruned and focused on males.
- Evaluate all pairs of genotypes for SNPs. For each SNP pair, analyze each of the 9 genotype pairs: **81,972,176,883** genotype pairs tested.
Distribute work over many threads (CPUs, up to 192 CPUs in new PCs).
For each genotype pair, X , make 2×2 table:
- Min. 20 occurrences of any genotype pair (support)
- Each table analyzed by Fisher test
- $p_{\text{obs}} = \min(\# \text{tests} \times p_{\text{nom}}, 1)$
- 69 genotype pairs significantly more frequent in cases than controls
- Genotypes \rightarrow variants \rightarrow genes: Network of 17 genes
- **Prediction**, classification: $c = 0 \rightarrow$ person with X must be a case!

Phenotype, Y	No. of individuals	
	With X	Without X
Affected, "case"	a	b
Unaffected, "control"	c	d

Ott "Genotype Patterns"

9

9

Prediction vs. Significance

- Presence of a given genotype pair, X , as an indicator of disease

Phenotype, Y	No. of individuals	
	With X	Without X
Affected, "case"	a	b
Unaffected, "control"	c	d

- Given a "case", what is the probability the test is significant?
Power = sensitivity = $a / (a + b)$
- Given presence of X in an individual, what is the probability that individual is a case? **Positive predictive value, PPV** = $a / (a + c)$, also called *confidence* in machine learning.
- Lo et al (2015) *Why significant variants aren't automatically good predictors*, PNAS 112 (45). DOI: 10.1073/pnas.1518285112

Ott "Genotype Patterns"

10

10

Cross-Validation

- Estimates of PPV for the same data that furnished the predictions \rightarrow tend to be too good
- Solution: Develop predictors in a set of data and apply the predictors to a new set of data.
- Same data: Build model in 90% of the data and apply resulting predictors to 10% of the data \rightarrow 10-fold **cross-validation**
- Better approach [1]: Leave-one-out method, *L1out*. Remove i -th individual from data and develop predictor \rightarrow apply to i -th individual. Do this for all individuals.
- Implement *L1out* (1) for *Gpairs* and (2) for polygenic risk score, *PRS*, as implemented in *plink* with the `--score` function.
- 1. Agresti (2019) *An introduction to categorical data analysis*. Wiley, Hoboken NJ

Ott "Genotype Patterns"

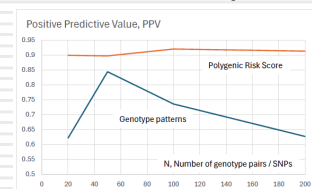
11

11

Decision Rules

DeWan et al, ...wet age-related macular degeneration. *Science*, 19 Oct 2006. DOI: 10.1126/science.1133807

- *Gpairs*: For a number N of best predicting genotype pairs, call an individual a "case" if she/he carries 20+ of such genotype pairs.
- *Polygenic Risk Score*: For N best predicting variants, call an individual a "case" if she/he has a score above the 95th percentile of controls.



Ott "Genotype Patterns"

12

12

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
Yale School of Public Health

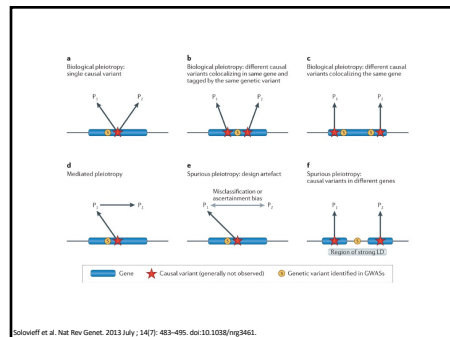
Yale SCHOOL OF PUBLIC HEALTH

1

Pleiotropy

- Phenomenon in which a genetic locus affects more than one trait or disease
- Molecular level
 - Single gene with multiple physiological function
 - Two domains of a single gene product with different functions and affecting multiple phenotypes
 - Gene product with a single function that affects multiple phenotypes acting in multiple tissues
- Statistical level
 - A locus displaying cross-phenotype associations is often considered pleiotropic
 - Can be at the variant, gene or region level

2



3

Early example of "pleiotropy"

Gregor Mendel documented one of the earliest examples of pleiotropy in his pea plant experiments

Violet flowers

- Seed coats = brown-grey
- Axils = red and spotted

White flowers

- Seed coats = white
- Axils = white and unspotted

Mendel J. G. 1866 Experiments in plant hybridization. Verhandlungen des naturforschenden Vereines in Brünn 4, 3-47 (in German).

4

Examples in humans

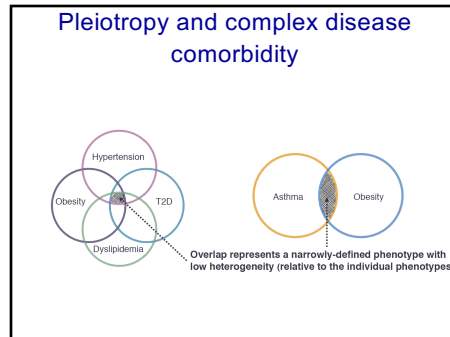
- Marfan syndrome
 - FBN1 (fibrillin-1)
 - thinness, joint hypermobility, limb elongation, lens dislocation, and increased susceptibility to heart disease.
- Holt-Oram syndrome,
 - TBX5 (transcription factor)
 - cardiac and limb defects
- Nijmegen breakage syndrome
 - NBS1 (DNA damage repair protein)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

5

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

6



7

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

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

Abundant Pleiotropy in Human Complex Diseases and Traits

Shanya Sivakumaran,^{1,6} Felix Agakov,^{1,5,6} Evropti Theodorou,^{1,6} James G. Prendergast,³ Lina Zgaga,^{1,4} Feri Manolio,³ Igor Rudan,¹ Paul McKie,¹ James F. Wilson,¹ and Harry Campbell^{1,*}

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

Disease Class	Genes			SNPs		
	Pleiotropic (%)	Nonpleiotropic (%)	p Value*	Pleiotropic (%)	Nonpleiotropic (%)	p Value*
All (comparison group)	233 (6.9)	1147 (83.1)	–	77 (4.0)	1610 (95.4)	–
Immune-mediated phenotypes	106 (37.7)	175 (62.3)	<0.0001	31 (8.3)	343 (91.7)	0.0066
Cancer	49 (24.8)	152 (75.2)	<0.0001	8 (4.8)	158 (95.2)	0.8256
Metabolic syndromes	79 (28.5)	199 (71.5)	<0.0001	30 (8.4)	327 (91.6)	0.0056

* Fisher's exact test p value.

10

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

11

A practitioners' guide for studying pleiotropy in genetic epi studies

doi:10.1093/aje/kwz028 (first of two)

Statistical Analysis of Multiple Phenotypes in Genetic Epidemiological Studies: From Cross-Phenotype Associations to Pleiotropy.

Bellizzi KS, Mora C, Dettler AC

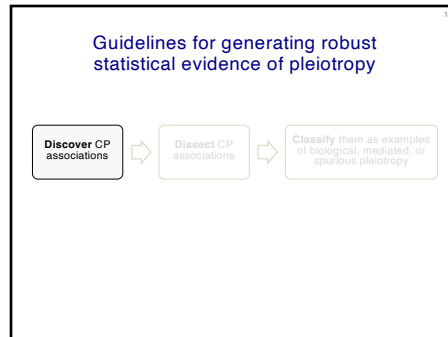
Abstract

In the context of genetics, pleiotropy refers to the phenomenon in which a single genetic locus affects more than one trait or disease. Genetic epidemiological studies have identified loci associated with multiple phenotypes, and these cross-phenotype associations are often incorrectly interpreted as examples of pleiotropy. Pleiotropy is only one possible explanation for cross-phenotype associations. Cross-phenotype associations may also arise due to issues related to study design, confounder bias, or non-genetic causal links between the phenotypes under analysis. Therefore, it is necessary to dissect cross-phenotype associations carefully to uncover true pleiotropic loci. In this review, we describe statistical methods that can be used to identify robust statistical evidence of pleiotropy. First, we provide an overview of univariate and multivariate methods for discovery of cross-phenotype associations and highlight important considerations for choosing among available methods. Then, we describe how to dissect cross-phenotype associations by using mediation analysis. Pleiotropic loci provide insights into the mechanistic underpinnings of disease comorbidity, and may serve as novel targets for interventions that simultaneously treat multiple diseases. Discovering between different types of cross-phenotype associations is necessary to realize the public health potential of pleiotropic loci.

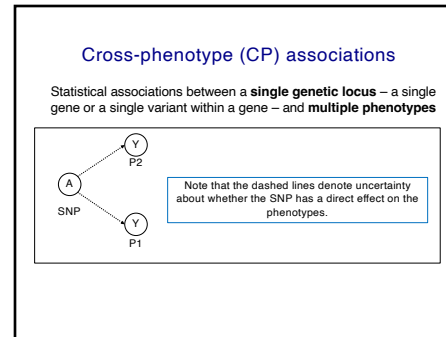
© The Author(s) 2017. Published by Oxford University Press on behalf of the Johns Hopkins Bloomberg School of Public Health. All rights reserved. For permissions, please email: journals.permissions@oup.com.

KEYWORDS: genetic epidemiology, mediation analysis, pleiotropy

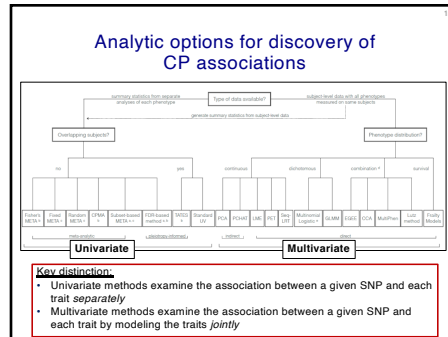
12



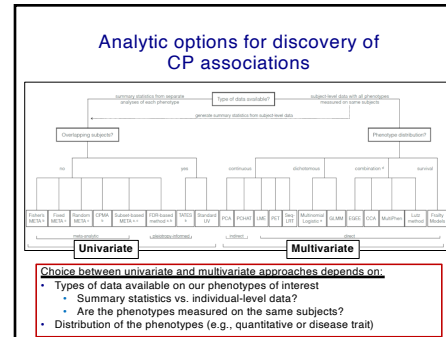
13



14



15



16

Univariate methods are by far the most commonly used to detect CP associations

- Univariate methods include (but are not limited to) the methods you've discussed in class so far:
 - allelic Chi-Square test
 - genotypic Chi-Square test
 - regression-based methods
- 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

17

Hypothetical example: Discovery of CP associations for hypertension and heart disease by using logistic regression

Step 1. Fit two univariate regression models within PLINK

$$E[\text{hypertension}] = \beta_0 + \beta_1 * SNP$$

$$E[\text{heart disease}] = \beta_0 + \beta_1 * SNP$$

Word of caution: The univariate tests of association should be **marginal** tests (conducted irrespectively of the second phenotype) NOT **conditional** tests (conducted on a subset defined based on absence/presence of the second phenotype). In this example, what that means is that the regression for hypertension should be fit on all subjects *irrespectively* of their heart disease status; and the regression for heart disease should be fit on all subjects *irrespectively* of their hypertension status. More on this later!

18

Hypothetical example: Discovery of CP associations for hypertension and heart disease by using logistic regression

Step 1. Fit two univariate regression models within PLINK

$$E[\text{hypertension}] = \beta_0 + \beta_1 * \text{SNP}$$

$$E[\text{heart disease}] = \beta_0 + \beta_1 * \text{SNP}$$

Step 2. For a given SNP, examine p-values for β_1 from each model.

- P-value for β_1 in hypertension model = 1.03×10^{-12}
- P-value for β_1 in heart disease model = 6.02×10^{-9}

Step 3. Declare CP associations at a given SNP, if the p-values for β_1 in each model surpass the study significance threshold.

- Assuming the standard GWAS significance threshold ($\alpha=5 \times 10^{-8}$), there is a statistically significant association with both hypertension and heart disease at this particular SNP. Therefore, we have sufficient statistical evidence to declare a CP association at this SNP.

19

Using multivariate methods to increase the power to detect cross-phenotype associations

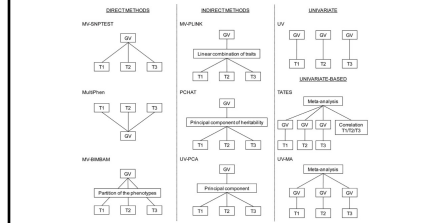
20

A Comparison of Multivariate Genome-Wide Association Methods

Tessal E. Galeano¹, Kristel van Steen¹, Lambertus A. M. Klenneman^{2,3}, Luc L. Jans^{1,4},
Eike M. Veenendaal^{1,5,6}

¹Department for Health Evidence, Radboud university medical center, Nijmegen, The Netherlands, ²Genetics and Health Unit, Murdoch University, University of Waikato, Hamilton, New Zealand, ³Department of Health, Behavior and Society, Johns Hopkins University, Baltimore, Maryland, ⁴Department of Genetic Epidemiology, Radboud university medical center, Nijmegen, ⁵Department of Molecular Biology and Genetics, Aarhus University, Aarhus, Denmark, ⁶Department of Human Genetics, Radboud university medical center, Nijmegen

PLoS ONE | www.plosone.org April 2014 | Volume 9 | Issue 4 | e95923



21

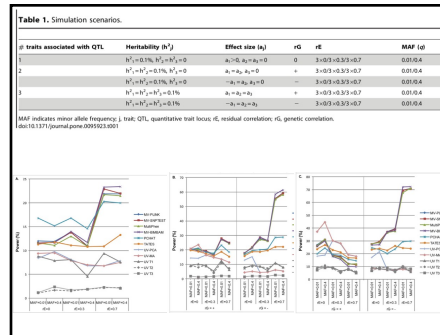
A comparison of univariate and multivariate GWAS methods for analysis of multiple dichotomous phenotypes

Yasmmyn D. Salinas¹, Andrew T. DeWan¹, and Zuoheng Wang²

¹Department of Chronic Disease Epidemiology, ²Department of Biostatistics, Yale School of Public Health, Yale University, 60 College St, New Haven, Connecticut, USA

Genet. Epidemiol. 41 (7), 689-689

23

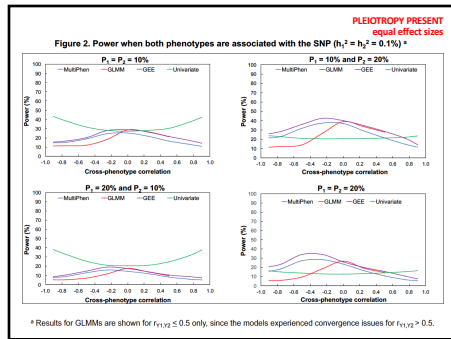


22

Simulation scenarios

# traits associated	h^2	$r_{1,2}$	P_1
1	$h_1^2 = 0.1, h_2^2 = 0\%$	[-0.9, 0.9]	$P_1 = P_2 = 10\%$
			$P_1 = P_2 = 20\%$
			$P_1 = 10\%, P_2 = 20\%$
2	$h_1^2 = h_2^2 = 0.1\%$	[-0.9, 0.9]	$P_1 = 20\%, P_2 = 10\%$
			$P_1 = P_2 = 20\%$
			$P_1 = 10\%, P_2 = 20\%$
2	$h_1^2 = 0.1\%, h_2^2 = 0.05\%$	[-0.9, 0.9]	$P_1 = P_2 = 10\%$
			$P_1 = P_2 = 20\%$
			$P_1 = 10\%, P_2 = 20\%$

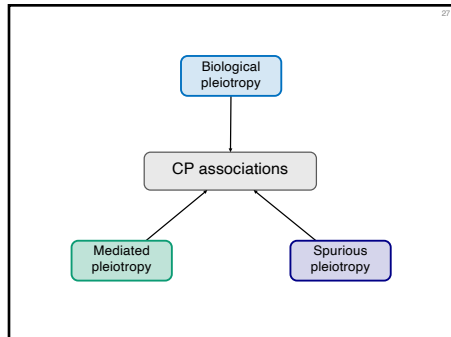
24



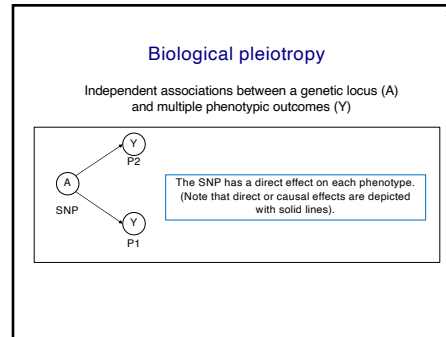
25

Problem: CP associations need not be indicative of pleiotropy

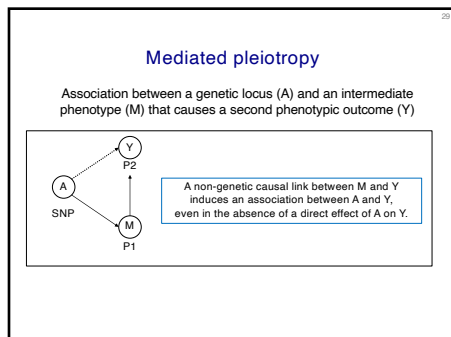
26



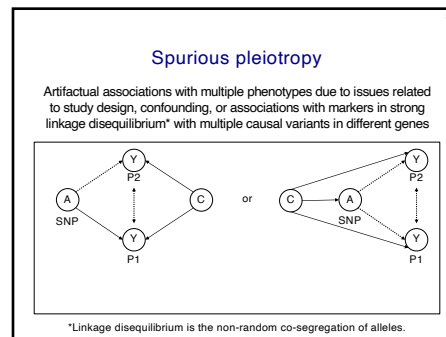
27



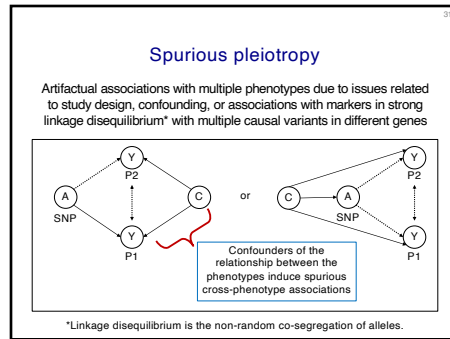
28



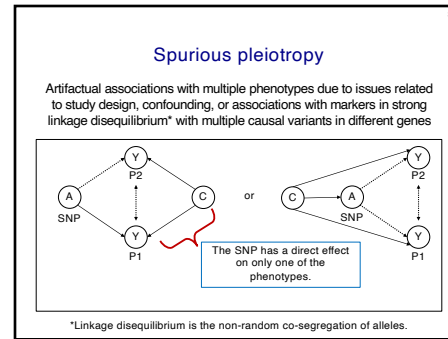
29



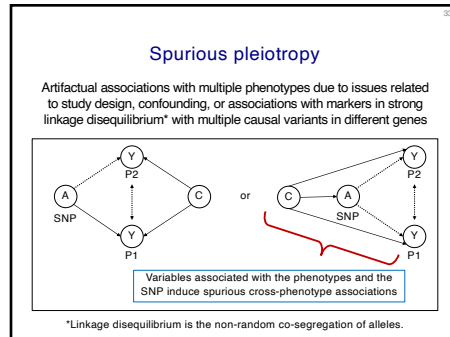
30



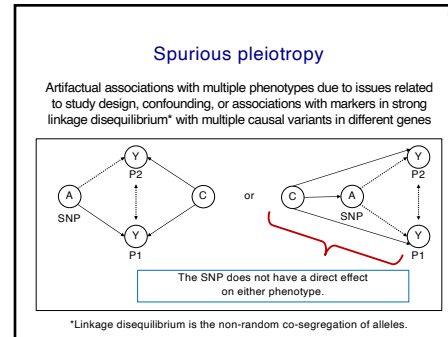
31



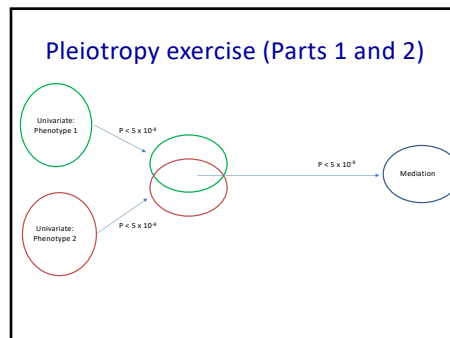
32



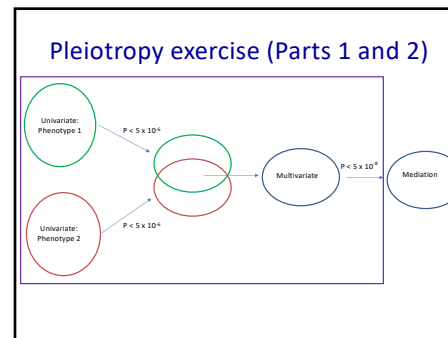
33



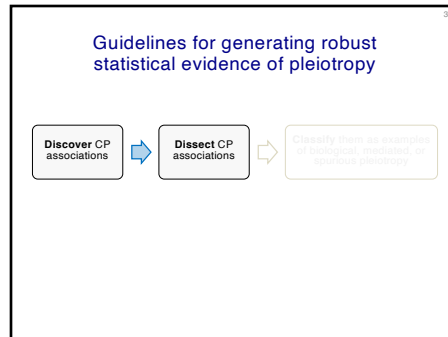
34



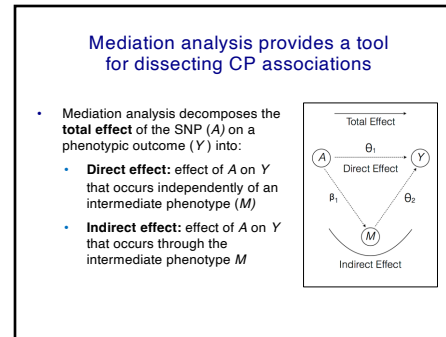
35



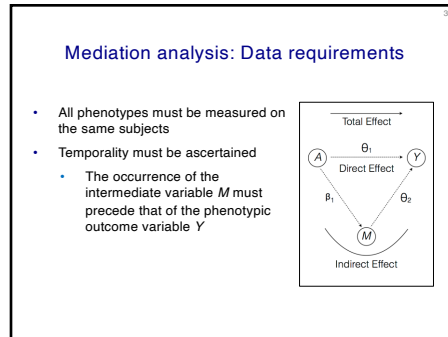
36



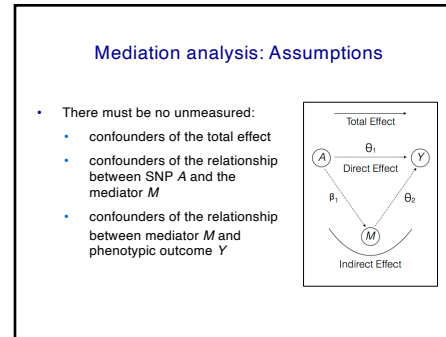
37



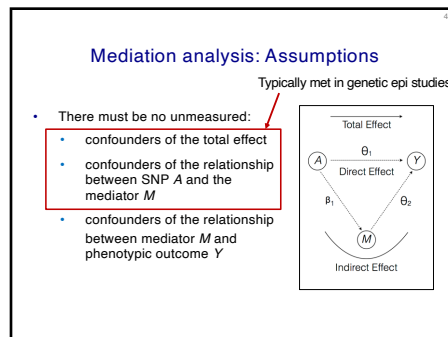
38



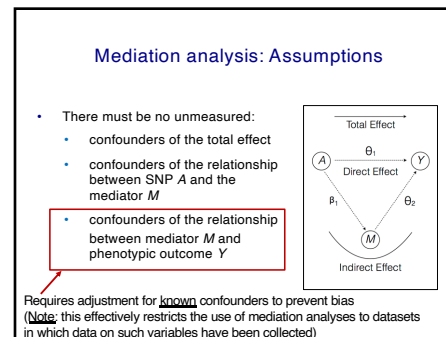
39



40



41



42

Mediation analysis: Regression-based approach

- Requires fitting two regression models, one for mediator M and one for phenotypic outcome Y :
 - $E[M | a, c] = \beta_0 + \beta_1 a + \beta_2' c$
 - $E[Y | a, m, c] = \theta_0 + \theta_1 a + \theta_2 m + \theta_4' c$

Assesses the effect of A on M , while controlling for measured confounders (C)

43

Mediation analysis: Regression-based approach

- Requires fitting two regression models, one for mediator M and one for phenotypic outcome Y :
 - $E[M | a, c] = \beta_0 + \beta_1 a + \beta_2' c$
 - $E[Y | a, m, c] = \theta_0 + \theta_1 a + \theta_2 m + \theta_4' c$

Assesses the effect of A on Y , while controlling for both M and C

44

Mediation analysis: Regression-based approach

- Requires fitting two regression models, one for mediator M and one for phenotypic outcome Y :
 - $E[M | a, c] = \beta_0 + \beta_1 a + \beta_2' c$
 - $E[Y | a, m, c] = \theta_0 + \theta_1 a + \theta_2 m + \theta_4' c$
- The parameter estimates from these models (namely β_1 , θ_1 , and θ_2) are used to estimate the direct and indirect effects

45

Guidelines for generating robust statistical evidence of pleiotropy

Discover CP associations → Dissect CP associations → Classify them as examples of biological, mediated, or spurious pleiotropy

46

Mediation analysis: Interpretation

- Mediated pleiotropy**
 - Complete mediation:** SNP A is associated with mediator M and the total effect of A on phenotypic outcome Y is equal to its indirect effect (i.e., the direct effect is equal to 0).
 - Incomplete mediation:** SNP A is associated with mediator M and A has both direct and indirect effects on phenotypic outcome Y (i.e., the total effect is equal to the sum of the direct and indirect effects)
- Biological pleiotropy**
 - SNP A is associated with mediator M , and the total effect of SNP A on phenotypic outcome Y is equal to its direct effect (i.e., the indirect effect is equal to 0)

47


Mediation analysis: Interpretation

- Mediated pleiotropy**
 - Complete mediation:** SNP A is associated with mediator M and the total effect of A on phenotypic outcome Y is equal to its indirect effect (i.e., the direct effect is equal to 0).
 - Biological pleiotropy**
 - SNP A is associated with mediator M , and the total effect of SNP A on phenotypic outcome Y is equal to its direct effect (i.e., the indirect effect is equal to 0)
 - Incomplete mediation:** SNP A is associated with mediator M and A has both direct and indirect effects on phenotypic outcome Y (i.e., the total effect is equal to the sum of the direct and indirect effects)

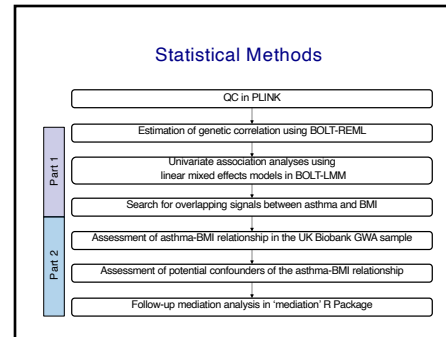
48

Phenotype definitions

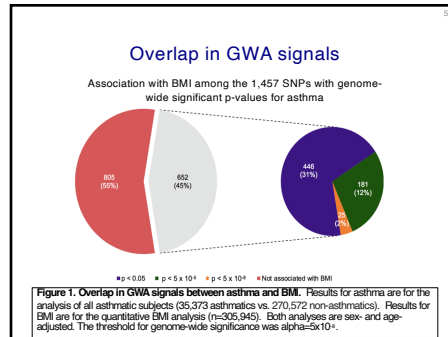
- BMI at baseline (kg/m²):
 - calculated based on height and weight measurements collected by trained UK Biobank staff at the recruitment sites
- Asthma diagnosed prior to baseline (yes/no):
 - ascertained via the question "Has a doctor ever told you that you had asthma?"
 - Note:** In mediation analyses,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



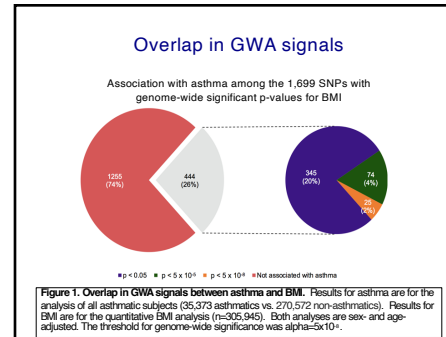
55



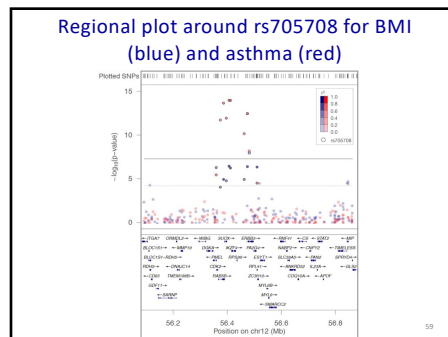
56



57



58



59

Cross-phenotype associations in 12q13.2

Table 2. Cross-phenotype associations in 12q13.2*

SNP	Gene	BP	Effect/reference allele	EAF	OR (95% CI)	P	OR (95% CI)	P
rs2094688	CDK2	56,366,321	G/A	0.3988	1.04 (1.02, 1.06)	3.50e-10†	20.06 (0.08, -0.04)	5.40e-10†
rs177814	EBF3	56,376,217	C/G	0.4217	1.06 (1.04, 1.08)	2.40e-10†	20.05 (0.07, -0.02)	7.50e-10†
rs705702	RSOX	56,390,636	G/A	0.3376	1.07 (1.05, 1.09)	3.10e-10†	20.05 (0.08, -0.03)	1.10e-10†
rs1087044†	RSOX	56,390,685	C/A	0.4279	1.06 (1.04, 1.08)	1.90e-10†	20.05 (0.07, -0.03)	1.60e-10†
rs170394	RSOX	56,412,487	G/T	0.3413	1.07 (1.05, 1.09)	1.90e-10†	20.06 (0.09, -0.04)	3.70e-10†
rs2494953	RSOX	56,416,925	C/A	0.3432	1.07 (1.05, 1.09)	1.90e-10†	20.06 (0.08, -0.04)	4.60e-10†
rs117799†	EBF3	56,470,225	C/T	0.4317	1.06 (1.04, 1.07)	8.80e-10†	20.05 (0.07, -0.03)	1.10e-10†
rs249239	EBF3	56,482,181	T/G	0.3470	1.07 (1.05, 1.09)	4.50e-10†	20.06 (0.08, -0.04)	4.20e-10†
rs705704	EBF3	56,484,013	A/G	0.4712	1.05 (1.03, 1.07)	7.20e-10†	20.06 (0.09, -0.04)	1.30e-10†
rs1171615†	EBF3	5651068	T/G	0.5109	1.04 (1.02, 1.06)	3.90e-10†	20.06 (0.08, -0.04)	4.50e-10†

Abbreviations: BP = base pair; BMI = body mass index; CI = confidence interval; EAF = effect allele frequency; OR = odds ratio; SNP = single nucleotide polymorphism

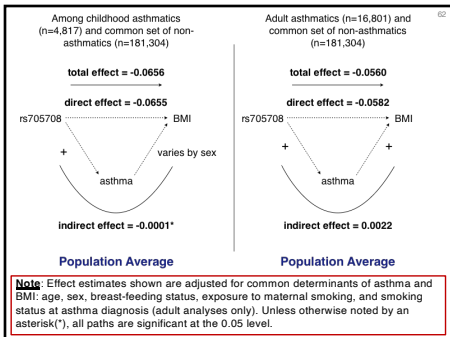
* Results shown for SNPs with $p < 5 \times 10^{-4}$ for asthma and $p < 0.05$ for BMI.
 † For non-top SNPs, the nearest gene is listed, with priority given to genes directly downstream of variant.
 ‡ P value from BOLT-LMM, derived using the standard "additive" mixed model.
 § P value from BOLT-LMM, derived using the Gaussian mixture model.

60

61

Decomposing the effect of rs705708 on BMI via mediation analysis

61



62

- 63
- ## Conclusions
- rs705708 has a positive direct effect on asthma
 - Stronger in magnitude for childhood asthma
 - rs705708 has a negative direct effect on BMI
 - Consistent in magnitude and direction in analyses including childhood vs. adult asthmatics
 - This suggests that locus 12q13.2, tagged by rs705708, has pleiotropic effects on asthma and BMI.

63

- 64
- ## Conclusions
- 12q13.2 is multigenic and our CP associations span genes *CDK2*, *RAB5*, *SUOX*, *IZK4*, *RPS26*, *ERBB3*, and *ESYT1*.
 - rs705708 is the top regional BMI signal and resides in *ERBB3*.
 - The top regional asthma signal, rs2456973, resides in *IZK4*.
 - While rs705708 and rs2456973 could be in LD with the same causative variant in either *ERBB3* or *IKZF4* or another gene in 12q13.2, it is also possible that each variant could tag a distinct, trait-specific causative variant in different genes.
 - Therefore, locus 12q13.2 displays pleiotropic effects on asthma and BMI, but this may not be an example of pleiotropy at the gene level (biological pleiotropy).

64

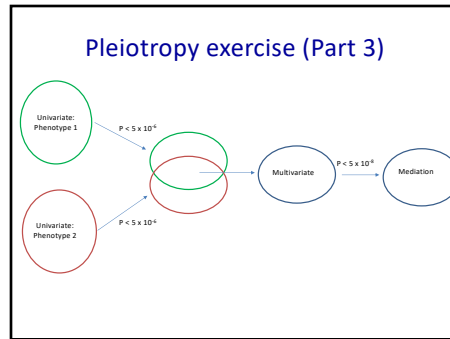
65

What if we expand this investigation to look at more phenotypes correlated with asthma?

65

- 66
- ## Asthma, T2D and anthropometric measures
- Obesity is a well-established risk factor for both asthma and T2D.
 - While highly correlated, waist circumference (WC) can provide distinct information on adiposity as it is a measure of visceral obesity, specifically WC adjusted for BMI. WC is often used in studies of chronic diseases.
 - Increased WC has been shown to be an additional risk factor for T2D and asthma even after adjusting for BMI.
 - Elevated blood glucose and T2D have been linked to increased risk of asthma in adults, and conversely, asthma has been associated with increased risk of developing T2D in adults.
 - Height is a highly heritable polygenic trait; there is evidence that shorter individuals have an increased risk for developing T2D and individuals with childhood onset asthma have shorter stature as adults compared to non-asthmatics.

66



79

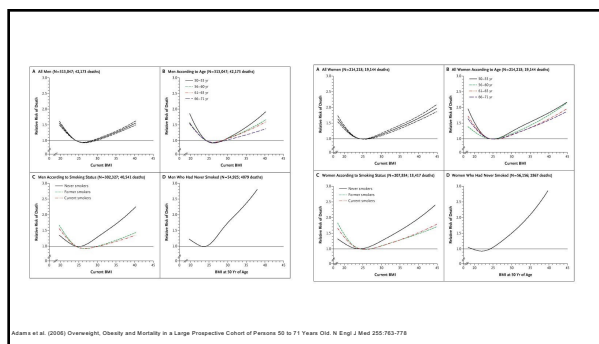
Yale

Mendelian randomization: An Introduction

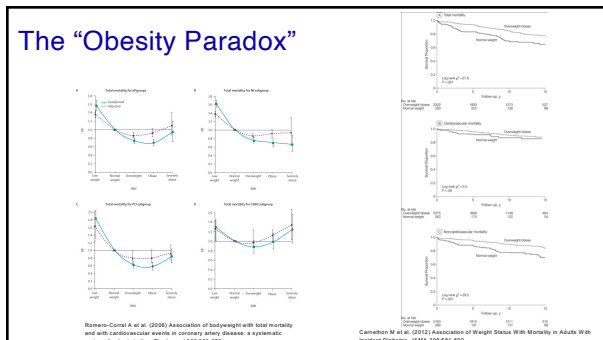
Andrew DeWan, PhD, MPH
Associate Professor of Epidemiology
Director, Yale Center for Perinatal, Pediatric and Environmental Epidemiology
Yale School of Public Health

Yale SCHOOL OF PUBLIC HEALTH

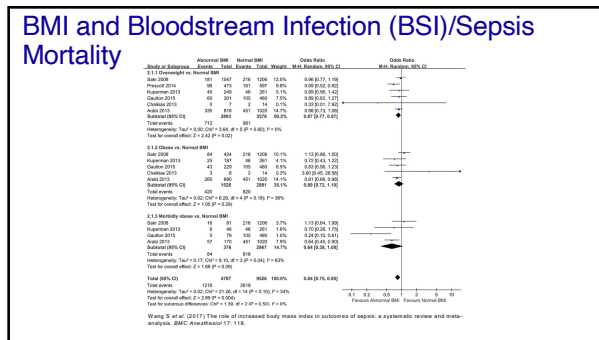
1



2



3



4

Table 2. Association of BMI and lifestyle factors with the risk of bloodstream infection among 64,037 participants in the HUNT Health Survey 1984-2001

Lifestyle variable	Age and sex adjusted				Age, sex, education and health-adjusted ^a			
	Person-years	BSI cases	RR	95% CI	Person-years	BSI cases	RR	95% CI
BMI, kg/m ²								
< 18.5	1279	15	1.00	1.00 (0.49, 2.04)	4462	9	1.00	1.00 (0.43, 2.33)
18.5-24.9	10340	107	1.00	1.00 (0.81, 1.22)	36631	368	1.00	1.00 (0.86, 1.16)
25.0-29.9	14822	107	1.00	1.00 (0.81, 1.22)	32093	308	1.00	1.00 (0.86, 1.17)
30.0-34.9	10992	126	1.20	1.12 (1.01, 1.31)	24952	274	1.20	1.12 (1.01, 1.31)
≥ 35.0	2044	19	1.47	1.36 (1.17, 1.57)	1737	18	1.47	1.36 (1.17, 1.57)
2-tailed P-value			0.00		0.00		0.00	
Smoking								
Never	36747	763	1.00	1.00 (0.90, 1.10)	17789	411	1.00	1.00 (0.89, 1.11)
Former	28820	262	1.00	1.00 (0.84, 1.18)	14854	84	1.00	1.00 (0.78, 1.27)
Current	23077	494	1.13	1.10 (1.04, 1.16)	20348	344	1.13	1.10 (1.04, 1.16)
Physical activity level ^b								
None	8144	109	1.00	1.00 (0.82, 1.21)	4696	176	1.00	1.00 (0.78, 1.27)
Light	23876	267	1.00	1.00 (0.91, 1.10)	11619	115	1.00	1.00 (0.84, 1.18)
Moderate	14899	67	1.00	1.00 (0.81, 1.22)	7474	67	1.00	1.00 (0.81, 1.22)
High	20886	213	1.24	1.18 (1.09, 1.28)	19774	217	1.24	1.18 (1.09, 1.28)
Alcohol intake ^c								
< 1 (g/day)	27422	366	1.00	1.00 (0.91, 1.10)	12619	261	1.00	1.00 (0.87, 1.14)
1-14 (g/day)	40376	403	1.00	1.00 (0.90, 1.10)	17742	111	1.00	1.00 (0.86, 1.16)
15-49 (g/day)	10621	111	1.00	1.00 (0.79, 1.26)	6759	122	1.00	1.00 (0.78, 1.27)
≥ 50 (g/day)	2012	12	1.73	1.58 (1.19, 2.16)	1346	14	1.73	1.58 (1.19, 2.16)

5

Areas of Concern (BMI/BSI as an example)

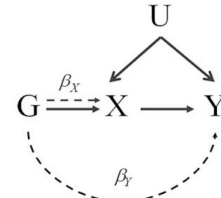
- Selection Bias: If obesity is associated with BSI risk, non-obese patients may have other characteristics that cause their BSI that in turn are more strongly associated with mortality
- Reverse Causation: if measured BMI is affected by BSI
- Confounding: if factors such as chronic diseases and smoking habits that affect both BMI and BSI mortality are not adequately adjusted

6

Mendelian randomization

- 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

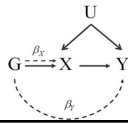
7



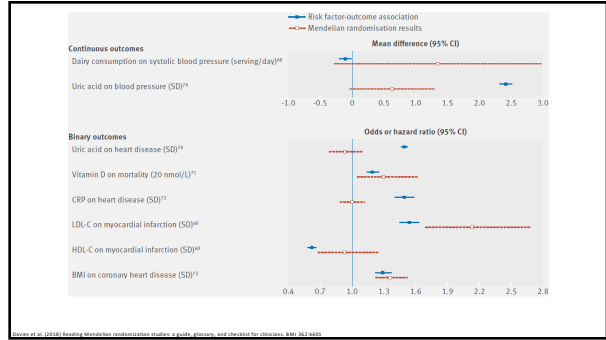
8

MR Assumptions

- The genetic instrument (G) is associated with the exposure (X)
- 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



10

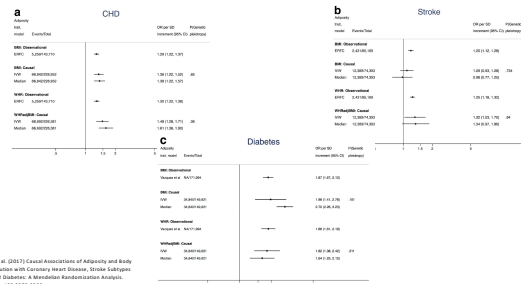
CRP and Heart Disease

Single nucleotide polymorphism	Allele frequency*	No of studies/cases/participants	Per allele higher concentration in CRP (95% CI), mg/L	Per allele higher concentration in CRP (95% CI), mg/L	Per allele risk ratio for CHD (95% CI)	Per allele risk ratio for CHD (95% CI)
rs3975077	0.06	10/15,133/96,807	0.23 (0.17 to 0.24)	0.23 (0.17 to 0.24)	0.93 (0.87 to 1.00)	0.93 (0.87 to 1.00)
rs1295	0.67	43/465,527/1,72,567	0.18 (0.16 to 0.20)	0.18 (0.16 to 0.20)	1.00 (0.98 to 1.02)	1.00 (0.98 to 1.02)
rs1130864	0.30	41/37,145/157,905	0.13 (0.12 to 0.15)	0.13 (0.12 to 0.15)	0.98 (0.96 to 1.00)	0.98 (0.96 to 1.00)
rs1809947	0.94	31/31,636/93,507	0.26 (0.23 to 0.29)	0.26 (0.23 to 0.29)	0.99 (0.94 to 1.03)	0.99 (0.94 to 1.03)

Category	Risk ratio* (95% CI) for CHD per 1 SD higher in CRP (mg/L)	Risk ratio* (95% CI) for CHD per 1 SD higher in CRP (mg/L)
Circulating usual concentrations of CRP	1.49 (1.40 to 1.59)	1.33 (1.23 to 1.43)
Adjusted for age, sex, and ethnicity		
Further adjusted†		
Genetically raised concentrations of CRPs	1.00 (0.90 to 1.13)	1.00 (0.89 to 1.12)
SNP analyses		
Haplotype analyses		

11

BMI and CHD/Stroke/Type 2 Diabetes



12

One-sample vs. two-sample designs

One-sample

- Genotype(s), risk factor and outcome all measured in the same set of study subjects
- Individual level data must be available

Two-sample

- Genotype(s) and risk factor measured in one set of study subjects and genotype(s) and outcome measured in a separate set of study subjects
- Can use summary statistics or individual level data

13

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

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

R²: proportion of the variance of the exposure explained by IV

N: sample size

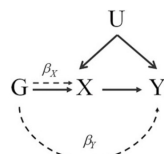
K: number of genetic variants

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

16

Pleiotropy

- Assumption that the IV is not associated with Y independently from X
- Presence of pleiotropy can bias the causal estimate
- Sensitivity analyses such as MR-Egger can be used to test whether or not the pleiotropy assumption has been violated



17

Testing MR: Wald Ratio

- Simple ratio of the effects of the instrument variable on the outcome over the instrument variable on the exposure
- Can be implemented in both one and two sample designs
 - One sample can use either a single variant or a GRS
 - Two sample design that uses multiple variants requires a method for combining Wald Ratios

$$\hat{\beta}_{IV} = \frac{\hat{\beta}_{ZY}}{\hat{\beta}_{ZX}}$$

18

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

Testing MR: Inverse variant weighted

- 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

For each variant calculate the Wald ratio:

$$\hat{\beta}_j = \frac{\hat{\gamma}_j}{\hat{\gamma}_j}$$

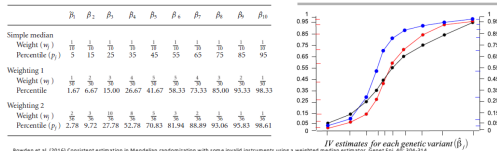
Combine into an overall estimate using a formula from meta-analysis literature:

$$\hat{\beta}_{IVW} = \frac{\sum_j \hat{\gamma}_j^2 \sigma_{Y_j}^{-2} \hat{\beta}_j}{\sum_j \hat{\gamma}_j^2 \sigma_{Y_j}^{-2}}$$

20

Testing MR: Weighted Median

- Calculate the Wald ratio for each instrument
- Select the median value according to the weighted method

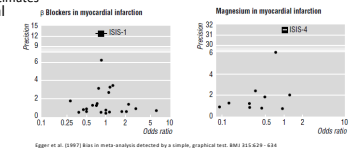


- Valid estimate when more than half of the genetic variants satisfy the IV assumptions
- No single IV contributes more than 50% of the weight

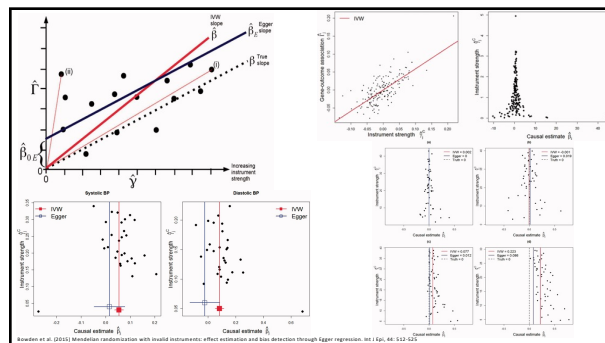
21

Testing MR: MR-Egger

- Provide a valid causal estimate in the presence of some violations of the MR assumptions (mainly pleiotropy)
- MR consisting of a single study with multiple IVs is analogous to a meta-analysis
- Bias resulting from pleiotropy is analogous to small study bias in meta-analysis
 - Small studies with less precise estimates tend to report larger estimates than big studies with more precise estimates
- Regress the standard normal deviate (odds ratio divided by its se) on the estimate's precision (inverse of the se)
 - Without bias, intercept = 0, and in the presence of bias the intercept is a measure of asymmetry



22



23

Databases and software

Data source	Description	Number of traits	Integrated with statistics package?
MR-Base	A curated database of genome-wide association study results with integrated R package for MR ²³	Over 1000	Yes
PhenoScanner	A curated database of genome-wide association study results with integrated R package for MR ²⁷	Over 500	Yes
GWAS catalog	Searchable database of genome-wide association study results ²⁸	Over 24,000	No

24

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^{4,5}, Ben M. Brumpton^{6,7,8}, Julie Paulsen⁹, Hallie C. Prescott^{10,11}, Randi M. Mørhus^{1,3}, Lise T. Gustad^{1,12}, Arne Mehl¹², Bjørn O. Asvold^{6,13}, Andrew T. DeWan¹⁴, Jan K. Damås^{1,14,15†}

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-Trøndelag 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

Table 1. Background characteristics.

Characteristic	Total population (n = 55,908)	BSI incidence (n = 2,547)	BSI death (n = 451)
Age (years) ^a	48.3 (26.5-62.3)	45.6 (32.9-71.4)	47.3 (37.7-74.5)
Male sex ^b	26,334 (47.1)	1,345 (52.8)	263 (58.3)
BMI (kg/m ²) ^c	26.3 (4.1)	27.7 (4.5)	27.9 (4.8)
Median follow-up time (years) ^d	21.1 (17.1-21.8)	13.8 (8.4-18.3)	13.3 (7.7-17.9)
Self-reported cancer ^e	3,995 (7.1)	144 (5.6)	24 (5.3)
Smoking ^f			
Never	23,394 (41.8)	876 (34.2)	156 (35.0)
Previous	15,133 (27.1)	493 (19.4)	84 (18.7)
Current	16,117 (29.4)	723 (28.0)	118 (26.9)
Physical activity ^g			
None	3,821 (7.0)	243 (11.9)	54 (15.4)
Slight	18,662 (33.3)	714 (27.9)	117 (26.3)
Moderate	17,167 (30.7)	693 (27.2)	116 (26.1)
High	13,810 (24.7)	397 (15.4)	64 (14.2)
Education ^h			
<9 years	19,033 (34.0)	1,305 (51.2)	240 (53.4)
10-12 years	23,468 (42.0)	762 (29.5)	125 (27.8)
≥13 years	10,402 (18.7)	274 (10.7)	43 (9.6)

BMI, body mass index; BSI, bloodstream infection. Data are presented as

^amean (standard deviation)

^bmedian (25th-75th percentiles), or

^cn (%). BSI incidence is based on first occurrence otherwise, last occurrence is used. Education defined as follows: <9 years ("primary school 7-10 years, continuation school, high school"), 10-12 years ("high school, intermediate school, vocational school, 1-2 years high school" and "university qualifying examination, junior college, A-levels"), and ≥13 years ("university or other post-secondary education, less than 4 years" and "university/college 4 years or more"). Activity defined as follows: none ("no light or vigorous activity"), slight ("<3 h light activity/week and no vigorous activity"), moderate ("≥3 h light activity/week or <1 h vigorous activity/week"), or high ("≥1 h vigorous activity/week")

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-Trøndelag 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 (Weng et al. (2018) *Mendelian randomization study of body mass index and risk of cardiovascular disease*. *PLoS One* 13(11): e0201111. doi:10.1371/journal.pone.0201111)
- 939 of 941 SNPs identified as associated with BMI ($p < 5 \times 10^{-8}$, 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)

29

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)

30

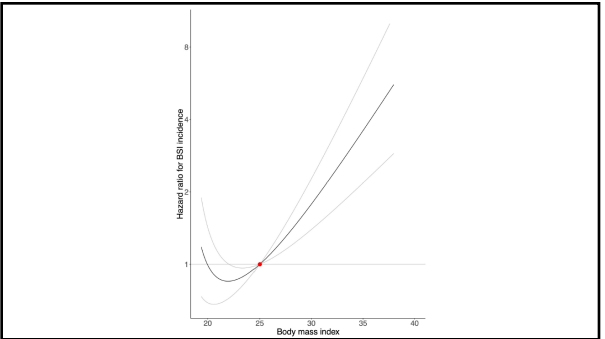
Table 1. Background characteristics.

Characteristic	Total population (n = 35,980)	BMI incidence (n = 2,547)	BMI death (n = 471)
Age (years) ^a	48.5 (16.2-82.3)	48.8 (15.2-82.4)	47.2 (15.1-82.3)
Male sex ^b	26,524 (47.1)	1,843 (52.3)	283 (58.3)
BMI (kg/m ²) ^c	26.5 (1.1)	27.7 (1.2)	27.5 (1.1)
Median follow-up time (years) ^d	21.1 (0.7-23.8)	13.8 (0.4-18.3)	13.3 (0.7-17.9)
Self-reported cancer ^e	1,999 (5.7)	1,448 (2.1)	24 (5.0)
Smoking ^f			
Never	25,399 (48.0)	878 (18.2)	136 (17.6)
Former	13,133 (27.6)	893 (18.4)	144 (17.4)
Current	16,117 (28.4)	773 (29.9)	118 (28.0)
Physical activity ^g			
None	3,821 (7.4)	2,633 (8.9)	54 (11.6)
Slight	15,962 (28.0)	734 (18.9)	117 (15.3)
Moderate	17,247 (28.6)	893 (18.4)	134 (16.3)
High	13,810 (27.4)	387 (19.4)	64 (18.2)
Education ^h			
<9 years	19,603 (55.7)	1,365 (53.4)	240 (58.4)
10-12 years	25,465 (48.0)	1,562 (24.0)	123 (16.6)
>13 years	10,912 (29.3)	274 (10.7)	41 (5.5)

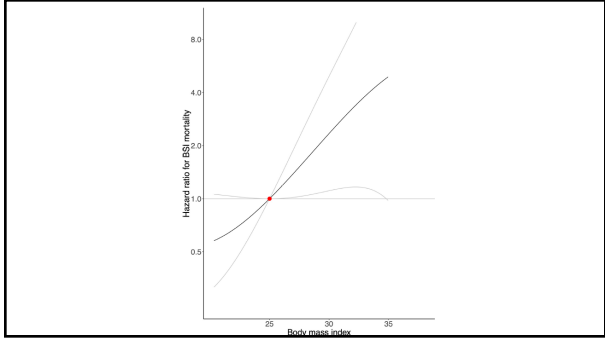
BMI, body mass index; BMI, bloodstream infection. Data are presented as mean (standard deviation).

^aMedian (IQR); 10-13 years (high school), 14-17 years (high school), 18-20 years (high school), 21-23 years (high school), and 24-27 years (university or other post-secondary education, less than 4 years) and 28-33 years (university or other post-secondary education, 4 years or more). ^bActivity defined as follows: none (no light or vigorous activity), slight (<1 h light activity/week and no vigorous activity), moderate (2-3 h light activity/week or 1 h vigorous activity/week), or high (>3 h vigorous activity/week).

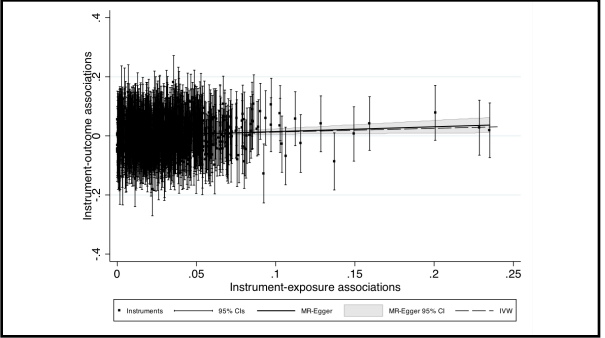
31



32



33



34

S5 Table. Mendelian randomization sensitivity analyses of linear association between body mass index and bloodstream infection mortality in the general population

	MR-Egger	IVW	P-value	Intercept	Lower	Upper	P-value
Over-sample							
MR-Egger random effects	1.18	1.04	1.23	0.002	1.00	0.99	1.00
IVW random effects	1.13	1.05	1.23	0.002	-	-	-
Median estimate - weighted	1.13	0.99	1.20	0.001	-	-	-
Five-sample							
MR-Egger random effects	1.08	0.95	1.18	0.070	1.00	0.99	1.01
IVW random effects	1.09	1.13	1.27	<0.001	-	-	-
Median estimate - weighted	1.10	1.10	1.20	0.002	-	-	-

MR-Egger random effects, MR-Egger random effects; IVW random effects, IVW random effects; Median estimate - weighted, median estimate - weighted; P-value, P-value; Intercept, intercept; Lower, lower bound of 95% CI; Upper, upper bound of 95% CI; P-value, P-value. The MR-Egger random effects, MR-Egger random effects, IVW random effects, IVW random effects, Median estimate - weighted, median estimate - weighted, P-value, P-value, Intercept, intercept, Lower, lower bound of 95% CI, Upper, upper bound of 95% CI, P-value, P-value. The MR-Egger random effects, MR-Egger random effects, IVW random effects, IVW random effects, Median estimate - weighted, median estimate - weighted, P-value, P-value, Intercept, intercept, Lower, lower bound of 95% CI, Upper, upper bound of 95% CI, P-value, P-value.

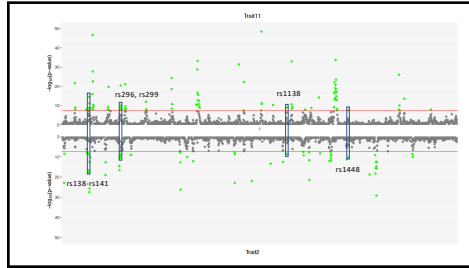
35

STROBE-MR: Guidelines for strengthening the reporting of Mendelian randomization studies

Authors (in alphabetical order):
 George Davey Smith, Neil M Davies, Niki Dimou, Matthias Egger, Valentina Gallo, Robert Golub, Julian PT Higgins, Claudia Langenberg, Elizabeth W Loder, J Brent Richards, Rebecca C Richmond, Veronika W Skrivankova, Sonja A Swanson, Nicholas J Timpson, Anne Tybjaerg-Hansen, Tyler J VanderWeele, Benjamin AR Woolf, James Yarmolinsky

PeerJ Preprints | <https://doi.org/10.7287/peerj.preprints.27857v1> | CC BY 4.0 Open Access | rec: 15 Jul 2019, publ: 15 Jul 2019

36



1

```

CHR F SNP BP P TOTAL NSIG SOS S01 S001 S0001
1 2 rs139 139000 2.86e-28 9 0 0 0 0 9
(INDEX) rs139 RB RQ ALLELES F P
rs137 -2 0.247 00/00 2 2.13e-17
rs138 -1 0.399 00/00 1 1.34e-09
rs139 0 -1 0.399 00/00 1 1.91e-18
rs139 0 1 0.000 0 1 1.15e-10
rs140 1 0.229 00/00 1 6.05e-12
rs140 1 0.229 00/00 1 1.85e-26
rs141 2 0.235 00/00 1 5.3e-09
rs141 2 0.235 00/00 1 5.39e-11
RANGE: chr11:37000..141000
SPAN: 4kb

```

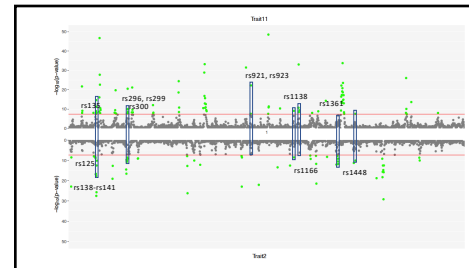
2

```

CHR F SNP BP P TOTAL NSIG SOS S01 S001 S0001
1 1 rs296 296000 1.15e-10 5 0 0 0 0 5
(INDEX) rs296 RB RQ ALLELES F P
rs295 -1 0.429 00/00 1 2.01e-08
rs296 0 1 00/00 2 2.6e-09
rs299 3 0.267 00/00 1 2.77e-09
rs299 3 0.267 00/00 2 7.29e-10
RANGE: chr12:95000..299000
SPAN: 4kb

```

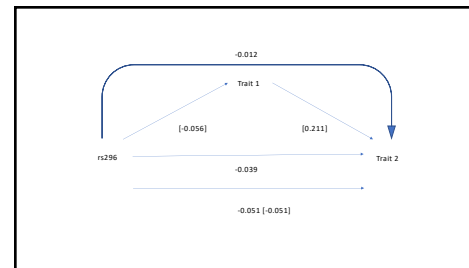
3



4

CHR	SNP	BP	P	TOTAL	NSIG	SOS	S01	S001	S0001
1	rs139	139000	2.86e-28	9	0	0	0	0	9
1	rs296	296000	1.15e-10	5	0	0	0	0	5
1	rs137	137000	2.13e-17	2	0	0	0	0	2
1	rs138	138000	1.34e-09	1	0	0	0	0	1
1	rs139	139000	1.91e-18	1	0	0	0	0	1
1	rs140	140000	6.05e-12	1	0	0	0	0	1
1	rs140	140000	1.85e-26	1	0	0	0	0	1
1	rs141	141000	5.3e-09	1	0	0	0	0	1
1	rs141	141000	5.39e-11	1	0	0	0	0	1
1	rs125	125000	7.3e-09	1	0	0	0	0	1
1	rs138-rs141	138000-141000	7.3e-09	1	0	0	0	0	1
1	rs296	296000	1.15e-10	5	0	0	0	0	5
1	rs299	299000	2.77e-09	1	0	0	0	0	1
1	rs300	300000	1.15e-10	5	0	0	0	0	5
1	rs321	321000	2.01e-08	1	0	0	0	0	1
1	rs923	923000	2.6e-09	2	0	0	0	0	2
1	rs1138	1138000	1.34e-09	1	0	0	0	0	1
1	rs336	336000	1.85e-26	1	0	0	0	0	1
1	rs1448	1448000	5.3e-09	1	0	0	0	0	1

5



6

Discovery and Refinement of Loci Associated with Lipid Levels

A full list of authors and affiliations appears at the end of the article.
 † These authors contributed equally to this work.

Abstract

Low-density lipoprotein (LDL) cholesterol, high-density lipoprotein (HDL) cholesterol, triglycerides, and total cholesterol are heritable, modifiable, risk factors for coronary artery disease. To identify new loci and refine known loci influencing these lipids, we examined 188,578 individuals using genome-wide and custom genotyping arrays. We identify and annotate 157 loci associated with lipid levels at $P < 5 \times 10^{-8}$, including 62 loci not previously associated with lipid levels in humans. Using dense genotyping in individuals of European, East Asian, South Asian, and African ancestry, we narrow association signals to 12 loci. We find that loci associated with blood lipids are often associated with cardiovascular and metabolic traits including coronary artery disease, type 2 diabetes, blood pressure, waist-hip ratio, and body mass index. Our results illustrate the value of genetic data from individuals of diverse ancestries and provide insights into biological mechanisms regulating blood lipids to guide future genetic, biological, and therapeutic research.

1

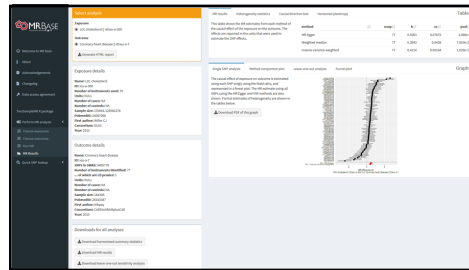
A comprehensive 1000 Genomes-based genome-wide association meta-analysis of coronary artery disease

A full list of authors and affiliations appears at the end of the article.
 † These authors contributed equally to this work.

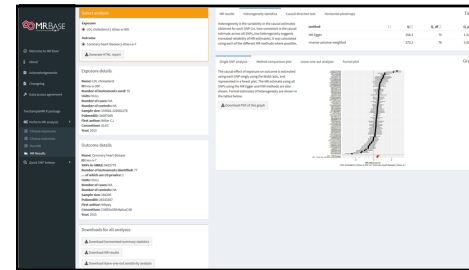
Abstract

Existing knowledge of genetic variants affecting risk of coronary artery disease (CAD) is largely based on genome-wide association studies (GWAS) analysis of common SNPs. Leveraging phased haplotypes from the 1000 Genomes Project, we report a GWAS meta-analysis of 185 thousand CAD cases and controls, interrogating 6.7 million common (MAF>0.05) as well as 2.7 million low frequency (0.005<MAF<0.05) variants. In addition to confirmation of most known CAD loci, we identified 10 novel loci, eight additive and two recessive, that contain candidate genes that newly implicate biological processes in vessel walls. We observed intra-locus allelic heterogeneity but little evidence of low frequency variants with larger effects and no evidence of synthetic association. Our analysis provides a comprehensive survey of the fine genetic architecture of CAD showing that genetic susceptibility to this common disease is largely determined by common SNPs of small effect size.

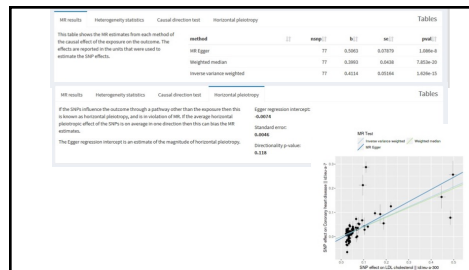
2



3



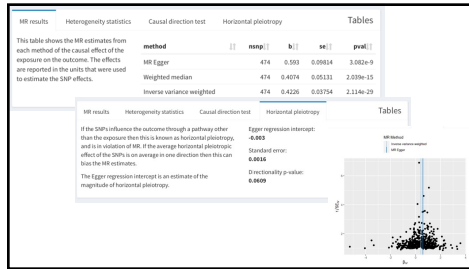
4



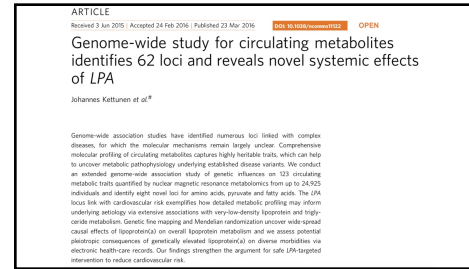
5



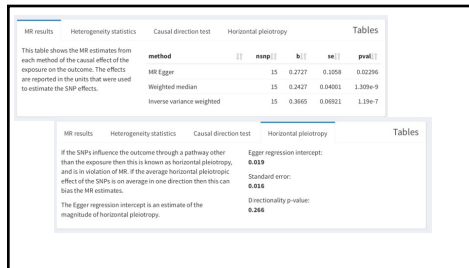
6



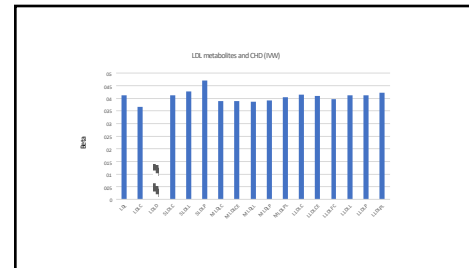
7



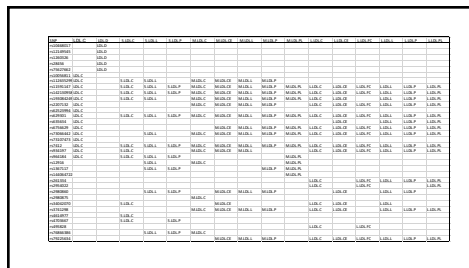
8



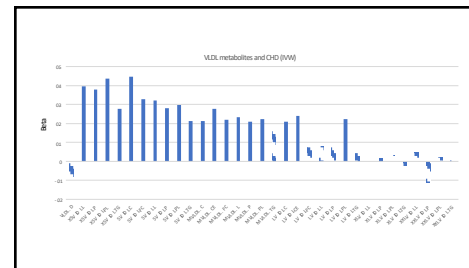
9



10



11




12



The Ethics and Regulation of Human Subjects Research

Wayne Patterson, PhD
Senior Consultant

1



The Nuremberg Code (1947)

Ten Basic Principles, including:


"The voluntary consent of the human subject is absolutely essential..."

"The experiment should be conducted so 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 continuation of the experiment is likely to result in injury, disability, or death to the experimental subject."



2


Tuskegee Study of Untreated Syphilis in the Negro Male (1932-1972)

3

National Research Act (1974)



Required the creation of the National Commission for the Protection of Human Subjects of Biomedical and Behavioral Research.



4

The Ethics of Conducting Research with Humans: The Belmont Report (1979)

- **Beneficence**
 - maximize benefits, minimize risks
- **Justice**
 - Who should bear the burdens of the research?
 - Who should benefit from results?
- **Respect for Persons**
 - Autonomy
 - Protect those with diminished autonomy

5


The Belmont Report was the basis for federal requirements of human research protections

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)

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




6

What is the Common Rule?

It is **the** Federal Policy for the Protection of Human Subjects

Originally promulgated in 1991, with no significant changes, until 1/21/19!

Rockefeller's Federal Wide Assurance (FWA) certifies compliance with this federal policy (for human research conducted or supported by Common Rule agencies...)




7

What's so Common about the Common Rule?

✓ 19 federal agencies follow the new Common Rule, e.g.,



- DHHS, including NIH (45 CFR 46, Subpart A)*
- DoD (32 CFR 219)
- NSF (45 CFR 690)
- Department of Energy (DoE) (10 CFR 745)
- Veterans Administration (38 CFR 16)
- Department of Education (DoEd) (34 CFR 97)

*FDA is within DHHS, but also has its own regulations
*DoJ has not signed on yet



8

First Question: Is your activity "human subjects research" (HSR)?





9

Specifically:

1. Is it HSR according to the Common Rule?
2. Is it HSR according to FDA?

(could be both!)




10

Start with the Common Rule

First assess:

Does the activity involve Research?




11

Common Rule Definition of Research:

"...a **systematic investigation**, including research development, testing and evaluation, **designed to develop or contribute to generalized knowledge...**"

(Both parts of the definition must be met)




12

Part I of the definition:
What's a Systematic Investigation?

an activity that involves a prospective plan which incorporates data collection, either quantitative and/or qualitative, and data analysis to answer a question

Does a case study involve a systematic investigation?



13


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 yet firmly established or accepted;

and

- ✓ the activity is not dependent on the unique characteristics of the target population or system in which it will be implemented




14

An activity is not likely to be generalizable if the intent is:

The evaluation or improvement of a process, practice, or program at the site where the activity is being conducted

Results only to be applied to populations, or inform practice within the target population or within the site where the activity is being conducted

Implementation and evaluation of an evidence-based practice, process, or program (is it functioning as intended within the site where the activity is being conducted or with the local target population)




15

If the activity IS research:
Does the research involve human subjects, according to the Common Rule?

A living individual about whom an investigator conducting research:


- (i) Obtains information or biospecimens through intervention or interaction with the individual, and uses, studies, or analyzes the information or biospecimens; or
- (ii) Obtains, uses, studies, analyzes, or generates identifiable private information or identifiable biospecimens.



16

Once you determine if the activity is or is not human subjects research according to the Common Rule...

You may still need to assess if the activity is human subjects research according to FDA




17

FDA Decisions

Does the activity evaluate an FDA-regulated test article (i.e., drug, biologic, device)?

Does the activity involve Human Subjects?
An individual who is, or becomes, a participant in research, either as a recipient of the test article or as a control. A subject may be either a healthy human or a patient. *Also included in the FDA human subject definition: The use of a biological specimen – even if de-identified from an individual used to test an investigational device*

Does the activity involve research (clinical investigation)?
Any experiment that involves a test article and one or more human subjects...



18

If the activity IS human subjects research, next question: Is it exempt from the federal regulations? *



*this does not mean exempt from Institutional review!




19

There are 6 HSR categories of research that are Exempt from IRB Review
Focus on: Exemption #4


Secondary research* for which consent is not required

*Secondary research only! (i.e., re-using identifiable information and/or identifiable biospecimens that were, or will be, are collected for another reason, e.g., clinical or research)



20


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:



21

Exemption 4(i)


The identifiable private information or identifiable biospecimens are publicly available;



22

Exemption 4(ii)


Identifiable private information...is recorded by the investigator in such a manner that the identity of the human subjects cannot readily be ascertained directly or through identifiers linked to the subject, the investigator does not contact the subjects, and the investigator will not re-identify subjects;



23

Exemption 4 (iii)


"The research involves only information collection and analysis involving the investigator's use of identifiable health information when that use is regulated under 45 CFR parts 160 AND 164, subparts A and E [HIPAA], for the purposes of "health care operations" or "research" as those terms are defined at 45 CFR 164.501 or "public health activities and purposes" as described under 45 CFR 164.512(b)"



24

Exemption 4 (iv)


The research is conducted by, or on behalf of, a Federal department or agency using government-generated or government-collected information obtained for non-research activities, if the research generates identifiable private information that is or will be maintained on information technology that is subject to and in compliance with section 208(b) of the E-Government Act of 2002, [44 U.S.C. 3501 note](#), if all of the identifiable private information collected, used, or generated as part of the activity will be maintained in systems of records subject to the Privacy Act of 1974, [5 U.S.C. 552a](#), and, if applicable, the information used in the research was collected subject to the Paperwork Reduction Act of 1995, [44 U.S.C. 3501 et seq.](#)



25

What are the ethical standards that should be considered for all exempt studies?

Criteria	Yes	No	NA
The research holds out no more than minimal risk to participants	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Selection of participants is equitable	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
If there is recording of identifiable information, there are adequate provisions to maintain the confidentiality of the data	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
If there are interactions with participants, there are adequate provisions to protect the privacy interests of participants	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
If there are interactions with participants, the consent process or information provided to potential subjects includes the following: <input type="checkbox"/> IRB—there are no interactions and no other need for consent	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
That the activity involves research	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
A description of the procedures	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
For Category 3 research that involves subject deception: A statement that subjects will be unaware of or misled regarding the nature or purposes of the research	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
That participation is voluntary	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Name and contact information for the researcher	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>



26

If the activity IS human subjects research, but does not qualify for exemption, it is HSR that is not exempt, i.e., it is subject to federal regulations governing human research protection...




...including review by a federally mandated Institutional Review Board (IRB)



27

Two Types of Non-Exempt Review

1. Expedited Review
2. Full Board Review




28

For a non-exempt study to qualify for Expedited (not full IRB Board) Review...

...The research must be all of the following:



- no greater than minimal risk
- not involve prisoners (per OHRP guidance)
- not be classified
- not involve identifiable data that would place subjects at risk of criminal or civil liability or be damaging to the subjects financial standing, employability, insurability, reputation, or be stigmatizing. If it could, reasonable protections must be in place so that risks related to invasion of privacy and breach of confidentiality are no greater than minimal, **and**

• Fit into one or more of these categories:
<https://www.hhs.gov/ohrp/regulations-and-policy/guidance/categories-of-research-expedited-review-procedure-1989/index.html>



29

If the nonexempt research doesn't qualify for expedited review, it must be reviewed at a convened IRB meeting.

30

Whether expedited or full board,
a study must meet federally-
defined criteria in order to be
approved

i.e.,

“The .111 Criteria”



31

§ 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:



32

1. Risks to subjects are minimized:

(i) By using procedures which are consistent with **sound research design** and which do not unnecessarily expose subjects to risk, and

(ii) Whenever appropriate, by using procedures already being performed on the subjects for diagnostic or treatment purposes



33

2. Risks to subjects are reasonable in relation to anticipated benefits, if any, to subjects, and the importance of the knowledge that may reasonably be expected to result



34

3. Selection of Subjects is Equitable

Consider:

- The setting in which the research will be conducted
- 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?



35

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?




36

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?



37

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)




38

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




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 re-identified?
- 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?
- Plans for data retention and destruction?



40

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 stored (University-provided database application like REDCap, IT file server, etc.)?
 - If there is 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 to authorized research personnel.



41

Data security, continued

- Transporting or transmitting the data
 - If any research data will be collected on a mobile device, such as an electronic tablet, cell phone, or wireless activity tracker, details are needed regarding the physical security of the device, electronic security, and how the transfer of data from device to research storage location will be securely accomplished.
 - If any research data will be directly entered/sent by subjects over the internet or via email, will a University-provided database application (like REDCap) be used, or is there an encrypted tunnel to the site/application?
- Access to the data
 - How will the investigators ensure only approved research personnel have access to the stored research data? Password-protected files, role-based security, etc.?
- Sharing of the data
 - Will data be transferred or disclosed to or from the University? Is a contract or data transfer agreement necessary? What (if any) identifiers will be included? How will the data be securely transferred or disclosed (University-approved secure file transfer, etc.)?




42

Using Social Media in your research

Recruitment

- Seek to normalize social media recruitment to the extent possible, drawing analogies to traditional recruitment efforts
- Ensure that the proposed online recruitment strategy complies with all applicable federal and state laws, e.g.
 - Recruitment advertisements
 - Web site "Terms of Use"
 - Tell potential subjects that information shared via social media is not secure.

https://catalyst.harvard.edu/pdf/regulatory/Social_Media_Guidance.pdf



43

Using Social Media in your research

Recruitment

- **Assure compliance between recruitment techniques and policies/terms of service of relevant websites.**
 - If a proposed technique conflicts with website policies and terms of service, request a written exception from the site, OR
 - Depending on IRB policy, provide a statement explaining why the recruitment strategy warrants approval without an explicit exception, to be evaluated by the IRB with input from institutional legal counsel.

https://catalyst.harvard.edu/pdf/regulatory/Social_Media_Guidance.pdf




44

Using Social Media in your research

Recruitment

- **Ensure that proposed social media recruitment strategies respect all relevant ethical norms, including:**
 - Proposed recruitment does not involve deception or fabrication of online identities
 - Proposed recruitment does not involve members of research team lurking or creeping social media sites in ways members are unaware of
 - Strategy must be sensitive to the privacy of potential participants and respectful of the norms of the community being recruited
 - Recruitment will not involve advancements or contact that could embarrass or stigmatize potential subjects

https://catalyst.harvard.edu/pdf/regulatory/Social_Media_Guidance.pdf




45

Using Social Media in your research

Recruitment

- **Enlist enrolled participants to facilitate introduction between members of their network and the research team. Ensure that consent will be obtained from current participants before they approach members of their online network for recruitment via their network or**
- **Ensure that a communication plan is in place for how the research team will handle online communication from enrolled participants that threatens the integrity of study**

https://catalyst.harvard.edu/pdf/regulatory/Social_Media_Guidance.pdf




46

Using Social Media in your research

Data source

- **A key issue in observational research using social media is whether the proposed project meets the criteria as human subjects research, and if so, what type of review is needed**
 - Identifiable/de-identified data
 - Minimal risk/greater than minimal risk



47

Using Social Media in your research

Data source

- **How is the data collected, transferred, etc.**
 - Specify if research data will be collected as part of the recruitment process via social media. If so, describe what data will be collected. If that data is of a sensitive or confidential nature, describe how that data will be transferred to secure institutional servers and how will it be protected upon receipt.




48

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?



49

That's it for the .111 criteria...
but that's not all!

Pregnant Women?
Subpart B of 45 CFR 46


Prisoners?
Subpart C of 45 CFR 46

Children?
Subpart D of 45 CFR 46

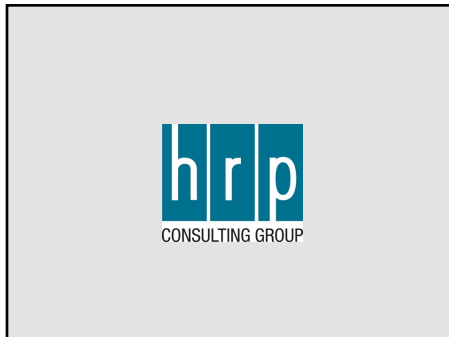
Department of Education (ED)?
Family Educational Rights and Privacy Act (FERPA) (34 CFR 99)
and the Protection of Pupil Rights Amendment (PPRA) (34 CFR 98)
See [resources provided by ED](#) when developing your research protocol

Investigational Drugs, biologics, devices?
FDA regulations at 21 CFR 50, 21 CFR 56, 21 CFR 312, 21 CFR 812

HIPAA?
45 CFR [Part 160](#) and Subparts A and E of [Part 164](#)




50



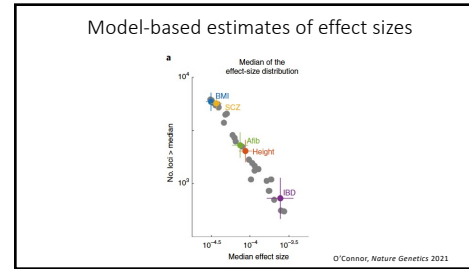
51

Genetic risk prediction

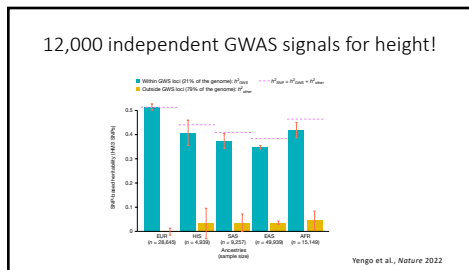


Genotype of an individual (Common SNPs) → Life-time risk of genetic disorders (Common complex genetic disorders)

1



2




3

Effect sizes of individual variants are very small

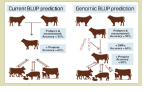
- Genotype at a single locus carries very little information about phenotype.
- It does not mean that one cannot predict phenotype from genotype.
- Accuracy (r^2) of an ideal genetic predictor equals heritability.

4

BLUP – Best Linear Unbiased Predictor



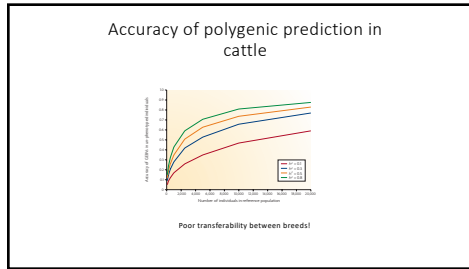
- Infinitesimal model
- Genetic effects are random
- Predict the expected genetic effect



5

$$\hat{g}_i = \mathbf{G}_i - (\mathbf{G} + \mathbf{I})^{-1}(\mathbf{y} - \mathbf{1}\hat{\mu})$$

6



7

Measuring risk of myocardial infarction

Coronary Risk Prediction in Adults (The Framingham Heart Study)

PETER W.F. WILSON, MD, WILLIAM J. CASTELLI, MD, and WILLIAM B. KANNEL, MD

The Framingham Heart Study, an ongoing prospective study of adult men and women, has shown that certain risk factors can lead to a possible 90% reduction of coronary artery disease. These factors include age, gender, total cholesterol level, high density lipoprotein cholesterol level, systolic blood pressure, cigarette smoking, glucose intolerance and cardiac enlargement. Risk prediction based on an individual's age and cholesterol level (total or other) may. Calculations and computer can be easily programmed using a microdotcom system.

Further risk allows calculation of the theoretical probability of cardiovascular events. These advance studies, based on population data, 2500 men and women participating in the Framingham study, indicate that 90% of coronary artery disease can be avoided by a 10% reduction in cholesterol level, 10% reduction in blood pressure, 10% reduction in cigarette smoking, 10% reduction in glucose intolerance and 10% reduction in cardiac enlargement. This study was published in *Journal of the American Medical Association*, 1976, 236: 692-698.

8

LDL levels and risk of disease

Annals of Internal Medicine ARTICLE

Nonoptimal Lipids Commonly Present in Young Adults and Coronary Calcium Later in Life: The CARDIA (Coronary Artery Risk Development in Young Adults) Study

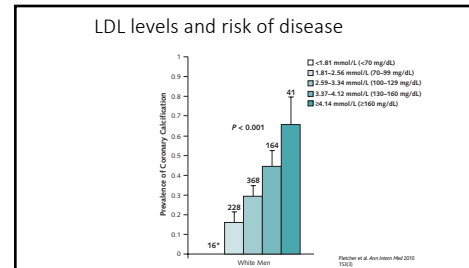
John A. Jacobson, MD, MPH, Rebecca Steinberger, PhD, MS, Wang Guo, PhD, Steve Kiehl, MD, MPH, Hong Guo, MS, Eric Vinogradov, PhD, and Barbara F. Kraljic, MD, PhD

~3500 subjects < 35 years old

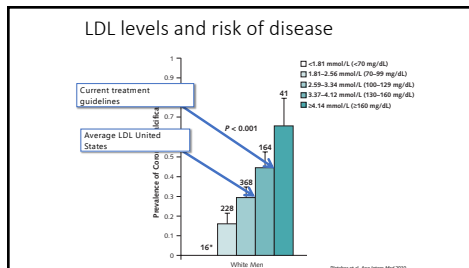
15-20 years

Archives of Internal Medicine 2008; 168: 1033-1039

9



10



11

Selecting populations for treatment


12

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

13

Applications in humans



Prediction of individual genetic risk to disease from genome-wide association studies
Wang et al. Nature Reviews Genetics 2015

LETTERS
Common polygenic variation contributes to risk of schizophrenia and bipolar disorder
Collier et al. Nature Genetics 2016

- LD-prune
- Exclude SNPs of very small effect

14

Extensions of BLUP – multiple variance scales and binary phenotypes

MultiBLUP:	Speed and Balding. <i>Genome Research</i> 2014
Bayesian analysis:	MacLeod et al. <i>Genetics</i> 2014
BSLMM:	Zhou et al. <i>PLOS Genetics</i> 2013
GeRSI:	Golan and Rosset. <i>AJHG</i> 2014

15

Methods that work with summary statistics

- Summary statistics are easily available
- Most methods require a separate small individual level dataset to tune parameters

16

LDpred – a Bayesian method using summary statistics

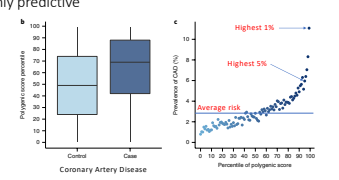
$$\beta_{i \sim M} \begin{cases} N\left(0, \frac{h^2}{M_p}\right) \text{ with probability } p \\ 0 \text{ with probability } (1 - p). \end{cases}$$

Vinkhuyzen et al. 2015

Also, check *BayesR*

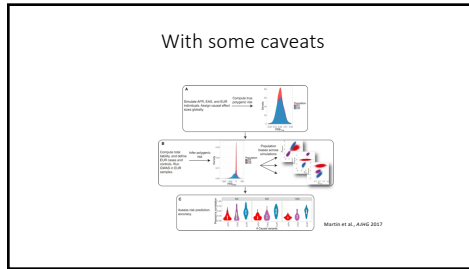
17

Extreme tails in the distributions of genetic risk scores are highly predictive

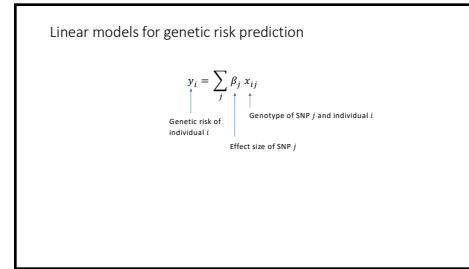


Khera et al. 2018

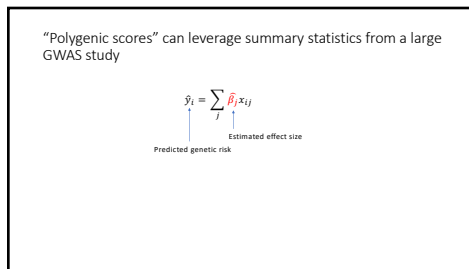
18



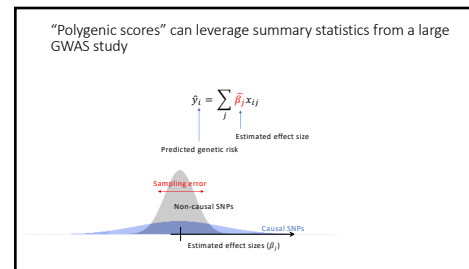
19



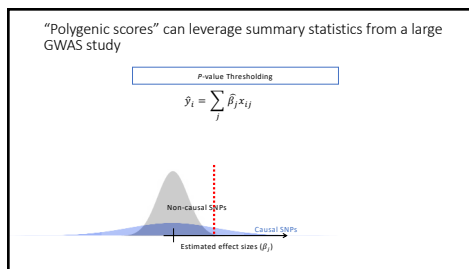
20



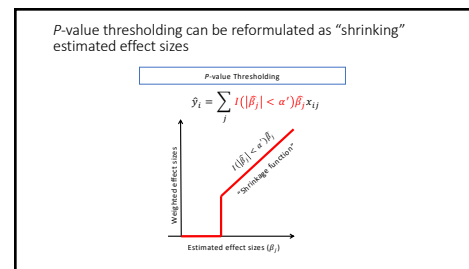
21



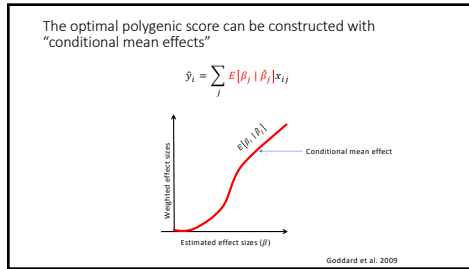
22



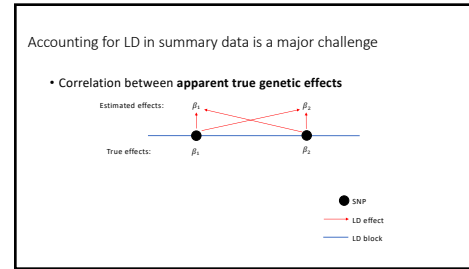
23



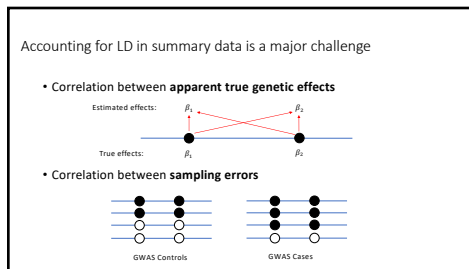
24



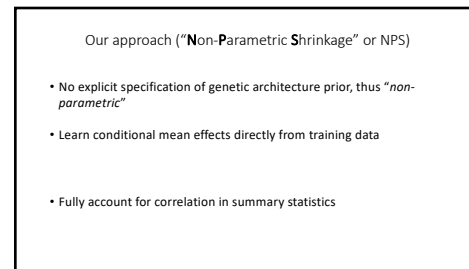
25



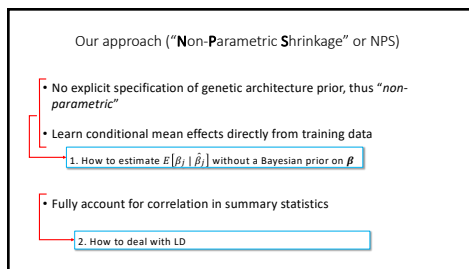
26



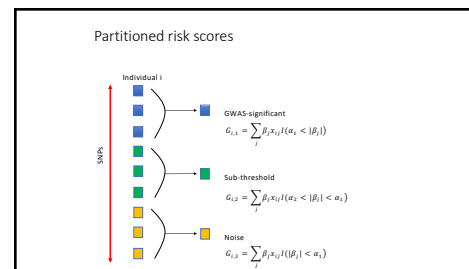
27



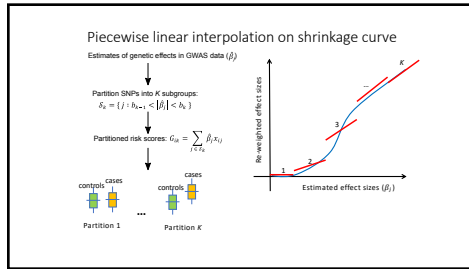
28



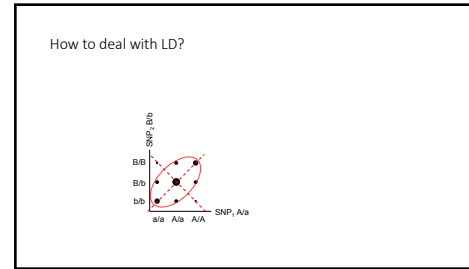
29



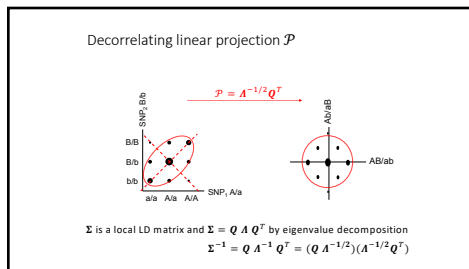
30



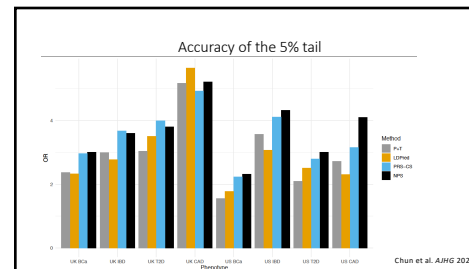
31



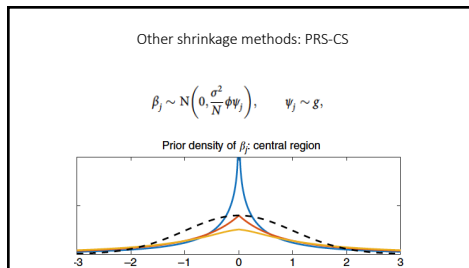
32



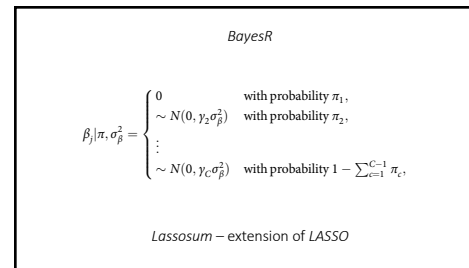
33



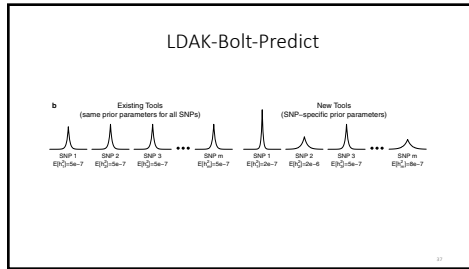
34



35



36

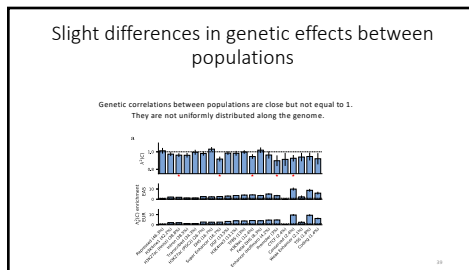


37

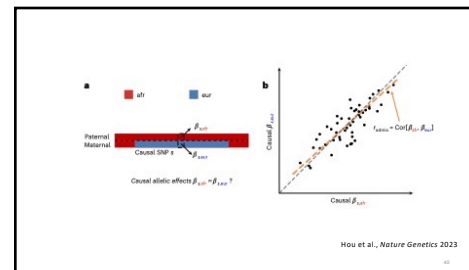
What makes PRS non-transferrable?

- Differences in allele frequencies between populations
- Differences in LD between populations
- Differences in effect sizes (although likely a minor contribution)

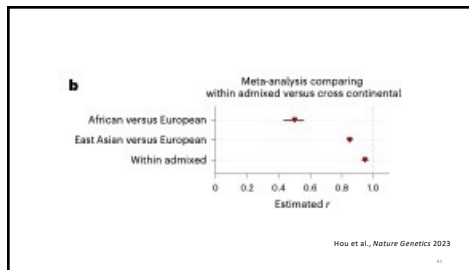
38



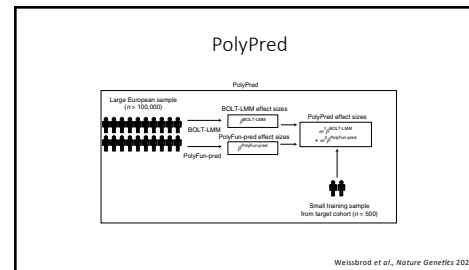
39



40

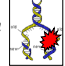





41



42

Forces responsible for genetic change



Mutation  μ
Selection  s
Drift  N_e
Population structure  F_{ST}

1

Mutations

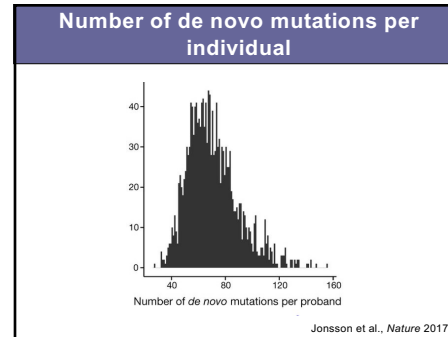
2

Mutation rate in humans and flies

 2.5×10^{-8} (Nachman & Crowell)
 1.8×10^{-8} (Kondrashov)





NGS estimates $\sim 1.2 \times 10^{-8}$ per nt changes genome
 ~ 70 per nt changes genome
 Other events: indels (10^{-9})
 repeat extensions/contractions (10^{-5})

3



4

Mutation rate is variable along the genome

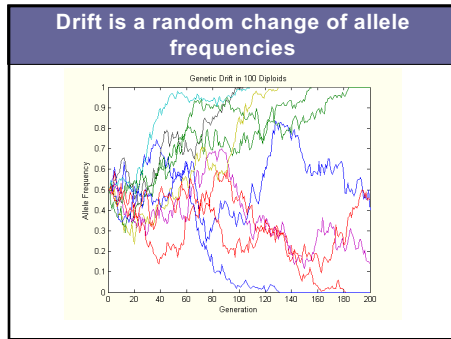
 Replication fidelity
 direct DNA damage
 DNA repair
 CpG deamination

Regional variation of mutation rate
Context dependence of mutation rate

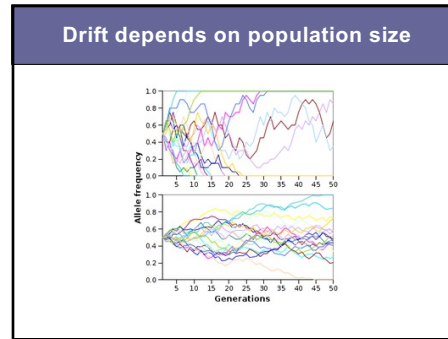
5

Genetic drift

6



7



8

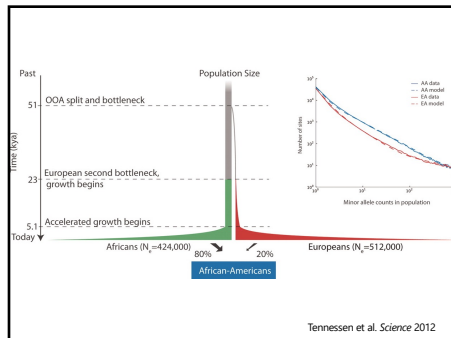
Effective population size

- In an idealized model, the intensity of genetic drift depends on population size (mean squared change in allele frequency is proportional to $1/N_e$)
- In more realistic situations, effective population size (N_e) is a parameter characterizing intensity of drift

9

Demographic history

10

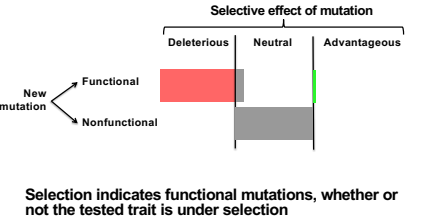


11

Selection

12

Most functional mutations are deleterious



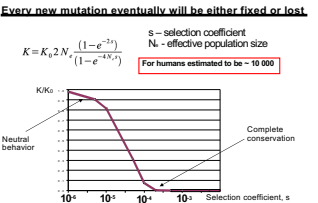
13

Selection coefficient

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

14

Conservation can be due to very weak selection!



15

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 ($\sim 300-400K$ years old)
- Protein coding regions are showing clear signs of selection (reduced diversity and excess of rare alleles)

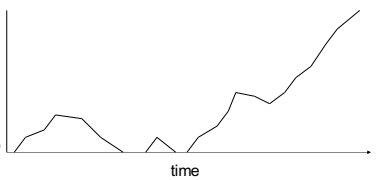
16

Methods of mathematical population genetics

17

Dynamic of allelic substitution

Mathematically, allele frequency change in a population follows a one-dimensional random walk



18

Diffusion approximation

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

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

19

Coalescent theory

Instead of modeling a population, we can model our sample

Time goes backwards !

20

Signatures of purifying selection

Reduced variation

Excess of rare alleles

21

Commonly used summary statistics to characterize variation

22

Number of segregating sites

```

. . . T C A A G T C A A G C G A T C A T G . . .
. . . T C A A G T C A A G C G A T C A G G . . .
. . . T C A G G T C A A G T G A T C A T G . . .
. . . T C A G G T C A A G T G A T C A T G . . .
. . . T C A A G T C A A G C G A T C A G G . . .
. . . T C A A G T C A A G C G A A C A G G . . .
    
```

k – number of sites variable in the sample
density of segregating sites is also frequently used
k is dominated by rare alleles
k strongly depend on sample size

23

Nucleotide diversity

$$\pi = \frac{2}{n(n-1)} \sum d_{ij} \quad d_{ij} - \text{number of nucleotide differences between sequences } i \text{ and } j$$

$$\pi = \frac{n}{(n-1)} \sum 2p_k (1 - p_k) \quad p_k - \text{allele frequency at site } k$$

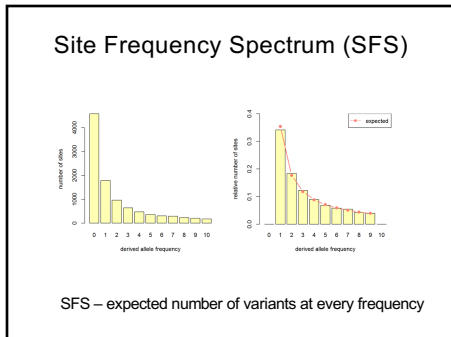
π – the average density of nucleotide differences between two sequences

π – per nucleotide heterozygosity

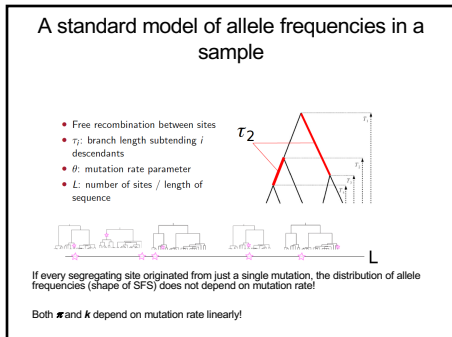
π is dominated by common alleles

π is independent of sample size

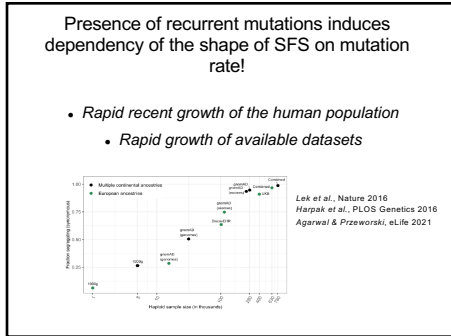
24



25



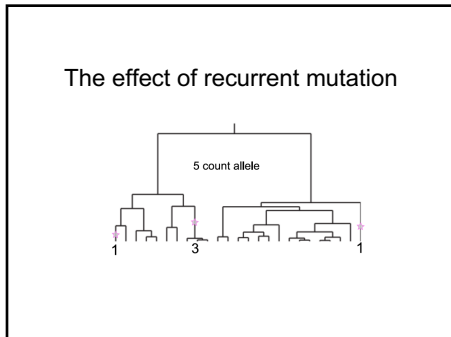
26



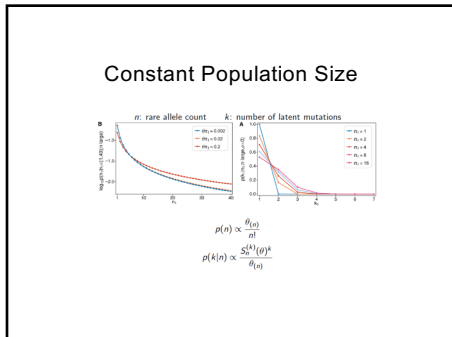
27



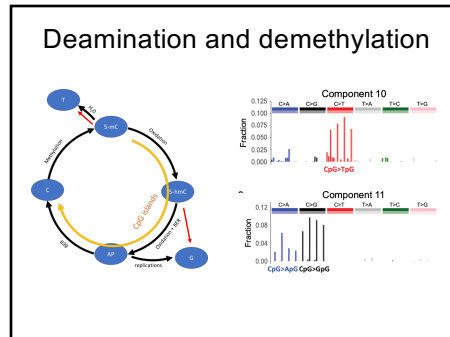
28



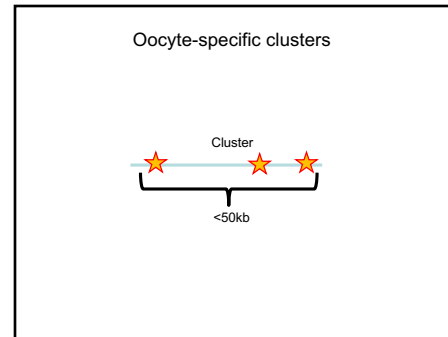
29



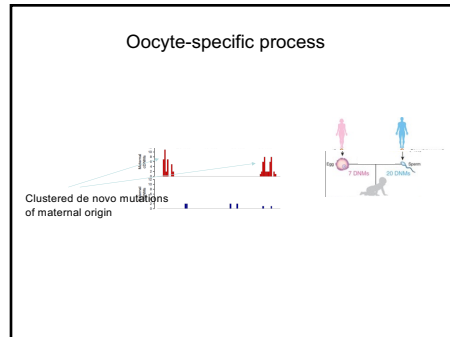
30



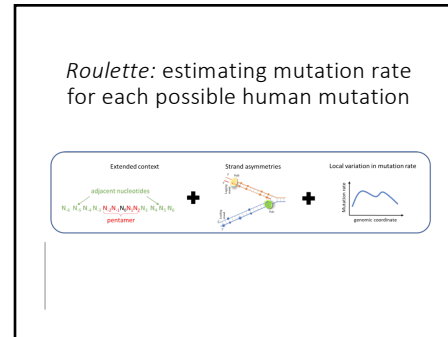
37



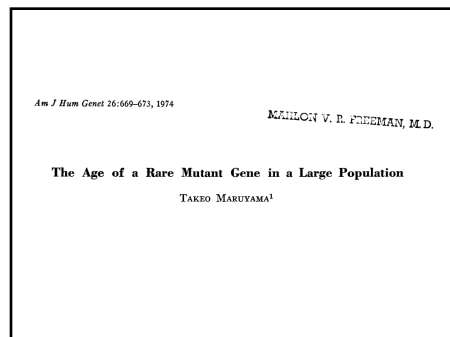
38



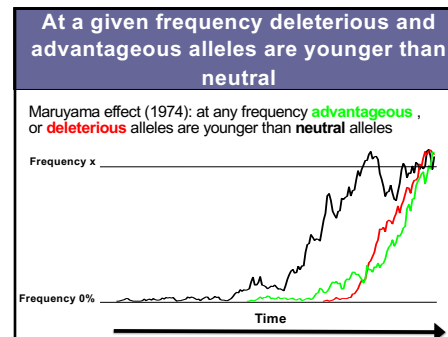
39



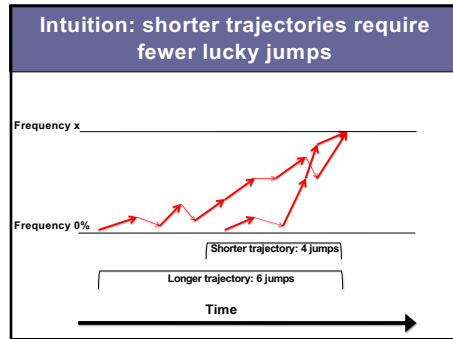
40



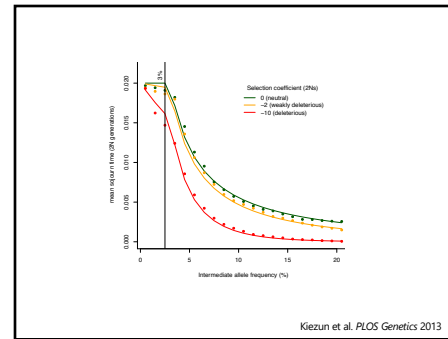
41



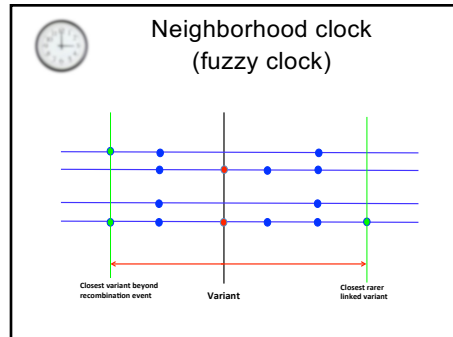
42



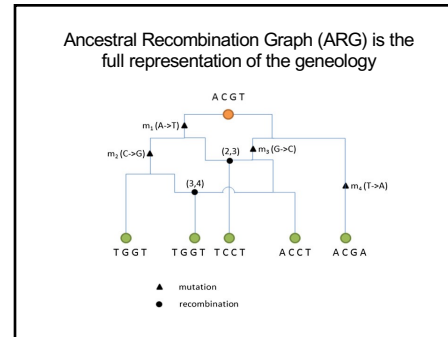
43



44



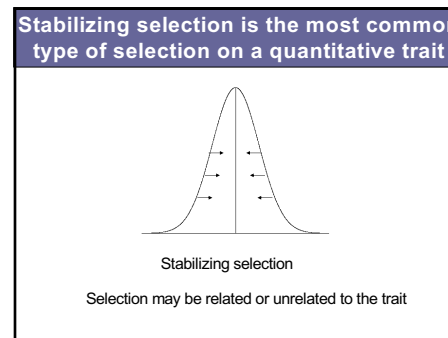
45



46



47



48

Technically, non-neutral genetic variation should not exist!

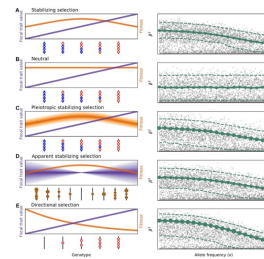
Forces to maintain variation:

Selection

Mutation

49

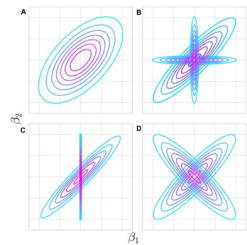
Possible theoretical models



Koch & Sunyaev *Front. Genet.* 2021

50

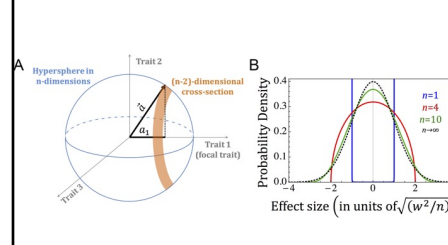
Shades of pleiotropy



Koch & Sunyaev *Front. Genet.* 2021

51

A highly pleiotropic model



Simons et al., *PLoS Biology* 2018

52

Functional annotation of genes and variants

1

Map variants onto genomic annotation

Watch for multiple transcripts!

Watch for conflicting annotations!

2

Nonsense 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

3

Selection inference from frequency of individual SNVs

$$\begin{aligned} & \text{Change in allele frequency} = \\ & = \text{Mutation} + \text{Selection} + \text{Drift} \end{aligned}$$

Of the order of 10^{-8} Demographic history Population structure

4

Focusing on rare deleterious PTVs

PTV – protein truncating variant
(a.k.a. nonsense)

Combine all PTVs per gene – we assume that they have identical effects

Consider each gene as a bi-allelic locus – PTV / no PTV

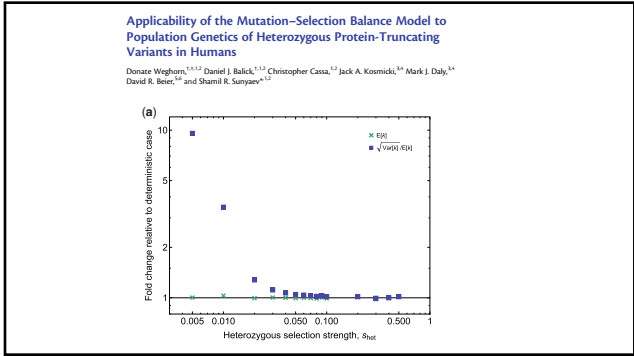
5

Selection inference using combined frequency of PTVs

$$\begin{aligned} & \text{Change in allele frequency} = \\ & = \text{Mutation} + \text{Selection} + \text{Drift} \end{aligned}$$

Assuming strong selection and a very large population, combined frequency of rare deleterious PTVs is expected to be Poisson distributed with $\lambda = U/hs$

6



7

Loss-of-function observed/expected upper bound fraction (LOEUF)

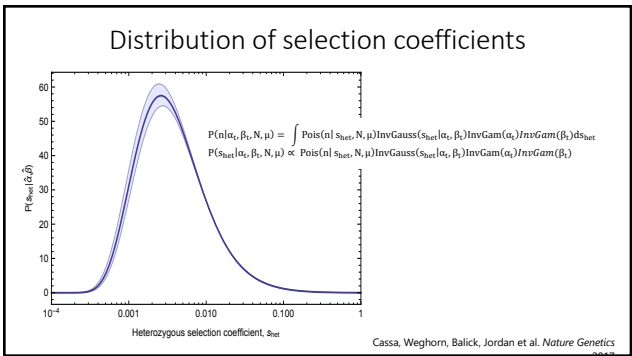
- LOEUF is based on the number of segregating sites as the statistic
- LOEUF assumes Poisson distribution for the number of segregating sites. It computes the expectation. The constraint metric is based on the Poisson likelihood ratio upper bound.

8

Treating combined PTVs as a bi-allelic locus

- We can use the total frequency of PTVs in the locus
- Theoretically, we can simply treat all PTV variation as a single bi-allelic locus with high mutation rate

9

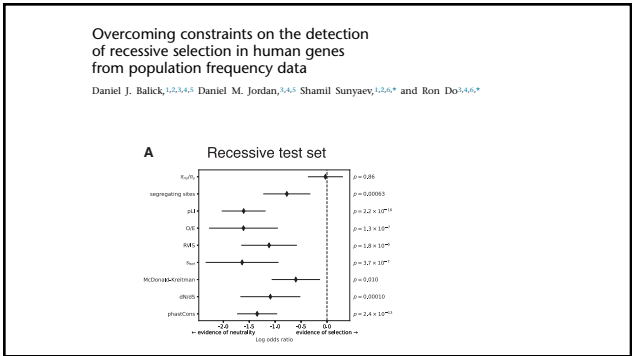


10

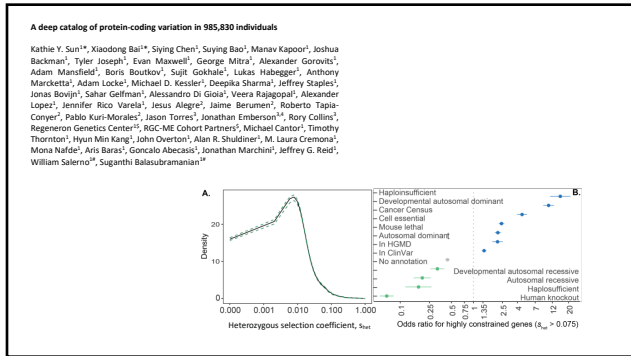
Distribution of selection coefficients

- 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

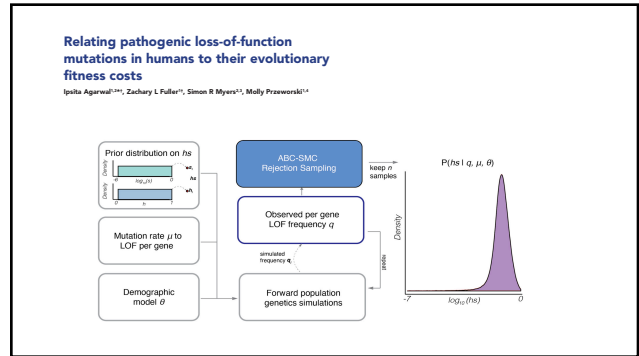
11



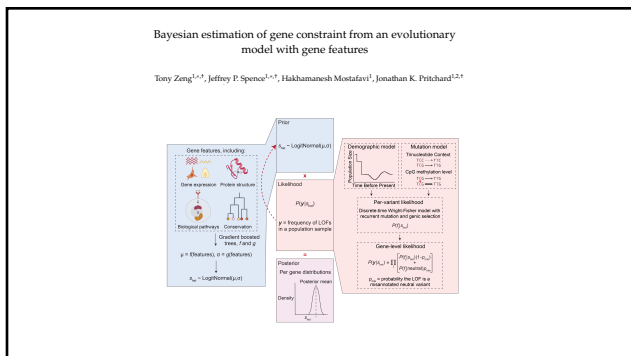
12



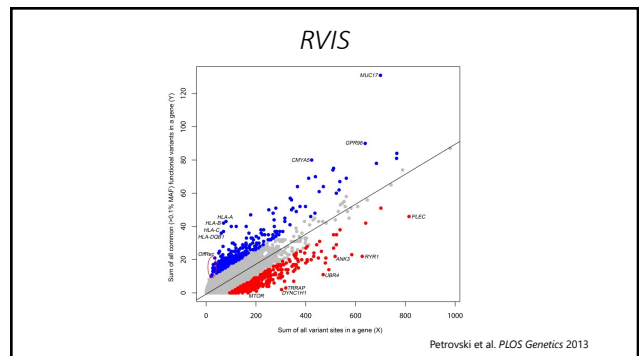
13



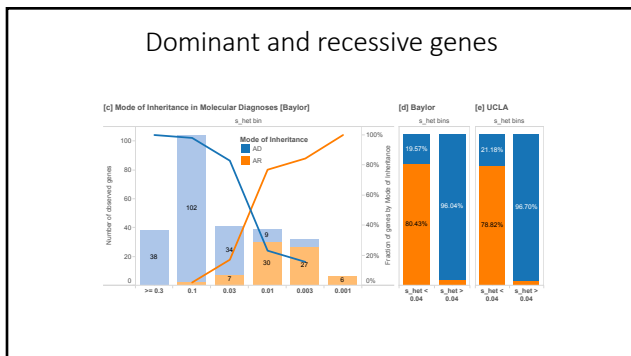
14



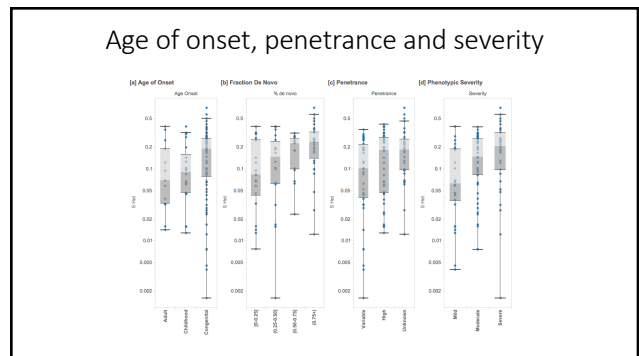
15



16

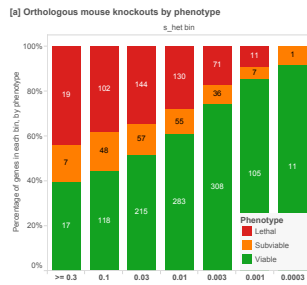


17



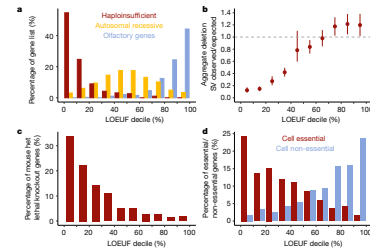
18

Concordance with the mouse knockout data



19

LOEUF (gnomAD)



20

Applications to Mendelian genetics – large cohorts make Mendelian genetics a data science

Article

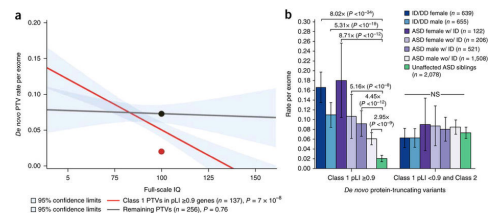
Evidence for 28 genetic disorders discovered by combining healthcare and research data

https://doi.org/10.1038/s41588-020-2832-5
 Received: 8 October 2019
 Accepted: 17 July 2020
 Published online: 14 October 2020
 Check for updates

DeNovoWEST – a method to identify significant recurrent *de novo* mutations controlling for mutation rate, weighting genes with *Sher* and weighting variants using variant effect predictors

21

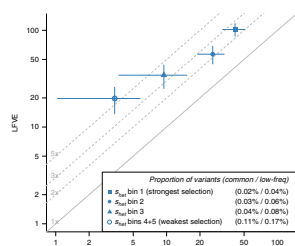
De Novo mutations in ASD



Kosmicki et al. *Nature Genetics* 2017

22

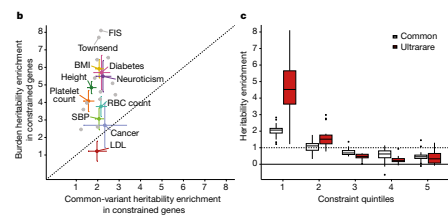
Heritability Enrichment



Gazal et al. *Nature Genetics* 2018

23

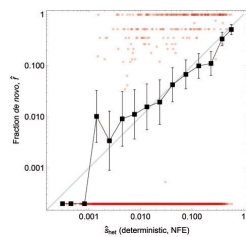
“Burden” heritability enrichment



Weiner, Nadig et al. *Nature* 2023

24

Selection in the present-day population

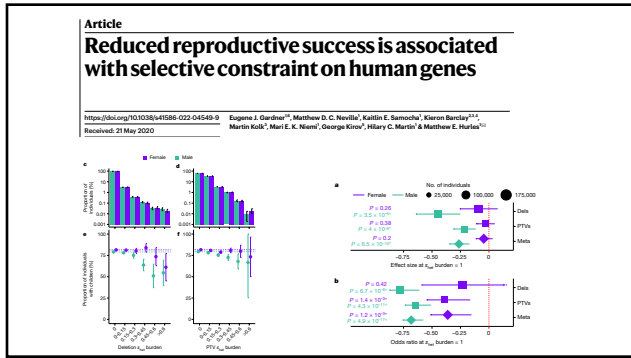


Fraction of *de novo* mutations (out of all variants) approximately equals selection coefficient.

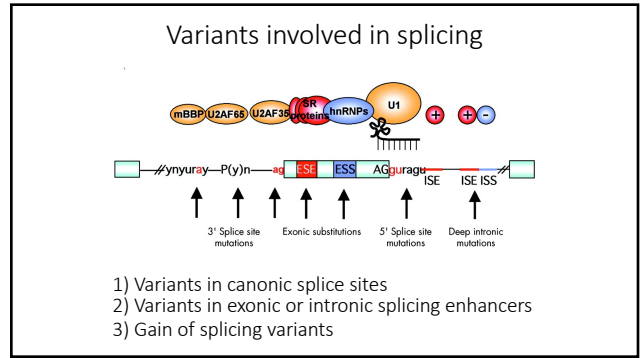
This result does not depend on phenotypic ascertainment.

Weghorn et al., *M&BE* 2019

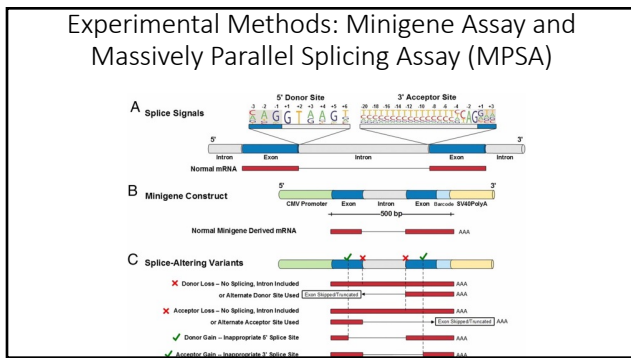
25



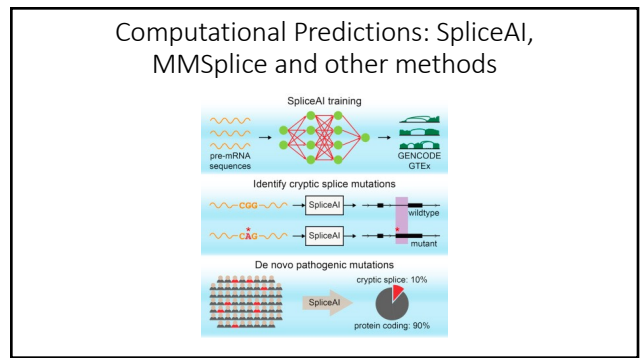
1



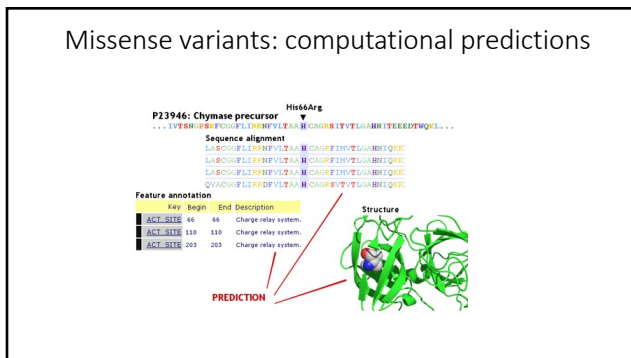
2



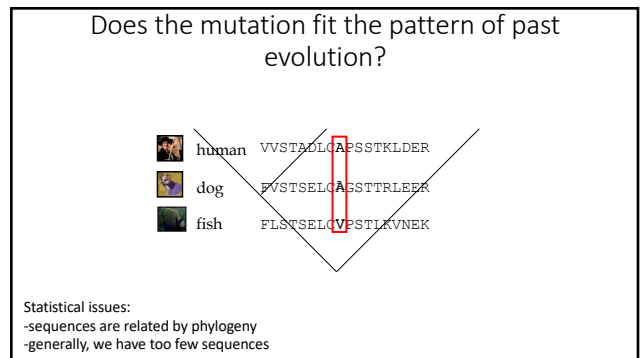
3



4



5



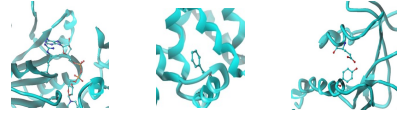
6

Does the mutation fit the pattern of past evolution?

- We assume a constant fitness landscape: what is good for fish is good for human!
- We can estimate whether the mutation fits the pattern of amino acid changes.
- We can also estimate rate of evolution at the amino acid site

7

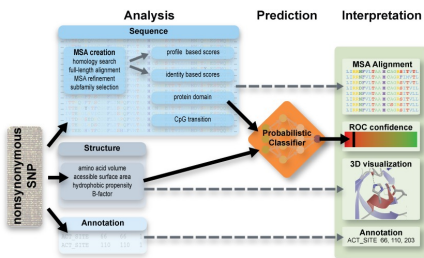
Protein structure view



- Most of pathogenic mutations are important for stability (good news?).
- $\Delta\Delta G$ is difficult to estimate.
- Unfolded protein response pathway has to be taken into account.
- Heuristic structural parameters help but less than comparative genomics.

8

PolyPhen2



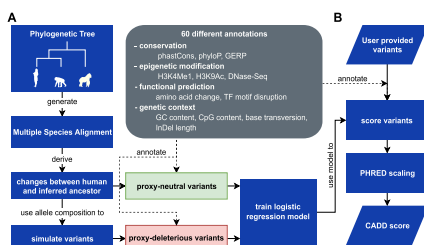
www.genetics.bwh.harvard.edu/pph2 Adzhubei, et al. Nature Methods 2010

9

SIFT is based on multiple sequence alignment

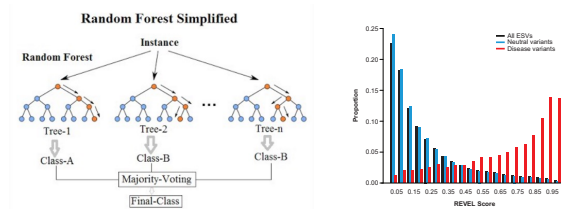
10

Umbrella methods - CADD



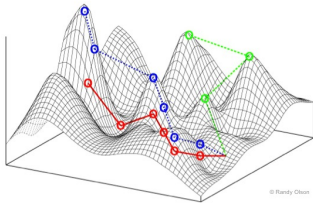
11

Umbrella methods - REVEL



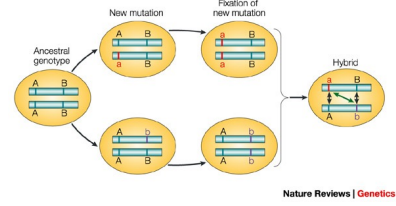
12

Ridges on the fitness landscape



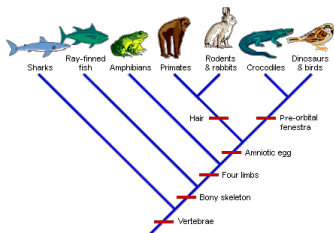
19

Dobzhansky-Muller incompatibility



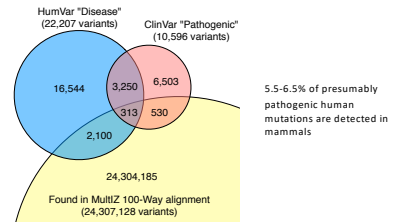
20

Looking at vertebrate species



21

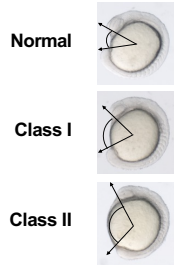
Many human pathogenic mutations are found in vertebrates



22

Zebrafish model

- Model of Bardet-Biedl Syndrome (obesity, renal failure, vision loss)
- Caused by defects in primary cilium
- Embryonic convergence / extension phenotype in zebrafish
- Easily scorable phenotype



Images: Phoebe

23

Testing double mutants

No injection		Human gene with disease mutant	
Knockdown		Double mutant (no suppression)	
Rescue with human gene		Double mutant (full suppression)	

Images: Phoebe

24

A newly identified gene

Clinical features

- Global developmental delay
- microcephaly
- feeding issues
- failure to thrive
- abnormal muscle tone
- low immunoglobulins
- frequent respiratory infections

Clinical testing

- normal female microarray
- metabolic testing – negative
- extensive genetic testing – negative

BTG2 <i>De novo</i>	TTN Compound het
NOS2 <i>De novo</i>	LAMA1 Compound het

Stephan Frangakis

25

The mutation is a reversal to the mammalian ancestral state

BTG2	R80	L128	Q140	V141	L142
<i>H. sapiens</i>	R	L	Q	V	L
<i>P. troglodytes</i>	•	•	•	•	•
<i>G. gorilla</i>	•	•	•	•	•
<i>M. musculus</i>	K	V	•	M	M
<i>R. norvegicus</i>	K	V	•	M	M
<i>H. glaber</i>	•	V	•	M	M
<i>S. domesticus</i>	K	V	•	M	M
<i>B. primigenius</i>	K	V	•	M	M
<i>E. ferus caballus</i>	K	V	•	M	M
<i>F. catus</i>	K	V	•	M	M
<i>C. lupus familiaris</i>	K	V	•	M	M
<i>D. novemcinctus</i>	K	V	•	M	M
<i>G. gallus</i>	K	P	•	M	M

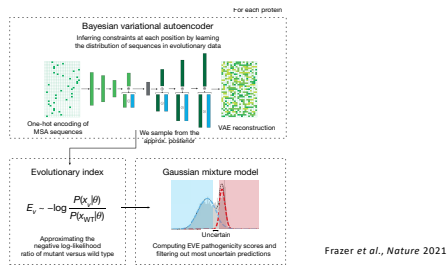
1

New methods directions

- Machine learning techniques have the potential to solve the epistasis problem
- Measures of population level constraint have the potential to solve the problem of distinguishing between strongly and weakly deleterious mutations.

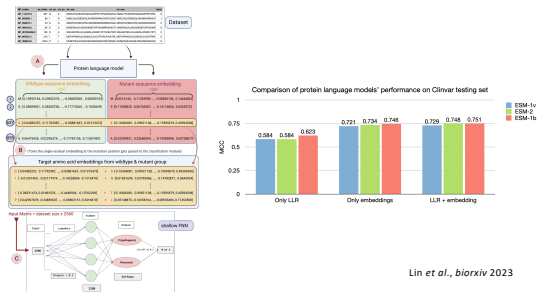
2

EVE – Variational Autoencoder



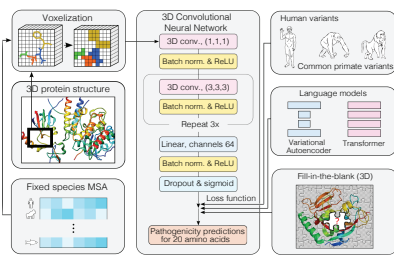
3

Large Language Models (VariPred)



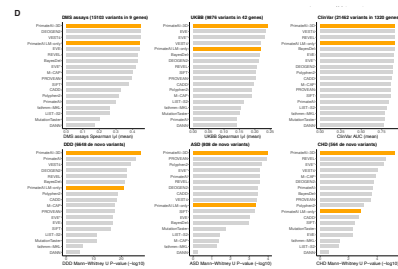
4

PrimateAI-3D



5

PrimateAI-3D



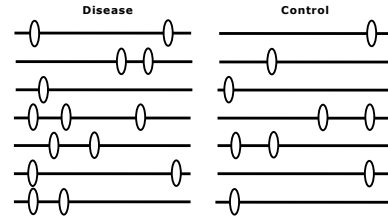
6

Applications

- Mendelian genetics
- Rare variant association studies

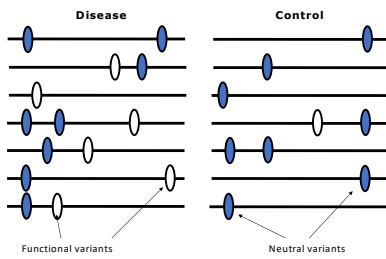
7

Rare variant collapsing study



8

Rare variant collapsing study



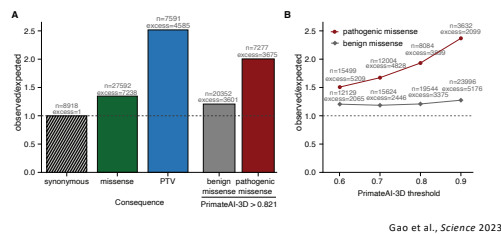
9

Predicting functional consequences increases power

- Inclusion of neutral variants reduces power of the test
- Combining variants with vastly different effect sizes reduces power of the test
- Most groups limit the tests to nonsense, splicing and missense variants that are predicted functional
- Assigning quantitative weights is probably a better approach, but nobody uses it in practice

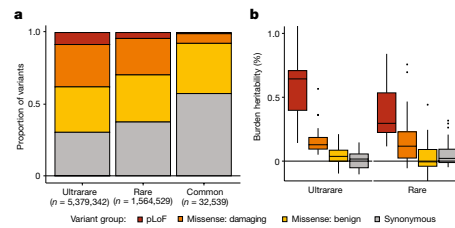
10

Damaging missense variants (as predicted by PrimateAI-3D) are enriched among de novo mutations in developmental disorders

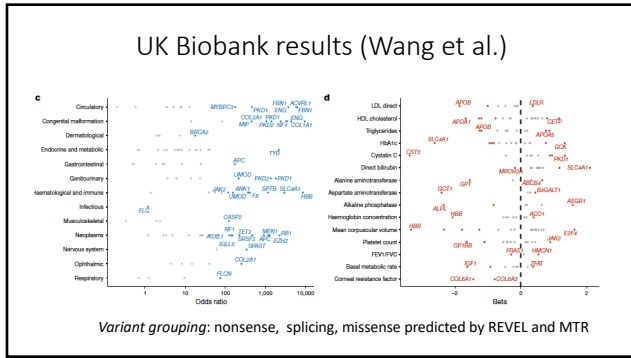


11

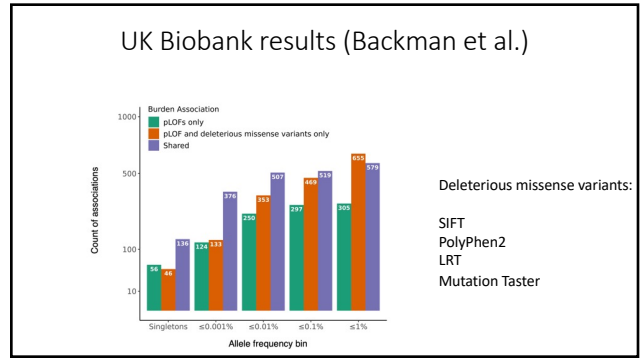
Burden heritability is significant for damaging missense variants (as predicted by PolyPhen2)



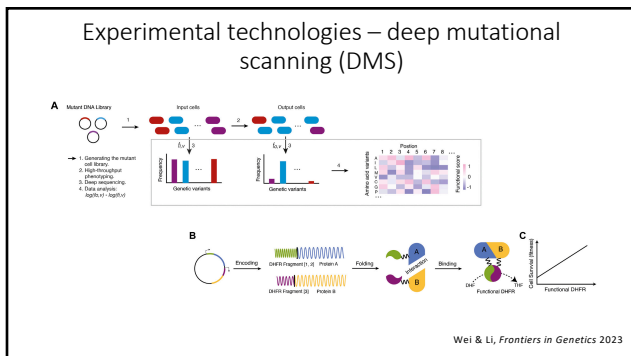
12



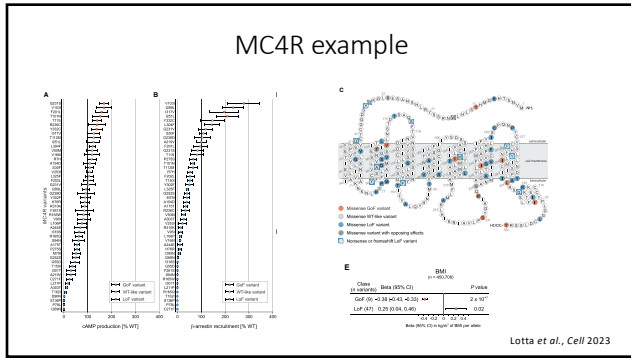
13



14



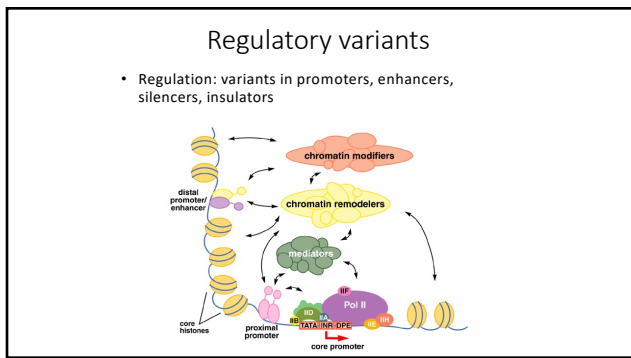
15



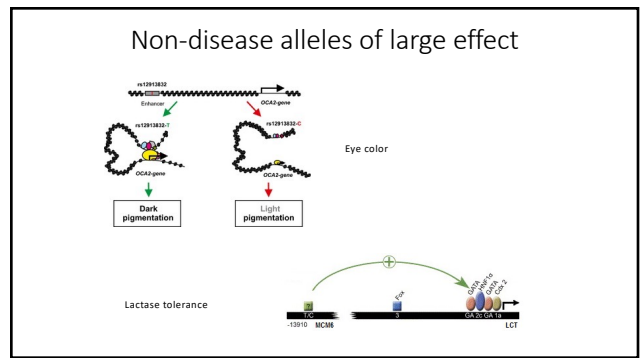
1

Non-coding variants

2



3



4

Ultraconserved elements

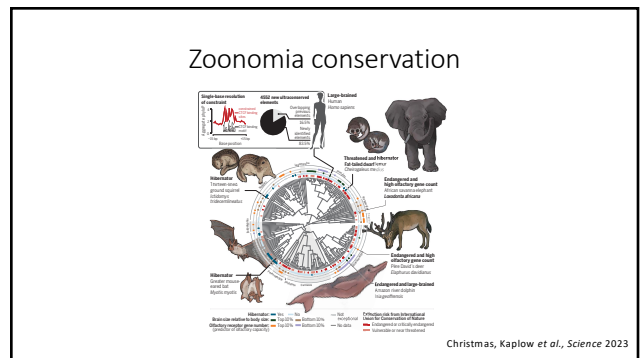
OPEN ACCESS | Study available on [PLOS](#)

Deletion of Ultraconserved Elements Yields Viable Mice

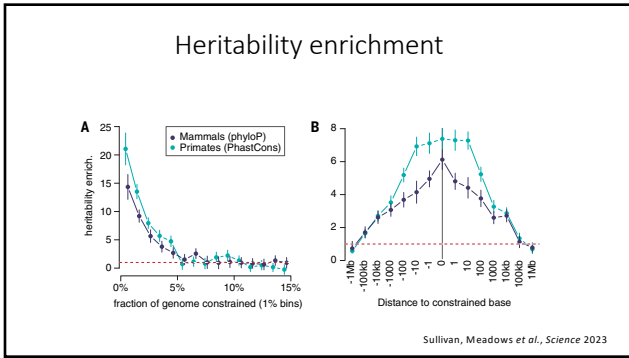
Hadar Ahissar^{1,2*}, Yuesen Zhu¹, Axel Visel¹, Amy Holt¹, Yuesen Afzal¹, Lea A. Pennacchio^{1,2}, Edward M. Rubin^{1,2*}

Abstract: Ultraconserved elements have been suggested to retain extended perfect sequence identity between the human, mouse, and rat genomes due to essential functional properties. To investigate the necessity of these elements in vivo, we removed four noncoding ultraconserved elements (ranging in length from 222 to 731 base pairs) from the mouse genome. To maximize the likelihood of observing a phenotype, we chose to delete elements that function as enhancers in a mouse transgenic assay and that are near genes that exhibit methylated phenotypes both when completely inactivated in the mouse and when their expression is altered due to other genomic modifications. Remarkably, all four resulting lines of mice lacking these ultraconserved elements were viable and fertile, and failed to reveal any critical abnormalities when assayed for a variety of phenotypes including growth, longevity, pathology, and metabolism. In addition, more targeted screens, initiated by the identification of deletions in mice in which genes in proximity to the investigated elements had been altered, also failed to reveal notable abnormalities. These results, while not inclusive of all the possible phenotypic impact of the deleted sequences, indicate that extreme sequence constraint does not necessarily reflect crucial functions required for viability.

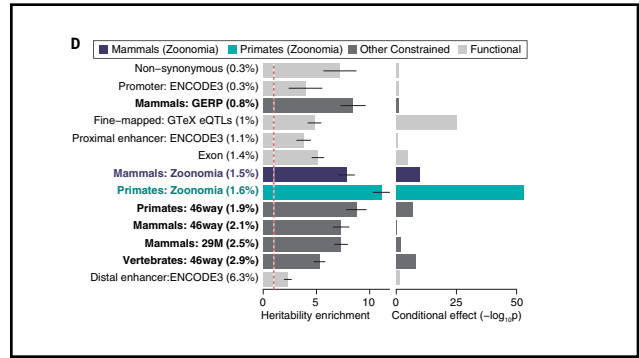
5



6



7



8

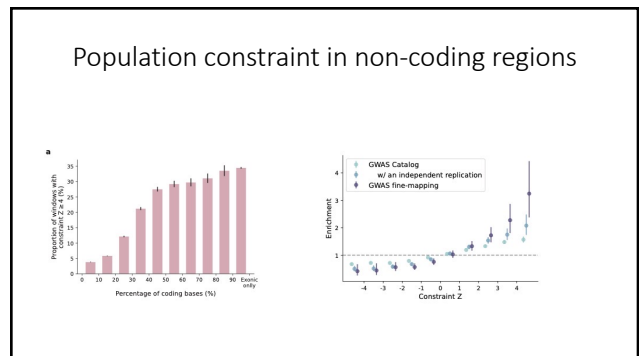
Population constraint in non-coding regions

Article
The sequences of 150,119 genomes in the UK Biobank
 A genome-wide mutational constraint map quantified from variation in 76,156 human genomes

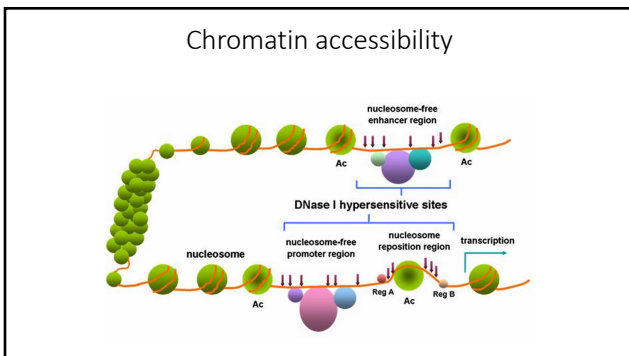
ARTICLE
 Extreme purifying selection against point mutations in the human genome

Noah Dahan^{1,3}, Mehreen R. Mughal^{1,3}, Ritika Ramanil¹, Yi-Fei Huang^{2,7} & Adam Sippel^{1,10}

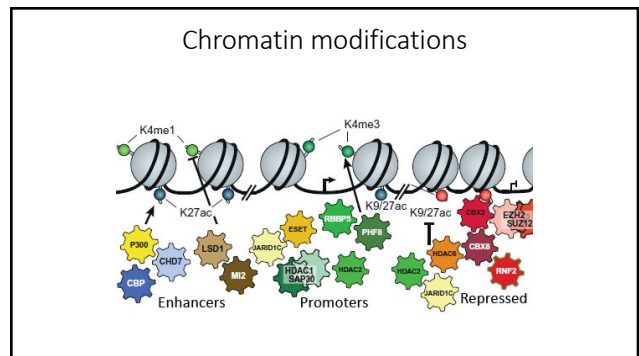
9



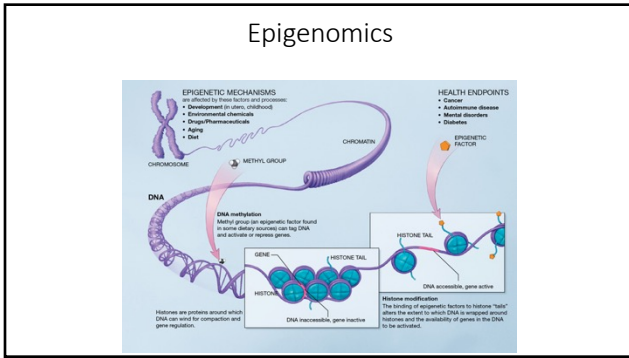
10



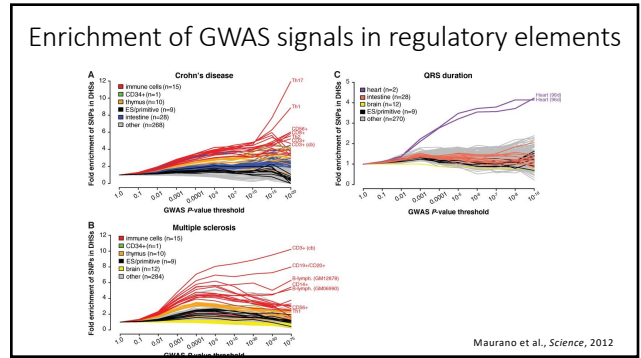
11



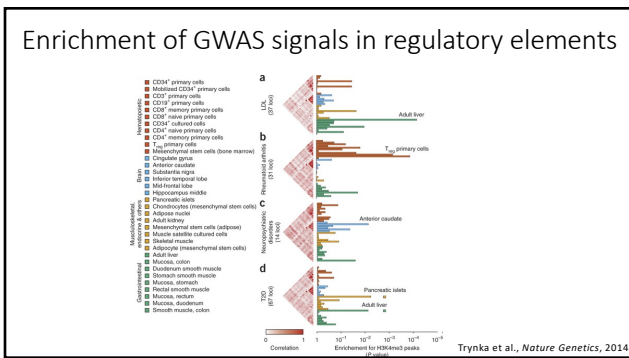
12



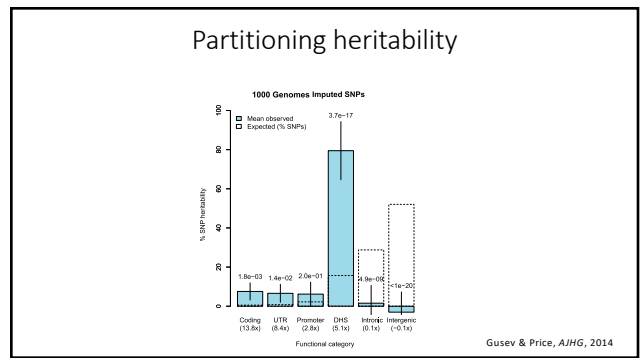
13



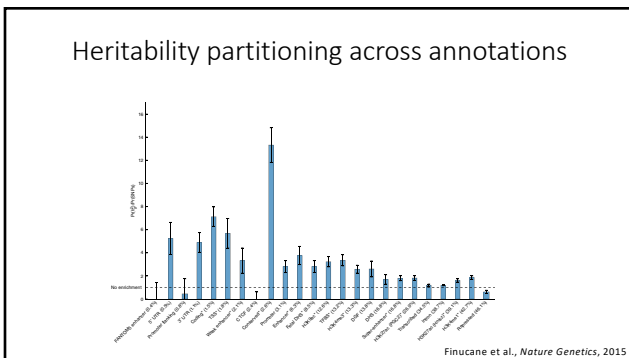
14



15



16



17

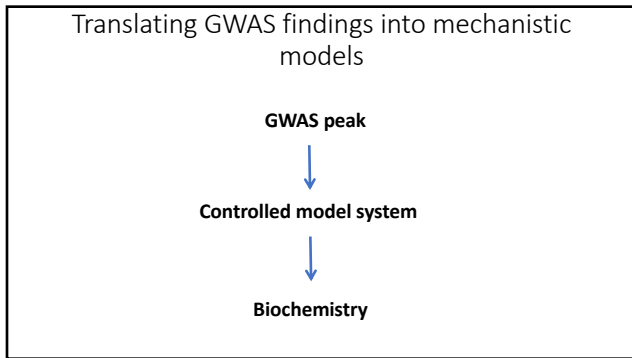
Application – function informed fine-mapping

Functionally informed fine-mapping and polygenic localization of complex trait heritability

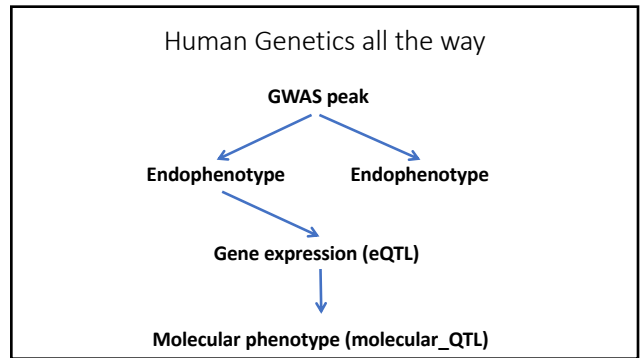
Omar Weissbrod^{1,2,3,4}, Farhad Hormozdiani^{1,2}, Christian Benner¹, Ran Gu¹, Jacob Ullrich^{1,2,4}, Steven Gazal^{1,2}, Armin P. Schoell¹, Boyce van de Geijn¹, Yuhai Ren^{1,2}, Carla Marques-Luna¹, Luke O'Connor¹, Matti Pirinen^{3,4,5}, Hilary K. Finucane^{3,4,5} and Alkes L. Price^{1,2,3,4,5}

- Estimate heritability enrichment and convert the estimates into prior probabilities
- Use these prior in fine-mapping (with SuSiE or FINEMAP)

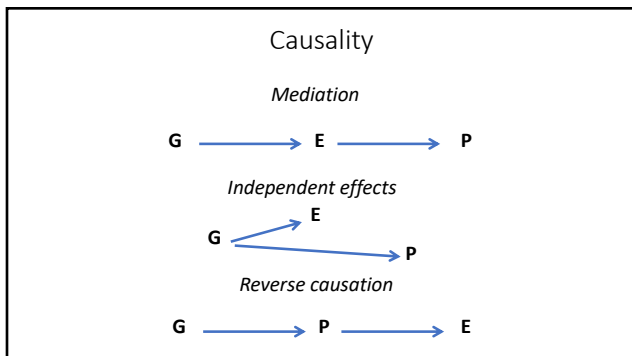
18



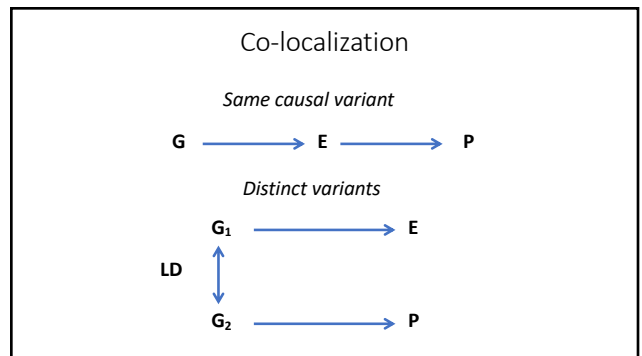
19



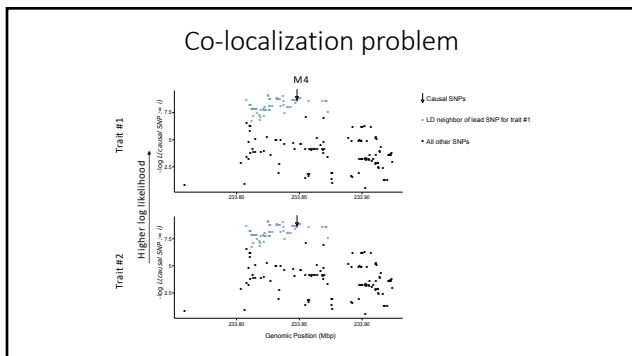
20



21



22



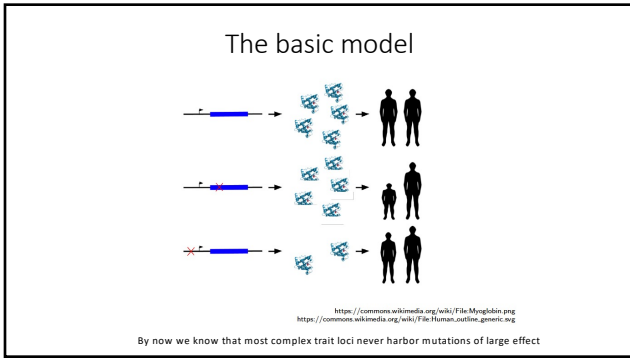
23

- ### Methods
- Coloc
 - eCAVIAR
 - JLIM

24

Genetic variants differ between Mendelian and complex traits

- | | |
|---|--|
| <ul style="list-style-type: none">• Complex trait variants <ul style="list-style-type: none">• Small effect size• Extremely large number of loci• Mostly non-coding (regulatory) | <ul style="list-style-type: none">• Mendelian & somatic cancer variants <ul style="list-style-type: none">• Large effect sizes• Small number of loci• Mostly coding• Are in “putatively causative” genes |
|---|--|



1

Hypothesis

- Most genes involved in Mendelian components of complex traits are also causative for cognate common forms.
- Variants involved in common forms alter regulatory sequence of these genes.
- This in turn induces changes in gene expression; regulatory variants are eQTLs.

2

Genes and phenotypes

(for complex traits, GWAS is restricted to non-coding variants)

Mend. trait	GWAS trait	Tissue
Breast cancer	Breast cancer	breast mammary tissue
Crohn's disease	Crohn's disease	small intestine terminal ileum colon sigmoid colon transverse
Dyslipidemia Hyperlipidemia Tanger's disease	HDL	liver adipose whole blood
Dwarfism	Height	skeletal muscle
Blood pressure	Blood pressure	heart atrial appendage kidney heart left ventricle
Dyslipidemia Hyperlipidemia	LDL	liver adipose tissue whole blood
Monogenic diabetes	Type II diabetes	pancreas skeletal muscle adipose whole blood
Ulcerative colitis	Ulcerative colitis	small intestine terminal ileum colon sigmoid colon transverse

Overall, 139 genes
 89 (64%) fall under a GWAS peak of a cognate complex trait

Examples include:
 LDL Receptor under a GWAS peak for LDL Cholesterol
 Estrogen receptor under a GWAS peak for breast cancer

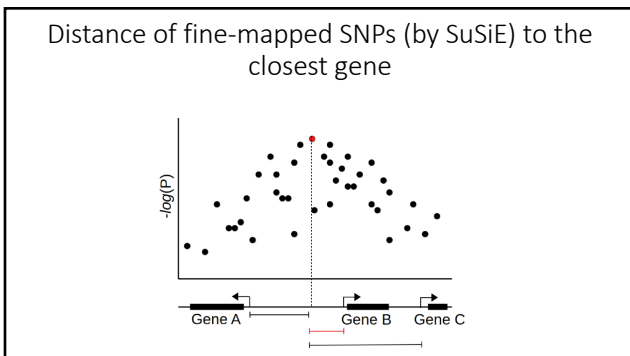
These genes are highly likely to mediate the effects of regulatory variants

3

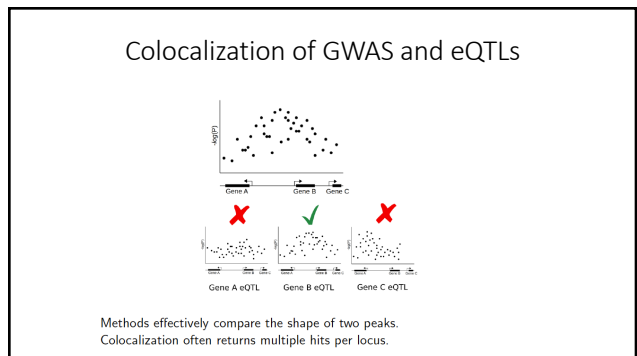
Statistical methods to locate the causative gene under GWAS peak

- Closest gene to peak
- Colocalization methods
 - ILIM
 - Coloc
 - eCAVIAR
- Transcriptome-wide association
 - FUSION
- Chromatin marks
 - Fine-mapping using SuSiE
 - Locate fine-mapped variants under chromatin modification peaks

4

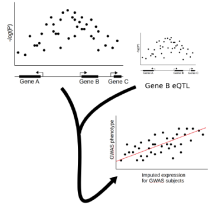


5



6

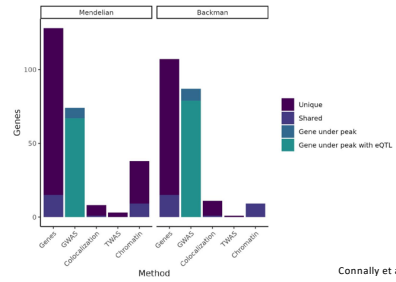
Transcriptome-wide association (TWAS)



TWAS often returns multiple hits per locus.

7

Results

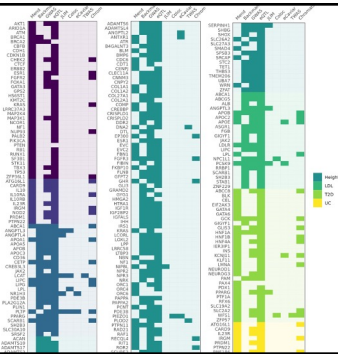


Connally et al., *eLife*, 2022

8

Our curated genes rarely colocalize

- This is true across all tested traits
- We also tried a chromatin method
 - It worked better
 - In large part because it favors the closest gene



9

But why?

Are eQTLs specific to...

- certain cell types?
- certain developmental stages?
- certain environmental conditions?

Are there inconsistent relationships...

- between gene expression and protein levels?
- between rate of transcription and gene expression?

10

I find it highly surprising that

- A context independent large change in expression of LDLR due to a nonsense mutation leads to a large phenotypic change
- A smaller change in expression does not affect LDL levels, while non-coding effect on LDLR does

11

Quantifying genetic effects on disease mediated by assayed gene expression levels

Douglas W. Yao^{1,2}, Luke J. O'Connor^{1,2,3}, Ailkes L. Price^{1,2,4} and Alexander Gusev^{1,2,4,5}

Feature Review

Where Are the Disease-Associated eQTLs?

Benjamin D. Umaru,^{1*} Alexis Battle,^{2,3*} and Yoav Gilad^{1,4*}

Limited overlap of eQTLs and GWAS hits due to systematic differences in discovery

Hakhamanesh Mostafaei^{1,3}, Jeffrey P. Spence¹, Sahin Nageci^{1,3}, Jonathan K. Pritchard^{1,3}

12

