Using genotypic data in our models

40[°] 00' 45'' N

Sarah Medland

SLP

Boulder 2024

Ecosystems (2003) 6: 31-45 DOI: 10.1007/s10021-002-0175-

Landscape Controls on Organic and Inorganic Nitrogen Leaching across an Alpine/Subalpine Ecotone, Green Lakes Valley, Colorado Front Range

ECOSYSTEMS

Overview

40[°] 00' 45''

Today I'll talk a bit about
How we get genotypic data
What we do to get it ready to use

Estimating relatedness

SLP

amblé síte

GWAS

PRS

40° 00' 45"

Locus/Variant/Marker

R

a given point in your genome

Sample site Stream

SLP

■ ≠ gene

40[°] 00' 45'' N

• Allele

I of the 2 copies of a variant at a

given locus

related concept – ambiguous snp

Etymology [edit]

The word "allele" is a short form of "allelomorph" ("other form", a word coined by British geneticists William Bateson and Edith Rebecca Saunders) in the 1900s,^{[7][8]} which was used in the early days of genetics to describe variant forms of a gene detected as different phenotypes. It derives from the Greek prefix ἀλληλο-, *allelo-*, meaning "mutual", "reciprocal", or "each other", which itself is related to the Greek adjective ἄλλος, *allos* (cognate with Latin *alius*), meaning "other".

SLP

ample site

n° 00' 45'

 Linkage disequilibrium/LD
 measure of whether an allele at one locus tends to be found more often with an allele at another locus.

related concept – LD block

40[°] 00' 45''

Terminology about locations within the genome
Base pair location/BP
Build

SLP

sample site

Centimorgan

strand

Obtaining genotypic data and getting it ready to use

Como cr

Sample site Stream

orest Cover

Kilont/sters

SLP

ARK-

40[°] 00' 45'' N

· NAV

How do we get genetic data?

n° 00' 45'

- 1. Recruit a large sample
 - From clinics
 - From the public
 - From an existing twin/cohort sample
 - Pay to access an existing sample
 - UK Biobank...

2. Collect information from the participants



40° 00' 45" N



Australian ASD and ADHD Study Investigating the impact of attention and behaviour on individuals and families

100%

Thinking back to when you were in primary school, do you remember having problems with paying attention or controlling your behaviour?

Survey Completion

O No

O Yes I had some problems, but I wasn't diagnosed with didn't have ADHD or ADD

O Yes I had some problems, and was diagnosed with ADHD or ADD

O I don't know



Contour Interval = 180m

SLP

3. Collect a DNA sample

How do we get genetic data?

40[°] 00' 45''



Kilonzeters Contour Interval = 180m

4. Extract the DNA

For non lab trained folk this webpage gives a lay overview that you might find helpful

https://learn.genetics.utah.edu/content/labs/extraction/howto/



How do we get genetic data?

II' 00' 45'

How do we get genetic data?

40[°] 00' 45"

5. Genotype the samples



Kilorzeters Contour Interval = 180m

5. Genotype the samples



IU, 00, 42,

- As DNA fragments pass over the BeadChip
- Each probe binds to a complementary sequence in the sample DNA, stopping one base before the locus of interest
- Single base extension that incorporates one of four labeled nucleotides
- When excited by a laser, the nucleotide label emits a signal
- The intensity of that signal conveys information about the allelic ratio at that locus

How do we get genetic data?

40 00 45"

5. Genotype the samples



What do we do to get it ready to use?

Quality Control

 Genotype level – allele frequency, missingness, Hardy-Weinberg Equilibrium (distribution of alleles)
 Sample level – missingness, heterozygosity, chromosomal distributions

Come back next year to learn how to do this

What do we do to get it ready to use?

10° 00' 45'

Imputation

- We pay to genotype ~.4-1M markers
- Through imputation we can get data for ~9M extra markers for free*

Reference set of haplotypes, for example, HapMap

04' 15" No

400011

4115

12 40° 00° 45" N

Genotype data with missing data at untyped SNPs (grey question marks)									ng tior	dat n m	ta a 1ari	at ks)				0 0 0 1 1 1 0 0 1	The reference haplotypes are used to impute alleles into the samples to cre imputed genotypes (orange)								ate	2						
1	?	?	?	1	?	1	?	0	2	2	?	?	2	?	0	1 1 1 0 0 1 0 0 1 1 1 0 1 1 1 0	1	1	1	1	1	2	1	0	0	2	2	0	2	2	2	0
0	?	?	?	2	?	2	?	0	2	2	?	?	2	?	о	7 1 7	0	0	1	0	2	2	2	0	0	2	2	2	2	2	2	0
1	?	?	?	2	?	2	?	0	2	1	?	?	2	?	0		1	1	1	1	2	2	2	0	0	2	1	1	2	2	2	0
1	· ?	2	· ?	2	· 2	1	· 2	1	2	2	2	· ?	2	2	0		1	1	2	0	2	2	1	0	1	2	2	1	2	2	2	0
2	•	•	•	2	•	2	•	1	2	1	•	•	2	•	0	↓ /	2	2	2	2	2	1	2	0	1	2	1	1	2	2	2	0
2	۲ ۵	÷	f D	4	۰ ۲	4	f D	1	4	÷	:	۲ ۲	4	:	0	Each sample is phased and the haplotypes	1	1	1	0	1	2	1	0	1	2	2	1	2	2	2	0
T	?	-	?	1	?	T	?	T	2	2	?	?	2	?	0	are modelled as a mosaic of those in the haplotype reference panel	1	1	2	1	2	1	2	0	0	2	1	1	1	2	1	1
1	?	?	?	2	?	2	?	0	2	1	?	?	2	?	1		2	2	2	1	1	1	1	0	1	2	1	0	1	2	1	1
2	?	?	?	1	?	1	?	1	2	1	?	?	2	?	1	0 ? ? ? 1 ? 1 ? 0 1 1 ? ? 1 ? 0	1	2	2	0	0	2	0	0	2	2	2	1	2	2	2	0
1	?	?	?	0	?	0	?	2	2	2	?	?	2	?	0	1 ? ? ? 1 ? 1 ? 0 1 1 ? ? 1 ? 0	-	2	2	~	0	2	0	~	2	2	2	+	2	2	4	Ŭ
																1 ? ? 1 ? 1 ? 0 1 0 ? ? 1 ? 0 1 ? ? ? 1 ? ? 1 ? ? 1 ? ? 1 ? ? 1 ? ? 1 ? ? 1 ? ? 1 ? ? 1 ? ? 1 ? ? 1 ? ? 1 ? ? 1 ? ? ? 0 1 ? ? ? ? 0 ? 1 1 ? ? 1 ? ? 0 1 ? ? ? ? 0 ? ? 1 1 ? ? 1 ? ? 0 ? 0 1 ?																

Contour Interval = 180m

Imputation reference sets

ff" 00° 45'

- Publicly Available References
 - HapMap
 - IKGP phase 3 version v5
- References only available via custom imputation servers
 - HRC 64,976 haplotypes 39,235,157 SNPs
 - CAPPA African American/Caribbean
 - Multi-ethnic HLA
 - Genome Asia Pilot GAsP
 - TopMed 97,256 haplotypes 308,107,085 SNPs (b38)



DIY – Use a cookbook!

http://genome.sph.umich.edu/wiki/Minimac3_Imputation_Cookbook OR http://genome.sph.umich.edu/wiki/IMPUTE2:_1000_Genomes_Imputation_Cookbook

- UMich Imputation Server
 - https://imputationserver.sph.umich.edu/
- Sanger Imputation Server
 - https://imputation.sanger.ac.uk/
- TOPMed Imputation Server

un° 0.01 4.51

https://imputation.biodatacatalyst.nhlbi.nih.gov/



2 commonly used genotype output formats

Hard call or best guess

Dosage data (most common – 1 number per SNP, 1-2)

##fileformat=VCFv4.1 ##filedate=2015.7.12 ##source=Minimac3 ##FORMAT=<ID=GT, Number=1, Type=String, Description="Genotype"> ##FORMAT=<ID=DS,Number=1,Type=Float,Description="Estimated Alternate Allele Dosage : [P(0/1)+2*P(1/1)]"> ##FORMAT=<ID=GP,Number=3,Type=Float,Description="Estimated Posterior Probabilities for Genotypes 0/0, 0/1 and 1/1 "> ##INFO=<ID=MAF,Number=1,Type=Float,Description="Estimated Alternate Allele Frequency"> ##INFO=<ID=R2,Number=1,Type=Float,Description="Estimated Imputation Accuracy"> ##INFO=<ID=ER2.Number=1 Type=Float.Description="Empirical (Leave-One-Out) R-square (available only for genotyped variants)"> #CHROM A0007 A0007 POS FORMAT A0001 A0001 A0004 A0004 A0009 A0009 A0010 A0003 A000: 10 27754636 10:27754636 PASS MAF=0.00032:R2=0.81788 0/0:0.000:1.000,0.000,0.000 0/0:0.000:1.000,0.000,0.000 GT:DS:GP 10 27754678 10:27754678 PASS MAF=0.00042:R2=0.77190 GT:DS:GP 0/0:0.000:1.000,0.000,0.000 0/0:0.000:1.000,0.000,0.000 10 27754849 10:27754849 PASS MAF=0.00001;R2=0.00262 GT:DS:GP 0/0:0.000:1.000,0.000,0.000 0/0:0.000:1.000,0.000,0.000 10 10:27754857 PASS 27754857 MAF=0.00120;R2=0.72916 GT:DS:GP 0/0:0.000:1.000.0.000.0.000 0/0:0.000:1.000.0.000.0.000 10 27754954 PASS 10:27754954 MAF=0.11410;R2=0.97841 GT:DS:GP 1/1:2.000:0.000,0.000,1.000 1/1:2.000:0.000,0.000,1.000 10 27755014 10:27755014 PASS MAF=0.00000;R2=0.00082 GT:DS:GP 0/0:0.000:1.000,0.000,0.000 0/0:0.000:1.000,0.000,0.000 10 27755016 10:27755016 PASS MAF=0.00003:R2=0.01909 GT:DS:GP 0/0:0.000:1.000,0.000,0.000 0/0:0.000:1.000,0.000,0.000 10 27755047 10:27755047 PASS MAF=0.02255;R2=0.87665 GT:DS:GP 0/0:0.000:1.000,0.000,0.000 0/0:0.000:1.000,0.000,0.000 10 27755175 10:27755175 PASS MAF=0.00004;R2=0.13821 GT:DS:GP 0/0:0.000:1.000,0.000,0.000 0/0:0.000:1.000,0.000,0.000 10 27755281 10:27755281 PASS MAF=0.00061;R2=0.86168 GT:DS:GP 0/0:0.000:1.000,0.000,0.000 0/0:0.000:1.000,0.000,0.000 10 27755330 10:27755330 PASS MAF=0.00273;R2=0.90295 0/0:0.000:1.000,0.000,0.000 GT:DS:GP 0/0:0.000:1.000,0.000,0.000 10 27755439 10:27755439 PASS MAF=0.00000;R2=0.00138 GT:DS:GP 0/0:0.000:1.000,0.000,0.000 0/0:0.000:1.000,0.000,0.000 10 27755489 10:27755489 PASS MAF=0.00003;R2=0.39172 GT:DS:GP 0/0:0.000:1.000.0.000.0.000 0/0:0.000:1.000.0.000.0.000

Not all markers are well imputed

Imputation quality evaluation

Minimac hides each of the genotyped SNPs in turn and then calculates 3 statistics:

IooRSQ - this is the estimated rsq for that SNP (as if SNP weren't typed).

empR - this is the empirical correlation between true and imputed genotypes for the SNP. If this is negative, the SNP alleles are probably flipped.

empRSQ - this is the actual R2 value, comparing imputed and true genotypes.

These statistics can be found in the *.info file

Be aware that, unfortunately, imputation quality statistics are not directly comparable between different imputation programs (MaCH/minimac vs. Impute vs. Beagle etc.).

SNP	A11	A12	Freq1	MAF	AvgCall	Rsq	Genotype	d	LooRsq	EmpR	EmpRsq	Dose1	Dose2	
1:10583	G	A	0.79288	0.20712	0.79288	-0.00000)	-	-	-	-	-	-	
1:10611	С	G	0.97889	0.02111	0.97889	0.00000	-	-	-	-	-	-		
1:13302	С	Т	0.86280	0.13720	0.86280	-0.00000)	-	-	-	-	-	-	
1:13327	G	С	0.96042	0.03958	0.96042	-0.00000)	-	-	-	-	-	-	
		-	-											
1:95207	182	Т	С	0.99547	0.00453	0.99547	0.10108	-	-	-	-	-	-	
1:95207	382	Т	Т	1.00000	0.00000	1.00000	0.00000	-	-	-	-	-	-	
1:95207	442	С	Т	0.62754	0.37246	0.99999	1.00507	Genotyp	ed	0.98810	0.99822	0.99645	0.99484	0.00421
1:95207	524	G	Α	0.78061	0.21939	1.00000	1.00511	Genotyp	ed	1.00059	1.00000	1.00000	0.99924	0.00083
1:95207	532:TG_T	R	D	0.78620	0.21380	0.99441	0.97729	-	-	-	-	-	-	
1:95207	558	С	Т	0.99399	0.00601	0.99399	0.05165	-	-	-	-	-	-	
1:95207	633	Α	С	0.93366	0.06634	0.99998	1.00482	Genotyp	ed	0.94847	0.99901	0.99802	0.99621	0.00372
1:95207	846	G	Т	0.98937	0.01063	0.98942	0.31316	-	-	-	-	-	-	

SLP

- Imputation accuracy is calculated differently by the two main imputation programs
- But is highly correlated and conceptually the same

The IMPUTE info measure I_A

This is based on measuring the relative statistical information about the population allele frequency, θ_j . If the G_{ij} 's were observed then the full data likelihood is given by

$$L(\theta_j) = \prod_{i=1}^{N} \theta_j^{G_{ij}} (1 - \theta_j)^{2 - G_{ij}}$$
(10)

For this likelihood the score and information are given by

$$U(\theta_j) = \frac{d \log L(\theta_j)}{d\theta_j} = \frac{X - 2N\theta_j}{\theta_j(1 - \theta_j)}$$
(11)

$$I(\theta_j) = \frac{-d^2 \log L(\theta_j)}{d\theta_j^2} = \frac{X}{\theta_j^2} + \frac{2N - X}{(1 - \theta_j)^2}$$
(12)

The IMPUTE info measure is based on the same idea used to calculate the SNPTEST information measure i.e. the ratio of the observed and complete information.

$$I_A = \frac{\mathbb{E}_{G,j}[I(\hat{\theta})] - V_G[U(\hat{\theta})]}{\mathbb{E}_{G,j}[I(\hat{\theta})]}$$
(13)

The MACH \hat{r}^2 measure

00" 45" N

This is the ratio of the empirically observed variance of the allele dosage to the expected binomial variance at Hardy-Weinberg equilibrium. At the jth SNP this is defined as

$$\hat{r}_{j}^{2} = \begin{cases} \frac{\sum_{i=1}^{N} e_{ij}^{2}}{N} - \left(\frac{\sum_{i=1}^{N} e_{ij}}{N}\right)^{2}}{2\hat{\theta}(1-\hat{\theta})} & \text{when } \hat{\theta} \in (0,1) \\ 1 & \text{when } \hat{\theta} = 0, \hat{\theta} = 1 \end{cases}$$
(1)

When all the genotypes are predicted with high certainty this ratio will be close to 1, although it can go above 1 (Figure 1). As the amount of uncertainty increases the allele dosages will tend to 2θ , the empirical variance will tend to 0 and so \hat{r}^2 tends to 0.

where the expectations are taken over the imputed genotype distribution and evaluated at the allele frequency estimate, $\hat{\theta}_j$. The exact terms are given by

$$\mathbb{E}_{G_{ij}}[I(\hat{\theta})] = \frac{2N}{\hat{\theta}(1-\hat{\theta})}$$

$$V_G[U(\hat{\theta})] = \frac{\sum_{i=1}^{N} (f_{ij} - e_{ij}^2)}{\hat{\theta}^2 (1-\hat{\theta})^2}$$
(14)
(14)
(15)

so that

$$I_{A} = \begin{cases} 1 - \frac{\sum_{i=1}^{N} (f_{ij} - e_{ij}^{2})}{2N\hat{\theta}(1 - \hat{\theta})} & \text{when } \hat{\theta} \in (0, 1) \\ 1 & \text{when } \hat{\theta} = 0, \hat{\theta} = 1. \end{cases}$$
(16)

So I_A is bounded above at 1 and will equal 0 when the sample mean variance of the imputed genotypes equals the variance you would expect if alleles where sampled with frequency $\hat{\theta}$.

Contour Interval = 100m



- After imputation you need to check that it worked and the data look ok
- Things to check
 - Plot r² across each chromosome look to see where it drops off
 - Plot MAF-reference MAF



Estimating relatedness

Como Cre

Sample site Stream

^zorest Céver

Kilometers ontour Interval

SLP

Slide Credit: Loci Yengo

40⁰ 00' 45" N

ARK-NAV

Genetic relationship matrices

ff" 00° 45

- Genetic relationship matrices (GRM) are important tools for estimating heritability
- Yesterday we used family level GRM based on expected relatedness



Genetic relationship matrices

- We can also calculate a GRM using SNP data.
- There are many ways to calculate a GRM using SNP data
- Common to use the standard estimator implemented in the software GCTA (but often calculated using Plink)

$$\hat{\pi}_{jk} = \frac{1}{m} \sum_{i} \frac{(x_{ij} - 2p_i)(x_{ik} - 2p_i)}{2p_i(1 - p_i)}$$

where, x_{ij} and x_{ik} are the minor allele count (x_{ij} , x_{ik} = 0,1 or 2) at SNP i for individuals j and k respectively, p_i the minor allele frequency (MAF) of SNP I and m the number of SNPs used to calculate the GRM.

Example of GRM between N=3 individuals m=1000 SNPs

[\$bash] zless myGRM.grm.gz 1 1 1000 0.99 1 2 1000 -0.01 1 3 1000 0.01 2 2 1000 1.03 2 3 1000 0.03 3 3 1000 1.01

SLP

<u>Am J Hum Genet.</u> 2011 Jan 7; 88(1): 76–82. doi: <u>10.1016/j.ajhg.2010.11.011</u>

40[°] 00' 45'' N ...

PMCID: PMC3014363 PMID: <u>21167468</u>



Jian Yang, 1,* S. Hong Lee, 1 Michael E. Goddard, 2,3 and Peter M. Visscher 1

40° 00' 45" N

The expectation (over a large sample of relatives) of the $\hat{\pi}_{jk}$ = the expected relatedness

NOT ST

Observed relatedness may be still vary within a type of pedigree relationship.



SLP

Genetic relationship matrices

- Later today we will use relatedness calculated from SNPs in an OpenMx model to estimate heritability
- On Thursday we will use relatedness calculated from SNPs to run Trio-GCTA

Analysing our genotypic data...

Como Cre

Sample site Stream

orest Cover

Kilom/sters

SLP

ARK-NAV

40° 00' 45" N

Association analyses

40° 00" 45" N



cases (n=1,000) people with heart disease



controls (n=1,000) people without heart disease

https://www.yourgenome.org/theme/genome-wide-association-studies/

controls



49% C 51%

SLP

1.0 Klionzeters Contour Interval = 180m

Association analyses

40° 00" 45" N



Genome-wide analyses of ADHD identify 27 risk loci, refine the genetic architecture and implicate several cognitive domains

Ditte Demontis ⁵², G. Bragi Walters, Georgios Athanasiadis, Raymond Walters, Karen Therrien, Trine Tollerup Nielsen, Leila Farajzadeh, Georgios Voloudakis, Jaroslav Bendl, Biau Zeng, Wen Zhang, Jakob Grove, Thomas D, Als, Jinjie Duan, F. Kyle Satterstrom, Jonas Bybjerg-Grauholm, Marie Bækved-Hansen Olafur O, Gudmundsson, Sigurdur H. Magnusson, Gisli Baldursson, Katrin Davidsdottir, Gyda S, Haraldsdottir, Esben Agerbo, Gabriel E. Hoffman, ADHD Working Group of the Psychiatric Genomics Consortium, iPSYCH-Broad Consortium, ... Anders D. Børglum ⁵² + Show authors

Contour Interv

Genome-wide Association Study (GWAS)

40 00 45" 1

16 -



SLP

Genome-wide analyses of ADHD identify 27 risk loci, refine the genetic architecture and implicate several cognitive domains

Ditte Demontis 🖾, G. Bragi Walters, Georgios Athanasiadis, Raymond Walters, Karen Therrien, Trine Tollerup Nielsen, Leila Farajzadeh, Georgios Voloudakis, Jaroslav Bendl, Biau Zeng, Wen Zhang, Jakob Grove, Thomas D. Als, Jinije Duan, F. Kyle Satterstrom, Jonas Bybierg-Grauholm, Marie Bækved-Hansen Olafur O, Gudmundsson, Sigurdur H, Magnusson, Gisli Baldursson, Katrin Davidsdottir, Gvda S, Haraldsdottir, Esben Agerbo, Gabriel E. Hoffman, ADHD Working Group of the Psychiatric Genomics Consortium, iPSYCH-Broad Consortium, ... Anders D. Børglum 🖾 🕇 Show authors

Nature Genetics 55, 198–208 (2023) Cite this article

16

15

14

13

12

11

9

8

6

5

4 3

40° 00' 45" N -

log10 (p) 10

38,899 people living with ADHD, 186,843 without

Mapping genomic loci implicates genes and synaptic biology in schizophrenia

Vassily Trubetskoy, Antonio F. Pardiñas, Ting Qi, Georgia Panagiotaropoulou, Swapnil Awasthi, Tim B. Bigdeli, Julien Bryois, Chia-Yen Chen, Charlotte A. Dennison, Lynsey S. Hall, Max Lam, Kyoko Watanabe, <u>Oleksandr Frei, Tian Ge, Janet C. Harwood, Frank Koopmans, Sigurdur Magnusson, Alexander L. Richards</u> Julia Sidorenko, Yang Wu, Jian Zeng, Jakob Grove, Minsoo Kim, Zhigiang Li, Indonesia Schizophrenia Consortium, PsychENCODE, Psychosis Endophenotypes International Consortium, The SynGO Consortium Schizophrenia Working Group of the Psychiatric Genomics Consortium + Show authors

Nature 604, 502–508 (2022) Cite this article

74,776 people living with Schizophrenia, 101,023 without



GWAS typically uses a significance threshold of 5x10⁻⁸

'00" 45" N

- This is based on the approximate number of independent tests conducted
 - Caveat: Many/most follow-up analyses use full distribution of effect estimates and don't restrict to significant loci

To account for multiple testing in genome-wide association studies (GWAS), a fixed P-value threshold of 5×10^{-8} is widely used to identify association between a common genetic variant and a trait of interest. Risch and Merikangas (1996) suggested this strict P-value threshold for studying the genetics of complex diseases due to the many false positive discoveries reported by candidate gene studies at that time. Later, the International HapMap Consortium (Altshuler and Donnelly 2005), Dudbridge and Gusnanto (2008), and Pe'er *et al.* (2008) independently suggested near-identical thresholds for common variant (minor allele frequency [MAF] >5%) GWAS. Each group of investigators sought to control the family-wise error rate

Zhongsheng Chen, Michael Boehnke, Xiaoquan Wen, Bhramar Mukherjee, Revisiting the genome-wide significance threshold for common variant GWAS, *G3 Genes*|*Genomes*|*Genetics*, Volume 11, Issue 2, February 2021, jkaa056, https://doi.org/10.1093/g3journal/jkaa056

Finding variants influencing the trait is important BUT

- The next steps to determine function and mechanism are \$\$\$\$\$
- Better to have *robust* findings that will replicate than find *more* variants

Publicly available summary statistics

un° 0.01 4.51

- Most GWAS efforts make their results publicly available
 - http://www.nealelab.is/uk-biobank
 - <u>https://pgc.unc.edu/for-</u> researchers/download-results/
 - https://enigma.ini.usc.edu/research/
 - download-enigma-gwas-results/

GWAS output = PRS input GWAS output = MR input GWAS output = LDscore input GWAS output = SNP h² input

. . .

Como Cre

Sample site Stream

orest Cover

Kilondeters

ARK

40⁰ 00' 45" N

-NAV

Polygenic Scores

 $40^{\circ} 00' 45''$

Perline will be talking more about this tomorrow morning

SLP

8am

Polygenic Scores

- Many names same concept
 - Polygenic Risk Scores (PRS)
 - Polygenic Scores (PGS)
 - Allelic Scores
 - Polygenic Index (PCI
- Polygenic risk score Weighted sum of alleles which quantify the effect of several genetic variants on an individual's phenotype.





Genotype=AG Genotype=GG 150 160 170 180 190 Height



In a new sample we would expect AG individuals to be on average 2cm taller than AA and 2cm shorter than GG















Complex traits are highly polygenic!

From above we can see there are many more genetic variants that contribute to the phenotype

Common variants typically have a small effect size (our example is an exaggeration for a common variant!). This would cause single-loci based prediction useless

We can combine the information we gain from several genetic variants to estimate an overall score and gain a better estimate of the trait. This is essentially what a PRS does

Polygenic Scores

- A couple of important gotchas
 - You need to make sure the weights are being applied to the right allele (ambiguous snps)
 - The individuals you are calculating the PRS for needs to be <u>completely</u> independent from the individuals in the GWAS

Polygenic Scores

40[°] 00' 45''

Numerous methods

Clumping and thresholding

SLP

Bayesian approaches

Article | Open Access | Published: 08 November 2019

Improved polygenic prediction by Bayesian multiple regression on summary statistics

Luke R. Lloyd-Jones ⊠, Jian Zeng ⊠, Julia Sidorenko, Loïc Yengo, Gerhard Moser, Kathryn E. Kemper, Huanwei Wang, Zhili Zheng, Reedik Magi, Tõnu Esko, Andres Metspalu, Naomi R. Wray, Michael E. Goddard, Jian Yang ⊠ & Peter M. Visscher ⊠

Nature Communications 10, Article number: 5086 (2019) Cite this article

- Combines a likelihood connecting the joint effects with GWAS summary statistics and a finite mixture of normal distribution priors for marker effects.
- Models the SNP effect sizes as a mixture of normal distributions with mean zero and different variances.
- Requires GWAS summary statistics with FREQ, BETA, SE and N; and an LD reference matrix



Lloyd-Jones, Jian Zeng, et al (2019)

LDpred2: better, faster, stronger 👌

Florian Privé 💌, Julyan Arbel, Bjarni J Vilhjálmsson 🐱

Bioinformatics, Volume 36, Issue 22-23, 1 December 2020, Pages 5424–5431,

https://doi.org/10.1093/bioinformatics/btaa1029

Published: 16 December 2020 Article history •

Addressed instability issues in LDpred providing a more stable workflow. Models long range LD such as that found near the HLA region.

Also derives an expectation of joint effects given are marginal effects and correlation between SNPs

 $\widehat{oldsymbol{\gamma}}_{ ext{joint}} = oldsymbol{S}^{-1}oldsymbol{R}^{-1}oldsymbol{S}\widehat{oldsymbol{\gamma}}_{ ext{marg}}$

Assumes:

$$\beta_j = S_{j,j} \gamma_j \sim \ \begin{cases} \mathscr{N}\left(0, \frac{h^2}{Mp}\right) & \text{with probability p,} \\ 0 & \text{otherwise,} \end{cases}$$



Bioinformatics, Volume 36, Issue 22-23, 1 December 2020, Pages 5424– 5431

With p= proportion of causal variants and h^2 estimated using Ldscore regression. Grid for p: p(1, 0.3, 0.1, 0.03, 0.01, 0.003 and 0.001).

Polygenic Scores

n° 00' 45'

Working across samples

- Replicating the finding
 - Requires the same snps, same weights etc
 - Restrict the pool of possible SNPs to those available across all cohorts
- Replicating the concept
 - Different clumping and/or weights



PRS – trait association

Think about your sample: > Is it a family based sample?

- Adjust for relatedness e.g. LMM
- > Is it homogeneous in terms of ancestry?
 -Always a good idea to adjust for genetic PCs
 > Does it match the GWAS ancestry?

Think about your trait:

- > Is it continuous linear regression
- > Binary logistic or probit regression
- > Ordinal cumulative linked mixed models
 > Always remember potential confounders of the trait and of the discovery GWAS

PGC-MDD1: N=18k max variance explained = 0.08%, p=0.018



PGC-MDD2: N=163k max variance explained =0.46%, p= 5.01e-08

Colodro-Conde L, Couvy-Duchesne B, et al, (2017) *Molecular Psychiatry*



C+T also allows us to explore the pattern of variance explained

Variance explained = partial R² for quantitative traits. Different ways of estimating it for binary traits

Applications

40° 00' 45'

- Quantify variance explained & explore architecture
- Risk stratification
 - (i.e. identifying people to later test for specific disease)
- Aid in clinical diagnosis
- Test for genetic overlap between traits
 - (e.g. does a Depression PRS predict cardiovascular disease?)
- Trait imputation when not measured
 - (obviously imperfect and dependent on heritability)
- Personalized treatment
 - (GWAS on treatment response are gaining power)
- Any hypothesis where you rely on a risk or liability
 - (e.g. GxE interactions)

Additional considerations in mental health

- When people present it is not usually a question of if someone has a future risk
 - Individuals or their families typically seek help
- Individuals often present with symptoms that might fit more than one diagnostic criterion
 - So the question is usually one of differential diagnosis

How is mental health different?

00'45

- Comorbidity (co-occurrence) is the rule than the exception
- Presentations and diagnoses are expected to change over time
 - This is ≠ misdiagnosis
 - PRS is static across the lifespan

Further complications...

10° 00' 45'

- Diagnoses are used for more than just treatment planning
 - Medico-legal contexts
 - Criminal, Civil and Family proceedings
 - Compensation
 - Access to support
 - Financial, Educational, Housing, Social

Further complications...

40° 00' 45'

- Because of this, diagnostic processes are ideally
 - Static
 - Reproducible
 - Easy to implement
 - Useable in low resource settings

SI P

Measurement invariant

- Wray NR, Goddard, ME, Visscher PM. Prediction of individual genetic risk to disease from genome-wide association studies. Genome Research. 2007; 7(10):1520-28.
- Evans DM, Visscher PM., Wray NR. <u>Harnessing the information contained within genome-wide association studies to improve individual prediction of complex disease risk</u>. Human Molecular Genetics. 2009; 18(18): 3525-3531.
- International Schizophrenia Consortium, Purcell SM, Wray NR, Stone JL, Visscher PM, O'Donovan MC, Sullivan PF, Sklar P. Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. Nature. 2009; 460(7256):748-52
- Evans DM, Brion MJ, Paternoster L, Kemp JP, McMahon G, Munafò M, Whitfield JB, Medland SE, Montgomery GW; GIANT Consortium; CRP Consortium; TAG Consortium, Timpson NJ, St Pourcain B, Lawlor DA, Martin NG, Dehghan A, Hirschhorn J, Smith GD. <u>Mining the human phenome using allelic scores that index biological intermediates</u>. PLoS Genet. 2013,9(10):e1003919.
- Dudbridge F. <u>Power and predictive accuracy of polygenic risk scores</u>. PLoS Genet. 2013 Mar;9(3):e1003348. Epub 2013 Mar 21. Erratum in: PLoS Genet. 2013;9(4). (Important discussion of power)
- Wray NR, Lee SH, Mehta D, Vinkhuyzen AA, Dudbridge F, Middeldorp CM. <u>Research review: Polygenic methods and their application to psychiatric traits.</u> J Child Psychol Psychiatry. 2014;55(10):1068-87. (Very good concrete description of the traditional methods).
- Wray NR, Yang J, Hayes BJ, Price AL, Goddard ME, Visscher PM. <u>Pitfalls of predicting complex traits from SNPs</u>. Nat Rev Genet. 2013;14(7):507-15. (Very good discussion of the complexities of interpretation).
- Witte JS, Visscher PM, Wray NR. <u>The contribution of genetic variants to disease depends on the ruler</u>. Nat Rev Genet. 2014;15(11):765-76. (Important in the understanding of the effects of ascertainment on PRS work).
- Shah S, Bonder MJ, Marioni RE, Zhu Z, McRae AF, Zhernakova A, Harris SE, Liewald D, Henders AK, Mendelson MM, Liu C, Joehanes R, Liang L; BIOS Consortium, Levy D, Martin NG, Starr JM, Wijmenga C, Wray NR, Yang J, Montgomery GW, Franke L, Deary IJ, Visscher PM. <u>Improving Phenotypic Prediction</u> by Combining Genetic and Epigenetic Associations. Am J Hum Genet. 2015; 97(1):75-85. (Important for the conceptualization of polygenicity)

