Using GREML to estimate SNP- heritability among 'unrelated' individuals

Matthew Keller

University of Colorado at Boulder

Three flavors of h²

Twin h²

- estimate of the narrow-sense (but in practice, broad-sense) h²
- Can be biased if strong assumptions are unmet
 SNP-h²
 - Estimate of the narrow-sense h² captured by causal variants tagged by SNPs in your analysis
 GWAS h²
 - Sum of r² from all independent genome-wide significant SNPs from a GWAS

Three types of missing h²

Almost always Twin $h^2 > SNP-h^2 > GWAS h^2$



Three types of missing h²

Almost always Twin $h^2 > SNP-h^2 > GWAS h^2$



Three types of missing h² Almost always Twin h² > SNP-h² > GWAS h²



Why care about SNP-h²?

- h_{snp}^2 should be unbiased by environmental factors that increase close relative similarity. In particular, doesn't rely on rMZ > rDZ due only to genetic differences (although still assumes no relationship between genetic & env. similarity among distant relatives)
- Estimating h²_{snp} from binned SNPs allows for estimates of relative importance of different SNP annotations. E.g., allows for estimates of allelic spectra (distribution of CV MAF)
- Can estimate r_g between low prevalence disorders that are impractical or impossible to estimate using twins/family designs
- As we continue to capture lower MAF SNPs through imputation or sequencing, GREML (but not LDSC) estimates of h_{snp}^2 approach full narrow-sense h^2 .

Multiple approaches to estimating h²_{snp}

- LD-score regression
- Least Squares Regression (Haseman-Elston)
- Mixed effects models (GREML):
 - Typical approach (GCTA assumptions)
 - Multi-GRM approaches
- Bayesian approaches

Multiple approaches to estimating h²_{snp}

LD-score regression

YESTERDAY

τοραγ

- Least Squares Regression (Haseman-Elston)
- Mixed effects models (GREML):
 - Typical approach (GCTA assumptions)
 - Multi-GRM approaches
- Bayesian approaches

pihat

$\hat{\pi} = E(IBD)$, usually genome-wide

- π̂ among close relatives captures long stretches of identical chromosomes, and estimate IBS at both common and rare alleles. Traditionally with close relatives, we know the expectation of this and use this (without variance) for modeling.
- π̂ among unrelateds (distant relatives) assumes base population is the current sample, and thus its expectation is 0. It is typically measured with SNPs, and so only captures IBS at measured SNPs and unmeasured SNPs in LD with measured SNPs. It can go negative (less related than average).

$\hat{\pi}$ = genome-wide mean correlation of SNP values between a pair of individuals *j*,*k*

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

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

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

$$\hat{\pi}_{jk} = \frac{1}{m} \sum_{i} (z_{ij})(z_{ik})$$

$$\hat{\pi}_{jk} = \frac{1}{m} \sum_{i} cor(x_{ij}, x_{ik})$$

H-E REGRESSION

$$\theta_{ij} = Z_i Z_j \leftarrow \cdots$$

product of centered scores (here, z-scores)

 $E[\theta_{ij}] = COV(Z_i, Z_j)$

$$E[\theta_{ij} \mid \hat{\pi}_{ij}] = \hat{\beta}_0 + \hat{\beta}_1 \hat{\pi}_{ij}$$

$$\hat{\beta}_1 = \hat{h}^2$$

$$\theta_{ij} = Z_i Z_j$$

$$E[\theta_{ij}] = COV(Z_i, Z_j)$$

$$E[\theta_{ij} \mid \hat{\pi}_{ij}] = \hat{\beta}_0 + \hat{\beta}_1 \hat{\pi}_{ij}$$

$$\hat{\beta}_1 = \hat{h}^2$$



$$\theta_{ij} = Z_i Z_j$$

$$E[\theta_{ij}] = COV(Z_i, Z_j)$$

$$E[\theta_{ij} \mid \hat{\pi}_{ij}] = \hat{\beta}_0 + \hat{\beta}_1 \hat{\pi}_{ij}$$

$$\hat{\beta}_1 = \hat{h}^2$$



$$\theta_{ij} = Z_i Z_j$$

$$E[\theta_{ij}] = COV(Z_i, Z_j)$$

$$E[\theta_{ij} \mid \hat{\pi}_{ij}] = \hat{\beta}_0 + \hat{\beta}_1 \hat{\pi}_{ij}$$

j

$$\hat{\beta}_1 = \hat{h}^2$$

(the slope of the regression is an estimate of h²)



$$\theta_{ij} = Z_i Z_j$$

$$E[\theta_{ij}] = COV(Z_i, Z_j)$$

$$E[\theta_{ij} \mid \hat{\pi}_{ij}] = \hat{\beta}_0 + \hat{\beta}_1 \hat{\pi}_{ij}$$

$$\hat{\beta}_1 = \hat{h}^2$$



$$\theta_{ij} = Z_i Z_j$$

$$E[\theta_{ij}] = COV(Z_i, Z_j)$$

$$E[\theta_{ij} \mid \hat{\pi}_{ij}] = \hat{\beta}_0 + \hat{\beta}_1 \hat{\pi}_{ij}$$

$$\hat{\beta}_1 = \hat{h}^2$$



$$\theta_{ij} = Z_i Z_j$$

$$E[\theta_{ij}] = COV(Z_i, Z_j)$$

$$E[\theta_{ij} \mid \hat{\pi}_{ij}] = \hat{\beta}_0 + \hat{\beta}_1 \hat{\pi}_{ij}$$

$$\hat{\beta}_1 = \hat{h}^2$$



Regression estimates of h²_{snp}

$$\theta_{ij} = Z_i Z_j$$

$$E[\theta_{ij}] = COV(Z_i, Z_j)$$

$$E[\theta_{ij} \mid \hat{\pi}_{ij}] = \hat{\beta}_0 + \hat{\beta}_1 \hat{\pi}_{ij}$$

$$\hat{\beta}_{\!1}=\hat{h}_{\,\rm snp}^2$$



GREML

genetics

2010

Common SNPs explain a large proportion of the heritability for human height

Jian Yang¹, Beben Benyamin¹, Brian P McEvoy¹, Scott Gordon¹, Anjali K Henders¹, Dale R Nyholt¹, Pamela A Madden², Andrew C Heath², Nicholas G Martin¹, Grant W Montgomery¹, Michael E Goddard³ & Peter M Visscher¹

GREML Model



GREML Model



We aren't interested in estimates of each u_i because such individual estimates are unreliable when m > n. Instead, estimate the <u>variance</u> of u_i .

GREML Model







REML find values of $\hat{\sigma}_A^2 \& \hat{\sigma}_e^2$ that maximizes the likelihood of the observed data. Intuitively, this makes the observed and implied var-covar matrices be as similar as possible.

Interpreting h² estimated from SNPs (h²_{snp})

- If close relatives included (e.g., sibs), $h_{snp}^2 \cong h^2$ estimated from a family-based method, because great influence of extreme pihats. Interpret h_{snp}^2 as from these designs.
 - If use 'unrelateds' (e.g., pihat < .05):
 - h² estimate 'uncontaminated' by shared environment and non-additive genetic effects
 - Does not rely on family/twin study assumptions
 - Evidence for h_{snp}^2 to degree similarity at SNPs corresponds to phenotypic similarity. Thus, h_{snp}^2 = proportion of V_P due to <u>CVs tagged by (in LD with)</u> <u>SNPs used in the GRM</u>.
 - Typically, h²_{snp} < h². It is the max r² possible from a PRS using those SNPs.

LD

Linkage disequilibrium (LD)

- Statistical association (e.g., r²) between 2 SNPs
- Typically arises from a mutation that occurs on a haplotype. This mutation will co-segregate with nearby SNPs. As it rises in frequency, so too will nearby SNPs.
- Decays as a function of number of recombination events that break the two SNPs apart, which is itself a function of:
 - Time (# generations) since the mutational event
 - Distance (cM) between the two SNPs
- SNPs can only predict SNPs that are similar in MAF. Rarerare or common-common. Rare-common is not possible.

How LD arises & decays



Bush & Moore, *PLoS Comp Bio*, 2012

SNPs can tag other nearby SNPs...



Bush & Moore, PLoS Comp Bio, 2012

LD drops as a function of distance



Dawson et al, *Nature*, 2002

...and high LD possible only if the two alleles are of similar frequencies.

Possible range of allele frequencies given LD (r²) between 2 SNPs

Allele frequency locus 2 (pB



Allele frequency at locus $1(p_A)$

Wray, TRHG, 2006

High LD possible only if the two alleles are of similar frequencies. Here A/B are major and a/b are minor alleles on haplotypes.





...and high LD possible only if the two alleles are of similar frequencies.

Possible range of allele frequencies given LD (r²) between 2 SNPs

Allele frequency locus 2 (pB



Allele frequency at locus $1(p_A)$

Wray, TRHG, 2006

...and high LD possible only if the two alleles are of similar frequencies...

Possible range of allele frequencies given LD (r²) between 2 SNPs

Allele frequency locus 2 (pB



Allele frequency at locus $1(p_A)$

Wray, TRHG, 2006
...and high LD possible only if the two alleles are of similar frequencies...

Possible range of allele frequencies given LD (r²) between 2 SNPs

Allele frequency locus 2 (pB



Allele frequency at locus $1(p_A)$

Wray, TRHG, 2006

...AND where the rare allele at SNP 1 is in PHASE with the rare allele at SNP 2



Why $h_{snp}^2 < h^2$ (almost always)

- Because we only estimate genetic variance from CVs in LD with the SNPs used in the analysis. Common CVs are in high LD with array/imputed SNPs, but this is less the case with rare CVs.
 - In particular:

$$\hat{h}_{snp}^2 \cong h^2 \frac{\overline{r}_{MQ}^2}{\overline{r}_{MM}^2}$$

where

 \overline{r}_{MQ}^2 is the average LD r² between CVs and SNPs

 \overline{r}_{MM}^2 is the average LD r² between SNPs and SNPs

THE END (extra slides after)

RUNNING GCTA

SNP QC

- Poor SNP calls can inflate SE and cause downward bias in h²_{snp}
- Clean data for
 - SNPs missing > ~.05
 - HWE p < 10e-6
 - − MAF < ~.01
 - Plate effects:
 - Remove plates with extreme average inbreeding coefficients or high average missingness

Individual QC

- Remove individuals <u>missing</u> > ~.02
- Remove <u>close relatives</u> (e.g., --grm-cutoff 0.05)
 - Correlation between pi-hats and shared environment can inflate h²_{snp} estimates
- Control for <u>stratification</u> (usually 5 to 20 PCs)
 - Different prevalence rates (or ascertainments)
 between populations can show up as h²_{snp}
- Control for <u>plates</u> and other technical artifacts
 - Be careful if cases & controls are not randomly placed on plates (can create upward bias in h²_{snp})

GCTA command & input



GCTA command & output

COMMAND:

gcta --grm-bin <path>/SNPs.rel05 --pheno <path>/pheno.txt --covar <path>/cov.txt --reml --out SNPgrm.randomCV

<u>OUTPUT:</u> cat SNPgrm.randomCV.hsq

Source	Variance	SE		
V(G)	0.300098	0.275857		
V(e)	1.730548	0.279257		
Vp	2.030646	0.049529		
V(G)/Vp	0.147785 0.135820			
logL	-2876.706			
logL0	-2877.338			
LRT	1.264			
Df	1			
Pval	0.1305			
Ν	3363			

GCTA command & output

COMMAND:

gcta --grm-bin <path>/SNPs.rel05 --pheno <path>/pheno.txt --covar <path>/cov.txt --reml --out SNPgrm.randomCV

<u>OUTPUT:</u> cat SNPgrm.randomCV.hsq

Source	Variance	SE	
V(G)	0.300098	0.275857	
V(e)	1.730548	0.279257	
Vp	2.030646	0.049529	
V(G)/Vp	0.147785	0.135820	$h^2_{\text{spp}} \in SE$
logL	-2876.706		shp & SE
logL0	-2877.338		
LRT	1.264		
Df	1		$0.147 - 1.96^{\circ} 0.134 = -0.12$
Pval	0.1305		0.147-1.96*0.134 = 0.41
N	3363		

GCTA command & output

COMMAND:

gcta --grm-bin <path>/SNPs.rel05 --pheno <path>/pheno.txt --covar <path>/cov.txt --reml --out SNPgrm.randomCV

<u>OUTPUT:</u> cat SNPgrm.randomCV.hsq

Source	Variance	SE	
V(G)	0.300098	0.275857	
V(e)	1.730548	0.279257	
Vp	2.030646	0.049529	
V(G)/Vp	0.147785	0.135820	
logL	-2876.706		
logL0	-2877.338		Likelihood Ratio Test
LRT	1.264	←	
Df	1		Testing if $V(G) > 0$
Pval	0.1305		$-2^{*}(-2877.3382876.706) = 1.26$
Ν	3363		χ^2 test, 1 df

LDAK

ARTICLES



2017

Reevaluation of SNP heritability in complex human traits

Doug Speed¹, Na Cai^{2,3}, the UCLEB Consortium⁴, Michael R Johnson⁵, Sergey Nejentsev⁶ & David J Balding^{1,7}

ARTICLE 2012

Improved Heritability Estimation from Genome-wide SNPs

Doug Speed,^{1,*} Gibran Hemani,² Michael R. Johnson,³ and David J. Balding¹

Changing GREML assumptions by weighting $\hat{\pi}$

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

• A more general form of this formula is:

$$\hat{\pi}_{jk} = \frac{1}{\sum W_i} \sum_i W_i (x_{ij} - 2p_i) (x_{ik} - 2p_i) [2p_i(1 - p_i)]^{\alpha})$$

Which reduces to the typical formulation above when: $W_i = 1 \forall i = 1 \dots m \& \alpha = -1$

The choice of W_i and α are arbitrary and depend on implicit assumptions about which types of SNPs tag CVs & CV effect sizes. If we heavily weight a certain type of SNP (e.g., those in genes), we assume such SNPs better tag CVs.

Typical (GCTA) assumptions implicit in $\hat{\pi}$

$$\hat{\pi}_{jk} = \frac{1}{\sum W_i} \sum_i W_i (x_{ij} - 2p_i) (x_{ik} - 2p_i) [2p_i(1 - p_i)]^{\alpha})$$

Assumptions

<u>Consequences</u>

 $W_i = 1 \forall i = 1 \dots m$ SNPs have equal weight, even if they are poorly imputed and redundantly tag the same CV

 $\alpha = -1$

Rarer SNPs (which tag rarer CVs) receive more weight, ostensibly due to NS. This means the variance explained per SNP is invariant across MAF:

$$G_{i} = (X_{i} - 2p_{i})[2p_{i}(1 - p_{i})]^{\alpha/2}$$

$$V[G_{i}] = [2p_{i}(1 - p_{i})]^{\alpha}V[(X_{i} - 2p_{i})]$$

$$V[G_{i}] = [2p_{i}(1 - p_{i})]^{\alpha}2p_{i}(1 - p_{i})$$

$$V[G_{i}] = [2p_{i}(1 - p_{i})]^{\alpha+1}$$

LDAK assumptions implicit in $\hat{\pi}$

$$\hat{\pi}_{jk} = \frac{1}{\sum W_i} \sum_i W_i (x_{ij} - 2p_i) (x_{ik} - 2p_i) [2p_i(1 - p_i)]^{\alpha})$$

Assumptions

<u>Consequences</u>

 $W_i = r_i w_i$ $w_i \cong \frac{1}{(1 + \sum r_{i,i'}^2)}$

Where r_i is the imputation INFO score and w_i is the LD score. High LD SNPs receive less weight, and poorly imputed SNPs receive less weight.

 $\alpha = -.25$

Lower MAF SNPs (which tag rarer CVs) receive less (vs. GCTA) weight. This means the variance explained per SNP increases with MAF:

$$V[G_i] = [2p_i(1-p_i)]^{\alpha+1}$$

Speed & Balding argued that LDAK weights are superior

- <u>Common sense:</u> Redundantly tagged CVs should not have higher effect sizes. Poorly imputed SNPs must tag CVs worse.
- <u>Model Fit</u>: log-likelihood from LDAK models was typically higher than log-likelihood from "GCTA" models
 - Moreover, h_{snp}^2 25-43% higher than GCTA models

Problems with LDAK approach

- Single GRM models depend heavily on assumptions and CV MAF matching the SNP MAF distribution
- Nothing about maximizing likelihoods ensures unbiasedness
- LD and imputation r² are highly positively related, but LDAK weights them oppositely. This gives extreme weight to a small number of unusual (well imputed, low LD, high MAF SNPs)



GREML-LDMS-I & -R

ANALYSIS

2018

LDMS-I

Comparison of methods that use whole genome data to estimate the heritability and genetic architecture of complex traits

nature

renetics

Luke M. Evans^{1*}, Rasool Tahmasbi¹, Scott I. Vrieze², Gonçalo R. Abecasis³, Sayantan Das³, Steven Gazal^{4*}, Douglas W. Bjelland¹, Teresa R. de Candia¹, Haplotype Reference Consortium⁶, Michael E. Goddard^{7,8}, Benjamin M. Neale⁵, Jian Yang⁹, Peter M. Visscher⁹ and Matthew C. Keller^{1,10*}

genetics

2015

Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index

Jian Yang^{1,2,24}, Andrew Bakshi¹, Zhihong Zhu¹, Gibran Hemani^{1,3}, Anna A E Vinkhuyzen¹, Sang Hong Lee^{1,4}, Matthew R Robinson¹, John R B Perry⁵, Ilja M Nolte⁶, Jana V van Vliet-Ostaptchouk^{6,7}, Harold Snieder⁶, The LifeLines Cohort Study⁸, Tonu Esko^{9–12}, Lili Milani⁹, Reedik Mägi⁹, Andres Metspalu^{9,13}, Anders Hamsten¹⁴, Patrik K E Magnusson¹⁵, Nancy L Pedersen¹⁵, Erik Ingelsson^{16,17}, Nicole Soranzo^{18,19}, Matthew C Keller^{20,21}, Naomi R Wray¹, Michael E Goddard^{22,23} & Peter M Visscher^{1,2,24}

LDMS-R

The LDMS Approach

- Single-GRM models are highly sensitive to assumptions about CV-SNP LD (e.g., that SNPs have same distribution as CVs) and CV effect size-MAF relationships. We don't want our estimates of genetic architecture to depend on our assumptions of genetic architecture.
- Moreover, even if we were to guess at these relationships perfectly for a trait, they are unlikely to hold across all traits.
- Akin to multiple regression, an alternative (LDMS) is to let the data tell us by fitting multiple GRMs, each with SNPs binned according to different MAF levels and LD levels
- Estimates associated with each GRM are free to soak up whatever variance is explained by those MAF/LD SNPs

LDMS justification

• Note that
$$\hat{h}_{snp}^2 \cong h^2 \frac{\overline{r}_{MQ}^2}{\overline{r}_{MM}^2}$$

• The range of MAF and range of LD will be smaller within a particular MAF/LD bin. As the MAF & LD range shrink for a given MAF/LD bin k of SNPs (M_k) and CVs (Q_k) ,

$$\frac{\overline{r}_{M_k Q_k}^2}{\overline{r}_{M_k M_k}^2} \to 1$$

and thus

$$\hat{h}_{snp,k}^2 \to h_k^2$$

LDMS-R vs. LDMS-I

- LDMS-R: Create 20 GRMs across 5 MAF bins (< .001, .001-.01, .05-.1, .1-.25, .25-.5) and 4 quartiles of LD scores within each bin, where SNPs take the average LD of SNPs in the surrounding ~ 200kb region.</p>
- However, SNPs with individually low LD that exist in regions of high LD explain more variation (Gazal et al, Nature Genetics, 2017)
- Thus, LDMS-I (unlike LDMS-R) uses each individual SNP's LD score for binning
- Because SEs tend to be ~2.5x larger than single-GRM estimates, both require large sample sizes (e.g., N > 30k) and therefore large amounts of RAM (e.g., >100 Gb)

RUNNING LDMS-I

Create LD quartiles

GCTA LD command: test.bed, test.bim, test.fam

gcta --bfile <path>/test --Id-score-region 200 --out LD.txt

SNP chr bp freq mean_rsq snp_num max_rsq ldscore_SNP ldscore_region rs4475691 1 836671 0.197698 0.000588093 999 0.216874 1.5875 2.75538 rs28705211 1 890368 0.278112 0.000573876 999 0.216874 1.5733 2.75538 rs9777703 1 918699 0.0301614 0.00131291 999 0.854464 2.31159 2.75538

Create LD quartiles in R:

LD <- read.table("LD.txt",header=T)

quants <- quantile(LD\$ldscore_SNP)</pre>

LD1 <- LD\$SNP[LD\$ldscore_SNP <= quants[2]]

write.table(LD1,"snp_group1.txt",row.names=F,col.names=F,quote=F)
<etc...>

Create GRMs in GCTA:

gcta --bfile <path>/test

--extract snps_group1.txt

--make-grm-bin

--out GRM.1

Run LDMS-I using GCTA



can spare

LDMS-I Output (3 GRM example)

TYPE: cat mgrm.randomCV.hsq

Source	Variance	SE
V(G1)	0.303900	0.184182
V(G2)	0.127654	0.309142
V(G3)	0.653199	0.328909
V(e)	0.926493	0.435653
Vp	2.011246	0.049641
V(G1)/Vp	0.151100	0.091277
V(G2)/Vp	0.063470	0.153765
V(G3)/Vp	0.324773	0.164408
logL	-2872.894	
N	3363	

 $h_{SNP} = 0.15 + 0.06 + 0.32 = 0.5391$

GREML vs. LDAK vs. LDMS-I



Comparison of methods that use whole genome data to estimate the heritability and genetic architecture of complex traits

Luke M. Evans¹*, Rasool Tahmasbi¹, Scott I. Vrieze², Gonçalo R. Abecasis³, Sayantan Das³, Steven Gazal⁴, Douglas W. Bjelland¹, Teresa R. de Candia¹, Haplotype Reference Consortium⁶, Michael E. Goddard^{7,8}, Benjamin M. Neale⁵, Jian Yang⁹, Peter M. Visscher⁹ and Matthew C. Keller^{1,10}*

- •We hope it's useful as a guide for best practices and proper interpretation of \hat{h}_{SNP}^2 .
- •We simulated 16 genetic architectures, 3 levels of stratification, and 3 SNP types (array, imputed, WGS) in order to compare \hat{h}_{SNP}^2 across 12 estimation methods (1728 different combos)
- •Here I highlight just a few of what I think are the most important points

Overview of Simulation Approach

- Genotypes from real WGS data (n=8k). Choose 1K rare (MAF < .0025) or common (MAF > .05) CVs.
- Pull out SNPs on UKB array & impute
- Vary 2 CV effect size dimensions $(\lambda_i = u_i [2pq]^{\alpha/2})$:
 - λ -LD (via u_i)
 - λ -MAF (via α)
- Compare \hat{h}_{SNP}^2 to true h^2 (=.50) across 3 methods on imputed data
- Repeat this 100 times for different sets of CVs; look at mean (to get bias) and SD (to get SE) \hat{h}_{SNP}^2

Simulation of phenotypes • CV effect size = $\lambda_i = u_i [2pq]^{\alpha/2}$



• Breeding values = $A_j = \sum_i \lambda_i x_{ij}$

• Phenotype values = $P_j = A_j + E_j$

3 Estimation Methods Compared
 <u>GREML-SC</u>: predictor is a single GRM (aka, "GCTA approach"). GRM built as usual from all imputed SNPs with MAC > 5 & imputation r² > .3

- <u>LDAK</u>: predictor is a single GRM from imputed SNPs and weighted by LD and imputation r².
 - GREML-LDMS-I: predictors are k = 8 GRMs created by binning imputed SNPs into 2 individual LD by 4 MAF categories. Within each bin, GRMs built as usual. $\hat{h}_{SNP}^2 = \Sigma(\hat{h}_{SNP_k}^2)$











LDAK results


LDAK results



LDAK results



LDAK results



GREML-LDMS-I results







Absolute Bias Across 4 Methods and hundreds of genetic architectures



Regarding LD-Score regression

- LD-score regression is robust to stratification and sample overlap. However:
 - it cannot estimate h² due to rare CVs, even when using imputed/WGS data
 - it is sensitive to assumptions about LD- λ
 - should provide a lower-bound of \hat{h}^2_{SNP} from other methods
- So long as genetic covariance is affected in the same way as genetic variances, estimates of genetic correlations should be OK.

Summary

- With datasets imputed to large WGS reference panels, \hat{h}^2_{SNP} can estimate full h^2 . It's important that we have unbiased estimators to know the true h^2 and for comparison to twin/family estimates (o/w things will get really confusing).
- Single-GRM approaches (incl. GREML-SC ("GCTA") and LDAK) are extremely sensitive to CV LD being similar to SNP LD across genome.
 - This is mostly influence by CV vs. SNP MAF, and also by assumptions of LD- λ relationship. MAF- λ less so.
- Binning SNPs by LD & MAF provides ~ unbiased estimates for the CVs tagged by SNPs used in analysis.
 - Even on well-imputed data, you'll still get an underestimate due to extremely rare variants

REAL TRAITS

LDMS-I on UKB phenotypes



STRATIFICATION & LONG-RANGE LD

Chance allele frequency differences b/w populations can induce long-range LD in stratified samples

Population 1

$$r^2 = \frac{(p_{AB} - p_A p_B)^2}{p_A p_a p_B p_b} = 0$$

Population 2

Α



а

$$r^2 = \frac{(p_{AB} - p_A p_B)^2}{p_A p_a p_B p_b} = 0$$

Chance allele frequency differences b/w populations can induce long-range LD in stratified samples



However, such "stratification-LD" is typically very small for pairs of common SNPs



But higher b/w rare (often ~ private) SNPs and common ancestry-informative SNPs



Effects of stratification on r²_{QM}/r²_{MM}

 In general, stratification inflates long-range r² between SNPs. However, within a given MAF bin, the ratio of r²_{QM}/r²_{MM} is ~ 1 because SNP-SNP & SNP-CV LDs are inflated similarly.

Effects of stratification on r²_{QM}/r²_{MM}

- In general, stratification inflates long-range r² between SNPs. However, within a given MAF bin, the ratio of r²_{QM}/r²_{MM} is ~ 1 because SNP-SNP & SNP-CV LDs are inflated similarly.
- However, across CVs and SNPs of different MAF, stratification induces differences in r²_{QM} & r²_{MM}. We observed:
- For rare CVs, $r_{QM}^2/r_{MM}^2 > 1$. Rare (ancestry specific) CVs are tagged by every common SNP that differs in allele frequency across ancestry (note $r_{QM}^2/r_{MM}^2 < 1$ in unstratified samples).
- For very common CVs, $r_{QM}^2/r_{MM}^2 \sim 1$. Very common CVs tend to have smaller MAF differences, and therefore less LD with common SNPs than typical between SNPs (note $r_{QM}^2/r_{MM}^2 > 1$ in unstratified samples).

This led to an opposite pattern of bias in stratified ("structured") samples when using single GRM GREML

Single GRM using WGS



Which once again was corrected by using LDMS GREML



ASSORTATIVE MATING

Positive primary phenotypic assortative mating (AM)

- AM: Assortment between mates leading to a correlation between phenotypic (and hence genetic) scores. Often conceptualized as mate choice based on similarity.
- Induces long-range (across chromosome) "directional"
 LD (δ) b/w CVs
 - $\delta = \text{covariance among CV effects; under positive AM,} \\ E[\delta] > 0; allelic effects in the same direction.$
 - Directional LD increases true $V_G \& h^2$ in the population.
 - This occurs for same reason the variance of a sum of positively correlated $X_i >$ variance of sum of independent X_i
 - For polygenic traits, the vast majority (>99%) of this increase is due to δ between different CVs, not to δ within CVs (homozygosity)

AM effects on pihat

- Assortment has ~ no influence on $\hat{\pi}_{jk}$
- Recall that $E[Z_j Z_k | \hat{\pi}_{jk}] = h^2 \hat{\pi}_{jk}$
- However, this is much different than the reverse conditional*: $E[\hat{\pi}_{jk}|Z_jZ_k] = \frac{rh^2}{m} < \frac{1}{m}$

where *r* is the mate correlation and *m* is the # CVs This is because δ *between* CVs. the major factor influencing h², plays no role in $\hat{\pi}_{jk}$ (or means in general)

$$\hat{\pi}_{jk} = \frac{1}{m} \sum_{i} cor(x_{ij}, x_{ik})$$

*Robinson et al., *Nature Human Behavior*, 2017 Yengo et al., *BioarXiv*, 2018

However, AM does bias h_{snp}^2 estimates

- AM typically leads to <u>upward</u> bias in estimates of equilibrium h_{snp}^2
- Occurs because AM creates positive covariances between CVs and these are correctly reflected in phenotypic covariances between individuals (product of means) but poorly reflected in pihat matrix (mean of products).
- Thus, variance of pihats is too small. Underestimated variance in a predictor leads to overestimates of the coefficients associated with that predictor.
- We derived this bias algebraically in HE regression estimates and confirmed it in simulation.
- REML also upwardly biased, but bias depends on ratio N/m.

Parameter *h*²_{snp}

- Define parameter h_{snp}^2 : proportion of phenotypic variance tagged by SNPs, accounting for their inter-correlations
- Equilibrium h_{snp}^2 : R² from linear model $Z \sim X_1 + X_2 + ... X_m$ for all *m* SNPs fit simultaneously as $n \rightarrow \infty$
- The parameter depends on how well CVs are tagged by SNPs (e.g., SNP density). Thus, it depends on the SNP chip and the population it is estimated in.

HE regression estimate of h_{snp}^2

 $E[Z_i Z_j] = COV(Z_i, Z_j)$

$$E[Z_i Z_j \mid \hat{\pi}_{ij}] = \hat{\beta}_0 + \hat{\beta}_1 \hat{\pi}_{ij}$$

$$\hat{\beta}_1 = \frac{COV(Z_i Z_j, \hat{\pi}_{ij})}{V(\hat{\pi}_{ij})} = \hat{h}_{snp}^2$$

HE regression estimate of h_{snp}^2

 $E[Z_i Z_j] = COV(Z_i, Z_j)$



We don't predict HE estimates to change as a function of m or N. GREML estimates are clearly a function of N/m, which occurs because when N>>m, the effects of each SNP are separable.

mean & 95% Cls of VA estimates (red or blue) or observed VA (green) by N/m. r(spouse)=.4, VA time 0 = 1



100

Predicted h²_{snp} biases assuming SNPs tag 50% of true VA



spousal correlation

SNP heritability

Predicted h²_{snp} biases assuming SNPs tag 50% of true VA



spousal correlation

SNP heritability

Potential degree of overestimation for various traits

Trait	r(spouse)	h ² _{ETFD} from literature	h ² _{snp} from literature	corrected h ² _{snp}	% Over- estimated
Extraversion	.01	.23	.15	.15	0
Neuroticism	.08	.24	.16	.155	.03
Height	.20	.70	.45	.39	.15
IQ	.35	.62	.35	.28	.25
Political Pref.	.48	.26	.18	.15	.20

Simulations



Rasool Tahmasbi

- Simulated populations under AM using GeneEvolve (Tahmasbi & Keller, 2016).
- CVs: 1000
- Heritability: 0.5
- Relative pruning: >.05
- Spousal phenotypic correlation: .4
- Took mean of 100 iterations

GeneEvolve Simulation Results



Heritability estimate

Sample Size

GeneEvolve Simulation Results



Heritability estimate

Sample Size

GeneEvolve Simulation Results



Heritability estimate

HE vs. GREML estimates

- For most realistic situations, m>>N, and thus GREML and HE estimates are similar: both over-estimate equilibrium h²_{snp}
- We can vary N (holding m constant) to see if AM is biasing estimates in real data
HE (blue) vs. REML (red) h_{snp}^2 estimates of <u>systolic BP</u> - UK Biobank

Error bars: 1 SEMs; GRM pruned for relatedness > .05; Covariates: sex, age, age-squared, PCs 1-4, townsend deprivation



HE (blue) vs. REML (red) h_{snp}^2 estimates of <u>height</u> - UK Biobank

Error bars: 1 SEMs; GRM pruned for relatedness > .05; Covariates: sex, age, age-squared, PCs 1-4



HE (blue) vs. REML (red) VA estimates of <u>fluid IQ</u> - UK Biobank

Error bars: 1 SEMs; GRM pruned for relatedness > .05; Covariates: sex, age, age–squared, PCs 1–4, townsend deprivation



Summary – bias due to AM

- AM creates upward biases in HE and REML h_{SNP}^2 estimates
 - We see evidence for this in UK Biobank data for height but not for fluid IQ
 - Natural selection creates negative LD among CVs. The combined effect of AM and NS could cancel each other out.
- Remaining issues:
 - Unsure how to account for the bias. LDMS GREML does not help.
 - Need to understand the effects of NS on h²_{snp}

Big picture: Using SNPs to estimate h²

- There has been a great deal of excitement about using SNPs to estimate h²
- Large sequence reference panels (TopMed) allow SNPs to be imputed down to MAF ~ .0001.
 - h_{snp}^2 will approach h^2
 - Also allows investigation of allelic spectra, and importance of biological/evolutionary annotations
 - By understanding true h², can begin understanding importance of familial environmental factors
- However, it is crucial to understand the factors that can bias these estimates
 - LDMS accounts for biases due to MAF & stratification
 - But not for biases caused by AM (and probably NS)

Acknowledgements

Collaborators

Consortia/Databases

Haplotype Reference Consortium UK Biobank

University of Queensland

Peter Visscher

Jian Yang

Naomi Wray

Mike Goddard

<u>CU</u>

Matt Jones

Broad Institute

Ben Neale

Funding

NIMH K01 MH085812 (Keller) NIMH R01 MH100141 (Keller)

<u>Postdoctoral</u> Fellows

Teresa de Candia Luke Evans Rasool Tahmasbi <u>Graduate</u> <u>Students</u> Emma Johnson Richard Border





Method	Description	Major assumptions	Simulation findings regarding ${\hat{h}}_{{ m SNP}}^2$	Computational issues
GREML-SC⁵	Often called the GCTA approach. Originally applied to common array SNPs only. Estimates $\hat{h}_{\rm SNP}^2$, the amount of h^2 caused by CVs tagged by SNPs used to create the GRM.	(i) Genetic similarity is uncorrelated with environmental similarity; (ii) an infinitesimal model; (iii) SNP effects are normally distributed, independent of LD, and inversely proportionate to MAF (α = -1).	Biased to the degree that the average LD among SNPs is different from the average LD between SNPs and CVs. This occurs in stratified samples and when MAF and LD distributions of SNPs do not match those of CVs.	Simple model tractable with large samples (>100,000).
GREML-MS ¹¹	The first multicomponent approach, usually applied by binning SNPs according to their MAF, annotation, or physical regions to explore genetic architecture.	Requires that the same assumptions of GREML-SC hold within each GRM.	Biased when CVs have generally higher or lower levels of LD than the SNPs used to make the GRM. Relatively large standard errors.	Run times and memory requirements higher than GREML-SC and increase as a function of the number of variance components estimated.
GREML-LDMS-R ⁷	A multicomponent approach that bins imputed SNPs by their MAF and regional LD.	Same as GREML-MS.	Use of regional LD scores can lead to biases when CVs have different LD on average compared to surrounding SNPs. Relatively large standard errors.	Same as GREML-MS.
GREML-LDMS-I	A multicomponent approach introduced here that bins imputed SNPs by their MAF and individual LD.	Same as GREML-MS.	Appears to be the least biased approach, even when traits have complex genetic architectures. Relatively large standard errors.	Same as GREML-MS.
LDAK-SC ^{15,20}	Introduced to account for redundant tagging of CVs by common SNPs. Recently modified to incorporate error due to imputation and to alter the MAF effect-size relationship.	Same as GREML-SC, except that allelic effects are a function of LD. Extended to assume that effects are also a function of imputation quality and weakly inversely proportionate to MAF ($\alpha = -0.25$).	Can correct for the overestimation observed in GREML-SC from redundant tagging of CVs, but otherwise about as biased as GREML-SC when assumptions are unmet, although the biases are sometimes in different directions.	Same as GREML-SC.
LDAK-MS ¹⁵	A multicomponent extension of LDAK-SC that bins SNPs by MAF.	Requires that the same assumptions of LDAK-SC hold within each GRM.	Less biased on average than LDAK-SC, but more biased than GREML-LDMS-I or -R). Relatively large standard errors.	Same as GREML-MS.
Threshold GRMs ²⁴	A multicomponent approach with two GRMs: the normal (unthresholded) GRM built from all SNPs and a second GRM with entries set to 0 if below a threshold. Conducted in samples that include close relatives.	Same as GREML-SC for the unthresholded GRM. Assumes no shared environmental influences among close relatives.	Estimates associated with unthresholded GRM similar to those of GREML-SC. When used in samples that include close relatives, the second GRM captures pedigree-associated variation but can be upwardly biased by shared environmental influences.	See GREML-SC.
LD score regression ¹⁹	Uses the slope from χ^2 (from GWAS) regressed on SNPs' LD scores to estimate the h^2 due to CVs in LD with common SNPs.	Infinitesimal model with allelic effects normally distributed.	Largely robust to confounding due to stratification and shared environmental influences. Estimates h^2 due to common CVs only, even when used on imputed or WGS data. Underestimates h^2 if the trait is not highly polygenic.	The most computationally efficient method of those compared and tractable for very large datasets.

Table 1 | Summary of commonly applied methods and a description of findings from simulations