

Documentation

Contact: Aysu Okbay <a.okbay@vu.nl>, Daniel Benjamin <daniel.benjamin@gmail.com>, Peter Visscher <peter.visscher@uq.edu.au>

Date: 8 March 2018

Introduction. Lee *et al.*¹ conducted genome-wide association analyses of four phenotypes: educational attainment (*EduYears*, $N = 1,131,881$), cognitive performance (*CP*, $N = 257,841$), self-reported math ability (*MA*, $N = 564,698$), and highest-level math class taken (*HM*, $N = 430,445$). In three separate hold-out cohorts, they analyzed the predictive power of polygenic scores derived from the genome-wide association study (GWAS) estimates. Due to IRB restrictions, it is not possible to release summary statistics for more than 10,000 single nucleotide polymorphisms (SNPs). Therefore, researchers with access to the individual-level genotype data cannot reproduce the polygenic scores from publicly available summary statistics (<https://www.thessgac.org/data>). As a partial remedy, we are releasing the polygenic scores directly to researchers (but due to the restrictions, we cannot release the underlying SNP-level weights themselves).

Scores for the Wisconsin Longitudinal Study (WLS) are available here:

“Lee_et_al_(2018)_PGS_WLS.txt”

If you use the scores, please cite:

Lee, J.J. *et al.* Gene discovery and polygenic prediction from a 1.1-million-person GWAS of educational attainment. Forthcoming in *Nat. Genet.* (2018).

The purpose of this document is to briefly describe the construction of the scores. For additional details, readers are referred to the Supplementary Note of Lee *et al.*¹, especially chapters 1 (details on association analyses and phenotypes) and 6 (prediction analyses).

Methodology. A polygenic score for an individual is defined as a weighted sum of a person’s genotypes at K SNPs,

$$\hat{g}_i = \sum_{j=1}^K x_{ij}w_j \quad (1)$$

Methodologies for PGS construction differ primarily across two dimensions: how to generate the weights w_j , and how to determine which K SNPs to include². Here, we use LDpred³, a Bayesian method that includes all measured SNPs and weights each SNP by (an approximation) to its conditional effect, given other SNPs. The theory underlying LDpred is derived assuming the variance-covariance matrix of the genotype data in the training sample is known and assuming some prior effect-size distribution. In practice, the matrix is not known but must be approximated using LD patterns from a reference sample. LDpred calculates posterior effect-size distributions

for the true effect sizes β (i.e., that are conditional on all other SNPs, unlike the GWAS estimates), and each SNP's weight is set equal to the mean of its (conditional) posterior effect-size distribution.

Genotype data and imputation. Genotype data from The Illumina Human Omni Express Bead Chip were available for 9,109 individuals and 713,014 variants. We imputed these genotypes against the Haplotype Reference Consortium (HRC) v1.1 European reference panel⁴ using the Michigan Imputation Server⁵. Prior to imputation, we identified the non-European individuals by plotting the principal components (PCs) of the covariance matrix of the individuals' genotype data⁶ together with the PCs of 1000 Genomes populations and visually inspecting the plots. We dropped the identified non-European individuals from the sample. Additionally, we excluded individuals that do not satisfy the following criteria: (i) genotype missingness rate is less than 0.05 in all chromosomes, (ii) there is no mismatch between surveyed sex and genetic sex, (iii) there is no mismatch between surveyed relationship data and genetic relatedness, (iv) the individual is not an outlier in terms of heterozygosity/homozygosity, and (v) the individual is not an ancestral outlier. We also dropped SNPs that have a call rate less than 0.95, Hardy-Weinberg exact test P -value less than 10^{-5} , or minor allele frequency less than 0.01. 607,469 autosomal SNPs and 8,527 individuals remained in the data.

Next, we checked the data against the HRC reference panel^a for consistency of strand, id names, positions, alleles, reference/alternative allele assignment, and allele frequency differences using version 4.2.5 of the HRC-1000G-check-bim.pl^b program. The program updates strand, position and reference/alternative allele assignment when possible. It removes a SNP if it has any of the following properties: (i) A/T or G/C alleles and a minor allele frequency greater than 0.4, (ii) alleles that do not match the HRC data, (iii) minor allele frequency discrepancy with the HRC data greater than 0.2, (iv) not available in the HRC data. After all checks, 604,710 SNPs remained which were taken forward for imputation. Genotype probabilities were imputed for 39,127,657 variants and 8,527 individuals.

Estimation of LD patterns. We estimated LD patterns using the imputed WLS genotype data. To obtain the LD reference data, we first converted the genotype probabilities for 38,909,200 biallelic SNPs to hard calls using Plink v1.9⁷. We restricted the set of genetic variants to 1,211,685 HapMap3⁸ SNPs, because these SNPs are generally well-imputed and provide good coverage of the genome in European-ancestry individuals. Next, we estimated a genetic relatedness matrix, restricting further to SNPs with minor allele frequency greater than 0.01. We dropped one individual from each of the 2,842 pairs of individuals with a genetic relatedness exceeding 0.025.

In order to make sure that there are no genetic outliers in the sample that can bias the LD estimates, we clustered the remaining 5,691 individuals based on identity-by-state distances in Plink v1.9⁷, again restricting to SNPs with minor allele frequency greater than 0.01. Plink reports a Z -score for each individual's IBS distance to his/her closest neighbor. We examined these Z -scores and marked an individual as genetic outlier if his/her Z -score was smaller than -5. We dropped these individuals and repeated the process, until no more individuals with a Z -score less than -5 remained in the data. The algorithm identified 6 outliers, which were then dropped from the data. In the final reference data, there were 5,685 individuals and 1,211,685 SNPs.

^a Site list was downloaded from <http://www.haplotype-reference-consortium.org/site>

^b Script available at <http://www.well.ox.ac.uk/~wrayner/tools/HRC-1000G-check-bim.v4.2.5.zip>

Weights. We provide two types of polygenic scores for educational attainment and cognitive performance based on different sets of summary statistics from Lee *et al.*¹: (i) a score based on standard GWAS summary statistics, which are the coefficient estimates from univariate GWAS of educational attainment and cognitive performance; and (ii) a score based on MTAG summary statistics, which are obtained from a multivariate analysis of educational attainment, cognitive performance, self-reported math ability and highest-level math class taken using the MTAG software tool⁹ (see below). For self-reported math ability and highest-level math class taken, we only provide the MTAG polygenic scores as univariate GWAS polygenic scores for these traits were not analyzed in Lee *et al.*¹. All weights were obtained from GWAS discovery samples that did not contain the WLS.

We adjusted the weights for linkage disequilibrium using the LDpred software tool³ and the reference genotype data whose construction is described above. The LD-adjusted univariate GWAS weights were obtained for the SNPs that are available in both the reference data and the standard GWAS summary statistics for the phenotype, and that pass the filters imposed by LDpred: (i) the variant has a minor allele frequency (MAF) greater than 1% in the reference data, (ii) the variant does not have ambiguous nucleotides, (iii) there is no mismatch between nucleotides in the summary statistics and reference data, and (iv) there is no high (>0.15) MAF discrepancy between summary statistics and validation sample. The LD-adjusted MTAG weights were further restricted to SNPs that are available in the GWAS summary statistics for all four phenotypes. Table 1 reports the number of SNPs used to construct each polygenic score. The posterior effect sizes were calculated assuming a fraction of causal SNPs equal to one and setting the LD window to $M/3000$, where M is the number of SNPs included in the score.

Table 1. Number of SNPs included in the polygenic scores.

	Educational attainment	Cognitive performance	Math ability	Highest-level math class taken
GWAS PGS	1,170,820	1,170,759	N/A	N/A
MTAG PGS	1,157,359	1,157,359	1,157,359	1,157,359

Polygenic scores. We calculated polygenic scores in Plink v1.9⁷ for 8,527 individuals, using genotype probabilities obtained from the HRC imputation and the LD-adjusted weights described above.

MTAG-based polygenic scores. MTAG⁹ is a method that uses GWAS summary statistics for a primary phenotype and for one or more secondary phenotypes to produce an updated set of summary statistics for the primary phenotype which, under certain assumptions, will be more precisely estimated than the input GWAS summary statistics.

There are costs and benefits to using an MTAG-based polygenic score. For instance, in all cases, MTAG-based polygenic scores will be more predictive of their corresponding phenotype in expectation. In some cases, however, MTAG can have a high false discovery rate (see

Supplementary Note section 1.4 of Turley *et al.*⁹), which may lead to spurious correlations between the MTAG-based polygenic score and other phenotypes.

We therefore offer the following recommendations. If in a regression, the dependent variable and the polygenic score correspond to the same phenotype, we recommend using the MTAG-based score. If the dependent variable and the polygenic score correspond to different phenotypes, but the coefficient of interest in the regression is not the coefficient associated with the polygenic score (e.g., if the polygenic score is only being used as a control variable in an experimental setting), then we also recommend using the MTAG-based polygenic score. Care should be taken when interpreting the coefficient of an MTAG-based polygenic score in this setting, however, since any observed association may be driven through channels involving the secondary phenotypes. This is especially true when the maxFDR is large (see Turley *et al.*⁹, Supplementary Note section 1.4). If researchers are interested in the coefficient on the polygenic score, they should either use GWAS-based scores, or justify why such channels would lead to negligible bias in their particular case.

Phenotypes. Table 2 summarizes the GWAS phenotype measures. For additional details, please see Supplementary Tables 17 and 40 in Lee *et al.*¹

Table 2. Phenotype measures

Phenotype	Measure
Educational attainment	For each component study, defined in accordance with the ISCED 1997 classification (UNESCO), leading to seven categories of educational attainment that are internationally comparable. The categories are translated into US years-of-schooling equivalents.
Cognitive performance	Meta-analysis of the following two studies: UK Biobank ¹⁰ : Standardized score on a test of verbal-numerical reasoning designed as a measure of fluid intelligence. COGENT ¹¹ : For each component study in the meta-analysis, the first unrotated principal component of performance on at least three neuropsychological tests (or at least two IQ-test subscales).
Self-reported math ability	How would you rate your mathematical ability? Very poor [0], Poor [1], About average [2], Good [3], Excellent [4]
Highest-level math class taken	Excluding statistics courses, what is the most advanced math class you have successfully completed? Pre-Algebra [1], Algebra [2], Geometry [3], Trigonometry [4], Pre-Calculus [5], Calculus [6], Vector Calculus [7], More than vector calculus [8]

Principal components. It is important to take a number of steps to minimize the risk that an observed association between the outcome of interest and the polygenic score is due to

unaccounted-for population stratification. A score is stratified if its distribution varies across members of different ancestry groups. Absence to control for differences in ancestry can severely bias estimates of effect sizes, since members of different groups may vary in the outcome of interest for environmental reasons¹². To reduce such concerns, we recommend controlling for the top 10 principal components (PCs) of the covariance matrix of the individuals' genotypic data⁶, which are included in "Lee_et_al_(2018)_PGS_WLS.txt". The principal components were obtained in Plink v1.9⁷ using SNPs with call rate greater than 0.99, minor allele frequency greater than 0.01, and imputation accuracy greater than 0.6. Prior to calculating the principal components, we excluded long-range LD regions on chromosomes 5 (44-51.5 Mb), 6 (25-33.5 Mb), 8 (8-12 Mb) and 11 (45-57 Mb). Remaining SNPs were LD-pruned ($R^2 < 0.1$ on a 1000kb window).

Variables. Table 3 provides a description of the variables included in "Lee_et_al_(2018)_PGS_WLS.txt".

Table 3. Description of variables

Variable	Description
<i>FID</i>	Family identifier
<i>IID</i>	Individual identifier
<i>PGS_EA3_GWAS</i>	Polygenic score for educational attainment, obtained using standard GWAS results
<i>PGS_CP_GWAS</i>	Polygenic score for cognitive performance, obtained using standard GWAS results
<i>PGS_EA3_MTAG</i>	Polygenic score for educational attainment, obtained using multivariate analysis of educational attainment, cognitive performance, self-reported math ability and highest-level math class taken
<i>PGS_CP_MTAG</i>	Polygenic score for cognitive performance, obtained using multivariate analysis of educational attainment, cognitive performance, self-reported math ability and highest-level math class taken
<i>PGS_MA_MTAG</i>	Polygenic score for self-reported math ability, obtained using multivariate analysis of educational attainment, cognitive performance, self-reported math ability and highest-level math class taken
<i>PGS_HM_GWAS</i>	Polygenic score for highest-level math class taken, obtained using multivariate analysis of educational attainment, cognitive performance, self-reported math ability and highest-level math class taken

References

1. Lee, J. J. *et al.* Gene discovery and polygenic prediction from a 1.1-million-person GWAS of educational attainment. *Nat. Genet.* (2018).
2. Wray, N. R. *et al.* Pitfalls of predicting complex traits from SNPs. *Nat. Rev. Genet.* **14**, 507–515 (2013).
3. Vilhjálmsson, B. J. *et al.* Modeling Linkage Disequilibrium Increases Accuracy of Polygenic Risk Scores. *Am. J. Hum. Genet.* **97**, 576–592 (2015).
4. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat. Genet.* **48**, 1279–1283 (2016).
5. Das, S. *et al.* Next-generation genotype imputation service and methods. *Nat. Genet.* **48**, 1284–1287 (2016).
6. Price, A. L. *et al.* The impact of divergence time on the nature of population structure: an example from Iceland. *PLoS Genet.* **5**, e1000505 (2009).
7. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* **4**, 7 (2015).
8. Altshuler, D. M., Gibbs, R. A. & Peltonen, L. Integrating common and rare genetic variation in diverse human populations. *Nature* **467**, 52–58 (2010).
9. Turley, P. *et al.* Multi-trait analysis of genome-wide association summary statistics using MTAG. *Nat. Genet.* (2018). doi:10.1038/s41588-017-0009-4
10. Bycroft, C. *et al.* Genome-wide genetic data on ~500,000 UK Biobank participants. *bioRxiv* 166298 (2017). doi:10.1101/166298
11. Trampush, J. W. *et al.* GWAS meta-analysis reveals novel loci and genetic correlates for general cognitive function: a report from the COGENT consortium. *Mol. Psychiatry* **22**, 336–345 (2017).
12. Hamer, D. & Sirota, L. Beware the chopsticks gene. *Mol. Psychiatry* **5**, 11–13 (2000).