

MIDUS Project 6 (Genetics):

Polygenic Risk Score (PRS) Documentation

For File:

M2MR_P6_PolygenicRiskScores_N2118_20211028

November 2021

NIA Grant #U19 AG051426 and post-doctoral support from a Templeton Foundation Grant. NIMH R01 MH124619 (Docherty, Shabalin), the Brain and Behavior Research Foundation (Shabalin, Docherty), the Simons Foundation (Docherty), the University of Utah EDGE Scholar Program, and the Huntsman Mental Health Institute

Polygenic Risk Score (PRS) Data

Overview

DNA was extracted from tissue samples (whole blood, saliva) obtained from MIDUS 2 (M2) and MIDUS Refresher 1 (MR1) Biomarker participants and then genotyped using the Illumina Omni Express array. The resultant SNPs were then processed via established protocols for genotype calling^A, ancestry measurement, and imputation of genotypes and an initial set of Polygenic Risk Scores (PRS) was computed using PRSice 2.0 [1]^B in 2019. In 2021 an additional individual PRS for Alzheimer's disease was computed using default parameters in PRS-CS. The purpose of this document is to provide information about the PRS, how they were generated, and considerations that must be made when using PRS in analyses. General information about the data file is provided first, and subsequent sections provide additional details about sample collection and processing as well as computation of the PRS.

Sample and Variables

In general, MIDUS data are processed and distributed separately by wave of data collection (M1, M2, Refresher1 etc.). At the time that DNA extraction was initiated, tissue samples were available from both M2 and MR1. SNP data from both waves were combined for efficiency of processing, so the PRS data file described here contains data for M2 and MR1 cases. Consistent with other MIDUS files, the file includes the full set of M2 (n=1255) and MR1 (n=863) biomarker cases for a total sample size of N=2118. The file includes a flag variable (see below) so that cases with PRS data may be easily identified.

MIDUS variable naming and coding conventions (seen Naming and Coding Conventions included with the Survey documentation) specify that the first 3-4 characters in a variable name indicate the wave (i.e. B for M2, RA for Refresher wave 1), the project, and the instrument or data type. Since the PRS data file includes data from both M2 and MR1, the convention is modified slightly such that each variable begins with BRA. The data were generated under the auspices of the Genetics project (P6) and are derived from DNA, thus the first 5 characters of the variables, except for MIDUS Administrative variables, begin with BRA6D. The remaining characters are determined by the data. For example, variables derived from principal components analysis include 'PC' while the PRS variables include 'PRS'.

The SPSS data file "M2MR_P6_PolygenicRiskScores_N2118_20211028.sav" contains 65 variables as follows:

- Administrative: there are 6 variables
 - M2MRID – contains the public identifiers for the MIDUS core sample and the MIDUS Refresher
 - SAMPLMAJ – the identifier created by the MIDUS Administrative Core to indicate the participants 'sample of origin' (e.g. MIDUS Refresher, Twin, etc.),
 - M2MRCASE – indicates if sample was obtained as part of the MIDUS 2 or MIDUS Refresher 1 data collection
 - GENCONSENT – flag variable indicating whether the participant consented to genetics or not.
 - BRA6DTISSUE – indicates the tissue sample source from which DNA was extracted. See below for information about sample collection.
 - BRA6DPAVAIL – flag variable indicating if PRS are available for a given case or not.
- Ancestry Principal Component Scores (PC): Genotype principal component analysis was performed as part of determination of ancestry. The top principal components (PCs) from that process (see details below) are included in the current data file (BRA6DPC1 to BRA6DPC5).
 - All five of the PC variables should be included in all PRS analyses. For an introduction to the interpretation of PC scores, see Novembre & Stephens (2008) [2].

^A "In next-generation sequencing (NGS) methods, a whole genome, or targeted regions of the genome, is randomly digested into small fragments (or short reads) that get sequenced and are then either aligned to a reference genome or assembled. Having aligned the fragments of one or more individuals to a reference genome, 'SNP calling' identifies variable sites, whereas 'genotype calling' determines the genotype for each individual at each site." (p. 443, Nielsen, Paul, Albrechtsen, & Song, 2011, <https://doi.org/10.1038/nrg2986>)

^B <https://choishinwan.github.io/PRSice/>

- Ancestry Scores: there are also five ancestry component scores (BRA6DPEUR -European, BRA6DPCEAS - East Asian, etc.). Ancestry component scores range from 0 to 1 and represent the probability of one's ancestry based on their genotype. See below for computation details.
- Polygenic Risk Scores (PRS) – there are 47 PRS variables generated by PRSice 2.0 [1]^c. The remaining characters in each variable name serve as a unique identifier for the phenotype (e.g. BRA6DPRSALZ “ALZ” for Alzheimer’s Disease, BRA6DPRSAOM “AOM” for Age of Menarche, etc.)
 - PRS are only included for individuals of predominately European descent, due to ancestry confounds present in non-European data. See details about PRS computation below. Importantly, as GWAS soon become available with multiple ancestries, PRS based on those GWAS will perform better across diverse ancestries.
 - Table 1 below contains a list of phenotypes for which PRS data are available and associated references for discovery genome-wide association studies.
- Polygenic Risk Score from PRS-CS – there is one variable, for Alzheimer’s Disease, generated using PRS-CS, a Bayesian method inferring posterior SNP effect sizes under continuous shrinkage (CS) priors (see below for more information). This variable appears in the data file after the above set of scores.
 - The new variable BRA6DPRSCSALZ, is the ‘score average’.
 - Note that the PRS for ALZ was derived using the default pi hat parameter in PRS-CS. For any PRS from smaller discovery GWAS, this parameter will need to be modified accordingly.

DNA Tissue Samples

Collection and storage of tissue samples for genetics was added to the MIDUS Biomarker protocol about halfway through the data collection period. Once that change was implemented, whole blood was stored for DNA extraction for the subset of MIDUS 2 participants and all the MIDUS Refresher participants who consented to genetic testing at the clinic visit. About half of the M2 participants completed the clinic visit before genetics was added to the protocol. At the end of the M2 data collection period individuals in this group were invited to provide saliva samples for genetic testing by mail using Oragene DNA Self-Collection Kits (OG-300). Saliva samples were obtained from about 70% of that group.

Analytic Guidelines

Background: PRS provide an index of an individual’s genetic propensity for a given phenotype based on common genetic variation. Put differently, PRSs provide “an individual-level genome-wide genetic proxy”¹ for a variable of interest. PRSs can be computed for any variable for which there has been a large discovery genome wide association study (e.g. height, weight, BMI, LDL, HDL, triglycerides, type-II diabetes, schizophrenia, subjective well-being, etc.) and can be used in a variety of applications. Summary data used to generate PRSs can be used to examine shared genetic covariation among two or more variables, epistasis or genome x genome interactions between two or more variables, and genome x environment interactions [3, 4]. PRS summary data are also commonly used as instrumental variables to infer quasi-causal relations between exposures and outcomes in observational studies, sometimes called Mendelian randomization studies [5, 6]. In all of these applications, the calculation of a polygenic risk score depends on estimates from a genome-wide association study (GWAS), specifically the effect sizes of single nucleotide polymorphisms (SNPs; β_i) that are statistically associated with the phenotype of interest. As a result, the predictive power and validity of polygenic risk score analyses depend on both the size and quality of the discovery GWAS and the degree those samples match the target sample in which the PRS analyses are performed.

IMPORTANT: PRS analyses are highly prone to ancestry bias, and the discovery GWAS from which SNP weights are obtained are almost universally based on participants of Non-Finnish, Northern European ancestry. Until non-European ancestry GWAS summary statistics are both available and reliable, PRS analyses should only be conducted with samples of European participants. **In addition, all PRS analysis**

^c <https://choishingwan.github.io/PRSice/>

requires inclusion of ancestry principal components as covariates, along with other critical covariates (age, sex, etc).

When conducting PRS analyses in MIDUS, it is also important to account for the nested structure of the data, as genotyped participants include twins (~ 100 pairs) that were raised in the same home. PRS analyses of twin data should typically only be considered in studies comparing sibling pairs. In other cases, one twin should be left out of analyses. PRS will be highly or completely overlapping across siblings due to the large amount of DNA shared and the way in which PRS aggregate common variation.

When conducting GxE and GxG studies, it is important to make sure that potential confounding variables are accounted for. “To properly control for confounders, researchers need to enter the covariate-by-environment and the covariate-by-gene interaction terms in the same model that tests the GxE term.” [7]. For example, if a researcher wanted to test a hypothesis about a particular polygenic risk score (PRS) interacting with a measured environment (E), and wanted to control for three potential confounders (cov1 - cov3), then the following regression equation should be estimated:

$$y = b_0 + b_1\text{PRS} + b_2E + b_3\text{cov1} + b_4\text{cov2} + b_5\text{cov3} + b_6\text{PRS}\times E + b_7\text{cov1}\times\text{PRS} + b_8\text{cov2}\times\text{PRS} + b_9\text{cov3}\times\text{PRS} + b_{10}\text{cov1}\times E + b_{11}\text{cov2}\times E + b_{12}\text{cov3}\times E + e$$

b_0 = intercept
 b_1 - b_{12} = slope
 e = residual error term

The same procedure should be followed to control for potential confounders when testing GxG interactions. When testing whether a particular variable statistically mediates or accounts for polygenic liability for a phenotype (i.e. when searching for endophenotypes), **it is important to bootstrap standard errors when estimating indirect effects**. This can be easily achieved using modern statistical software, including SPSS, SAS, R, and Mplus. Finally, when testing multiple PRS → phenotype associations, **it is important to adjust p-value thresholds for multiple comparisons** using a Bonferroni correction ($p < .05/\text{total number of comparisons}$). In some cases, FDR correction may be used. Technical details about steps for generating the PRSs are provided in the following section.

Generating Polygenic Risk Scores

Genotyping Quality Control: Genotyping quality control (QC) for MIDUS samples was performed on approximately 500k common variants using single nucleotide polymorphism (SNP) clustering in the Illumina Genome Studio software^D. Samples were assayed on 2 different versions of an Illumina OmniExpress array (1.0 and 1.1) and, consequently, SNP calling using Genome Studio required separate processing from different arrays. The data from the different arrays had previously been processed jointly, resulting in major bias and a loss of quality. To correct this, each batch was re-processed separately with Genome Studio, using the appropriate clustering files for each respective genotyping array. This significantly increased SNP call rate and improved other QC metrics.

Ancestry Analysis: Ancestry of the MIDUS participants was estimated using Admixture software^E with a 1000 Genomes data (Phase 3) reference [8] All 5 super-populations^F were used as a basis for estimation.

^D <https://www.illumina.com/techniques/microarrays/array-data-analysis-experimentaldesign/genomestudio.html>

^E <http://software.genetics.ucla.edu/admixture/admixture-manual.pdf>

^F 26 different populations from around the world comprise the 1000 Genomes reference data, and these populations are divided into 5 “super-populations”: African, Ad Mixed American, East Asian, European, and South Asian. For a list of populations and their descriptions, see: <http://www.internationalgenome.org/category/population/>

To calculate ancestry component scores, genotype principal components analysis (PCA) was performed on MIDUS genotypes after linkage disequilibrium (LD) pruning SNPs at a 0.2 R^2 threshold. The PCA was run using RaMWAS [9], a Bioconductor package which comprises a complete toolset for QC, GWAS and methylome-wide association studies. RaMWAS includes functions for PCA for capturing batch effects and detection of outliers, association analysis while correcting for top PCs and covariates, creation of QQ-plots and Manhattan plots, and annotation of significant results. In total, five ancestry component scores were calculated: European (EUR), East Asian (EAS), Ad-mixed American (AMR), Southeast Asian (SAS), and African (AFR).

Inclusion Criteria: To date, discovery GWASs have focused almost exclusively on participants of Non-Finnish Northern European ancestry. Consequently, the estimated effect sizes of individual SNPs are only known for individuals of European ancestry, and, as a result, the calculation of polygenic scores are only valid for participants of predominantly European ancestry. Therefore, to exclude ancestrally heterogeneous samples from the data, the top principal components (PCs), defined as those components which accounted for > 0.1% of the genotype variance, $n_{pc} = 4$, were used to establish PC centroid limits centered around 1000 Genomes CEU data such that 99% of the CEU data fell within the limits. Only MIDUS samples also falling within these limits ($N = 1309$) were considered ancestrally homogenous and, thus, were included in polygenic risk scoring.

The Illumina OmniExpress arrays tag a sufficient number of variants on the X and Y chromosomes to determine biological sex (e.g. 17,707 SNPs on X chromosome and 1,367 on Y for array v. 1.1). Samples were excluded if self-reported sex did not match biological sex as determined by genotype ($N = 13$), as this indicates either invalid self-reports, genotyping errors, or accidental I.D. swaps. In sum, after filtering out samples that did not pass ancestry- and sex-checks, PRSs were calculated for a final sample of $N = 1296$ participants.

Genotype Imputation: After MIDUS genotype samples were filtered via inclusion criteria, genotypes were imputed to approximately 80 million SNPs using minimac3 [10] and Eagle [11] with the 1000 Genomes reference panel [8] on the Michigan Imputation Server [10]. SNPs with ambiguous strand orientation, >5% missing calls, or Hardy-Weinberg equilibrium $p < 0.001$ were excluded prior to imputation. After imputation, SNPs with minor allele frequency below 0.01 or an average call rate (AvgCall) below 0.9 were excluded, resulting in ~8 million common variants with sufficient minor allele frequency and call rate for inclusion. This variant calling pipeline follows the established best practice methods from the Broad Institute, incorporated in the freely available Genome Analysis Toolkit (GATK) [12]. All genomic data were handled using Plink 1.9 [13-15].

Polygenic Risk Scoring:

To be more precise, a PRS may be defined “as a single value estimate of an individual’s propensity to a phenotype, calculated by computing the sum of risk alleles corresponding to a phenotype of interest in each individual, weighted by the effect size estimate of the most powerful GWAS on the phenotype”¹. PRS is traditionally calculated as $PRS_k = \sum_i \beta_i SNP_{ik}$, where PRS for individual k in the target sample is calculated by the summation of each SNP (measured for both the person k and passing a set association threshold in the discovery GWAS) multiplied by the effect size, β , of that SNP in the discovery GWAS. PRSs are then typically transformed to a standardized metric ($M = 0$, $SD = 1$).

PRSice 2.0 [1]^G was used to calculate polygenic risk scores. A default threshold of $p = 1.0$ (including all SNPs of infinitesimal effect) was used to calculate scores. A variety of thresholds can be used, and some researchers prefer a threshold of $p = 0.5$. Thresholds are not arbitrary. We adopt a more liberal threshold of $p = 1.0$, due to the diverse array of medical and psychiatric phenotypes, each with varying, undefined levels of genetic complexity and a tendency toward massive polygenicity. Stricter thresholds (0.3, 0.5) may be needed for any given PRS, particularly if there is precedence in the literature for such a prediction. If multiple thresholds for PRS are used, analytic p-values must be corrected for multiple comparison. Approximately 4000 phenotypes are currently available for PRS calculation, but for the MIDUS samples, only phenotypes with adequate discovery GWAS sample size ($N \sim 10,000$) were selected for polygenic risk scoring. This helps to minimize the standard error of the SNP weights used in the calculation of PRSs and increases predictive power. See Table 1 below for a list of phenotypes for which PRS data is available and associated references for discovery genome-wide association studies. Note, because this field is changing rapidly, researchers should use these references as a starting point and regularly check the literature for more recent advances.

PRS-CS was used to calculate polygenic risk score of Alzheimer Disease. PRS-CS is a Bayesian method inferring posterior SNP effect sizes under continuous shrinkage (CS) priors. We used the European LD reference panel constructed using the UK Biobank data. For better stability of the results, we ran a Markov chain Monte Carlo procedure in PRS-CS for 2,000 iterations (default is 1,000), burning in the first 1,000 iterations (default is 500). Other PRS-CS parameters were set to the default values. Note that unlike PRSice, PRS-CS does not have a p-value threshold parameter.

^G <https://choishingwan.github.io/PRSice/>

Table. Phenotypes and Samples Sizes from Discovery GWAS

Phenotype	N cases	N controls	N total
ADHD Females [16]	30,961	938,949	969,910
ADHD Males [16]	68,694	969,312	1,038,006
Adventurousness [17]	-	-	557,923
Age at Menarche [18]	-	-	370,000
Alzheimer's Disease [19]	25,580	48,466	74,046
Alzheimer's Disease [PRS-CS] [20]	90,338	1,036,225	1,126,563
Anorexia Nervosa [21]	3,495	10,982	14,477
Anxiety [22]	-	-	18,186
Asthma [23]	10,365	16,110	26,475
Autism Spectrum Disorder [24]	18,381	27,969	46,350
Automobile Speeding Propensity [17]	-	-	404,291
Bipolar Disorder [25]	29,764	169,118	198,882
Birth Height [26]	-	-	28,495
Birth Weight [27]	-	-	153,781
Body Mass Index [28]	-	-	234,069
Child IQ [29]	-	-	17,989
Cigarettes Per Day [30]	-	-	73,853
Cognitive Performance [31]	-	-	257,841
College Education [32]	20,040	75,387	95,427
Coronary Artery Disease [33]	22,233	64,762	86,995
Crohn's Disease [34]	6,299	15,148	21,447
Cross-Disorder [35]	33,332	27,888	61,220
Depressive Symptoms [36]	-	-	180,866
Disinhibition [17]	-	-	315,894
Drinks Per Week [17]	-	-	414,343
Educational Attainment [31]	-	-	766,344
Ever Smoked [17]	-	-	518,633
Extraversion [37]	-	-	63,030
General Risk Tolerance [17]	-	-	557,923
HDL Cholesterol [38]	-	-	99,900
Height [39]	-	-	253,288
Infant Head Circumference [40]	-	-	10,768
Intracranial Volume [41]	-	-	30,717
LDL Cholesterol [38]	-	-	95,454
Major Depressive Disorder [42]	59,851	113,154	173,005
Neuroticism [43]	-	-	168,105
Number of Sexual Partners [17]	-	-	370,711
Obsessive-Compulsive Disorder [44]	2,688	7,037	9,725
Post-Traumatic Stress Disorder [45]	2,488	7,466	9,954
Rheumatoid Arthritis [46]	5,539	20,169	25,708
Risk Tolerance [17]	-	-	939,908

Schizophrenia [47]	36,989	113,075	150,064
Subjective Well-Being [36]	-	-	298,420
Suicide Death*			
Total Cholesterol [37]	-	-	100,184
Triglycerides [37]	-	-	96,598
Type 2 Diabetes [48]	62,892	596,424	659,316
Ulcerative Colitis [34]	7,211	20,783	27,994
Waist-to-Hip Ratio [49]	-	-	142,762

*Contact lab for archival citation

References

1. Euesden, J., C.M. Lewis, and P.F. O'Reilly, *PRSice: Polygenic Risk Score software*. Bioinformatics, 2015. **31**(9): p. 1466-1468.
2. Novembre, J. and M. Stephens, *Interpreting principal component analyses of spatial population genetic variation*. Nature Genetics, 2008. **40**(5): p. 646-649.
3. Choi, S.W., T.S.-H. Mak, and P.F. O'Reilly, *Tutorial: A guide to performing polygenic risk score analyses*. Nature Protocols, 2020. **15**(9): p. 2759-2772.
4. Dudbridge, F., *Power and predictive accuracy of polygenic risk scores*. PLOS Genetics, 2013. **9**(3): p. Article e1003348.
5. Bowden, J., G. Davey Smith, and S. Burgess, *Mendelian randomization with invalid instruments: Effect estimation and bias detection through Egger regression*. International Journal of Epidemiology, 2015. **44**(2): p. 512-525.
6. Davey Smith, G. and G. Hemani, *Mendelian randomization: Genetic anchors for causal inference in epidemiological studies*. Human Molecular Genetics, 2014. **23**(R1): p. R89-R98.
7. Keller, M.C., *Gene \times environment interaction studies have not properly controlled for potential confounders: The problem and the (simple) solution*. Biological Psychiatry, 2014. **75**(1): p. 18-24.
8. Auton, A., et al., *A global reference for human genetic variation*. Nature, 2015. **526**(7571): p. 68-74.
9. Shabalin, A.A., et al., *RaMWAS: Fast methylome-wide association study pipeline for enrichment platforms*. Bioinformatics, 2018. **34**(13): p. 2283-2285.
10. Das, S., et al., *Next-generation genotype imputation service and methods*. Nature Genetics, 2016. **48**(10): p. 1284-1287.
11. Loh, P.-R., et al., *Reference-based phasing using the Haplotype Reference Consortium panel*. Nature Genetics, 2016. **48**(11): p. 1443-1448.
12. McKenna, A., et al., *The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data*. Genome Research, 2010. **20**(9): p. 1297-303.
13. Chang, C.C., et al., *Second-generation PLINK: Rising to the challenge of larger and richer datasets*. GigaScience, 2015. **4**(1): p. Article s13742-015-0047-8.
14. Purcell, S. and C. Chang, *PLINK 1.9*, www.Cog-genomics.Org/plink/1.9/. n.d.
15. Purcell, S., et al., *PLINK: A tool set for whole-genome association and population-based linkage analyses*. American Journal of Human Genetics, 2007. **81**(3): p. 559-575.
16. Martin, J., et al., *A genetic investigation of sex bias in the prevalence of attention-deficit/hyperactivity disorder*. Biological Psychiatry, 2018. **83**(12): p. 1044-1053.
17. Karlsson Linnér, R., et al., *Genome-wide association analyses of risk tolerance and risky behaviors in over 1 million individuals identify hundreds of loci and shared genetic influences*. Nature Genetics, 2019. **51**(2): p. 245-257.
18. Day, F.R., et al., *Genomic analyses identify hundreds of variants associated with age at menarche and support a role for puberty timing in cancer risk*. Nature Genetics, 2017. **49**(6): p. 834-841.

19. Lambert, J.-C., et al., *Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer's disease*. Nature Genetics, 2013. **45**(12): p. 1452-1458.
20. Wightman, D.P., et al., *A genome-wide association study with 1,126,563 individuals identifies new risk loci for Alzheimer's disease*. Nature Genetics, 2021. **53**(9): p. 1276-1282.
21. Duncan, L., et al., *Significant locus and metabolic genetic correlations revealed in genome-wide association study of anorexia nervosa*. American Journal of Psychiatry, 2017. **174**(9): p. 850-858.
22. Otowa, T., et al., *Meta-analysis of genome-wide association studies of anxiety disorders*. Molecular Psychiatry, 2016. **21**(10): p. 1391-1399.
23. Moffatt, M.F., et al., *A large-scale, consortium-based genomewide association study of asthma*. New England Journal of Medicine, 2010. **363**(13): p. 1211-1221.
24. Grove, J., et al., *Identification of common genetic risk variants for autism spectrum disorder*. Nature Genetics, 2019. **51**(3): p. 431-444.
25. Stahl, E.A., et al., *Genome-wide association study identifies 30 loci associated with bipolar disorder*. Nature Genetics, 2019. **51**(5): p. 793-803.
26. van der Valk, R.J., et al., *A novel common variant in DCST2 is associated with length in early life and height in adulthood*. Human Molecular Genetics, 2015. **24**(4): p. 1155-68.
27. Horikoshi, M., et al., *Genome-wide associations for birth weight and correlations with adult disease*. Nature, 2016. **538**(7624): p. 248-252.
28. Locke, A.E., et al., *Genetic studies of body mass index yield new insights for obesity biology*. Nature, 2015. **518**(7538): p. 197-206.
29. Benyamin, B., et al., *Childhood intelligence is heritable, highly polygenic and associated with FBNPIL*. Molecular Psychiatry, 2014. **19**(2): p. 253-8.
30. Furberg, H., et al., *Genome-wide meta-analyses identify multiple loci associated with smoking behavior*. Nature Genetics, 2010. **42**(5): p. 441-447.
31. Lee, J.J., et al., *Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals*. Nature Genetics, 2018. **50**(8): p. 1112-1121.
32. Rietveld, C.A., et al., *GWAS of 126,559 individuals identifies genetic variants associated with educational attainment*. Science, 2013. **340**(6139): p. 1467-1471.
33. Schunkert, H., et al., *Large-scale association analysis identifies 13 new susceptibility loci for coronary artery disease*. Nature Genetics, 2011. **43**(4): p. 333-338.
34. Jostins, L., et al., *Host-microbe interactions have shaped the genetic architecture of inflammatory bowel disease*. Nature, 2012. **491**(7422): p. 119-24.
35. Cross-Disorder Group of the Psychiatric Genomics Consortium, *Identification of risk loci with shared effects on five major psychiatric disorders: A genome-wide analysis*. Lancet, 2013. **381**(9875): p. 1371-1379.
36. Okbay, A., et al., *Genetic variants associated with subjective well-being, depressive symptoms, and neuroticism identified through genome-wide analyses*. Nature Genetics, 2016. **48**(6): p. 624-633.
37. van den Berg, S.M., et al., *Meta-analysis of genome-wide association studies for extraversion: Findings from the Genetics of Personality Consortium*. Behavior Genetics, 2016. **46**(2): p. 170-82.
38. Teslovich, T.M., et al., *Biological, clinical and population relevance of 95 loci for blood lipids*. Nature, 2010. **466**(7307): p. 707-13.

39. Wood, A.R., et al., *Defining the role of common variation in the genomic and biological architecture of adult human height*. Nat Genet, 2014. **46**(11): p. 1173-86.
40. Taal, H.R., et al., *Common variants at 12q15 and 12q24 are associated with infant head circumference*. Nature Genetics, 2012. **44**(5): p. 532-538.
41. Hibar, D.P., et al., *Common genetic variants influence human subcortical brain structures*. Nature, 2015. **520**(7546): p. 224-229.
42. Wray, N.R., et al., *Genome-wide association analyses identify 44 risk variants and refine the genetic architecture of major depression*. Nature Genetics, 2018. **50**(5): p. 668-681.
43. Turley, P., et al., *Multi-trait analysis of genome-wide association summary statistics using MTAG*. Nature Genetics, 2018. **50**(2): p. 229-237.
44. International Obsessive Compulsive Disorder Foundation Genetics Collaborative and O. C. D. Collaborative Genetics Association Studies, *Revealing the complex genetic architecture of obsessive-compulsive disorder using meta-analysis*. Molecular Psychiatry, 2018. **23**(5): p. 1181-1188.
45. Duncan, L.E., et al., *Largest GWAS of PTSD (n=20 070) yields genetic overlap with schizophrenia and sex differences in heritability*. Molecular Psychiatry, 2018. **23**(3): p. 666-673.
46. Stahl, E.A., et al., *Genome-wide association study meta-analysis identifies seven new rheumatoid arthritis risk loci*. Nature Genetics, 2010. **42**(6): p. 508-514.
47. Ripke, S., et al., *Biological insights from 108 schizophrenia-associated genetic loci*. Nature, 2014. **511**(7510): p. 421-427.
48. Xue, A., et al., *Genome-wide association analyses identify 143 risk variants and putative regulatory mechanisms for type 2 diabetes*. Nature Communications, 2018. **9**(1): p. Article 2941.
49. Shungin, D., et al., *New genetic loci link adipose and insulin biology to body fat distribution*. Nature, 2015. **518**(7538): p. 187-196.