# Genetic nature or genetic nurture? Introducing social genetic parameters to quantify bias in polygenic score analyses

Sam Trejo & Benjamin W. Domingue

Published online: 18 Dec 2019.

Submit your article to this journal

View related articles

View Crossmark data

Routledge
Taylor & Francis Group

Check for updates

# Genetic nature or genetic nurture? Introducing social genetic parameters to quantify bias in polygenic score analyses

Sam Trejo and Benjamin W. Domingue

Graduate School of Education, Stanford University, Stanford, CA, USA

### ABSTRACT

Results from a genome-wide association study (GWAS) can be used to generate a polygenic score (PGS), an individual-level measure summarizing identified genetic influence on a trait dispersed across the genome. For complex, behavioral traits, the association between an individual's PGS and their phenotype may contain bias (from geographic, ancestral, and/or socioeconomic confounding) alongside the causal effect of the individual's genes. We formalize the introduction of a different source of bias in regression models using PGSs: the effects of parental genes on offspring outcomes, known as genetic nurture. GWAS do not discriminate between the various pathways through which genes become associated with outcomes, meaning existing PGSs capture both direct genetic effects and genetic nurture effects. We construct a theoretical model for genetic effects and show that the presence of genetic nurture biases PGS coefficients from both naïve OLS (between-family) and family fixed effects (within-family) regressions. This bias is in opposite directions; while naïve OLS estimates are biased away from zero, family fixed effects estimates are biased toward zero. We quantify this bias using two novel parameters: (1) the genetic correlation between the direct and nurture effects and (2) the ratio of the SNP heritabilities for the direct and nurture effects.

## Introduction

### Genomics & the Social Sciences

Spurred by the plummeting cost of DNA sequencing and technological developments in processing large amounts of genetic data, researchers have made great strides in connecting genes to biological and social outcomes in a replicable manner. The key tool is the genome-wide association study (GWAS); a GWAS uses genotype and phenotype data from many individuals to probe the relationship between a given trait and thousands of regions of the genome (Pearson and Manolio 2008). GWAS are conducted on a wide variety of outcomes, ranging from proximal, biological phenotypes, such as blood pressure (Giri et al. 2019) and height (Yengo et al. 2018a), to distal, behavioral phenotypes, such as depression (Hyde et al. 2016; Okbay et al. 2016) and educational attainment (Lee et al. 2018).

Findings from GWAS are often used to generate a predictor – a polygenic score (PGS) – meant to summarize an individual's genetic predisposition for a given trait. PGSs offer great promise to social scientists interested in incorporating genes into

biosocial models of human behavior (Belsky and Israel 2014). In the short term, PGSs may be used as control variables in studies of environmental effects (Rietveld et al. 2013), used in gene–environment interaction studies to probe whether genetic effects are environmentally contingent (Barcellos, Carvalho, and Turley 2018; Papageorge and Thom 2017; Trejo et al. 2018), and used to better understand how genetic factors influence developmental processes (Belsky et al. 2013, 2016). In the long run, PGSs might be used to identify those who would benefit most from early medical or educational interventions (Torkamani, Wineinger, and Topol 2018) i.e., for a developmental disorder like dyslexia.

## The Problem of Confounding

A point of emphasis is that the same technique, GWAS, is being used to map the genetic architecture of a diverse set of phenotypes. It is not obvious that the methodology used to identify the underlying genetics of proximal, biological phenotypes can be deployed without side effect to interrogate the genetics of complex, socially contextualized phenotypes. Especially in the case of traits like depression and educational attainment, it is critical that existing GWAS results be interpreted cautiously (Martschenko, Trejo, and Domingue 2019); while PGSs have been shown to predict complex phenotypes, they capture a broad range of information and therefore the associations between an individual's PGS and their downstream outcomes cannot be readily interpreted as the causal effect of genes. An individual's genome contains fine-grain information about their place in the intricate structure of a population (Hamer and Sirota 2000; Novembre et al. 2008), meaning that GWASs for complex traits may inadvertently identify genes related to confounding environmental variables such as ancestry, geography, or socioeconomic status.

Recent work in human studies has begun to elucidate a novel source of confounding: social genetic effects (Domingue and Belsky 2017). Social genetic effects, also known as indirect genetic effects, are defined as the influence of one organism's genotype on a different organism's phenotype. The idea of social genetic effects originated in evolutionary theory (Moore, Brodie, and Wolf 1997; Wolf et al. 1998), and social genetic effects have been observed in animal populations (Baud et al. 2018; Bergsma et al. 2008; Canario, Lundeheim, and Bijma 2017; Petfield et al. 2005). Social science is now beginning to study such effects in human populations; examples include among social peers (Domingue et al. 2018; Sotoudeh, Mullan, and Conley 2019), sibling pairs (Cawley et al. 2017; Kong et al. 2018), and parents and their children (Armstrong-Carter et al. 2019; Bates et al. 2018; Kong et al. 2018; Wertz et al. 2018). Genetic nurturance refers to the social genetic effect that parents have on their children. The existence of within-family social genetic effects, like genetic nurture effects, complicates attempts to derive causal estimates from GWAS.

For recent breakthroughs in the genetic architecture of complex traits to provide novel insights to researchers in the biomedical and social sciences, the relationships discovered in a GWAS must mostly reflect causal relationships between an individual's genes and their phenotype. If, for example, the genes identified for a complex trait predict it only through spurious correlation, PGSs will provide little use toward broadening our understanding of genetic and environmental influences. Thus, validating PGSs within families is vitally important for sifting out causation from correlation among the genetics identified in GWASs of complex traits (Belsky et al. 2018; Domingue et al. 2015; Lee et al. 2018; Rietveld et al. 2014). Environmental differences are muted between siblings and, conditional on parental genotype,

child genotype is randomly assigned through a process known as genetic recombination (Conley and Fletcher 2017). This makes family fixed effect regression models that compare genetic differences in siblings to phenotypic differences in siblings the gold standard for testing and understanding whether genes are causally related to downstream outcomes. Within-family research designs, however, are not without their own complications. Genetic nurture may lead to bias in estimates derived from within-family studies, though the extent of this bias has not yet been explored.

## Accounting for Genetic Nurture

In this paper, we describe how genetic nurture influences PGS construction and can introduce bias into within-family and between-family regression analyses using PGSs. We construct a theoretical model for additive genetic effects and show that, unlike other sources of bias in PGSs, the presence of genetic nurture can bias PGS coefficients from both naïve OLS (between-family) regressions and family fixed effects (within-family) regressions. We quantify the magnitude of this bias for a given trait using two novel parameters: (1) the genetic correlation between the direct and nurture effects and (2) the ratio of the SNP heritabilities for the direct and nurture effects. Bias is in opposite directions; whereas naïve OLS estimates are biased upwards, family fixed effects estimates are biased downwards. These findings highlight a shortcoming of existing PGSs and have important implications for the use and interpretation of research designs using PGSs for traits where genetic nurture is a relevant causal pathway.

## Empirical Motivation

### Empirical Model

We motivate our theoretical framework by first considering the empirical specifications used in recent work (Belsky et al. 2018; Domingue et al. 2015; Lee et al. 2018). Consider the following two models relating an individual's PGS constructed from recent GWAS results ($\widehat{PGS'}^{D}_{ij}$) to their outcome ($Y_{ij}$):

$$\text{Model 1}: Y_{ij} = \hat{\psi}_0 + \hat{\psi}_1 \widehat{PGS'}^{D}_{ij} + \boldsymbol{X_{ij}}\widehat{\boldsymbol{\Theta}} + \epsilon_{ij}$$
$$\text{Model 2}: Y_{ij} = \hat{\pi}_0 + \hat{\pi}_1 \widehat{PGS'}^{D}_{ij} + \boldsymbol{X_{ij}}\widehat{\boldsymbol{\Theta}} + \Gamma_j + \epsilon_{ij}$$

(2a.i)

$\widehat{PGS'}^{D}_{ij}$ : Normalized PGS constructed from the observed linear relationship between genotype and outcome
$Y_{ij}$ : Outcome for individual $i$ in family $j$
$\Gamma_j$ : Family $j$ fixed effect
$\boldsymbol{X_{ij}}$ : $12 \times 1$ vector of covariates comprised of sex, age, and the first 10 principal components of genotype

Model 1 treats individuals as though they are unrelated whereas Model 2 compares siblings using a family fixed effect. Thus, Model 1 leverages covariation in $Y_{ij}$ and $\widehat{PGS'}^{D}_{ij}$ between individuals from different families while Model 2 compares individuals in the same family. In effect, Model 2 asks whether sibling differences in PGS translate into sibling differences in the outcome.

## Unresolved Questions

We consider a brief empirical example to motivate scrutiny of between- versus within-family findings. Table 1 displays results from Model 1 and Model 2 using data from the National Longitudinal Study of Adolescent to Adult Health (Harris et al. 2019) for six phenotypes: educational attainment (Lee et al. 2018), cognitive ability (Lee et al. 2018), depressive symptoms (Turley et al. 2018), birth weight (Warrington et al. 2019), body mass index (Locke et al. 2015), and height (Wood et al. 2014). We further discuss the Add Health data and PGS construction in Sections A1 and A2 of the Appendix.

In Table 1, a one standard deviation increase in the educational attainment PGS is associated with an additional 0.8 year of schooling between-families ($\hat{\psi}_1$) but less than half of that within-families ($\hat{\pi}_1$). If we compare the six phenotypes, the relative size of $\hat{\psi}_1$ and $\hat{\pi}_1$ (as captured by their ratio) varies dramatically. For years of schooling and cognitive performance, bootstrapped $p$-values show that the differences seen within and between-family are statistically significant (i.e. $\hat{\psi}_1 \neq \hat{\pi}_1$). These findings are consistent with those from other data; using data from the United Kingdom, similar analyses found that PGS coefficients for cognitive traits were on average 60% greater between families than within-families (Selzam et al. 2019). In both cases, there was no evidence for differences between within- and between-family results for non-cognitive traits.

What drives differences between $\hat{\psi}_1$ and $\hat{\pi}_1$? One possibility is that the between-family models are confounded while the within-family models capture the true causal effects of the PGS. Alternatively, it may be that some of the processes captured by GWAS function differently within-families versus between-families (genetic nurturance, sibling spillovers, niche formation, etc.). Answering this question requires a formal treatment so as to parse differences between $\hat{\psi}_1$ and $\hat{\pi}_1$ across phenotypes. Our theoretical model, developed below, suggests that bias in both $\hat{\psi}_1$ and $\hat{\pi}_1$ may depend in part on two novel parameters: (1) the

**Table 1.** The association between polygenic score and observed trait for six phenotypes, within-families and between-families.

| | $\hat{\psi}_1$ | $\hat{\pi}_1$ | $\frac{\hat{\pi}_1}{\hat{\psi}_1}$ | $p(\hat{\psi}_1 = \hat{\pi}_1)$ |
|---|---|---|---|---|
| Years of Schooling | 0.81** | 0.35** | 0.44 | <0.01 |
| Cognitive Ability | 3.16** | 1.70** | 0.54 | 0.02 |
| CESD Depression Index | 0.13** | 0.06 | 0.42 | 0.25 |
| Birth Weight | 3.10** | 3.26* | 1.05 | 0.91 |
| Body Mass Index | 2.01** | 2.34** | 1.17 | 0.59 |
| Height | 2.50** | 2.51** | 1.00 | 0.97 |

* 0.05 ** 0.01 . All models control for sex, age, and the first 10 principal components of individual genotype. All models use only individuals of European ancestry. Models without individual-fixed effects use a sample of unrelated individuals, whereas the family-fixed effect models use a sample of sibling pairs. The sample of unrelated individuals contains one randomly selected sibling from each pair. All polygenic scores are standardized within sample to be mean 0 and standard deviation 1. Cognitive ability is measured through the Peabody Picture Vocabulary Test during Wave 1 of Add Health, when respondents were approximately 16 years old. Birth weight is retrospectively reported by respondents' parents during Wave 1 of Add Health. Years of schooling, CESD depression index, body mass index, and height are measured during Wave 4 of Add Health, when respondents were approximately 28 years old. Height is reported in centimeters, birth weight is reported in ounces, and cognitive ability is reported in IQ score points. The CESD depression index is normalized within sample to be mean 0 and standard deviation 1. A traditional regression table is available in Section A12 of the Appendix. $P$-values for the test that $\hat{\psi}_1 = \hat{\pi}_1$ were calculated through bootstrap resampling with replacement simulated with 1000 repetitions.

underlying genetic correlation of direct and nurture effects and (2) the ratio of the SNP heritabilities for direct effects and nurture effects.

## Theoretical Model

### *Direct Genetic Effects and Genetic Nurture Effects*

Historically, biosocial analyses have modeled complex traits as a function of both direct genetic effects and environment influences on an individual. Motivated by recent work highlighting the relevance of genetic nurture effects (Bates et al. 2018; Belsky et al. 2018; Kong et al. 2018; Wertz et al. 2018), we extend this model to include the genes of an individual's parent. Thus, we assume the outcome $Y_{ij}$ is a function of individual $i$'s genotype, the genotypes of the parents in family $j$, and distinct individual-level and family-level environments. We choose to have a common effect of parental genetics at a given locus, instead of separate maternal and paternal effects, given the current lack of strong empirical evidence of differences across parents (Kong et al. 2018). Note also that environmental components are defined as the strictly non-genetic sources of variation in $Y_{ij}$. In other words, an environment influences a child's outcome irrespective of their or their parents' genotype. To the extent that family-level features of a child's environment are a result of their parents' genotype (i.e. some portion of a family's socioeconomic status), they are captured by terms involving parental genotype.

$$Y_{ij} = \beta_0 + f(G_{ij}) + f(G_j) + f(E_{ij}) + f(E_j) + \epsilon_{ij} \qquad (3a.i)$$

$f(G_{ij})$ : Effect of $i$'s genome on $Y_{ij}$
$f(G_j)$ : Effect of family $j$'s genome on $Y_{ij}$
$f(E_{ij})$ : Effect of $i$'s environment on $Y_{ij}$
$f(E_j)$ : Effect of family $j$'s environment on $Y_{ij}$

We make three important assumptions to simplify the exposition of this model: no gene–environment interaction, no gene–environment correlation, and no assortative mating (see Section A3 of the Appendix for additional details regarding these assumptions). We consider the likely implications of violations for our results in the Discussion. While our model is simplistic, it illustrates the key empirical phenomenon of interest; these higher-order features of the real world should not change the key implications derived from our model.

### *True Polygenic Scores*

Complex, behavioral traits are associated with many genes across the genome that simultaneously produce very small effects (Chabris et al. 2015; Visscher et al. 2017). To increase statistical power and simplify computation, researchers often summarize the relevant genetics of individual $i$ into a single linear predictor called a PGS (Dudbridge 2013). This has become a widely utilized technique (Duncan et al. 2019) and standard usage relies on the assumptions that genetic effects are linear and additive. Recent meta-analyses of twin studies support the linear, additive model for genetic effects (Polderman et al. 2015). For the remainder of the paper, we approximate $f(G_{ij})$ and $f(G_j)$ in our theoretical model (3a.i) using PGSs that summarize information over $n$ independent genetic loci.

$$f\left(G_{ij}\right) \approx \sum_{z=1}^{n} \alpha^{z} g_{ij}^{z} = PGS_{ij}^{D}$$

$$f\left(G_{j}\right) \approx \sum_{z=1}^{n} \delta^{z} g_{j}^{z} = PGS_{j}^{N}$$

(3b.i)

$\alpha^{z}$ : True causal effect of a one allele change at $i$'s genetic locus $z$ on $Y_{ij}$

$\delta^{z}$ : True causal effect of a one allele change at either parent in family $j$'s genetic locus $z$ on $Y_{ij}$

$g_{ij}^{z}$ : Total number of risk alleles at $i$'s genetic locus $z$ $(0, \ 1, \text{or } 2)$

$g_{j}^{z}$ : Total number of parental risk alleles in family $j$ at genetic locus $z$ $(0, \ 1, \ 2, \ 3, \text{or } 4)$

$PGS_{ij}^{D}$: PGS constructed from the true casual linear effect of $i$'s genes on $Y_{ij}$

$PGS_{j}^{N}$: PGS constructed from the true casual linear effect of the parents in family $j$'s genes on $i$'s $Y_{ij}$

Notice that $\boldsymbol{\alpha}$ is the vector of causal allelic weights used to construct the true, underlying PGS for direct genetic effects. In the same vein, $\boldsymbol{\delta}$ is the vector of causal allelic weights used to construct the true, underlying PGS for genetic nurture effects. Note that both $\boldsymbol{\alpha}$ and $\boldsymbol{\delta}$ are structural parameters that are never empirically observed.

$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha^{1} \\ \alpha^{2} \\ \vdots \\ \alpha^{n} \end{bmatrix} \qquad \boldsymbol{\delta} = \begin{bmatrix} \delta^{1} \\ \delta^{2} \\ \vdots \\ \delta^{n} \end{bmatrix}$$

(3b.ii)

$\boldsymbol{\alpha}$ : $n \times 1$ vector of the true causal effects of a one allele change at $i$'s genetic locus $z$ on $Y_{ij}$

$\boldsymbol{\delta}$ : $n \times 1$ vector of the true causal effects of a one allele change at a parent in family $j$'s genetic locus $z$ on $Y_{ij}$

We can now rewrite our theoretical model (3a.i) using PGSs.

$$Y_{ij} = \beta_{0} + \beta_{1} PGS_{ij}^{D} + \beta_{2} PGS_{j}^{N} + f\left(E_{ij}\right) + f\left(E_{j}\right)$$

(3b.iii)

Notice that, because $PGS_{ij}^{D}$ and $PGS_{j}^{N}$ are not standardized, an individual's value for $PGS_{ij}^{D}$ and $PGS_{j}^{N}$ represents the true effect of their genes and their parents' genes, respectively, on $Y_{ij}$ in the units of $Y_{ij}$. Thus, $\beta_{1}$ and $\beta_{2}$ are both equal to 1 by construction and Equation 3b.iii can equivalently be written as:

$$Y_{ij} = \beta_{0} + PGS_{ij}^{D} + PGS_{j}^{N} + f\left(E_{ij}\right) + f\left(E_{j}\right)$$

(3b.iv)

## Transmitted Genetic Nurture Alleles

The presence of social genetic effects, such as genetic nurture effects, will only bias GWAS estimates of direct genetic effects when a social or biological process induces a correlation between the genetics of an individual and the genetics of his or her relevant social relationships. In the case of genetic nurture effects, biological recombination acts as such a process;

children randomly inherit half of each parent's genome, leading to a mechanical correlation between parental genetics and child genetics. To capture the portion of the genetic nurture PGS that was transmitted to individual $i$ in family $j$ from their parents, we introduce a third PGS parameter, $PGS_{ij}^{N}$, that is absent from our formal model that determines outcomes.

$$\sum_{z=1}^{n} \delta^z g_{ij}^z = PGS_{ij}^{N} \tag{3c.i}$$

$PGS_{ij}^{N}$: PGS constructed from the true causal linear effect of $i$'s genes on $i$'s child's $Y$

$PGS_{ij}^{N}$ is constructed using aspects of both $PGS_{ij}^{D}$ and $PGS_{j}^{N}$, the two PGSs corresponding to the two causal sources of genetic effects present in our theoretical model. Like $PGS_{ij}^{D}$, $PGS_{ij}^{N}$ is constructed using $g_{ij}$ (as opposed to $g_j$) and therefore varies within-families. However, like $PGS_{j}^{N}$, $PGS_{ij}^{N}$ is constructed using the allelic weights $\boldsymbol{\delta}$, which correspond to genetic nurture effects (as opposed to direct genetic effects).

The relationship $PGS_{ij}^{N}$ and $PGS_{j}^{N}$ hinges on the relationship between $g_{ij}$ and $g_j$. Because the alleles transmitted from parent $g_j$ to child $g_{ij}$ are determined stochastically through genetic recombination, we can compute the correlation between $g_{ij}$ and $g_j$. This, in turn, leads us to the correlation between the relevant polygenic scores (see Section A4 of the Appendix).

$$\rho_{PGS_{ij}^{N},\, PGS_{j}^{N}} = \frac{\sqrt{2}}{2} \tag{3c.ii}$$

Note that this quantity does not vary between traits; this is due to the fact that same vector of allelic weights is used to construct both $PGS_{ij}^{N}$ and $PGS_{j}^{N}$ and differences are due entirely to the trait-independent recombination of parental alleles.

### The Relationship between Direct Genetic Effects and Genetic Nurture Effects

While $PGS_{ij}^{N}$ is absent from our underlying theoretical model, it provides the link between $PGS_{ij}^{D}$ and $PGS_{j}^{N}$ that allows for the presence of genetic nurture effects to distort GWAS results and subsequent PGS analysis. Thus, crucial elements of our theoretical model are relationships between $PGS_{ij}^{D}$ and $PGS_{ij}^{N}$ and between $PGS_{ij}^{N}$ and $PGS_{j}^{N}$. As we have seen above, $PGS_{ij}^{N}$ and $PGS_{j}^{N}$ have a mechanical correlation that does not vary between traits as it depends only on genotype (and not allelic weights). However, any differences between $PGS_{ij}^{D}$ and $PGS_{ij}^{N}$ are a result of differences between allelic weights used to construct each PGS ($\boldsymbol{\alpha}$ and $\boldsymbol{\delta}$, respectively). Because $\boldsymbol{\alpha}$ and $\boldsymbol{\delta}$ vary across traits, relationship between $PGS_{ij}^{D}$ and $PGS_{ij}^{N}$ also vary between traits. We term the correlation between $PGS_{ij}^{D}$ and $PGS_{ij}^{N}$ for a given trait the *direct-nurture genetic correlation*.

$$\rho_g = \frac{cov\left(PGS_{ij}^{D}, PGS_{ij}^{N}\right)}{var\left(PGS_{ij}^{D}\right)^{\frac{1}{2}} var\left(PGS_{ij}^{N}\right)^{\frac{1}{2}}} \tag{3d.i}$$

$\rho_g$ : Correlation between genetic nurture effects and direct genetic effects

In general, we expect the direct-nurture genetic correlation to be positive. To describe the relevant intuition, consider the case of educational attainment. Presumably, some genetic pathways that lead a parent to create an environment conducive to their child succeeding in school will also have impacted the parent's schooling. However, there may also exist genetic pathways that contribute to a parent's ability to create a positive educational environment for their children that do not influence the amount of educational attainment that the parent themselves received. In a related vein, recent evidence (Warrington et al. 2019, Figure, 2) suggests that the maternal genes that predict fetal birthweight also predict an individual's educational attainment whereas the fetal genes that predict birthweight are unrelated to those genes that predict educational attainment. Based on this logic, we expect direct genetic effects and genetic nurture effects to have a correlation of less than 1. While this value could in fact be negative, existing evidence suggests that $\rho_g$ is typically positive (Armstrong-Carter et al. 2019; Bates et al. 2018; Belsky et al. 2018; Kong et al. 2018; Wertz et al. 2018).

## *Underlying Allelic Weights*

We now turn to a discussion of $\alpha^z$ and $\delta^z$. In the course of this discussion, we will show that $\rho_g$ from Equation 3d.i is a critical structural parameter of our theoretical model. We first discuss $\alpha^z$. Allelic weights $\alpha^z$ ($z \in \{1 \ldots n\}$) are taken from a distribution with variance $\sigma$. Without loss of generality, we define a "risk" allele at each genetic locus $z$ such that this distribution has mean 0. Turning to $\delta^z$, we assume that $\delta$ are drawn from an identical distribution as $\alpha$, except with variance $\frac{\lambda^2}{4}\sigma$. Note that this is related to the variance of the $\alpha$ distribution but scaled by a parameter $\lambda$; this allows for the average effects sizes to differ between direct genetic effects and genetic nurture effects. Finally, we assume that genetic nurture effects and genetic nature effects have a similar genetic architecture such that $\alpha$ are $\delta$ distributed identically across genetic loci with respect to the mean and variance of the risk allele frequencies. See Section A6 of the Appendix for a formal treatment of this assumption.

The assumptions of our theoretical model entail that $\lambda$ represents the ratio of the SNP heritabilities of genetic nurture effect and direct genetic effects (see Section A7 of the Appendix). We call this value the *direct-nurture heritability ratio*.

$$\lambda = \frac{2 \sum_{z=1}^{n} (\delta^z \overline{g_{ij}^z})}{\sum_{z=1}^{n} (\alpha^z \overline{g_{ij}^z})} = \frac{h_N^2}{h_D^2} \tag{3e.i}$$

$\lambda$ : Ratio of the SNP heritabilities of genetic nurture effects and direct genetic effects
$h_D^2$ : SNP heritability for direct genetic effects of $Y_{ij}$
$h_N^2$ : SNP heritability for genetic nurture effects of $Y_{ij}$

Without loss of generality, we normalize all variables such that the variance of $PGS_{ij}^D$ is equal to one.

$$var\left(PGS_{ij}^D\right) = 1 \tag{3e.ii}$$

We can now use our theoretical model to derive the variance of the remaining PGSs as a function of the direct-nurture heritability ratio (see Sections A8 and A9 of the Appendix).

$$var\left(PGS_{ij}^{N}\right) = \frac{\lambda^2}{4}$$

$$var\left(PGS_{j}^{N}\right) = \frac{\lambda^2}{2} \tag{3e.iii}$$

Using these derivations, we can now gain insight into the implications of genetic nurture effects for GWAS and PGSs.

## Analytic Results

### *Observed PGS*

Up until this point, all our work has been theoretical; we have defined the functional form of a set of causal relationships between underlying parameters of interest which are difficult to observe directly. We now transport our theoretical model into the real world and consider its implications for GWAS and subsequent PGSs. In reality, we observe not $\boldsymbol{\alpha}$ but $\widehat{\boldsymbol{\alpha}}$. We then use this observed $\widehat{\boldsymbol{\alpha}}$ to construct not $PGS_{ij}^{D}$ but $\widehat{PGS_{ij}^{D}}$, which, as we will see, contains information from both $PGS_{ij}^{D}$ and $PGS_{ij}^{N}$.

We obtain $\widehat{\boldsymbol{\alpha}}$ by fitting the following regression for $n$ SNPs via GWAS.

$$Y_{ij} = \hat{\beta}_0^z + \hat{\alpha}^z g_{ij}^z + \boldsymbol{X_{ij}}\widehat{\boldsymbol{\Theta}} + \epsilon_{ij} \tag{4a.i}$$

$\hat{\alpha}^z$ : Allelic weight from the observed linear relationship of a one allele change at $i'$s $z^{th}$ gene and $Y_{ij}$

We can now plug in from our theoretical model (3b.iv) to derive how the estimator for each allelic weight is impacted by omitted variable bias (Wooldridge 2015, Chapter, 3).

$$\hat{\alpha}^z = \frac{var\left(g_{ij}^z\right)\alpha^z + cov\left(g_{ij}^z, g_j^z\right)\delta^z}{var\left(g_{ij}^z\right)} \tag{4a.ii}$$

To further analyze this expectation, we can separate $g_j^z$ into the sum of alleles that were transmitted to $i$ and the alleles that were not transmitted.

$$g_j^z = g_{ij}^z + f_{ij}^z \tag{4a.iii}$$

$f_{ij}^z$ : Total number of risk alleles at the parents in family $j'$s genetic locus $z$ that were **not** transmitted to $i$ $(0, 1, \text{or } 2)$

Thus yielding:

$$\hat{\alpha}^z = \frac{var\left(g_{ij}^z\right)\alpha^z + cov\left(g_{ij}^z, g_{ij}^z + f_{ij}^z\right)\delta^z}{var\left(g_{ij}^z\right)}$$

$$\hat{\alpha}^z = \frac{var\left(g_{ij}^z\right)\alpha^z + cov\left(g_{ij}^z, g_{ij}^z\right)\delta^z + cov\left(g_{ij}^z, f_{ij}^z\right)\delta^z}{var\left(g_{ij}^z\right)}$$

(4a.iv)

In the absence of assortative mating, transmitted alleles are uncorrelated with non-transmitted allele, meaning that $cov\left(g_{ij}^z, f_{ij}^z\right) = 0$.

$$\hat{\alpha}^z = \frac{var\left(g_{ij}^z\right)\alpha^z + cov\left(g_{ij}^z, g_{ij}^z\right)\delta^z}{var\left(g_{ij}^z\right)}$$

$$\hat{\alpha}^z = \frac{var\left(g_{ij}^z\right)\alpha^z + var\left(g_{ij}^z\right)\delta^z}{var\left(g_{ij}^z\right)}$$

(4a.v)

$$E[\hat{\alpha}^z] = E\left[\frac{var\left(g_{ij}^z\right)\alpha^z + var\left(g_{ij}^z\right)\delta^z}{var\left(g_{ij}^z\right)}\right]$$

$$E[\hat{\alpha}^z] = \alpha^z + \delta^z$$

We can already see that the estimated allelic weights $\widehat{\boldsymbol{\alpha}}$ intermingle both $\boldsymbol{\alpha}$ and $\boldsymbol{\delta}$. As a result, the quantity typically used for analysis, $\widehat{PGS}_{ij}^{D}$, will contain information about both the direct genetic effects and the genetic nurture effects (4a.vi). We can see this by using $\widehat{\boldsymbol{\alpha}}$ to construct our observed direct genetic PGS, $\widehat{PGS}_{ij}^{D}$.

$$\widehat{PGS}_{ij}^{D} = \sum_{z=1}^{n} \hat{\alpha}^z g_{ij}^z$$

$$E\left[\widehat{PGS}_{ij}^{D}\right] = E\left[\sum_{z=1}^{n} \hat{\alpha}^z g_{ij}^z\right]$$

(4a.vi)

$$E\left[\widehat{PGS}_{ij}^{D}\right] = \sum_{z=1}^{n} (\alpha^z + \delta^z) g_{ij}^z$$

$$E\left[\widehat{PGS}_{ij}^{D}\right] = PGS_{ij}^{D} + PGS_{ij}^{N}$$

Finally, constructed PGSs are typically normalized within sample, so we convert $\widehat{PGS}_{ij}^{D}$ to $\widehat{PGS}_{ij}^{\prime D}$.

$$\widehat{PGS}_{ij}^{\prime D} = \frac{\widehat{PGS}_{ij}^{D} - \overline{\widehat{PGS}_{ij}^{D}}}{var\left(\widehat{PGS}_{ij}^{D}\right)^{\frac{1}{2}}}$$

$$\mathrm{E}\left[\widehat{PGS}'^{\mathrm{D}}_{ij}\right] = \frac{\left(PGS^{\mathrm{D}}_{ij} + PGS^{\mathrm{N}}_{ij}\right) - \left(\overline{PGS^{\mathrm{D}}_{ij}} + \overline{PGS^{\mathrm{N}}_{ij}}\right)}{var\left(PGS^{\mathrm{D}}_{ij} + PGS^{\mathrm{N}}_{ij}\right)^{\frac{1}{2}}} \tag{4a.vii}$$

Thus, we have shown that the PGS derived from GWAS has information from both an individual's PGS for direct genetic effects and genetic nurture effects.

### *Between-Family Analyses*

Our theoretical model suggests that PGSs constructed from GWAS estimates capture both the direct genetic effects and the genetic nurture effects of a given allele. We now explore the implications of this result for analyses using such PGSs, beginning with the between-family analysis (Model 1). We use the derived $\widehat{PGS}'^{\mathrm{D}}_{ij}$ to calculate the expected bias in $\hat{\psi}_1$ from Model 1,

$$\text{Model 1}: Y_{ij} = \hat{\psi}_0 + \hat{\psi}_1 \widehat{PGS}'^{\mathrm{D}}_{ij} + \boldsymbol{X_{ij}\hat{\Theta}} + \hat{\epsilon}_{ij} \tag{4b.i}$$

We will see that the inclusion of genetic nurture effects biases direct genetic effect coefficients away from zero. Recall that (due to fact that $PGS^{\mathrm{D}}_{ij}$ is unstandardized in our theoretical model (3b.iv)) $\psi_1$ and $\pi_1$ are equal to one. Thus, the expected values of $\mathrm{E}\left[\hat{\psi}_1\right]$ and $\mathrm{E}[\hat{\pi}_1]$ represent the expected inflation or deflation of estimated PGS coefficients (i.e. bias) in between-family and within-family analyses, respectively.

$$\hat{\psi}_1 = \frac{cov\left(\widehat{PGS}'^{\mathrm{D}}_{ij}, Y_{ij}\right)}{var\left(\widehat{PGS}'^{\mathrm{D}}_{ij}\right)} \tag{4b.ii}$$

Note that $cov\left(\widehat{PGS}'^{\mathrm{D}}_{ij}, Y_{ij}\right)$ is given by true causal relationships from our theoretical model (3b.iv).

$$\hat{\psi}_1 = \frac{cov\left(\widehat{PGS}'^{\mathrm{D}}_{ij}, PGS^{\mathrm{D}}_{ij}\right)\beta_1 + cov\left(\widehat{PGS}'^{\mathrm{D}}_{ij}, PGS^{\mathrm{N}}_{j}\right)\beta_2}{var\left(\widehat{PGS}'^{\mathrm{D}}_{ij}\right)} \tag{4b.iii}$$

$\beta_1$ and $\beta_2$ from our theoretical model are equal to 1 by construction and therefore fall away.

$$\hat{\psi}_1 = \frac{cov\left(\widehat{PGS}'^{\mathrm{D}}_{ij}, PGS^{\mathrm{D}}_{ij}\right) + cov\left(\widehat{PGS}'^{\mathrm{D}}_{ij}, PGS^{\mathrm{N}}_{j}\right)}{var\left(\widehat{PGS}'^{\mathrm{D}}_{ij}\right)} \tag{4b.iv}$$

We solve (see Section A9 of the Appendix for details) to obtain the magnitude of the bias in $\hat{\psi}_1$.

$$\mathrm{E}\left[\hat{\psi}_1\right] = \sqrt{1 + \lambda\rho_g + \frac{\lambda^2}{4}} \tag{4b.v}$$

Thus, estimates for $\hat{\psi}_1$ will be biased upwards by a factor of $\sqrt{1 + \lambda\rho_g + \frac{\lambda^2}{4}}$. We will unpack this quantity further in the discussion but note that it depends on unobserved parameters $\lambda$ and $\rho_g$.

### Within-Family Analyses

Let us now turn to the within-family analysis (Model 2),

$$\text{Model 2}: Y_{ij} = \hat{\pi}_0 + \hat{\pi}_1 \widehat{PGS}'^D_{ij} + \boldsymbol{X_{ij}}\widehat{\boldsymbol{\Theta}} + \Gamma_j + \hat{\epsilon}_{ij}$$

We will see that the inclusion of genetic nurture effects biases direct genetic effect coefficients toward zero. We first translate Model 2 to an equivalent model based on differences (Wooldridge 2015, Chapter 14):

$$\left(Y_{1j} - Y_{0j}\right) = \hat{\pi}_1\left(\widehat{PGS}'^D_{1j} - \widehat{PGS}'^D_{0j}\right) + \left(\hat{\epsilon}_{1j} - \hat{\epsilon}_{0j}\right)$$
$$\Delta^1_0 Y_{ij} = \hat{\pi}_1 \Delta^1_0 \widehat{PGS}'^D_{ij} + \Delta^1_0 \hat{\epsilon}_{ij}$$

(4c.i)

$PGS^N_{0j}$ : $PGS^N_{ij}$ of sibling 0 in family $j$
$PGS^N_{1j}$ : $PGS^N_{ij}$ of sibling 1 in family $j$

As before, we begin by deriving the expected value of $\hat{\pi}_1$.

$$\hat{\pi}_1 = \frac{cov\left(\Delta^1_0 \widehat{PGS}'^D_{ij}, \Delta^1_0 Y_{ij}\right)}{var\left(\Delta^1_0 \widehat{PGS}'^D_{ij}\right)}$$

(4c.ii)

$cov\left(\Delta^1_0 \widehat{PGS}'^D_{ij}, \Delta^1_0 Y_{ij}\right)$ is given by true causal relationships from our theoretical model.

$$\hat{\pi}_1 = \frac{cov\left(\Delta^1_0 \widehat{PGS}'^D_{ij}, \Delta^1_0 PGS^D_{ij}\right)\beta_1 + cov\left(\Delta^1_0 \widehat{PGS}'^D_{ij}, \Delta^1_0 PGS^N_j\right)\beta_2}{var\left(\Delta^1_0 \widehat{PGS}'^D_{ij}\right)}$$

(4c.iii)

Again, $\beta_1$ and $\beta_2$ from our theoretical model are equal to 1 by construction and therefore fall away.

$$\hat{\pi}_1 = \frac{cov\left(\Delta^1_0 \widehat{PGS}'^D_{ij}, \Delta^1_0 PGS^D_{ij}\right) + cov\left(\Delta^1_0 \widehat{PGS}'^D_{ij}, \Delta^1_0 PGS^N_j\right)}{var\left(\Delta^1_0 \widehat{PGS}'^D_{ij}\right)}$$

(4c.iv)

Notice that between siblings there is no variation in family genetic nurturing environment, meaning that $\Delta^1_0 PGS^N_j = 0$.

$$\hat{\pi}_1 = \frac{cov\left(\Delta^1_0 \widehat{PGS}'^D_{ij}, \Delta^1_0 PGS^D_{ij}\right)}{var\left(\Delta^1_0 \widehat{PGS}'^D_{ij}\right)}$$

(4c.v)

Now we just solve (see Section A11 of the Appendix for details) to obtain the magnitude of the bias in $\hat{\pi}_1$.

$$E[\hat{\pi}_1] = \frac{1 + \frac{\lambda \rho_g}{2}}{\sqrt{1 + \lambda \rho_g + \frac{\lambda^2}{4}}} \tag{4c.vi}$$

Recall that, because $PGS_{ij}^{D}$ is unstandardized in our theoretical model (3b.iv), the expected value of $\hat{\pi}_1$ represents inflation or deflation of PGS coefficient estimates and is itself interpretable as a measure of bias. Thus, we have shown that our observed estimates for $\hat{\psi}_1$ will be biased downwards by a factor of $\frac{1 + \frac{\lambda \rho_g}{2}}{\sqrt{1 + \lambda \rho_g + \frac{\lambda^2}{4}}}$. We discuss this quantity further in the discussion.

## Discussion

### Bias

Our theoretical model illustrates that, in the presence of genetic nurture (i.e. $h_N^2 \neq 0$), regression analyses using PGSs to estimate the effects of an individual's genetics on their outcomes will suffer from bias. Between-family OLS models will be biased upwards by a factor of $\sqrt{1 + \lambda \rho_g + \frac{\lambda^2}{4}}$ while within-family fixed effect models will be biased downwards by a factor of $\frac{1 + \frac{\lambda \rho_g}{2}}{\sqrt{1 + \lambda \rho_g + \frac{\lambda^2}{4}}}$. Figure 1 plots the bias in both OLS and family fixed effect regressions using PGSs as a function of various direct-nurture genetic correlations and direct-nurture heritability ratios.

The absence of bias is represented in Figure 1 by the horizontal gray line. For any given trait, bias is always more extreme in the between-family models. The magnitude of bias is a function of two parameters: $\rho_g$, the trait's direct-nurture genetic correlation, and $\lambda$, the trait's direct-nurture heritability ratio. As $\lambda$ increases, so does the magnitude of the bias. However, $\rho_g$ has opposing effects on the bias within and between families; as $\rho_g$ increases, we note more less within-family bias toward zero more between-family bias away from zero. Thus, taken together, $\hat{\psi}_1$ and $\hat{\pi}_1$ can provide a useful set of upper and lower bounds of the true relationship between a PGS and an outcome.

There is an intuitive interpretation of the trends presented in Figure 1. In the between-family estimates, the inclusion of genetic nurture effects in $\widehat{PGS}_{ij}^{D}$ leads to an upward bias as it is capturing both differences in genetic composition between individuals *and* differences in the family environments between individuals that result from differences in their parents' genetic composition. The larger that $\rho_g$ is, the greater extent to which an individual with a beneficial allele for educational attainment reaps the reward from the same genes twice; first when their parents provide a more nurturing environment, and second when they themselves inherit the beneficial allele.

On the other hand, in family fixed effects models (within-family), the inclusion of genetic nurture in the observed PGS leads to downward bias when $\rho_g$ is less than unity. This is because there are no differences in genetic nurture systematically driving
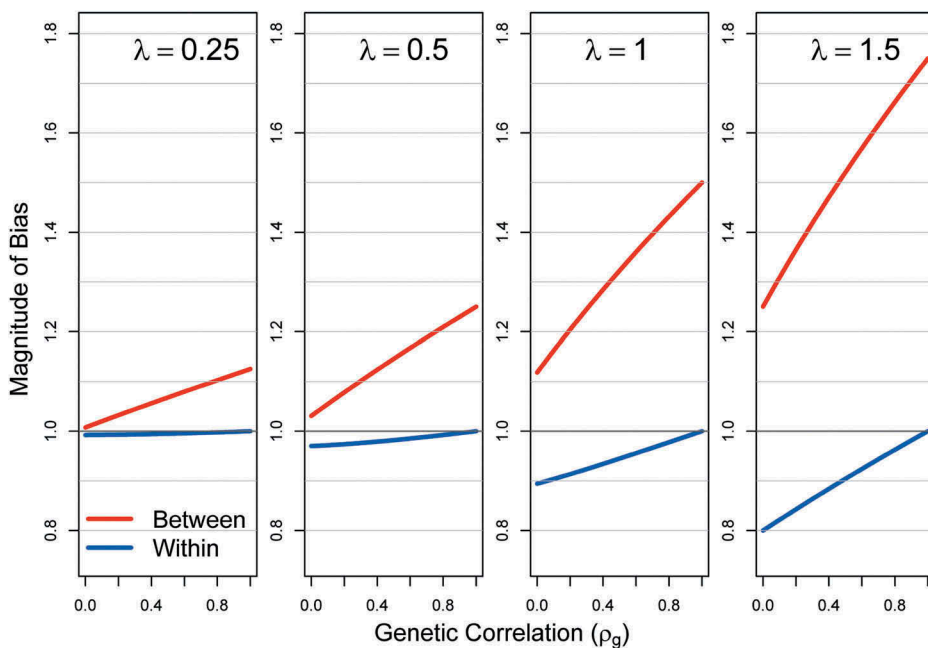
**Figure 1.** Bias due to genetic nurture in within- and between-family regressions using polygenic scores. Gray line at y = 1 represents no bias. $\rho_g$ is the direct-nurture genetic correlation and $\lambda$ is the direct-nurture heritability ratio. Results are derived analytically from a theoretical model.

differences in educational outcomes between siblings. Regardless of their own $PGS_{ij}^{D}$, siblings have identical $PGS_{j}^{N}$ because they are in the same family $j$. Thus, any extent that genetic nurture effects causes $\widehat{a}$ to diverge from $a$ amounts to measurement error in the allelic weights and causes downward attenuation bias.

The insights from our theoretical model offer a partial explanation of the large differences in $\frac{\hat{\pi}_1}{\hat{\psi}_1}$ observed across the various phenotypes considered in Table 1. For traits like years of schooling, cognitive ability, and (more speculatively) depression, where genetic nurture effects likely play an important role (i.e. large $\lambda$), we see large differences between $\hat{\psi}_1$ and $\hat{\pi}_1$ (i.e. $\frac{\hat{\pi}_1}{\hat{\psi}_1} < 1$). On the other side of the coin, for traits like body mass index and height, where most of the genetic contribution is likely to be direct (i.e. small $\lambda$), we see almost no difference between $\hat{\psi}_1$ and $\hat{\pi}_1$ (i.e. $\frac{\hat{\pi}_1}{\hat{\psi}_1} \approx 1$). These results coincide with the predictions of our theoretical model. Further, our model suggests that the differences in $\frac{\hat{\pi}_1}{\hat{\psi}_1}$ observed between years of schooling, cognitive ability, and depression may be a function of differences in $\rho_g$ between the traits.

Recall that these theoretical results are based on several simplifying assumptions: no gene–environment interaction, no gene–environment correlation, and no assortative mating. These three assumptions are unlikely to hold for most complex traits studied of interest to the social and biological sciences. Nonetheless, we can use the results from this simplified model to begin to probe how such violations might influence our results. Consider the case of positive genetic assortative mating, which exists for many of the

traits considered in this paper (Yengo et al. 2018b). A correlation between maternal and paternal genetics induces a positive correlation between transmitted and non-transmitted parental alleles shown in (4a.iv). This is because non-transmitted maternal alleles would be correlated with transmitted paternal alleles and vice versa. Thus, $cov\left(g_{ij}^z, f_{ij}^z\right)$ is no longer equal to zero and does not fall out of our equation. In such a case, the expected value of a GWAS allelic weight, shown in (4a.v), would include information about both non-transmitted genetic nurture effects in addition to the direct genetic effects and the transmitted genetic nurture effects. Thus, the bias documented in PGS analyses will increase as positive genetic assortative mating increases.

Next, we consider the case of gene-environment correlation. Notice that the bias that genetic nurture may cause in GWAS and PGS results from a special case of gene-environment correlation (i.e. a gene–environment correlation that mechanically exists due to the correlation of genetics between parents and their children induced through genetic inheritance). Thus, in our theoretical model's specification, gene-environment correlation only exists when environmental features relevant to an outcome are correlated with an individual's genetics *after* accounting for their parents' genetics. When such a gene-environment correlation exists, GWAS results and PGSs become biased in exactly the same way as they do from genetic nurture effects alone. For example, the existence of positive gene-environment correlation effectively increases $\lambda$ (as additional outcome variance explained by a non-genetic nurture environmental component), therein increasing the magnitude of the bias.

Finally, the case of gene–environment interaction is difficult to consider more generally, as it would vary as a function of the magnitude, direction, and the pathways of the interaction. Thus, the effects of gene–environment interactions on how genetic nurture effects influence GWAS and PGSs remain uncertain.

### *Direct-Nurture Genetic Correlation and Heritability Ratio*

Thousands of GWAS have been conducted in the last decade (Mills and Rahal 2019; Visscher et al. 2017). Nonetheless, to our knowledge, few or no GWAS has been conducted in human populations that independently identifies the direct genetic effects and the genetic nurture effects for a complex trait (a notable exception is GWAS on maternal influences on child birthweight (Beaumont et al. 2018; Warrington et al. 2019), though this indirect genetic affect may not be social in nature). Thus, for virtually all complex traits, little is known about the parameters of interest identified in our models: the direct-nurture genetic correlation and the direct-nurture heritability ratio ($\rho_g$ and $\lambda$). Critically, the two parameters are readily estimable with existing data and methods.

A better understanding of $\rho_g$ and $\lambda$ would offer value beyond aiding in the comparison of results from within-family and between-family regressions. The genetic pathways discovered in GWASs and summarized in PGSs offer researchers a puzzle to unpack (Freese 2018). Understanding why some phenotypes have strong versus weak genetic nurture effects, or why a phenotype's direct genetic effects and genetic nurture effects are more versus less linked, could help researchers glean insight into the underlying mechanisms at play.

Moreover, it would be interesting to understand how $\rho_g$ and $\lambda$ are influenced by the social environment. For example, while gene–environment interaction studies have been

the conventional way to understand how the environment moderates the influence of genetics, recent work has proposed a genetic correlation–environment interaction study (Wedow et al. 2018). In a genetic correlation–environment interaction study framework, the social environment can transform the genetic link between two traits. Exploring how the environment shapes $\rho_g$ would be a special case of a genetic correlation–environment interaction study where the two traits influence the same phenotype (directly and socially). Social policymakers might prefer a low $\rho_g$ for valued life outcomes like educational attainment to reduce the accumulation of inequality across generations.

While which specific social, physical, or economic factors moderate $\rho_g$ and $\lambda$ for various traits remains to be explored empirically, there may exist *a priori* reasons to suspect certain environmental modifiers. For example, imagine that individual variation in height is a function of both direct genetic effects that shape physiological development and genetic nurture effects that influence access to socioeconomic and nutritional resources during childhood (thereby reducing the likelihood of stunting). If there is a large casual effect of height on socioeconomic status, we would expect individuals with a greater genetic predisposition for height (i.e. a high $\widehat{PGS}_{ij}^{D}$) to be more likely to attain a high socioeconomic position where their children to have access to nutritional and health resources (i.e. a high $\widehat{PGS}_{ij}^{N}$), resulting in a positive $\rho_g$. However, this relationship could be modified by environmental features; if, for instance, the causal effect of height on social status is due to labor force discrimination, outlawing the use of height for employment decisions would uncouple $\widehat{PGS}_{ij}^{D}$ and $\widehat{PGS}_{ij}^{N}$ and reduce $\rho_g$. Alternatively, a social policy that provides adequate healthcare and nutrition to all children could effectively eliminate stunting and undo the relationship between parental genetics for height altogether, forcing both $\rho_g$ and $\lambda$ to zero.

## Implications for the Use of Polygenic Scores

While within-family analyses have demonstrated that many PGSs do have a significant causal signal for direct genetic effects, comparing the results from within-family analyses to results from between-family analyses is complicated by the presence of genetic nurture effects. To what extent do existing PGSs capture direct genetic effects, genetic nurture effects, and socioeconomic or geographic confounding? Until we better understand $\rho_g$ and $\lambda$ for a wide variety of traits, our ability to use within-family analyses to validate between-family discoveries will be limited. Analyses using PGSs should be interpreted accordingly.

Even in the absence of confounding due to population stratification, the observation that existing PGSs likely have genetic nurture components complicates their use and interpretation. Say, for instance, a researcher wonders whether there exists moderation of the association between an individual's PGS and their educational attainment as a function of school-level socioeconomic status (Trejo et al. 2018). Because a component of the PGS is capturing the benefit of having a parent with higher educational attainment and (in turn a higher socioeconomic status), any detected gene–environment interaction lacks a clear interpretation. It might be the case that household socioeconomic status interacts with school-level socioeconomic status in shaping educational attainment, or alternatively, it could be that an individual's genetic

composition interacts with school-level socioeconomic status in shaping educational attainment. These two different results have very different theoretical and practical implications but are indistinguishable in analyses using existing PGSs, which contain both direct genetic effects and genetic nurture effects.

## *Future Research*

The models constructed in this paper highlight key areas for future research in the field of social science genomics. Across a range of complex phenotypes, there is much work to be done toward separating out the genetics nurture effects from direct genetic effects. Utilizing random cage-mate assignment in mice, a recent study in mice conducted a social genetic effects GWAS and direct genetic effects GWAS in parallel and identified statistically significant genome-wide social genetic effect loci for 16 phenotypes (Baud et al. 2018). For these 16 phenotypes, the mean social-direct genetic correlation was 0.53 and the mean social-direct heritability ratio was 1.29. Crucially, social genetic effects arise from partially different loci as direct genetic effects and can have effects of differing magnitudes or directions at the same loci.

Unfortunately, social relationships are often not randomly assigned in human populations. Nonetheless, existing GWAS methods could be modified to use dyads of parents and their children. For example, a genetic nurture effects GWAS could be conducted by controlling for child's genetics in a GWAS of parental genetics on child phenotype. Alternatively, a GWAS conducted using variation only amongst sibling pairs would provide information on direct genetic pathways untainted by confounding or genetic nurture effects. Results from such a sibling GWAS might then be used to back out information about the genetic nurture effects from existing GWAS results of unrelated individuals. Nonetheless, obtaining large samples of parent-child or sibling pairs might prove challenging for many complex phenotypes. If statistical power is a problem, methods such as LD score regression (Bulik-Sullivan et al. 2015) could be used with smaller samples to identify estimates of $\rho_g$ and $\lambda$. Indeed, having estimates of $\rho_g$ and $\lambda$ for a trait would allow researchers to correct for bias in between-family and within-family analyses that use PGSs by dividing observed regression coefficients by the quantities displayed in Equation 4b.v and Equation 4c.vi, respectively.

It is also possible that the direct effects and genetic nurture effects are not independent. Imagine, for example, that parents with a higher educational attainment PGSs tend to invest more heavily in their children with higher PGSs than do parents with lower PGSs (differential investment in low birth weight children has been observed by across socio-economic lines (Hsin 2012; Restrepo 2016)). In such a case, the effects of parental genetics on a child would vary as a function as their child's genetics. Research designs that investigate the existence of interaction between direct genetic effects and genetic nurture effects may prove a fruitful avenue for future inquiry.

Finally, work should be done to extend our framework to other social genetic effects within families, such as between sibling pairs. We chose to start with genetic nurture effects because, unlike sibling effects, genetic nurture effects are likely to be unidirectional, with the causal effect pointing from parent to child, making them more straightforward to model (effects between siblings are likely reciprocal) (Kong et al. 2018) and because parental effects generalize to all families, not just those with multiple children. Nonetheless, social genetic effects between siblings may also complicate the interpretation of GWAS results and warrant attention.

## Conclusion

In summary, we formalize a theoretical model for additive direct genetic effects and genetic nurture effects and show that, unlike bias from other confounders, the presence of genetic nurture can bias coefficients from between-family and within-family regressions using PGSs. While within-family analyses that compare siblings using family fixed effects are considered the gold-standard, they are not without their own complications. Even if we were able to run a GWAS on an infinitely large sample using the methods in practice today, the presence of confounding genetic nurture effects would mean that it would be impossible to obtain precise estimates of the causal effect of an individual's genes on their life outcomes. Until GWAS can be conducted controlling for parental genotype, models using PGSs may be biased by genetic nurture effects. Obtaining estimates of our two novel social genetic parameters, the direct-nurture genetic correlation and the direct-nurture heritability ratio, may allow researchers to correct for such a bias.

## ORCID

Sam Trejo http://orcid.org/0000-0002-9880-5354
Benjamin W. Domingue http://orcid.org/0000-0002-3894-9049

# References

Armstrong-Carter, E., S. Trejo, L. Hill, K. Crossley, D. Mason, and B. W. Domingue. 2019. The earliest origins of genetic nurture: Prenatal environment mediates the association between maternal genetics and child development. *Psyarxiv*. doi:10.22201/fq.18708404e.2004.3.66178.

Barcellos, S. H., L. S. Carvalho, and P. Turley. 2018. Education can reduce health differences related to genetic risk of obesity. *Proceedings of the National Academy of Sciences*. doi:10.1007/BF00871674.

Bates, T. C., B. S. Maher, S. E. Medland, K. McAloney, M. J. Wright, N. K. Hansell, K. S. Kendler, N. G. Martin, and N. A. Gillespie. 2018. The nature of nurture: Using a virtual-parent design to test parenting effects on children's educational attainment in genotyped families. *Twin Research and Human Genetics*. doi:10.1017/thg.2018.11.

Baud, A., F. P. Casale, J. Nicod, and O. Stegle. 2018. Genome-wide association study of social genetic effects on 170 phenotypes in laboratory mice. *BioRxiv*. doi:10.1101/302349.

Beaumont, R. N., N. M. Warrington, A. Cavadino, J. Tyrrell, M. Nodzenski, M. Horikoshi, F. Geller, R. Myhre, R. C. Richmond, L. Paternoster, et al. 2018. Genome-wide association study of offspring birth weight in 86 577 women identifies five novel loci and highlights maternal genetic effects that are independent of fetal genetics. *Human Molecular Genetics*. doi:10.1093/hmg/ddx429.

Belsky, D. W., B. W. Domingue, R. Wedow, L. Arseneault, J. D. Boardman, A. Caspi, D. Conley, J. M. Fletcher, J. Freese, P. Herd, et al. 2018. Genetic analysis of social-class mobility in five longitudinal studies. *Proceedings of the National Academy of Sciences*. doi:10.1073/pnas.1801238115.

Belsky, D. W., and S. Israel. 2014. Integrating genetics and social science: Genetic risk scores. *Biodemography and Social Biology*. doi:10.1080/19485565.2014.946591.

Belsky, D. W., T. E. Moffitt, D. L. Corcoran, B. Domingue, H. Harrington, S. Hogan, R. Houts, S. Ramrakha, K. Sugden, B. S. Williams, et al. 2016. The genetics of success: How single-nucleotide polymorphisms associated with educational attainment relate to life-course development. *Psychological Science* 27:957–72. doi:10.1177/0956797616643070.

Belsky, D. W., T. E. Moffitt, T. B. Baker, A. K. Biddle, J. P. Evans, H. L. Harrington, R. Houts, M. Meier, K. Sugden, B. Williams, et al. 2013. Polygenic risk and the developmental progression to heavy, persistent smoking and nicotine dependence: Evidence from a 4-decade longitudinal study. *JAMA Psychiatry (Chicago, ill.)*. doi:10.1001/jamapsychiatry.2013.736.

Bergsma, R., E. Kanis, E. F. Knol, and P. Bijma. 2008. The contribution of social effects to heritable variation in finishing traits of domestic pigs (Sus Scrofa). *Genetics*. doi:10.1534/genetics.107.084236.

Bulik-Sullivan, B. K., P.-R. Loh, H. K. Finucane, S. Ripke, J. Yang, N. Patterson, M. J. Daly, A. L. Price, and B. M. Neale,, . 2015. LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nature Genetics* 47 (3):291–95. Nature Publishing Group. doi:10.1038/ng.3211.

Canario, L., N. Lundeheim, and P. Bijma. 2017. The early-life environment of a pig shapes the phenotypes of its social partners in adulthood. *Heredity*. doi:10.1038/hdy.2017.3.

Cawley, J., E. Han, J. (June) Kim, and E. C. Norton. 2017. Testing for peer effects using genetic data. *NBER Working Paper No. 23719*. doi:10.3386/w23719.

Chabris, C. F., J. J. Lee, D. Cesarini, D. J. Benjamin, and D. I. Laibson. 2015. The fourth law of behavior genetics. *Current Directions in Psychological Science* 24 (4):304–12. doi:10.1177/0963721415580430.

Conley, D., and J. Fletcher. 2017. *The genome factor what the social genomics revolution reveals about ourselves, our history, and the future*. Princeton, NJ: Princeton University Press.

Domingue, B. W., and D. W. Belsky. 2017. The social genome: current findings and implications for the study of human genetics. *PLoS Genetics* 13:3. doi:10.1371/journal.pgen.1006615.

Domingue, B. W., D. W. Belsky, D. Conley, K. M. Harris, and J. D. Boardman. 2015. Polygenic influence on educational attainment: New evidence from the national longitudinal study of adolescent to adult health. *AERA Open* 1 (3):1–13. doi:10.1177/2332858415599972.

Domingue, B. W., D. W. Belsky, J. M. Fletcher, D. Conley, J. D. Boardman, and K. M. Harris. 2018. The social genome of friends and schoolmates in the national longitudinal study of adolescent to adult health. *Proceedings of the National Academy of Sciences* 201711803. doi:10.1073/pnas.1711803115.

Dudbridge, F. 2013. Power and predictive accuracy of polygenic risk scores. *PLoS Genetics* 9:3. doi:10.1371/journal.pgen.1003348.

Duncan, L., H. Shen, B. Gelaye, K. Ressler, M. Feldman, R. Peterson, and B. Domingue. 2019. Analysis of polygenic score usage and performance in diverse human populations. *Nature Communications*. doi:10.1101/398396.

Freese, J. 2018. The arrival of social science genomics. *Contemporary Sociology*. doi:10.1177/0094306118792214a.

Giri, A., J. N. Hellwege, J. M. Keaton, J. Park, C. Qiu, H. R. Warren, E. S. Torstenson, C. P. Kovesdy, Y. V. Sun, O. D. Wilson, et al. 2019. Trans-ethnic association study of blood pressure determinants in over 750,000 individuals. *Nature Genetics* 51 (1):51. Nature Publishing Group. doi:10.1038/s41588-018-0303-9.

Hamer, D. H., and L. Sirota. 2000. Beware the chopsticks gene. *Molecular Psychiatry*. doi:10.1038/sj.mp.4000662.

Harris, K. M., C. T. Halpern, E. A. Whitsel, J. M. Hussey, L. A. Killeya-jones, J. Tabor, and S. C. Dean. 2019. Cohort profile: The national longitudinal study of adolescent to adult health (add health). *International Journal of Epidemiology* 1–12. doi:10.1093/ije/dyz115.

Hsin, A. 2012. Is biology destiny? Birth weight and differential parental treatment. *Demography* 1385–405. doi:10.1007/s13524-012-0123-y.

Hyde, C. L., M. W. Nagle, C. Tian, X. Chen, S. A. Paciga, J. R. Wendland, J. Y. Tung, D. A. Hinds, R. H. Perlis, and A. R. Winslow. 2016. Identification of 15 genetic loci associated with risk of major depression in individuals of european descent. *Nature Genetics*. doi:10.1038/ng.3623.

Kong, A., G. Thorleifsson, M. L. Frigge, B. J. Vilhjalmsson, A. I. Young, T. E. Thorgeirsson, S. Benonisdottir, A. Oddsson, B. V. Halldorsson, G. Masson, et al. 2018. The nature of nurture: effects of parental genotypes. *Science* 359 (6374):424–28. doi:10.1126/science.aan6877.

Lee, J. J., R. Wedow, A. Okbay, E. Kong, O. Maghzian, M. Zacher, T. A. Nguyen-Viet, P. Bowers, J. Sidorenko, R. Karlsson Linnér, et al. 2018. Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals. *Nature Genetics*. doi:10.1038/s41588-018-0147-3.

Locke, A. E., S. I. Bratati Kahali, A. E. Berndt, T. H. Justice, F. R. Pers, C. P. Day, C. Powell, S. Vedantam, M. L. Buchkovich, J. Yang, et al. 2015. Genetic studies of body mass index yield new insights for obesity biology. *Nature*. doi:10.1038/nature14177.

Martin, A. R., C. R. Gignoux, R. K. Walters, G. L. Wojcik, B. M. Neale, S. Gravel, M. J. Daly, C. D. Bustamante, and E. E. Kenny. 2017. Human demographic history impacts genetic risk prediction across diverse populations. *The American Journal of Human Genetics* 100 (4):635–649.

Martschenko, D., S. Trejo, and B. W. Domingue. 2019. Genetics and education: Recent developments in the context of an ugly history and an uncertain future. *AERA Open* 5 (1):1–15. doi:10.1093/ije/dyx041.

Mills, M. C., and C. Rahal. 2019. A scientometric review of genome-wide association studies. *Communications Biology* 2 (1):9. Springer US. doi:10.1038/s42003-018-0261-x.

Moore, A. J., E. D. Brodie, and J. B. Wolf. 1997. Interacting phenotypes and the evolutionary process: I. Direct and indirect genetic effects of social interactions. *Evolution*. doi:10.2307/2411187.

Novembre, J., T. Johnson, K. Bryc, A. R. Zoltán Kutalik, A. A. Boyko, A. Indap, A. Indap, K. S. King, S. Bergmann, M. R. Nelson, et al. 2008. Genes mirror geography within Europe. *Nature*. doi:10.1038/nature07331.

Okbay, A., M. L. Bart, J. Baselmans, E. De Neve, M. G. Patrick Turley, M. Nivard, S. Alan Fontana, F. W. Meddens, C. A. Rietveld, J. Derringer, et al. 2016. Genetic variants associated with subjective well-being, depressive symptoms, and neuroticism identified through genome-wide analyses. *Nature Genetics*. doi:10.1038/ng.3552.

Papageorge, N., and K. Thom. 2017. Genes, education, and labor market outcomes: Evidence from the health and retirement study. *Working Paper* 055654. doi:10.2139/ssrn.2982606.

Pearson, T. A., and T. A. Manolio. 2008. How to interpret a genome-wide association study. *JAMA : the Journal of the American Medical Association* 299 (11):1335–44. doi:10.1001/jama.299.11.1335.

Petfield, D., S. F. Chenoweth, H. D. Rundle, and M. W. Blows. 2005. Genetic variance in female condition predicts indirect genetic variance in male sexual display traits. *Proceedings of the National Academy of Sciences of the United States of America*. doi:10.1073/pnas.0409378102.

Polderman, T. J. C., B. Benyamin, C. A. de Leeuw, P. F. Sullivan, A. van Bochoven, P. M. Visscher, and D. Posthuma. 2015. Meta-analysis of the heritability of human traits based on fifty years of twin studies. *Nature Genetics* 47 (7):702–09. doi:10.1038/ng.3285.

Restrepo, B. J. 2016. Parental investment responses to a low birth weight outcome: Who compensates and who reinforces ? *Journal of Population Economics* 969–89. doi:10.1007/s00148-016-0590-3.

Rietveld, C. A., D. Conley, N. Eriksson, T. Esko, S. E. Medland, A. A. E. Vinkhuyzen, J. Yang, J. D. Boardman, C. F. Chabris, C. T. Dawes, et al. 2014. Replicability and robustness of genome-wide-association studies for behavioral traits. *Psychological Science* 25 (11):1975–86. doi:10.1177/0956797614545132.

Rietveld, C. A., S. E. Medland, J. Derringer, J. Yang, N. W. Tõnu Esko, H. Martin, J. Westra, K. Shakhbazov, A. Abdellaoui, A. Agrawal, et al. 2013. GWAS of 126,559 individuals identifies genetic variants associated with educational attainment. *Science*. doi:10.1126/science.1235488.

Selzam, S., S. J. Ritchie, J.-B. Pingault, C. A. Reynolds, and F. Paul. 2019. Comparing within- and between-family polygenic score prediction authors. *The American Journal of Human Genetics* 105:351–63.

Sotoudeh, R., K. Mullan, and D. Conley. 2019. Effects of the peer metagenomic environment on smoking behavior. *Proceedings of the National Academy of Sciences* 116 (33). doi:10.1073/pnas.1806901116.

Torkamani, A., N. Wineinger, and E. Topol. 2018. The personal and clinical utility of polygenic risk scores. *Nature Reviews Genetics* 19:581–90. September Springer US. doi:10.1038/s41576-018-0018-x.

Trejo, S., D. Belsky, J. Boardman, J. Freese, K. Harris, P. Herd, K. Sicinski, and B. Domingue. 2018. Schools as moderators of genetic associations with life course attainments: Evidence from the WLS and add heath. *Sociological Science*. doi:10.15195/v5.a22.

Turley, P., R. K. Walters, O. Maghzian, J. J. Aysu Okbay, M. Lee, A. Fontana, T. A. Nguyen-Viet, R. Wedow, M. Zacher, N. A. Furlotte, et al. 2018. Multi-trait analysis of genome-wide association summary statistics using MTAG. *Nature Genetics*. doi:10.1038/s41588-017-0009-4.

Visscher, P. M., N. R. Wray, Q. Zhang, P. Sklar, M. I. McCarthy, M. A. Brown, and J. Yang. 2017. 10 years of GWAS discovery: Biology, function, and translation. *American Journal of Human Genetics*. doi:10.1016/j.ajhg.2017.06.005.

Warrington, N. M., R. N. Beaumont, M. Horikoshi, F. R. Day, and O. Helgeland. 2019. Maternal and fetal genetic effects on birth weight and their relevance to cardio-metabolic risk factors. *Nature Genetics*. no. Ld. doi:10.1038/s41588-019-0403-1.

Wedow, R., M. Zacher, B. M. Huibregtse, K. M. Harris, B. W. Domingue, and J. D. Boardman. 2018. Education, smoking, and cohort change: forwarding a multidimensional theory of the environmental moderation of genetic effects. *American Sociological Review*. doi:10.1177/0003122418785368.

Wertz, J., T. E. Moffitt, J. Agnew-Blais,L. Arseneault, D. W. Belsky, D. L. Corcoran, R. Houts, T. Matthews, J. A. Prinz, L. S. Richmond-Rakerd, et al. 2018. Using DNA from mothers and children to study parental investment in children's educational attainment. *BioRxiv*.

Wolf, J. B., E. D. Brodie, J. M. Cheverud, A. J. Moore, and M. J. Wade. 1998. evolutionary consequences of indirect genetic effects. *Trends in Ecology and Evolution*. doi:10.1016/S0169-5347(97)01233-0.

Wood, A. R., T. Esko, J. Yang, S. Vedantam, T. H. Pers, S. Gustafsson, A. Y. Chu, K. Estrada, J. Luan, Z. Kutalik, et al. 2014. Defining the role of common variation in the genomic and biological architecture of adult human height. *Nature Genetics*. doi:10.1038/ng.3097.

Wooldridge, J. M. 2015. *Introductory econometrics: A modern approach*. Mason, OH: South-Western.

Yengo, L., J. Julia Sidorenko, K. E. Kemper, Z. Zheng, A. R. Wood, M. N. Weedon, T. M. Frayling, J. Hirschhorn, A. J. Yang, and P. M. Visscher. 2018a. Meta-analysis of genome-wide association studies for height and body mass index in □700000 individuals of European ancestry. *Human Molecular Genetics* 27:3641–49. doi:10.1093/hmg/ddy271.

Yengo, L., M. R. Robinson, M. C. Keller, K. E. Kemper, Y. Yang, M. Trzaskowski, J. Gratten, P. Turley, D. Cesarini, D. J. Benjamin, et al. 2018b. Imprint of assortative mating on the human genome. *Nature Human Behavior* 2:1–41.

# Appendix

## A1. Add Health

The National Longitudinal Study of Adolescent to Adult Health (Add Health) is a nationally representative cohort drawn from a probability sample of 80 high schools and 52 middle schools in roughly 80 US communities, and representative of schools in the United States in 1994–95 with respect to region, urban setting, school size, school type, and race or ethnic background. About 15,000 Add Health respondents (or 96%) consented to genotyping during the Wave 4 interview in 2008–09 for purposes of approved Add Health Wave 4 research. Of those who consented to genotyping, ~12,000 (or 80%) agreed to have their DNA archived for future testing. DNA extraction and genotyping were conducted on this archive sample using two platforms (Illumina Omni1 for siblings and the Illumina Omni2.5 for unrelated individuals). Quality control procedures were performed on the genetic samples collected yielding genetic data from ~10,000 individuals on 609,130 overlapping SNPs.

We focus our analyses on a set of genetically homogenous respondents of European descent. We restrict our sample to only European ancestry individuals because differences in linkage disequilibrium and allele frequencies across ancestral groups complicate the interpretation of PGS–phenotype associations (Martin et al. 2017). Although we recognize the importance of research in more diverse samples, Add Health does not contain large enough samples of any ancestral group other than European ancestry to conduct well-powered within-family analyses. Fortunately, the theoretical results regarding how genetic nurturance can induce bias in GWASs and the PGSs constructed using their summary statistics, which we emphasize is the core contribution of our paper, applies to genetic analyses conducted in all ancestry groups (although the relevant underlying $\lambda$ and $\rho_g$ parameters may vary between such groups due to environmental differences). Add Health also contains a variety of data on students' academic performance, personal characteristics measured in adolescence (cognitive ability, personality characteristics, professional aspirations, physical health, and functioning, etc.).

## A2. European Ancestry Identification and Polygenic Score Construction

To identify a sample of European-ancestry respondents, we calculated the first two principal components of participants with known ancestries from the 1000 Genomes Project. We then projected Add Health individuals onto those principal components, obtaining the loadings of each Add Health respondent on each PC. We use these loadings to assign each individual to 1 of 5 super-populations in the 1000 Genomes data: European, African, East Asian, South Asian and Admixed and restrict our sample to only individuals of European ancestry. In this European ancestry sub-sample, we calculated new PCs to use as controls in our polygenic score regression analyses.

Polygenic scores (PGSs) were created using SNPs in the Add Health genetic database that were matched to SNPs with reported results in a GWAS. We also removed all SNPs where the risk allele identified via GWAS could not be readily identified in the Add Health genetic database. For each SNP, a loading was calculated as the number of trait-associated alleles multiplied by the effect size

estimated in the original GWAS. SNPs with relatively large $p$-values will have small effects (and thus be down-weighted in creating the composite), so we do not impose a $p$-value threshold. Loadings were summed across the SNP set to calculate the polygenic score. The scores are standardized within sample to have a mean of 0 and standard deviation of 1. PGS generated from GWAS that included Add Health was constructed from summary statistics with Add Health removed. Links to the summary statistics used to construct each score are provided below:

Body mass index: https://portals.broadinstitute.org/collaboration/giant/images/c/c8/Meta-analysis_Locke_et_al%2BUKBiobank_2018_UPDATED.txt.gz

Height: https://portals.broadinstitute.org/collaboration/giant/images/6/63/Meta-analysis_Wood_et_al%2BUKBiobank_2018.txt.gz

Own birth weight: http://mccarthy.well.ox.ac.uk/publications/2019/EggBirthWeight_NatureGenetics/Fetal_BW_European_meta.NG2019.txt.gz

Educational attainment*: http://ssgac.org/documents/MTAG_EA.to10K.txt

Cognitive performance*: http://ssgac.org/documents/MTAG_CP.to10K.txt

Depressive symptoms*: http://ssgac.org/documents/MTAG_DEP_CLUMPED.to10K.txt

*Note that the Social Science Genomics Association Consortium phenotypes (i.e. educational attainment, cognitive performance, and depressive symptoms) used 23andMe data in their published GWAS. While our polygenic scores were generated using the full set of summary statistics, the publicly available data (linked to above) contain only the 10K SNPs with the lowest $p$-values in compliance with 23andMe's data-sharing policies.

## A3. Simplifying Assumptions of Structural Model

We make several simplifications to simplify the exposition of the structural model. In defining our structural model, we have already assumed away gene–environment interactions (notice the lack of an interaction term between genes and environment). We further assume no gene-environment correlation (net of the correlation mechanically induced by the combination of genetic nurturance and genetic inheritance).

$$cov\left(PGS_{ij}^{D}, E_{ij}\right) = cov\left(PGS_{j}^{N}, E_{j}\right) = 0$$

Finally, we assume no assortative mating. This means that the polygenic scores of parents are uncorrelated.

$$PGS_{j}^{N} = \sum_{z=1}^{n} \delta^z g_{j}^{z} = \sum_{z=1}^{n} \delta^z \left(g_{j}^{z_m} + g_{j}^{z_p}\right) = \sum_{z=1}^{n} \delta^z g_{j}^{z_m} + \sum_{z=1}^{n} \delta^z g_{j}^{z_p} = PGS_{j}^{N_m} + PGS_{j}^{N_p}$$

$$var\left(PGS_{j}^{N}\right) = var\left(PGS_{j}^{N_m} + PGS_{j}^{N_p}\right) = var\left(PGS_{j}^{N_m}\right) + var\left(PGS_{j}^{N_p}\right) = 2var\left(PGS_{ij}^{N}\right)$$

$g_{j}^{z_m}$ : Total number of major alleles at the mother in family $j$'s genetic locus $z$ $(0, 1, \text{or } 2)$
$g_{j}^{z_p}$ : Total number of major alleles at the father in family $j$'s genetic locus $z$ $(0, 1, \text{or } 2)$
$PGS_{j}^{D_m}$ : PGS constructed from the true causal linear effect of the mother in family $j$'s genes on $Y_{ij}$

$PGS_{j}^{D_p}$ : PGS constructed from the true causal linear effect of the father in family $j$'s genes on $Y_{ij}$

This also entails that a child's PGS is related to their parent's and sibling's PGS only directly through the genes they receive through the process of recombination. Thus, the correlation between the same PGSs for all parent-child and sibling-sibling pairs is .5.

$$\rho_{PGS_{0j}^{D},PGS_{1j}^{D}} = \rho_{PGS_{j}^{D_m},PGS_{0j}^{D}} = \rho_{PGS_{j}^{D_m},PGS_{1j}^{D}} = \rho_{PGS_{j}^{N_p},PGS_{0j}^{D}} = \rho_{PGS_{j}^{N_p},PGS_{1j}^{D}} = .5$$

$$\rho_{PGS_{0j}^{N},PGS_{1j}^{N}} = \rho_{PGS_{j}^{N_m},PGS_{0j}^{N}} = \rho_{PGS_{j}^{N_m},PGS_{1j}^{N}} = \rho_{PGS_{j}^{N_p},PGS_{0j}^{N}} = \rho_{PGS_{j}^{N_p},PGS_{1j}^{N}} = .5$$

## A4. Correlation between $PGS_{ij}^{N}$ and $PGS_{j}^{N}$

$$\rho_{PGS_{ij}^{N},PGS_{j}^{N}} = \frac{cov\left(PGS_{ij}^{N}, PGS_{j}^{N}\right)}{var\left(PGS_{ij}^{N}\right)^{\frac{1}{2}} var\left(PGS_{j}^{N}\right)^{\frac{1}{2}}}$$

$$\rho_{PGS_{ij}^{N},PGS_{j}^{N}} = \frac{cov\left(PGS_{ij}^{N}, PGS_{j}^{N_m} + PGS_{j}^{N_p}\right)}{var\left(PGS_{ij}^{N}\right)^{\frac{1}{2}} var\left(PGS_{j}^{N_m} + PGS_{j}^{N_p}\right)^{\frac{1}{2}}}$$

$$\rho_{PGS_{ij}^{N},PGS_{j}^{N}} = \frac{cov\left(PGS_{ij}^{N}, PGS_{j}^{N_m}\right) + cov\left(PGS_{ij}^{N}, PGS_{j}^{N_p}\right)}{var\left(PGS_{ij}^{N}\right)^{\frac{1}{2}} \left[var\left(PGS_{j}^{N_m}\right) + var\left(PGS_{j}^{N_p}\right) + 2cov\left(PGS_{j}^{N_m}, PGS_{j}^{N_p}\right)\right]^{\frac{1}{2}}}$$

$$\rho_{PGS_{ij}^{D},PGS_{j}^{N}} = \frac{\rho_{PGS_{ij}^{D}, PGS_{j}^{N_m}} var\left(PGS_{ij}^{D}\right)^{\frac{1}{2}} \left(PGS_{j}^{N_m}\right)^{\frac{1}{2}} + \rho_{PGS_{ij}^{D}, PGS_{j}^{N_p}} var\left(PGS_{ij}^{N}\right)^{\frac{1}{2}} \left(PGS_{j}^{N_p}\right)^{\frac{1}{2}}}{var\left(PGS_{ij}^{D}\right)^{\frac{1}{2}} \left[var\left(PGS_{j}^{N_m}\right) + var\left(PGS_{j}^{N_p}\right) + 0\right]^{\frac{1}{2}}}$$

$$\rho_{PGS_{ij}^{N},PGS_{j}^{N}} = \frac{.5\lambda^2 + .5\lambda^2}{\lambda[\lambda^2 + \lambda^2]^{\frac{1}{2}}}$$

$$\rho_{PGS_{ij}^{N},PGS_{j}^{N}} = \frac{\sqrt{2}}{2}$$

## A5. Correlation between $PGS_{ij}^{D}$ and $PGS_{j}^{N}$

$$\rho_{PGS_{ij}^{D},PGS_{j}^{N}} = \frac{cov\left(PGS_{ij}^{D}, PGS_{j}^{N}\right)}{var\left(PGS_{ij}^{D}\right)^{\frac{1}{2}} var\left(PGS_{j}^{N}\right)^{\frac{1}{2}}}$$

$$\rho_{PGS_{ij}^{D},PGS_{j}^{N}} = \frac{cov\left(PGS_{ij}^{D}, PGS_{j}^{N_m} + PGS_{j}^{N_p}\right)}{var\left(PGS_{ij}^{D}\right)^{\frac{1}{2}} var\left(PGS_{j}^{N_m} + PGS_{j}^{N_p}\right)^{\frac{1}{2}}}$$

$$\rho_{PGS_{ij}^{D},PGS_{j}^{N}} = \frac{cov\left(PGS_{ij}^{D}, PGS_{j}^{N_m}\right) + cov\left(PGS_{ij}^{D}, PGS_{j}^{N_p}\right)}{var\left(PGS_{ij}^{D}\right)^{\frac{1}{2}} [var\left(PGS_{j}^{N_m}\right) + var\left(PGS_{j}^{N_p}\right) + 2cov\left(PGS_{j}^{N_m}, PGS_{j}^{N_p}\right)]^{\frac{1}{2}}}$$

$$\rho_{PGS_{ij}^{D},PGS_{j}^{N}} = \frac{\rho_{PGS_{ij}^{D}, PGS_{j}^{N_m}} var\left(PGS_{ij}^{D}\right)^{\frac{1}{2}} \left(PGS_{j}^{N_m}\right)^{\frac{1}{2}} + \rho_{PGS_{ij}^{D}, PGS_{j}^{N_p}} var\left(PGS_{ij}^{N}\right)^{\frac{1}{2}} \left(PGS_{j}^{N_p}\right)^{\frac{1}{2}}}{var\left(PGS_{ij}^{D}\right)^{\frac{1}{2}} \left[var\left(PGS_{j}^{N_m}\right) + var\left(PGS_{j}^{N_p}\right) + 0\right]^{\frac{1}{2}}}$$

$$\rho_{PGS_{ij}^{D},PGS_{j}^{N}} = \frac{.5\rho\lambda + .5\rho\lambda}{[\lambda^2 + \lambda^2]^{\frac{1}{2}}}$$

$$\rho_{PGS_{ij}^{D},PGS_{j}^{N}} = \frac{\sqrt{2}}{2}\rho$$

## A6. Underlying Allelic Weights

Specifically, we assume:

$$\rho_{\boldsymbol{\alpha},\overline{\boldsymbol{g}}_{ij}} = \rho_{\boldsymbol{\delta},\overline{\boldsymbol{g}}_{ij}}$$

$$\rho_{\boldsymbol{\alpha}^2,\, var\left(\boldsymbol{\alpha},\overline{\boldsymbol{g}}_{ij}\right)} = \rho_{\boldsymbol{\delta}^2,\, var\left(g_{ij}\right)}$$

$\overline{g}_{ij}^z$ : Population mean risk allele frequency count at genetic locus $z$

$var\left(g_{ij}^z\right)$ : Population variance of risk allele frequency at genetic locus $z$

$\overline{\boldsymbol{g}}_{ij}$: An $n \times 1$ vector of the population mean risk allele frequency for all genetic loci $z$

$\boldsymbol{var}(\boldsymbol{g}_{ij})$: An $n \times 1$ vector of the population variation of risk allele frequency for all genetic loci $z$

$\rho_{\alpha}$ : Correlation between vectors comprised of $\alpha^z$ and $\overline{g}_{ij}^z$ for all genetic loci $z$

$\rho_{\boldsymbol{\delta},\overline{\boldsymbol{g}}_{ij}}$ : Correlation between vectors comprised of $\delta^z$ and $\overline{g}_{ij}^z$ for all genetic loci $z$

$\rho_{\boldsymbol{\alpha}^2,\, var\left(\boldsymbol{\alpha},\overline{\boldsymbol{g}}_{ij}\right)}$ : Correlation between vectors comprised of $(\alpha^z)^2$ and $var\left(g_{ij}^z\right)$ for all genetic loci $z$

$\rho_{\boldsymbol{\delta}^2,\, var\left(g_{ij}\right)}$ : Correlation between vectors comprised of $(\delta^z)^2$ and $var\left(g_{ij}^z\right)$ for all genetic loci $z$

## A7. λ is the Direct-Nurture Heritability Ratio

$$\rho_{\boldsymbol{\alpha},\overline{\boldsymbol{g}}_{ij}} = \rho_{\boldsymbol{\delta},\overline{\boldsymbol{g}}_{ij}}$$

$$\frac{cov\left(\boldsymbol{\alpha},\overline{\boldsymbol{g}}_{ij}\right)}{var(\boldsymbol{\delta})^{\frac{1}{2}}var\left(\overline{\boldsymbol{g}}_{ij}\right)^{\frac{1}{2}}} = \frac{cov\left(\boldsymbol{\delta},\overline{\boldsymbol{g}}_{ij}\right)}{var(\boldsymbol{\delta})^{\frac{1}{2}}var\left(\overline{\boldsymbol{g}}_{ij}\right)^{\frac{1}{2}}}$$

$$\frac{cov\left(\boldsymbol{\alpha},\overline{\boldsymbol{g}}_{ij}\right)}{var(\boldsymbol{\alpha})^{\frac{1}{2}}} = \frac{cov\left(\boldsymbol{\delta},\overline{\boldsymbol{g}}_{ij}\right)}{var(\boldsymbol{\delta})^{\frac{1}{2}}}$$

$$\frac{var(\boldsymbol{\delta})^{\frac{1}{2}}}{var(\boldsymbol{\alpha})^{\frac{1}{2}}}cov\left(\boldsymbol{\alpha},\overline{\boldsymbol{g}}_{ij}\right) = cov\left(\boldsymbol{\delta},\overline{\boldsymbol{g}}_{ij}\right)$$

$$\frac{var(\boldsymbol{\delta})^{\frac{1}{2}}}{var(\boldsymbol{\alpha})^{\frac{1}{2}}}[\sum_{z=1}^{n}\left(\alpha^z\overline{g}_{ij}^z\right) - \overline{\alpha}\overline{g}_{ij}^n] = \sum_{z=1}^{n}(\delta^z\overline{g}_{ij}^z) - \overline{\delta}\overline{g}_{ij}^n$$

$$\frac{var(\boldsymbol{\delta})^{\frac{1}{2}}}{var(\boldsymbol{\alpha})^{\frac{1}{2}}}[\sum_{z=1}^{n}\left(\alpha^z\overline{g}_{ij}^z\right) - 0] = \sum_{z=1}^{n}(\delta^z\overline{g}_{ij}^z) - 0$$

$$\frac{var(\boldsymbol{\delta})^{\frac{1}{2}}}{var(\boldsymbol{\alpha})^{\frac{1}{2}}} = \frac{\sum_{z=1}^{n}(\delta^z\overline{g}_{ij}^z)}{\sum_{z=1}^{n}(\alpha^z\overline{g}_{ij}^z)}$$

$$\frac{\left(\frac{\lambda^2\sigma}{4}\right)^{\frac{1}{2}}}{(\sigma)^{\frac{1}{2}}} = \frac{\sum_{z=1}^{n}(\delta^z\overline{g}_{ij}^z)}{\sum_{z=1}^{n}(\alpha^z\overline{g}_{ij}^z)}$$

$$\lambda = \frac{\sum_{z=1}^{n}(\delta^z\overline{g}_{ij}^z)}{\sum_{z=1}^{n}(\alpha^z\overline{g}_{ij}^z)}$$

## A8. Variance of $PGS_{ij}^{N}$ is $\frac{\lambda^2}{4}$

$$\rho_{\alpha^2, var(g_{ij})} = \rho_{\delta^2, var(g_{ij})}$$

$$\frac{cov\left(\alpha^2, var\left(g_{ij}\right)\right)}{var\left(\delta^2\right)^{\frac{1}{2}} var\left(var\left(g_{ij}\right)\right)^{\frac{1}{2}}} = \frac{cov\left(\delta^2, var\left(g_{ij}\right)\right)}{var\left(\delta^2\right)^{\frac{1}{2}} var\left(var\left(g_{ij}\right)\right)^{\frac{1}{2}}}$$

$$\frac{cov\left(\alpha^2, var\left(g_{ij}\right)\right)}{var\left(\alpha^2\right)^{\frac{1}{2}}} = \frac{cov\left(\delta^2, var\left(g_{ij}\right)\right)}{var\left(\delta^2\right)^{\frac{1}{2}}}$$

$$\frac{var\left(\delta^2\right)^{\frac{1}{2}}}{var\left(\alpha^2\right)^{\frac{1}{2}}} cov\left(\alpha^2, var\left(g_{ij}\right)\right) = cov\left(\delta^2, var\left(g_{ij}\right)\right)$$

$$\frac{var\left(\delta^2\right)^{\frac{1}{2}}}{var\left(\alpha^2\right)^{\frac{1}{2}}} \sum_{z=1}^{n} (\alpha^z)^2 var\left(g_{ij}^z\right) - (\bar{\alpha}^z)^2 \overline{var}\left(g_{ij}^z\right) = \sum_{z=1}^{n} (\delta^z)^2 var\left(g_{ij}^z\right) - (\delta^z)^2 \overline{var}\left(g_{ij}^z\right)$$

$$\frac{var\left(\delta^2\right)^{\frac{1}{2}}}{var\left(\alpha^2\right)^{\frac{1}{2}}} \sum_{z=1}^{n} (\alpha^z)^2 var\left(g_{ij}^z\right) - 0 = \sum_{z=1}^{n} (\delta^z)^2 var\left(g_{ij}^z\right) - 0$$

$$\frac{var\left(\delta^2\right)^{\frac{1}{2}}}{var\left(\alpha^2\right)^{\frac{1}{2}}} \sum_{z=1}^{n} (\alpha^z)^2 var\left(g_{ij}^z\right) = \sum_{z=1}^{n} (\delta^z)^2 var\left(g_{ij}^z\right)$$

$$\frac{var\left(\delta^2\right)^{\frac{1}{2}}}{var\left(\alpha^2\right)^{\frac{1}{2}}} \sum_{z=1}^{n} var\left(\alpha^z g_{ij}^z\right) = \sum_{z=1}^{n} var\left(\delta^z g_{ij}^z\right)$$

$$\frac{var\left(\delta^2\right)^{\frac{1}{2}}}{var\left(\alpha^2\right)^{\frac{1}{2}}} var\left(PGS_{ij}^{D}\right) = var\left(PGS_{ij}^{N}\right)$$

$$\frac{var\left(\delta^2\right)^{\frac{1}{2}}}{var\left(\alpha^2\right)^{\frac{1}{2}}} = var\left(PGS_{ij}^{N}\right)$$

$$\frac{\left(\frac{\lambda^4 \sigma^4}{16}\right)^{\frac{1}{2}}}{(\sigma^4)^{\frac{1}{2}}} = var\left(PGS_{ij}^{N}\right)$$

$$var\left(PGS_{ij}^{N}\right) = \frac{\lambda^2}{4}$$

## A9. Variance of $PGS_{j}^{N}$ is $\frac{\lambda^2}{2}$

$$var\left(PGS_{j}^{N}\right) = var\left(PGS_{j}^{N_m} + PGS_{j}^{N_p}\right)$$

$$var\left(PGS_{j}^{N}\right) = var\left(PGS_{j}^{N_m}\right) + var\left(PGS_{j}^{N_p}\right) + 2cov\left(PGS_{j}^{N_m}, PGS_{j}^{N_p}\right)$$

$$var\left(PGS_{j}^{N}\right) = var\left(PGS_{ij}^{N}\right) + var\left(PGS_{ij}^{N}\right) + 2cov\left(PGS_{j}^{N_m}, PGS_{j}^{N_p}\right)$$

$$var\left(PGS_{j}^{N}\right) = \frac{\lambda^2}{4} + \frac{\lambda^2}{4} + 0$$

$$var\left(PGS_j^N\right) = \frac{\lambda^2}{2}$$

## A10. Between-Family Analyses

$$\hat{\psi}_1 = \frac{cov\left(\widehat{PGS}'^D_{ij}, Y_{ij}\right)}{var\left(\widehat{PGS}'^D_{ij}\right)}$$

$$\hat{\psi}_1 = \frac{cov\left(\widehat{PGS}'^D_{ij}, PGS^D_{ij}\right)\beta_1 + cov\left(\widehat{PGS}'^D_{ij}, PGS^N_j\right)\beta_2}{var\left(\widehat{PGS}'^D_{ij}\right)}$$

$$\hat{\psi}_1 = \frac{cov\left(\widehat{PGS}'^D_{ij}, PGS^D_{ij}\right) + cov\left(\widehat{PGS}'^D_{ij}, PGS^N_j\right)}{var\left(\widehat{PGS}'^D_{ij}\right)}$$

$$\hat{\psi}_1 = cov\left(\frac{\widehat{PGS}'^D_{ij} - \overline{\widehat{PGS}'^D_{ij}}}{var\left(\widehat{PGS}'^D_{ij}\right)^{\frac{1}{2}}}, PGS^D_{ij}\right) + cov\left(\frac{\widehat{PGS}'^D_{ij} - \overline{\widehat{PGS}'^D_{ij}}}{var\left(\widehat{PGS}'^D_{ij}\right)^{\frac{1}{2}}}, PGS^N_j\right)$$

$$\hat{\psi}_1 = \frac{cov\left(\widehat{PGS}^D_{ij}, PGS^D_{ij}\right) + cov\left(\widehat{PGS}^D_{ij}, PGS^N_j\right)}{var\left(\widehat{PGS}^D_{ij}\right)^{\frac{1}{2}}}$$

$$E\left[\hat{\psi}_1\right] = E\left[\frac{cov\left(\widehat{PGS}'^D_{ij}, PGS^D_{ij}\right) + cov\left(\widehat{PGS}'^D_{ij}, PGS^N_j\right)}{var\left(\widehat{PGS}'^D_{ij}\right)^{\frac{1}{2}}}\right]$$

$$E\left[\hat{\psi}_1\right] = \frac{cov\left(PGS^D_{ij} + PGS^N_{ij}, PGS^D_{ij}\right) + cov\left(PGS^D_{ij} + PGS^N_{ij}, PGS^N_j\right)}{var\left(PGS^D_{ij} + PGS^N_{ij}\right)^{\frac{1}{2}}}$$

$$E\left[\hat{\psi}_1\right] = \frac{cov\left(PGS^D_{ij}, PGS^D_{ij}\right) + cov\left(PGS^N_{ij}, PGS^D_{ij}\right) + cov\left(PGS^D_{ij}, PGS^N_j\right) + cov\left(PGS^N_{ij}, PGS^N_j\right)}{var\left(PGS^D_{ij} + PGS^N_{ij}\right)^{\frac{1}{2}}}$$

$$E\left[\hat{\psi}_1\right] = \frac{1 + \frac{\lambda\rho_g}{2} + \frac{\lambda\rho_g}{2} + \frac{\lambda^2}{4}}{\left(1 + \lambda\rho_g + \frac{\lambda^2}{4}\right)^{\frac{1}{2}}}$$

$$E\left[\hat{\psi}_1\right] = \frac{1 + \lambda\rho_g + \frac{\lambda^2}{4}}{\left(1 + \lambda\rho_g + \frac{\lambda^2}{4}\right)^{\frac{1}{2}}}$$

$$E\left[\hat{\psi}_1\right] = \sqrt{1 + \lambda\rho_g + \frac{\lambda^2}{4}}$$

## A11. Within-Family Analyses

$$\hat{\pi}_1 = \frac{cov\left(\Delta_0^1 \widehat{PGS}'^{\mathrm{D}}_{ij}, \Delta_0^1 Y_{ij}\right)}{var\left(\Delta_0^1 \widehat{PGS}'^{\mathrm{D}}_{ij}\right)}$$

$$\hat{\pi}_1 = \frac{cov\left(\Delta_0^1 \widehat{PGS}'^{\mathrm{D}}_{ij}, \Delta_0^1 PGS^{\mathrm{D}}_{ij}\right)\beta_1 + cov\left(\Delta_0^1 \widehat{PGS}'^{\mathrm{D}}_{ij}, \Delta_0^1 PGS^{\mathrm{N}}_{j}\right)\beta_2}{var\left(\Delta_0^1 \widehat{PGS}'^{\mathrm{D}}_{ij}\right)}$$

$$\hat{\pi}_1 = \frac{cov\left(\Delta_0^1 \widehat{PGS}'^{\mathrm{D}}_{ij}, \Delta_0^1 PGS^{\mathrm{D}}_{ij}\right) + cov\left(\Delta_0^1 \widehat{PGS}'^{\mathrm{D}}_{ij}, \Delta_0^1 PGS^{\mathrm{N}}_{j}\right)}{var\left(\Delta_0^1 \widehat{PGS}'^{\mathrm{D}}_{ij}\right)}$$

$$\hat{\pi}_1 = \frac{cov\left(\Delta_0^1 \widehat{PGS}'^{\mathrm{D}}_{ij}, \Delta_0^1 PGS^{\mathrm{D}}_{ij}\right)}{var\left(\Delta_0^1 \widehat{PGS}'^{\mathrm{D}}_{ij}\right)}$$

$$\hat{\pi}_1 = cov\left(\frac{\Delta_0^1 \widehat{PGS}^{\mathrm{D}}_{ij} - \Delta_0^1 \overline{\widehat{PGS}^{\mathrm{D}}_{ij}}}{var\left(\Delta_0^1 \widehat{PGS}^{\mathrm{D}}_{ij}\right)^{\frac{1}{2}}}, \Delta_0^1 PGS^{\mathrm{D}}_{ij}\right)$$

$$\hat{\pi}_1 = \frac{cov\left(\Delta_0^1 \widehat{PGS}^{\mathrm{D}}_{ij}, \Delta_0^1 PGS^{\mathrm{D}}_{ij}\right)}{var\left(\Delta_0^1 \widehat{PGS}^{\mathrm{D}}_{ij}\right)^{\frac{1}{2}}}$$

$$\mathrm{E}[\hat{\pi}_1] = \mathrm{E}\left[\frac{cov\left(\Delta_0^1 \widehat{PGS}^{\mathrm{D}}_{ij}, \Delta_0^1 PGS^{\mathrm{D}}_{ij}\right)}{var\left(\Delta_0^1 \widehat{PGS}^{\mathrm{D}}_{ij}\right)^{\frac{1}{2}}}\right]$$

$$\mathrm{E}[\hat{\pi}_1] = \frac{cov\left(\Delta_0^1 PGS^{\mathrm{D}}_{ij} + \Delta_0^1 PGS^{\mathrm{N}}_{ij}, \Delta_0^1 PGS^{\mathrm{D}}_{ij}\right)}{var\left(\Delta_0^1 PGS^{\mathrm{D}}_{ij} + \Delta_0^1 PGS^{\mathrm{N}}_{ij}\right)^{\frac{1}{2}}}$$

$$\mathrm{E}[\hat{\pi}_1] = \frac{cov\left(\Delta_0^1 PGS^{\mathrm{D}}_{ij}, \Delta_0^1 PGS^{\mathrm{D}}_{ij}\right) + cov\left(\Delta_0^1 PGS^{\mathrm{N}}_{ij}, \Delta_0^1 PGS^{\mathrm{D}}_{ij}\right)}{var\left(\Delta_0^1 PGS^{\mathrm{D}}_{ij} + \Delta_0^1 PGS^{\mathrm{N}}_{ij}\right)^{\frac{1}{2}}}$$

$$\mathrm{E}[\hat{\pi}_1] = \frac{var\left(\Delta_0^1 PGS^{\mathrm{D}}_{ij}\right) + cov\left(\Delta_0^1 PGS^{\mathrm{N}}_{ij}, \Delta_0^1 PGS^{\mathrm{D}}_{ij}\right)}{var\left(\Delta_0^1 PGS^{\mathrm{D}}_{ij} + \Delta_0^1 PGS^{\mathrm{N}}_{ij}\right)^{\frac{1}{2}}}$$

$$\mathrm{E}[\hat{\pi}_1] = \frac{1 + \frac{\lambda\rho_g}{2}}{\left(1 + \lambda\rho_g + \frac{\lambda^2}{4}\right)^{\frac{1}{2}}}$$

## A12. Empirical Results (Full Regression Table)

**Table A1.** The relationship between polygenic score and observed trait for six phenotypes in Add Health, within families and between families.

| | Years of Schooling | | Cognitive Ability | | CESD Depression Index | | Birth Weight | | Body Mass Index | | Height | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Family Fixed Effect | | X | | X | | X | | X | | X | | X |
| Educational Attainment PGS | 0.808** (0.0275) | 0.354** (0.132) | | | | | | | | | | |
| Cognitive Ability PGS | | | 3.173** (0.158) | 1.699** (0.629) | | | | | | | | |
| Depression PGS | | | | | 0.130** (0.0132) | 0.0556 (0.0650) | | | | | | |
| Birth Weight PGS | | | | | | | 3.071** (0.317) | 3.256* (1.360) | | | | |
| Body Mass Index PGS | | | | | | | | | 1.999** (0.0963) | 2.344** (0.542) | | |
| Height PGS | | | | | | | | | | | 2.496** (0.0932) | 2.508** (0.407) |
| r2 | 0.173 | 0.756 | 0.0836 | 0.736 | 0.0361 | 0.598 | 0.0353 | 0.816 | 0.0837 | 0.695 | 0.583 | 0.864 |
| N | 5323 | 734 | 5087 | 702 | 5323 | 734 | 4524 | 647 | 5269 | 727 | 5299 | 733 |

+ 0.10 * 0.05 ** 0.01 . All models control for sex, age, and the first 10 principal components of individual genotype. All models use only individuals of European ancestry. Models without individual-fixed effects use a sample of unrelated individuals, whereas the family-fixed effect models use a sample of sibling pairs. The sample of unrelated individuals contains one randomly selected sibling from each pair. All polygenic scores are standardized within sample to be mean 0 and standard deviation 1. Cognitive ability is measured through the Peabody Picture Vocabulary Test during Wave 1 of Add Health, when respondents were approximately 16 years old. Birth weight is retrospectively reported by respondents' parents during Wave 1 of Add Health. Years of schooling, CESD depression index, body mass index, and height are measured during Wave 4 of Add Health, when respondents were approximately 28 years old. Height is reported in centimeters, birth weight is reported in ounces, and cognitive ability is reported in IQ score points. The CESD depression index is normalized within sample to be mean 0 and standard deviation 1.