

# Gene-environment interactions with essential heterogeneity \*

Johannes Hollenbach

RWI &  
Paderborn University

Hendrik Schmitz

RWI &  
Paderborn University

Matthias Westphal

University of Hagen &  
RWI

July 2025

## Abstract

We study gene-environment interactions in the context of the long-term effect of education on cognition. Our central parameter of interest is the interaction effect between endogenous education and a predetermined measure of genetic endowment. Education is instrumented by a reform that raised compulsory schooling in England. We use this setting to show that two-stage least squares (2SLS) estimates of interaction effects can be misleading when there is essential heterogeneity (e.g., selection into gains) and complier status depends on the interaction variable. The 2SLS estimator cannot disentangle interaction effects from shifts in complier groups. Estimating marginal treatment effects addresses this problem by fixing the underlying population and unobserved heterogeneity. The results show complementarities between education and genetic predisposition in determining later-life memory, our measure of cognition. The marginal treatment effect estimates reveal a substantially larger gene-environment interaction, exceeding the 2SLS estimate by a factor of at least 4.7.

**JEL Classification:** C31, J14, J24

**Keywords:** Two-stage least squares estimation, marginal treatment effects, gene-environment interactions, cognitive decline

---

Johannes Hollenbach: RWI - Leibniz Institute for Economic Research, Hohenzollernstr. 1-3, 45128 Essen, Germany ([johannes.hollenbach@rwi-essen.de](mailto:johannes.hollenbach@rwi-essen.de)); Hendrik Schmitz: Paderborn University, Warburger Str. 100, 33098 Paderborn, Germany ([hendrik.schmitz@uni-paderborn.de](mailto:hendrik.schmitz@uni-paderborn.de)); Corresponding author: Matthias Westphal: FernUniversität in Hagen, Universitätsstraße 47, 58097 Hagen, Germany ([matthias.westphal@fernuni-hagen.de](mailto:matthias.westphal@fernuni-hagen.de)).

\*We thank the editor and four anonymous referees, as well as Silvia Barcellos, Pietro Biroli, Leandro Carvalho, Jason Fletcher, Lauren Schmitz, and Kevin Thom for their excellent comments and suggestions, which substantially improved the paper. We are also grateful to Souvik Banerjee, Martin Fischer, Hendrik Jürges, and Kristina Strohmaier for their thoughtful feedback as discussants. We further thank participants of the Initiative in Social Genomics group meetings at UW–Madison, the Applied Micro Workshop at UW–Milwaukee, the CINCH-dggö Academy in Health Economics in Essen, the EuHEA PhD Conferences in Bologna and Lucerne, as well as the Brown Bag Seminars in Wuppertal and at RWI Essen for their comments, questions, and lively exchange. Financial support from Deutsche Forschungsgemeinschaft (DFG, project number 437564156) is gratefully acknowledged. This paper uses data from the English Longitudinal Study of Ageing (ELSA). ELSA is funded by the National Institute on Aging (R01AG017644), and by UK Government Departments coordinated by the National Institute for Health and Care Research (NIHR).

# 1 Introduction

The recent availability of genetic data has revived the old debate in the social sciences about nature versus nurture in determining success over the life course (see, e.g., [Behrman and Taubman, 1989](#); [Björklund and Salvanes, 2011](#); [Plug and Vijverberg, 2003](#)). The focus is on estimating gene-environment ( $G \times E$ ) interactions to assess how the effects of environmental exposures or individual decisions vary by genetic endowment. These interaction models are typically specified as

$$Y_i = \beta_0 + \beta_1 E_i + \beta_2 G_i + \beta_3 G_i \times E_i + X_i' \gamma + \varepsilon_i, \quad (1)$$

where  $Y_i$  denotes a (long-run) outcome of interest,  $E_i$  represents an endogenous environmental exposure or individual decision,  $G_i$  captures a pre-determined genetic endowment, and  $X_i$  is a vector of control variables. Recent studies have focused on the causal identification of  $\beta_1$ ,  $\beta_2$ , and  $\beta_3$  by instrumenting  $E_i$  and  $G_i \times E_i$  and by removing factors correlated with the environment from  $G_i$  (see, e.g., [Barcellos et al., 2018, 2021, 2025](#); [Pereira et al., 2022](#); [Schmitz and Conley, 2017](#)). As an alternative to instrumenting  $E_i$ , some studies directly estimate interactions of  $G_i$  with a plausibly exogenous variable  $Z_i$ , which we refer to as  $G \times Z$  ([Ahlskog et al., 2024](#); [van den Berg et al., 2023a,b](#)).

The focus of our paper is on the estimation of the interaction coefficient,  $\beta_3$ , which is the central parameter in the gene-environment literature. In its intended interpretation, it measures how the causal effect of the environment varies with genetic endowment, all else being equal. However, as we demonstrate, the commonly used two-stage least squares (2SLS) or reduced-form approaches may not provide a reliable estimate of this effect, even with a valid instrumental variable. This is the case when two conditions hold simultaneously: First, compliers to the instrument for  $E_i$  have different unobserved characteristics between different values of  $G_i$ . Second, the (individual) treatment effects of  $E_i$  on  $Y_i$  exhibit essential heterogeneity. This occurs when the propensity to take the treatment correlates with the unobserved effect heterogeneity ([Heckman et al., 2006](#)). A prominent example of essential heterogeneity is self-selection into treatment based on unobserved gains. These conditions frequently occur in real-world settings that are investigated with causal methods. As a result, 2SLS conflates two different changes when estimating the  $G \times E$  coefficient: first, how the local average treatment effect (LATE) of  $E_i$  on  $Y_i$  changes with  $G_i$ , which is the interaction effect of interest. Second, how the complier subpopulation of this LATE shifts as  $G_i$  varies.

In this paper, we (1) comprehensively describe the problem, (2) propose a solution, and (3) apply it to a real-world setting. Using a numerical example, we show that relying on 2SLS estimates of  $\beta_3$  to provide evidence on how genes and the environment interact can be misleading in a setting with essential heterogeneity and a substantial gradient in the

first-stage coefficients across different  $G_i$ . In our simulation example, the 2SLS coefficient even has the opposite sign of the actual interaction effect. We propose a solution that maintains a fixed underlying population when comparing the effect of  $E_i$  on  $Y_i$  for different values of  $G_i$ . Estimating marginal treatment effects (MTEs) offers a suitable approach to achieve this (Heckman and Vytlačil, 2005). We apply this method to the long-term effect of education  $E_i$  on cognition in later life  $Y_i$  using data from the English Longitudinal Study of Aging (ELSA). We select our sample around the pivotal cohort of a compulsory schooling reform, which extended the minimum school-leaving age from 14 to 15 for individuals born after 1933. Our measure of cognition is the word recall test, a widely used indicator that has been shown to predict cognitive decline (Apolinario et al., 2016; Tsoi et al., 2017). We use data from six waves between 2002 and 2012 when individuals in our sample were between 65 and 80 years old. To measure genetic endowment, we use a polygenic index (PGI), a summary measure that predicts individual-level educational attainment based on the aggregated effects of many DNA differences between individuals. When estimating MTEs, we rely on a recently developed partial identification method by Mogstad et al. (2018), also used by Rose and Shem-Tov (2021).

Our paper makes three main contributions to the literature. The first is purely pedagogical. While it is well-documented that selection into gains poses problems for 2SLS in general (see, e.g. Heckman and Vytlačil, 2005), the problem that interaction effects are difficult to interpret in this setting still deserves attention. We aim to provide an accessible and intuitive presentation of the problem. The problem we describe is not limited to the gene-environment literature and is, in principle, relevant to any interaction effect between an endogenous (and instrumented) treatment variable and observable characteristics. In Appendix F we cover other settings where researchers are interested in effect heterogeneity by observables and where the same problems might occur, asking for the same kind of solution. A more important contribution is to provide a transparent and easy-to-implement solution by using marginal treatment effects in this setting. Of course, other ways exist to separate possibly correlated observed and unobserved effect heterogeneity. Kline and Walters (2019) discuss the general equivalence between instrumental variable methods and control functions (Blundell and Powell, 2003; Imbens and Newey, 2009). In control function approaches, the essential heterogeneity is absorbed by a control variable (which might incorporate instrumental variables and functional form assumptions). Although certainly possible, we are not aware of a control function approach in the setting of endogenous interaction terms. Arold et al. (2025) use a control function approach in a  $G \times E$ -study, but they employ the approach by Altonji and Mansfield (2018) and base their control function on group-level averages of observed characteristics without using instruments. The implicit assumption is that the control function is modeled correctly (which may be a stronger assumption than the restrictions we impose). The inability of 2SLS to incorporate unobservable differences between complier groups that could (partially) explain gene-

environment interactions is also mentioned in [Barcellos et al. \(2021\)](#). They find differences in returns to schooling between individuals with different genetic endowments and use a linear MTE estimation to check whether unobservable factors can explain these disparities, which is not the case in their study. Our third contribution is a substantive one to the literature on gene-environment interactions, a dynamic field with numerous recent papers in areas related to ours. We are unaware of any study estimating the causal effects of education and its interaction with genetic makeup on memory in later life. [Banks and Mazzonna \(2012\)](#) study the effect of the same reform we do on memory, but without looking at gene-environment interactions. [Ding et al. \(2019\)](#) study the relationship between genes/educational attainment and recall using data from the Health and Retirement Study (HRS), but do not use exogenous variation in education. [Anderson et al. \(2020\)](#) estimate a positive bidirectional relationship between educational attainment and intelligence using genetic variants as instruments. [Schmitz and Conley \(2017\)](#) study whether the effect of the Vietnam War draft lottery on schooling outcomes differs by a genetic predisposition for education. Going beyond educational outcomes, [Barcellos et al. \(2018\)](#) estimate whether genetic predisposition to obesity moderates the effect of education on health using the UK compulsory schooling reform for the 1957 birth cohort as an exogenous variation and a different data set. [Ahlskog et al. \(2024\)](#) estimate reduced-form interactions between compulsory schooling exposure in Sweden and a set of different PGIs (i.e., they focus on  $G \times Z$ ) on different outcomes. They find significant interactions for two outcomes (wages and educational attainment), both with the PGIs for educational attainment that we also use in this paper. However, one drawback of the focus on  $G \times Z$  is that reduced-form regressions do not differentiate between an implicit first-stage gradient in  $G_i$  and heterogeneous direct effects of  $Z_i$  on the outcome along  $G_i$  (i.e., through  $E_i$ ). Besides [Schmitz and Conley \(2017\)](#) and [Barcellos et al. \(2018\)](#), the earliest study in economics on how education can compensate for the effects of genetic differences is probably [Papageorge and Thom \(2020\)](#), who study the impact on labor market outcomes.

Our results are as follows: Applying a benchmark 2SLS estimator, we find a zero effect of education on recall for individuals in the lowest quintile of  $G_i$ , that is, those with the lowest genetic propensity for education. On average, moving to a higher quintile of  $G_i$  goes along with an increase in the effect of  $E_i$  on  $Y_i$  by an insignificant 0.1 words correctly recalled. Using marginal treatment effects, the interaction effect is much larger than when estimating with 2SLS: The average linearized interaction effect across all quintiles of  $G_i$  indicates that the effect of  $E_i$  on  $Y_i$  increases by 0.46–0.47 words per quintile. This corresponds to roughly 10–15 percent of the standard deviation of the outcome variable. While education does not improve memory in the group with the lowest genetic endowment, it increases word recall by about 1.8 words in the highest quintile of  $G_i$  compared to the lowest. 2SLS would considerably underestimate this gene-environment complementarity. In our application, genetic endowment is correlated with the complier status: The share of

compliers to the education reform is highest in the lowest quintile of  $G_i$  (65 percent) and monotonically decreases to 36 percent in the highest quintile. Moreover, there is evidence of selection into gains. Overall, the 1947 UK compulsory schooling reform has increased schooling, especially for those with lower genetic propensity for schooling (first stage results). However, these individuals have no returns to schooling in terms of cognition. Instead, significant returns are seen for those with a higher genetic propensity.

The paper proceeds as follows: Section 2 describes the institutional setting of our application and the data used. Section 3 presents 2SLS estimates of gene-environment interactions in our application. Section 4 outlines the challenges in identifying the gene-environment interplay from an econometric perspective and presents our suggested solution. Section 5 gives an overview of the partial identification approach to estimate MTEs and presents our main results. Section 6 concludes.

## 2 Institutional Setting and Data

### 2.1 Compulsory schooling reform in the UK

In our application, we exploit exogenous variation from a compulsory schooling reform in the UK. Based on the Education Act of 1944, two reforms were introduced to raise the minimum school-leaving age in England, Scotland, and Wales. We use the first reform, which took effect on April 1, 1947.<sup>1</sup> This reform raised the minimum age for leaving school from 14 to 15. Given that students in the UK typically entered school at the age of 5, the 1947 reform effectively extended compulsory education from nine to ten years. The first birth cohorts to be affected by this change, i.e., the first to be required to attend school for an additional year (the “pivotal cohort”), were those born in April 1933. This particular reform from 1947 has served as exogenous variation for compulsory schooling in studies on the effect of education on wages (Clark and Royer, 2013; Devereux and Hart, 2010; Harmon and Walker, 1995; Oreopoulos, 2006), other labor market outcomes (Clark, 2023), health (Clark and Royer, 2013; Jürges et al., 2013; Powdthavee, 2010; Silles, 2009), health knowledge (Johnston et al., 2015), mortality (Clark and Royer, 2013; Gathmann et al., 2015), and cognitive abilities (Banks and Mazzonna, 2012).

To demonstrate the strong response to the compulsory schooling reform from 1947, Figure 1 shows aggregated cohort-level data from ELSA. It depicts the share of individuals with different levels of schooling by birth cohort. The pivotal cohorts of both compulsory

---

<sup>1</sup>The second part was enacted much later, in 1972, raising the school-leaving age to 16. Since we are interested in studying memory in old age, we use only the 1947 reform. Cohorts affected by the second reform in 1972 are, for the most part, still too young at the time of data collection for the English Longitudinal Study of Ageing, our data source.

schooling reforms are marked with vertical lines. The highest line (circle markers) shows how the 1947 reform caused a significant increase in the share of students leaving school at age 15 or later from about 40% to almost 100%. The middle line (diamond markers) shows how the second reform in 1972 lead to a still remarkable but comparably smaller increase in the share of leaving school at 16 or later from 75% to about 90%. The lowest line (triangle markers) can be read as a placebo test, showing the general trend in increased years of schooling but no discontinuity at the two reform cut-offs (Clark and Royer, 2013).

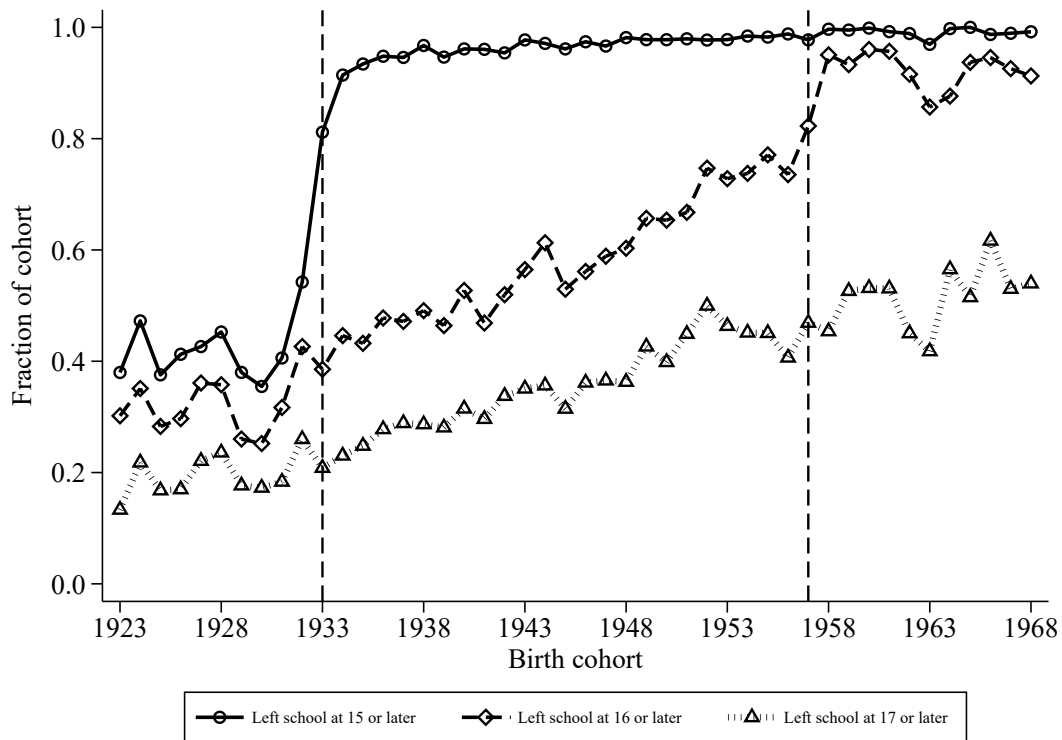


Figure 1: Education by birth cohort

*Notes:* This figure illustrates the shares of students leaving school at 15 or later, 16 or later, and 17 or later changed over birth cohorts and how these shares were affected by two compulsory schooling reforms in England using data from ELSA waves 1–9 without the sample selection described in Chapter 2.2. Vertical dashed lines indicate the first affected birth cohorts of two school-leaving age increases. The three groups are not mutually exclusive and do not add up to 100%. The illustration is adapted from Clark and Royer (2013) to fit our definition of educational attainment.

Besides the high compliance rates, Figure 1 also reveals noncompliance. That is, despite being disallowed by the new compulsory schooling age, some individuals reportedly leave school at the age of 14 after the reform. According to Clark and Royer (2013), who studied the reforms extensively, this noncompliance is primarily due to individuals born in the summer months, who turned 15 before the start of the next school year.

## 2.2 Sample and Variables

### Sample

We use data from the English Longitudinal Study of Ageing (ELSA), a large representative



microdata set providing information on health and other socioeconomic characteristics of individuals aged 50 and over in England (Banks et al., 2023). ELSA was launched in 2002 and is conducted every two years. It currently comprises eleven waves of interviews.<sup>2</sup> We use individuals aged 65–80 from waves 1–6 of ELSA. Data collection for wave 6 took place in 2012 and 2013 when individuals born in 1933 – our cutoff – turned 80. Thus, starting with wave 7, only individuals born after the cutoff can theoretically enter the sample. We exclude the 1933 birth cohort because we lack information on birth month and cannot accurately assign this cohort to pre- or post-reform status (the cutoff is April 1933). We also restrict the data to birth cohorts ten years before and after the reform cut-off. Finally, for our main analysis, we need to limit the data to individuals for whom genetic data is available. This reduces the number of individuals by about 50 percent and may introduce a selection bias if the compulsory schooling reform affects the willingness to be genotyped. We find that the sample is selective regarding the outcome variable: Individuals who consent to be genotyped have higher recall scores on average (see Table A.2 in the Appendix). However, we do not find evidence of a statistically significant effect of the compulsory schooling reform on the probability of being genotyped (see Table B.2 in the Appendix). Similarly, the willingness to be genotyped does not interact with the impact of compulsory schooling on the probability of going to school until at least the age of 15. In our preferred estimation sample, we use all available observations per individual.<sup>3</sup> In doing so, we do not assess effects at a single point in time, but implicitly receive average effects over multiple years. The robustness of this choice concerning panel attrition and alternative samples is addressed in Section 3.2. This sample comprises 11,027 observations from 3,009 individuals born between 1923 and 1943, who are observed between 2002 and 2013.

## Cognition

Cognitive abilities – as a broad concept – include “the ability to reason, plan, solve problems, think abstractly, comprehend complex ideas, learn quickly and learn from experience. It is not merely book learning, a narrow academic skill, or test-taking smarts. Rather, it reflects a broader and deeper capability for comprehending our surroundings – ‘catching on,’ ‘making sense’ of things, or ‘figuring out’ what to do.” (Gottfredson, 1997). The sum of these abilities is called intelligence (Schiele and Schmitz, 2023). A wide range of cognitive tests measure different aspects of cognitive abilities to accommodate this multifaceted notion. ELSA offers several measures, including cognitive capacity, temporal orientation, literacy, and numerical ability. We use test scores from the word recall test in which an interviewer reads ten words to the respondent, who is then asked to recall

---

<sup>2</sup>For details of the ELSA sampling procedure, questionnaire content, and fieldwork methodology, see Steptoe et al., 2013.

<sup>3</sup>199 individuals are observed once, 471 twice, 856 three times, 550 four times, 466 five times, and 467 in all six waves.

as many words as possible. This test is administered twice: immediately after the words are read (immediate recall) and five minutes later (delayed recall). The scores from both instances are added together to yield a total recall score, which can range from 0 to 20. It serves as a measure of episodic memory, susceptible to aging (Rohwedder and Willis, 2010). Episodic memory is considered a component of fluid intelligence, reflecting the innate cognitive ability to store and retrieve information. It is distinct from crystallized intelligence that people acquire over a lifetime (using their fluid intelligence). Word recall has been shown to predict cognitive decline (Bruno et al., 2013; Tsoi et al., 2017) and is an important part of measures for (mild) cognitive impairment (Apolinario et al., 2016; Cadar et al., 2020). Furthermore, it is widely used in economics as a reliable and accessible measure of cognitive functioning (see e.g., Banks and Mazzonna 2012; Bonsang et al. 2012; Christelis et al. 2010; Schiele and Schmitz 2023). In our estimation sample, the total recall score, our dependent variable, has a mean of 9.67 correctly recalled words (out of 20) with a standard deviation (SD) of 3.37 words (see Table 1).

## Education

ELSA does not provide information on respondents' years of education, but on the age at which they completed their continuous full-time education. However, the data is aggregated at the low (finished age 14 or earlier) and high (finished age 19 or later) ends. Our treatment variable  $E_i$  is a binary variable equal to one if the individual has left school at 15 or later, and zero otherwise. By design, and as observable in Figure 1, the proportion of individuals having left school at 15 or later (i.e., having stayed in school for at least ten years) is affected by the 1947 education reform that raised the minimum school-leaving age from 14 to 15.

Education is assessed retrospectively, and thus potentially affected by recall bias, a common concern in older age samples. Yet, respondents may better be able to recall the year of school completion (especially so close after the end of World War II) than general years of education. Moreover, Figure 1 (and also the subsequent regression analyses) match remarkably well with the corresponding estimates of Clark and Royer (2013), who use a survey collected from 1991 to 2004 – more than one decade before our estimation sample.<sup>4</sup> Hence, we believe that our education information is unlikely to be significantly affected by recall bias.

## Genes

We use an Educational Attainment Polygenic Index provided by ELSA and based on Lee et al. (2018) to measure genetic makeup. This indicator predicts educational attainment

---

<sup>4</sup>The corresponding first stage coefficients are 0.445 versus 0.479 – a difference smaller than our standard error.



based on differences in genetic variants across individuals. The education PGI we use explains 11-13 percent of the variation in educational attainment in the original discovery sample (Lee et al., 2018). An individual's PGI represents their genetic propensity (or genetic risk – depending on the application) for a particular trait – not just according to one genetic marker, but over many genetic variants. The PGI we are using thus represents individual genetic propensity for educational attainment. For a more detailed explanation of polygenic indices and their construction, see Appendix C. The PGI is normally distributed. Individuals whose genetic endowment puts them on the left side of this distribution have a lower genetic propensity to pursue education; individuals on the right side have a higher propensity. As the choice equations in Section 4 emphasize, this propensity is not deterministic. Individuals with a high PGI are not necessarily highly educated, and highly educated individuals do not necessarily have a high PGI for educational attainment. Our analysis uses the quintiles of this index, yielding five equally sized groups.

When estimating gene-environment interactions, researchers often use a polygenic index of the outcome they are investigating since it is the obvious choice and will produce an effect. However, the choice is not set in stone. As Biroli et al. (2025) point out, “any PGI could be used if warranted by theory or for empirical reasons”. We target the PGI towards the environmental variable (education) by using an education PGI, and the outcome we investigate is memory. Education PGIs are associated with several different outcomes besides educational attainment: wealth at retirement (Barth et al., 2020), labor market earnings (Papageorge and Thom, 2020) and socioeconomic success (Belsky et al., 2018). In our setting, we can use the education PGI to demonstrate heterogeneous responses to the education reform by the relevant part of the genetic endowment. At the same time, the effect of education on cognition likely varies with genetic propensity for education – which is what we want to estimate.

We include the first ten principal components of the genetic data as controls, which are summary scores of the overall variation of the genetic data in ELSA, condensed into a smaller number of dimensions. They reflect population stratification, i.e., different frequencies of genetic variants among subpopulations that could be responsible for spurious correlations with outcomes of interest. This could occur if both an outcome and certain genetic variants are more common in one stratum of the population than in another and they do not mate randomly (Barth et al., 2022). Price et al. (2006) show that including principal components as control variables can mitigate these confounding effects. Therefore, adding principal components as controls has become a convention in gene-environment studies (see, e.g., Barcellos et al., 2018; Barth et al., 2020, 2022; Biroli et al., 2025; Pereira et al., 2022). Both principal components and polygenic indices combine information from differences of genetic variants across the population. However, they serve different purposes: principal components capture overall genetic similarity and

population structure, while PGIs predict specific traits, such as educational attainment, based on gene-outcome associations.<sup>5</sup>

While predetermined at conception, the effect of an individual’s genetic endowment is not entirely exogenous. Genetic makeup is fully inherited from the parents, whose own genetic endowment also influences the family environment in which the children are raised. This environment, in turn, partially determines later-life outcomes, creating a correlation between the child’s genetic endowment and their developmental context (Biroli et al., 2025; Houmark et al., 2024). We are interested in the differential effect of  $E_i$  (for which we have a valid instrument) by educational attainment PGI. While based on predetermined genetic variants, the PGI also reflects this correlation. Houmark et al. (2024) show that these family genetic correlations are due to family characteristics and can be effectively and almost entirely accounted for by parental education. We follow them as well as Barth et al. (2020), Barth et al. (2022), and Papageorge and Thom (2020) and add parental education as additional controls. ELSA includes information on the age at which a respondent’s mother and father left school, truncated at both ends (age 14 or under and age 19 or over). The vast majority (about 60%) of mothers and fathers left school at age 14 or before. Therefore, we condense the information on education into a categorical variable with three permutations: One if both parents have no or low education (i.e., left school at age 14 or earlier), one if at least one parent stayed in school beyond age 14, and one if information on parental education is missing. We have missing information for 988 individuals in our sample. Since we are running local estimations, we choose not to drop them.

## 2.3 Descriptive Statistics

Table 1 shows descriptive statistics of our main sample of individuals for whom genetic information is available as well as of “treatment” ( $E_i = 1$ ) and “control” ( $E_i = 0$ ) groups separately. Overall, about three-quarters of the observations in the sample are in the treatment group; 66 percent were born in 1933 or later, and 52 percent are female. The treatment group scores significantly higher in recall than the control group. More educated individuals ( $E_i = 1$ ) exhibit a more favorable genetic endowment (significantly less observations in the first and more in the top quintile). Unsurprisingly, individuals in the treatment group are, on average, younger since they are more likely to be born after the compulsory schooling reform. Table A.1 extends the statistics. Table A.3 in the appendix shows the sample means by quintiles of the education PGI. Instrument assignment, age, and proportion of women do not vary across quintiles of the education PGI. However, individuals in higher quintiles perform better on the recall test. The difference between an average person in the lowest PGI quintile and an average person in the highest quintile is

---

<sup>5</sup>Nevertheless, we show in a robustness check that including principal components does not drive our results (see Table 5).

1.33 words, a sizable difference compared to the overall mean of 9.67. Not surprisingly, the probability of having more schooling is also higher in higher education PGI quintiles.

Table 1: Descriptive statistics

	Main sample	By $E_i$		
	Mean (SD)	$E_i=1$	$E_i=0$	Difference (SE)
<i>Outcome <math>Y_i</math></i>				
Recall score	9.67 (3.37)	10.11	8.08	-2.03 (0.07)***
<i>Treatment <math>E_i</math></i>				
Left school $\geq 15$	0.78 (0.41)	1.00	0.00	-1.00 (0.00)
<i>Polygenic index <math>G_i</math></i>				
1st PGI quintile	0.20 (0.40)	0.18	0.25	0.07 (0.01)***
2nd PGI quintile	0.19 (0.40)	0.19	0.21	0.02 (0.01)**
3rd PGI quintile	0.20 (0.40)	0.21	0.19	-0.02 (0.01)**
4th PGI quintile	0.21 (0.41)	0.21	0.20	-0.01 (0.01)
5th PGI quintile	0.20 (0.40)	0.22	0.15	-0.07 (0.01)***
<i>Instrument <math>Z_i</math></i>				
Born 1933 or later	0.66 (0.47)	0.82	0.13	-0.69 (0.01)***
<i>Selected Controls (for a complete list, see Table A.1)</i>				
Female	0.52 (0.50)	0.52	0.50	-0.02 (0.01)**
Birth year	1934.89 (5.00)	1936.29	1929.92	-6.37 (0.10)***
Parental education:				
Missing	0.25 (0.43)	0.20	0.41	0.21 (0.01)***
Both left school $\leq 14$	0.57 (0.49)	0.58	0.55	-0.03 (0.01)**
At least one left school $\geq 15$	0.18 (0.39)	0.22	0.04	-0.18 (0.01)***
Observations	11,027	8,590	2,437	

Notes: This table presents descriptive statistics using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. The categories for parental education include: Missing information of at least one parent, both parents left full-time education at age 14 or before or have no education, and at least one parent stayed in school until age 15 or longer. We include mean and standard deviation of the main sample as well as means by  $E_i$ , the difference of means and standard errors of a t-test for equality of means. \* $p < 0.1$ , \*\* $p < 0.05$ , and \*\*\* $p < 0.01$ .

### 3 Benchmark 2SLS estimation

#### 3.1 Empirical Strategy

We start by estimating the gene-environment interactions using “conventional” methods. Since education is a choice variable, an OLS regression will yield biased estimates. We estimate the following 2SLS regression:

$$E_i = \pi_0 + \pi_1 G_i + \pi_2 Z_i + \pi_3 G_i \times Z_i + X' \gamma + f(t) + u_i \quad (2)$$

$$Y_i = \beta_0 + \beta_1 G_i + \beta_2 \widehat{E_i} + \beta_3 \widehat{G_i \times E_i} + X' \delta + f(t) + \varepsilon_i \quad (3)$$

Eq. (2) is the first stage, where we regress education  $E_i$  on our instrument  $Z_i$ , genetic predisposition  $G_i$  and the interaction of  $G_i$  and  $Z_i$ .<sup>6</sup> Eq. (3) shows the second stage. Here, we regress the outcome variable  $Y_i$  (the total recall score for individual  $i$ ) on the predicted values  $\widehat{E}_i$  from the first stage,  $G_i$  and the predicted values  $\widehat{G_i \times E_i}$ . In both stages, we add the same controls  $X_i$  which include an indicator variable for sex, the first ten principal components of the genetic data (see Section 2.2 for a description), as well as the ten interactions of the principal components with the instrument, fixed effects for survey wave, as well as a categorical variable for parental education. Furthermore,  $f(t)$ , is a function that captures a linear cohort trend and its interaction with the instrument  $Z_i$ . This specification estimates a fuzzy regression discontinuity model with the re-centered distance to the reform cohort of 1933 (the cohort trend) as the running variable. Finally,  $u_i$  and  $\varepsilon_i$  capture all unobserved factors that affect outcome variables in their respective stages. We cluster standard errors at the individual level in all analyses.

Besides the potential problems due to essential heterogeneity, this specification linearizes the  $G \times E$  effect (and the effect of  $G$  itself). This may also mask potentially interesting non-linearities. To be more flexible, we extend our analysis by fully saturating our specification using information on the quintiles of the education polygenic index. These effects compare better to our MTE approach because we directly estimate effects by quintiles. Accordingly, we estimate the following adapted model:

$$E_i = \sum_{g=1}^5 \left[ \pi_{g,0}^f \mathbb{1}[G_i = g] + \pi_{g,\Delta}^f \mathbb{1}[G_i = g] \times Z \right] + X' \gamma^f + f(t) + \omega_i \quad (4)$$

$$Y_i = \sum_{g=1}^5 \beta_{g,0}^f \mathbb{1}[G_i = g] + \beta_{1,1}^f \widehat{E}_i + \sum_{g=2}^5 \beta_{g,1}^f \widehat{[G_i = g] \times E_i} + X' \delta^f + f(t) + \eta_i \quad (5)$$

This is the equivalent of the 2SLS model described above in Eqs. (2) and (3), but with sets of indicator variables for the five quintiles of the PGI ( $G_i = g$  with  $g \in \{1, 2, 3, 4, 5\}$ ). To distinguish the coefficients from the base model, we add the superscript  $f$ . While in the baseline first stage (Eq. 2),  $\pi_2$  informs about the share of compliers in the data, the  $\pi_{1,\Delta}^f$  to  $\pi_{5,\Delta}^f$  of Eq. (4) inform about the share of compliers by PGI quintile. In the second stage (Eq. 5), we include  $\widehat{E}_i$  as the reference category that captures the local average treatment effect for the lowest quintile ( $\beta_{1,1}^f$ ). The coefficients  $\beta_{2,1}^f$  to  $\beta_{5,1}^f$  inform about gene-environment interactions relative to the lowest quintile.

---

<sup>6</sup>Note that there are technically two first stages, one with the dependent variable  $E_i$  and one with the dependent variable  $G_i \times E_i$ . Depending on how  $G_i$  is included, there are more. With  $G_i$  as quintiles of the PGI, there are six first stages. For the sake of simplicity, we only show one of them here.

## 3.2 Assumptions

We need to assume that the compulsory schooling reform is a valid instrument to identify the causal effects of extending schooling beyond the age of 14. Specifically, the 1933 birth cohort cutoff must be exogenous to the individuals in our sample. This is plausible given that the reform was announced in 1944 and the sample does not suffer from selective non-response or attrition (discussed below). Additionally, we assume that only compulsory schooling changes discontinuously for individuals born after April 1933, without other factors changing simultaneously (the exclusion restriction). Finally, we assume that no individual leaves school earlier because of the reform (the monotonicity assumption).

The exclusion restriction deserves the most discussion, for instance, as spillovers might exist and because two significant events occurred around the time our sample cohorts were born: the Great Depression and World War II. Although individuals may have experienced rationing or evacuations, those on either side of the 1933 cutoff were affected similarly (Clark and Royer, 2013). Furthermore, the compulsory schooling reform may also have increased the general quality of schooling, affecting not only compliers, but also generating spillover effects to always-takers. However, as Clark (2023) documents for the 1947 UK reform, nearly all compliers attended lower-track schools that ended at the minimum leaving age. This makes it unlikely that spillovers to non-complying groups exist. The lower-track schools emphasized practical education, exhibited lower quality (e.g., class size and teacher qualifications), which did not change due to the reform (as resources adjusted to increased enrollment, see Clark and Royer, 2013). These facts suggest that the reform did not affect the quality of schooling (not even for compliers). Thus, we can interpret our treatment effects in terms of years of schooling (as is commonly done in this literature). Clark (2023) also finds that the reform did not raise the probability of students receiving formal academic or vocational qualifications. Nevertheless, as Clark and Royer (2013) note, citing official reports from the period, “the extra year created by the 1947 change introduced some students to more advanced materials and helped other students master more basic material,” suggesting a natural progression in curricula rather than an overhaul.

Additionally, panel attrition may be a concern in older-age samples. If the education reform affected survival or survey participation, it could lead to a disproportionate representation of healthier, more educated individuals among respondents. This selective attrition could bias the estimates if the instrument indirectly influences the sample’s composition through differential attrition at older ages. Clark and Royer (2013) comprehensively investigated the effect of the 1947 reform on mortality and reported no or negligible effects. Nevertheless, we test for differential panel attrition in our sample by filling in the observations for each individual where necessary from the first wave they were observed in until wave 6, and creating an attrition indicator if they did not respond (for whatever reason) in a

subsequent wave. We then regress this indicator on the instrument to assess whether the compulsory schooling reform predicts survey non-response. Table B.1 presents the results. The estimate is negligibly small and not statistically significant. One major difference of our sample compared to related studies (e.g., Banks and Mazzonna, 2012) is that we exclude individuals who did not provide genetic information to ELSA. Therefore, we also examine whether the probability of sharing genetic information jumps discontinuously at the cutoff (see Table B.2). We find that this is not the case. We conclude that, although panel attrition is generally a concern, it is not related to the schooling reform in our sample.

### 3.3 Results

#### OLS, reduced-form, and second-stage results

Table 2 presents the OLS, reduced form, and 2SLS regression results (Eq. 5) in Columns (1), (2), and (3), respectively. Panel A includes controls for each quintile – but no interaction of  $G$  with  $E$  or  $Z$ . Panel B adds these interactions. Finally, we use the standardized PGI as a continuous interaction variable to show linear effects in Panel C. Without interactions, we see a considerable correlation: Individuals who left school at age 15 or later recall about 1.1 words more later in life. Using the compulsory schooling reform as exogenous variation, however, the reduced form and the 2SLS estimates suggest that the causal effect, if there is any, is considerably smaller. We find that individuals who would have dropped out at age 14, but had to stay in school for at least one additional year, recall about 0.15 words more. This is a small and statistically insignificant effect. This effect is also more negligible compared to Banks and Mazzonna (2012), who find a relevant and significant impact on the recall score in their preferred specification. However, as we use an unrestricted sample regarding the school-leaving age, control for  $G$  (including its principal components) and parental education, and use additional individuals from later waves, our estimates may be more conservative.

The OLS coefficients of  $G_i = 2$  to  $G_i = 5$  suggest that, in general, only individuals in the highest PGI quintile score statistically significantly higher on the recall test relative to individuals in the lowest – about 0.77 words higher than individuals in the lowest quintile. This positive relationship between an education PGI and cognitive performance is also documented by Jeong et al. (2024), who use data from the US Health and Retirement Study. The interactions in Panel B use the first quintile ( $G = 1$ ) as a reference category, so that the remaining interaction coefficients are interpreted as the additional effect of higher quintiles, relative to the first one. An additional year of education ( $E_i$ ) is associated with an increase of about 0.62 words later in life for individuals in the lowest PGI quintile. The markups on this association for individuals in the four higher PGI quintiles ( $E_i \times (G_i = g)$ ) are positive across all quintiles. However, their magnitude varies. For individuals in the third quintile, an additional year of schooling results in about 0.45 more words recalled



than for individuals in the first quintile. For the fourth quintile, this relative premium versus the lowest quintile increases to 0.76. Except for the interaction coefficient of the fourth quintile, they are not statistically significant, but importantly, all indicate a positive gene-environment interaction. All in all, this suggests that the association between genes, more education, and memory are mutually reinforcing. When we include the standardized PGI as a continuous variable (in a separate regression, shown in Panel C), its interaction coefficient suggests that a one standard deviation increase in PGI is associated with an additional rise in recall score by 0.18 words. However, our OLS results only represent correlations.

Reduced-form estimates that regress the instrument  $Z_i$  and its interaction with  $G_i$  directly on recall are reported in Column (2). Our 2SLS estimates are reported in Column (3). The coefficient of the reduced form without interacting with  $G_i$  is almost zero (Panel A). Nevertheless, when considering gene-instrument interactions (Panel B), we see positive effects for all quintiles except the lowest. However, only effects for quintiles two and five have a relevant size, and none of the interaction coefficients are statistically significant. Additionally, we visualize the reduced form alongside the first stage using sample means for each birth cohort in Figure B.1 in the Appendix. The corresponding coefficients are reported in Table B.3. When considering the raw means like this, there are apparent and large differences by  $G_i$  in the first stage that discussed in the subsequent section, but little to no differences in the reduced form. The 2SLS regression for education finds a small, positive, but not statistically significant effect of more education on later-life recall when not interacting with education PGI (Panel A). Furthermore, there is a zero effect of an additional year of schooling on recall for those in the lowest PGI quintile (Panel B) and a positive estimate for individuals in the upper quintiles. The standard errors are large, so we cannot be certain that these interactions differ from zero. The linear interaction effect using a standardized PGI (Panel C) is also close to zero. Based on these results, we would conclude that, after resolving the endogeneity problem with  $E_i$  by instrumenting – if anything – there may only be a small positive interaction effect that cannot be precisely estimated. The cognitive returns to education are likely not much higher for individuals with higher genetic endowment. However, consistent with the problem outlined in Section 4, recall that when comparing the effects between two quintiles, the complier group also changes, which may offset the small and monotonic gene-environment interaction.

### First-stage results

We report the coefficients of  $Z_i$  by  $G_i$ , that is,  $\pi_{1,\Delta}^f$  to  $\pi_{5,\Delta}^f$  of Eq. (4) in Figure 2.<sup>7</sup> It shows that overall, there is a large share of compliers to the education reform in the data. However, it varies substantially over the quintiles of the PGI. Complier share monotonically decreases

---

<sup>7</sup>Regression results are reported in Table B.3 in the Appendix.

Table 2: OLS, reduced-form and 2SLS estimates

	Dependent variable – total recall score					
	OLS (1)		Reduced form (2)		2SLS (3)	
<i>Panel A: baseline estimate, w/o interactions</i>						
$E_i$	1.099	(0.138)***			0.154	(0.423)
$Z_i$			0.075	(0.209)		
<i>Panel B: Including nonlinear interactions with PGI quintiles</i>						
$E_i$	0.620	(0.245)***			−0.021	(0.449)
$Z_i$			−0.022	(0.294)		
$G_i = 1$	reference category		reference category		reference category	
$G_i = 2$	0.187	(0.269)	0.355	(0.257)	0.239	(0.397)
$G_i = 3$	0.195	(0.274)	0.558	(0.265)**	0.524	(0.451)
$G_i = 4$	0.177	(0.270)	0.802	(0.246)***	0.784	(0.448)*
$G_i = 5$	0.774	(0.314)**	1.084	(0.266)***	0.871	(0.591)
$E_i \times (G_i = 1)$	reference category				reference category	
$E_i \times (G_i = 2)$	0.349	(0.325)			0.314	(0.497)
$E_i \times (G_i = 3)$	0.452	(0.331)			0.091	(0.552)
$E_i \times (G_i = 4)$	0.759	(0.328)**			0.049	(0.577)
$E_i \times (G_i = 5)$	0.433	(0.364)			0.394	(0.698)
$Z_i \times (G_i = 1)$			reference category			
$Z_i \times (G_i = 2)$			0.185	(0.317)		
$Z_i \times (G_i = 3)$			0.053	(0.327)		
$Z_i \times (G_i = 4)$			0.028	(0.316)		
$Z_i \times (G_i = 5)$			0.167	(0.330)		
<i>Panel C: Including linear interaction with continuous PGI</i>						
$E_i \times G_i$	0.179	(0.111)			0.011	(0.209)
$Z_i \times G_i$			−0.008	(0.102)		
Controls	Yes		Yes		Yes	
Observations	11,027		11,027		11,027	

Notes: This table presents OLS, reduced-form and 2SLS estimates of the effect of staying in school until at least age 15 ( $E_i$ ), an education PGI ( $G_i$ ) and their gene-environment interaction ( $G \times E$ ) on recall later in life using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. In panel A, we show estimates of education (respectively, the instrument  $Z_i$  – born in 1933 or later) on recall without interacting with genetic endowment. For estimates in Panel B, we use quintiles of the polygenic index and estimate non-linear interaction effects. Panel C shows estimates of a linear effect when including the standardized PGI as a continuous variable. Coefficients in all panels are obtained from separate regressions. Controls in each case include a linear cohort trend, its interaction with the instrument, sex, survey wave fixed effects, parental education, the first ten principal components of the genetic data as well as interactions of each principal component with the instrument. Standard errors clustered at the individual level shown are in parentheses. \* $p < 0.1$ , \*\* $p < 0.05$ , and \*\*\* $p < 0.01$ .

along the PGI. In the lowest quintile ( $G_i = 1$ ), 65 percent of all individuals increased their length of education due to the reform. The share of compliers reduces to a still sizable 36 percent in the highest quintile. The compulsory schooling reform had a more substantial impact on individuals in the lowest quintiles of the PGI, who are disadvantaged in terms of the genetic endowment that predicts education. Therefore, the reform was likely effective in targeting disadvantaged children. It drastically increased their probability of staying

in school until at least age 15. Our estimates in Table 2 suggest that the reform may not have successfully reduced differences in the cognitive returns to education but may have increased them. This finding is consistent with Barcellos et al. (2021), who document that the UK’s 1972 compulsory schooling reform reduced disparities in education and qualifications between children from different backgrounds but ultimately increased differences in socioeconomic status.

Our first-stage results document that the complier status and the genetic type correlate substantially. We will now demonstrate that this finding provides a necessary condition for a 2SLS estimate of the  $G_i \times E_i$  coefficient to be problematic in our setting – it may not have a well-defined causal interpretation.

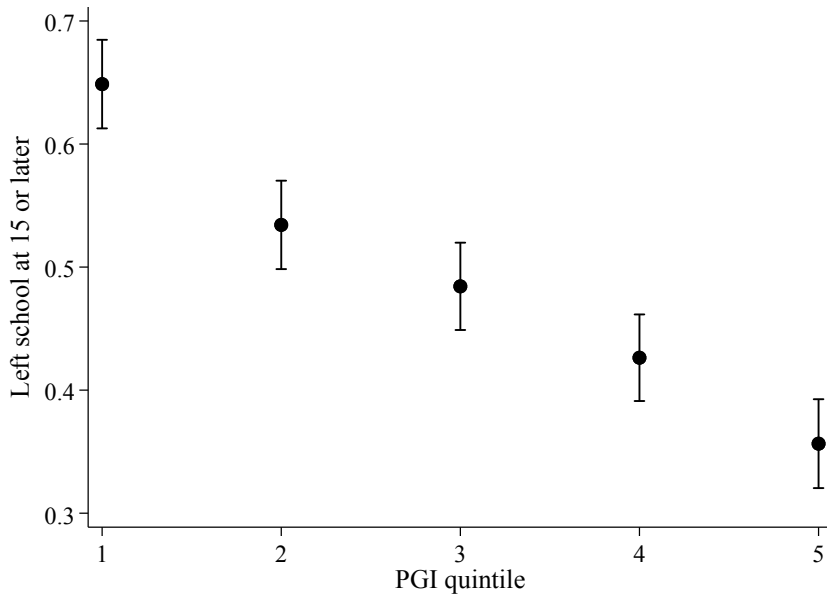


Figure 2: Strength of the first stage by quintiles of the polygenic index

Notes: This figure shows the complier shares by PGI quintile using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. The shares correspond to  $\pi_{1,\Delta}^f$  to  $\pi_{5,\Delta}^f$ , the five estimated first-stage coefficients of Eq. (4). We add 95% confidence intervals. The point estimates and their standard errors are reported in Table B.3.

## 4 Potential identification problems of interaction effects

### 4.1 The problem

We are interested in the effect of a particular environment or life decision (here, education),  $E_i$ , on an outcome  $Y_i$  (here, old-age recall ability) and how this effect interacts with genetic endowment  $G_i$ . For simplicity, first assume that  $E_i$  and  $G_i$  are binary variables. Each individual has four potential outcomes,  $Y_i^j(G_i)$ ,  $j \in \{0, 1\}$ ,  $G_i \in \{0, 1\}$ , where  $j$  denotes educational status and  $G_i$  the attributes of genetic endowment. For example,  $Y_i^1(0)$  is the

potential outcome with a high educational level ( $E_i = 1$ ) and a low genetic propensity ( $G_i = 0$ ). However, only one of the four is realized and observed by the researcher. The observation rule is

$$\begin{aligned}
Y_i &= E_i \cdot G_i \cdot Y_i^1(1) + E_i \cdot (1 - G_i) \cdot Y_i^1(0) \\
&\quad + (1 - E_i) \cdot G_i \cdot Y_i^0(1) + (1 - E_i) \cdot (1 - G_i) \cdot Y_i^0(0) \\
&= Y_i^0(0) \\
&\quad + (Y_i^1(0) - Y_i^0(0)) E_i \\
&\quad + (Y_i^0(1) - Y_i^0(0)) G_i \\
&\quad + (Y_i^1(1) - Y_i^0(1) - Y_i^1(0) - Y_i^0(0)) G_i \times E_i
\end{aligned}$$

The second equality shows how the observation rule corresponds to the  $G \times E$ -workhorse model in Eq. (1) which, for ease of exposition suppresses that the slope coefficients might be individual-specific, allowing for heterogeneous treatment effects:

$$Y_i = \beta_0 + \beta_1 E_i + \beta_2 G_i + \beta_3 G_i \times E_i + \varepsilon_i \quad (6)$$

The gene-environment interaction effect is calculated as  $Y_i^1(1) - Y_i^0(1) - (Y_i^1(0) - Y_i^0(0))$ , that is, the difference in the effect of  $E_i$  on  $Y_i$  when  $G_i = 1$  (which is  $Y_i^1(1) - Y_i^0(1)$ ) and the effect of  $E_i$  on  $Y_i$  when  $G_i = 0$  (which is  $Y_i^1(0) - Y_i^0(0)$ ).

Assume that  $G_i$  is pre-determined while  $E_i$  is a choice variable and, therefore, endogenous.<sup>8</sup> Further assume that we have a binary instrument  $Z_i$  that fulfills the classic LATE assumptions (Imbens and Angrist, 1994). Expressing the gene-environment regression equation as two separate regressions for  $G_i = 0$  and  $G_i = 1$  yields

$$\begin{aligned}
Y_i &= \beta_0 + \beta_1 E_i + \varepsilon_i & \text{for } G_i = 0 \\
Y_i &= (\beta_0 + \beta_2) + (\beta_1 + \beta_3) E_i + \varepsilon_i & \text{for } G_i = 1
\end{aligned}$$

---

<sup>8</sup>The extension of our framework to an endogenous  $G_i$  entails the same kind of problems. Our proposed solution applies to this case but is not straightforward in applications as it requires an instrumental variable for  $G_i$ . In Schmitz and Westphal (2025), we apply marginal treatment effect (MTE) estimation with two endogenous variables in a different context, namely causal mediation analysis. However, the estimation of interaction effects with two endogenous variables is beyond the scope of this paper.

In Wald notation, separately estimating 2SLS regressions for  $G_i = 0$  and  $G_i = 1$  yields the following estimates:

$$\begin{aligned}\hat{\beta}_1 &= \frac{\mathbb{E}[Y_i|Z_i = 1, G_i = 0] - \mathbb{E}[Y_i|Z_i = 0, G_i = 0]}{\mathbb{E}[E_i|Z_i = 1, G_i = 0] - \mathbb{E}[E_i|Z_i = 0, G_i = 0]} && \text{for } G_i = 0 \\ \hat{\beta}_1 + \hat{\beta}_3 &= \frac{\mathbb{E}[Y_i|Z_i = 1, G_i = 1] - \mathbb{E}[Y_i|Z_i = 0, G_i = 1]}{\mathbb{E}[E_i|Z_i = 1, G_i = 1] - \mathbb{E}[E_i|Z_i = 0, G_i = 1]} && \text{for } G_i = 1\end{aligned}$$

Using the LATE theorem (Imbens and Angrist, 1994) which states that 2SLS estimates are average treatment effects for the compliers, we can rewrite these expressions as:

$$\begin{aligned}\hat{\beta}_1 &= \mathbb{E}[Y_i^1(0) - Y_i^0(0)|C(G_i = 0)] \\ \hat{\beta}_1 + \hat{\beta}_3 &= \mathbb{E}[Y_i^1(1) - Y_i^0(1)|C(G_i = 1)]\end{aligned}$$

where  $C(G_i = 0)$  stands for compliers within the group  $G_i = 0$  and  $C(G_i = 1)$  for compliers within the group  $G_i = 1$ . The mechanics of the LATE require that the group-specific effects ( $\hat{\beta}_1$  and  $\hat{\beta}_1 + \hat{\beta}_3$ ) are average treatment effects for the  $G_i$ -specific compliers. Without further covariates, the joint interaction regression of Eq. (6) (where  $Z_i$  and  $Z_i \times G_i$  are used as instrumental variables for  $E_i$  and  $E_i \times G_i$ ) yields the same results as the two separate estimations and delivers the estimate

$$\hat{\beta}_3 = (\hat{\beta}_1 + \hat{\beta}_3) - \hat{\beta}_1 = \mathbb{E}[Y_i^1(1) - Y_i^0(1)|C(G_i = 1)] - \mathbb{E}[Y_i^1(0) - Y_i^0(0)|C(G_i = 0)] \quad (7)$$

This shows that the 2SLS estimate of the interaction coefficient puts together two effects of two potentially different groups: the effect of  $E_i$  on  $Y_i$  given that  $G_i = 1$  in the group  $C(G_i = 1)$  and the effect of  $E_i$  on  $Y_i$  given that  $G_i = 0$  in the group  $C(G_i = 0)$ . Hence, an estimated interaction effect via 2SLS could come from two sources: actual differences in the effect of  $E_i$  on  $Y_i$  by realization of  $G_i$  and/or differences in these effects between the groups  $C(G_i = 1)$  and  $C(G_i = 0)$ .

Put differently, this implies that a 2SLS estimation does not yield a well-defined interaction effect if two conditions hold simultaneously: First, the individual response to the instrument depends on  $G_i$ . Second, the (individual) treatment effects of  $E_i$  on  $Y_i$  exhibit essential heterogeneity. This occurs when the propensity to take the treatment correlates with the unobserved effect heterogeneity (Heckman et al., 2006). As a result, 2SLS estimates may differ in size from true interaction effects.

A simple simulation model visualizes this potential problem. The model and its parameterization are outlined in Appendix D. Set up as an illustrative example, Figure 3 shows the average effects of  $E_i$  on  $Y_i$  (depending on  $G_i$ ) for four groups in the simulated data. Group 1 on the left are always-takers (AT), irrespective of their realization of  $G_i$ . This is because

their gains from  $E_i$  are so large that they choose more education regardless of  $Z_i$  and  $G_i$ . The example also produces individuals that are always-takers when  $G_i = 1$  but compliers when  $G_i = 0$  (Group 2), compliers when  $G_i = 1$  and never-takers (NT) when  $G_i = 0$  (Group 3) and never-takers, irrespective of  $G_i$  (Group 4). Absent simulated data, many of the effects depicted in Figure 3 are unobserved by the researcher. We sort these four groups on the horizontal axis according to their willingness to take education  $E_i$ . Those on the left are most willing, and those on the right are least willing. The blue triangles show  $E[Y_i^1(1) - Y_i^0(1)]$ , the first part of the interaction effect. The red circles show  $E[Y_i^1(0) - Y_i^0(0)]$ , the second part. Thus, the interaction effect for each group is the difference between their blue triangle and red circle. Our data-generating process is set up so that the interaction effect equals 1.5 for each individual and, consequently, for each group. However, as per Eq. (7), 2SLS calculates it as the difference between the filled blue triangle and the filled red circle, that is  $E[Y_i^1(1) - Y_i^0(1)|C(G_i = 1)] - E[Y_i^1(0) - Y_i^0(0)|C(G_i = 0)] = 0.2 - 1.4 = -1.2$ . Not only is the estimate different in magnitude, but because of how our example is set up, it is even negative while the true interaction effect is positive.

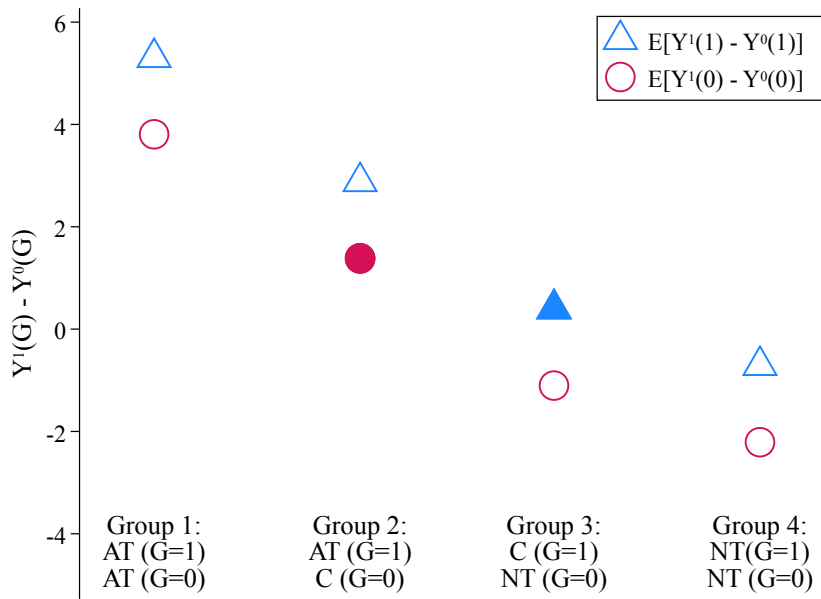


Figure 3: Effects of  $E_i$  on  $Y_i$  by  $G_i$  and complier type in the simulation model

Notes: This figure visualizes stylized potential outcomes from our simulation model. Potential outcomes, their differences and resulting treatment effects are defined by the data-generating process outlined in Appendix D using generated data.

The figure also illustrates the conditions under which 2SLS does not fail: This occurs when either (i) the circles and triangles are on horizontal lines, meaning both complier groups have the same effect of  $E_i$  on  $Y_i$ , and/or (ii)  $G_i$  does not affect the complier status of individuals, meaning groups with  $C(G_i = 1)$  and  $C(G_i = 0)$  do not differ on average. This is the case when  $G_i$  does not affect  $E_i$ . We will elaborate on this point in the following sections.



## 4.2 A solution

We suggest going beyond estimating the two points that form the 2SLS estimate. Instead, we propose to estimate the MTE curve (see, e.g., [Heckman and Vytlacil, 2005](#)) by genetic endowment  $G_i$ . The MTE framework expands the discrete points from Figure 3 to continuous functions on the unit interval. An introduction and more formal account of MTEs are presented in Section 5.1. Exemplary stylized MTE curves from simulated data are shown in Figure 4.

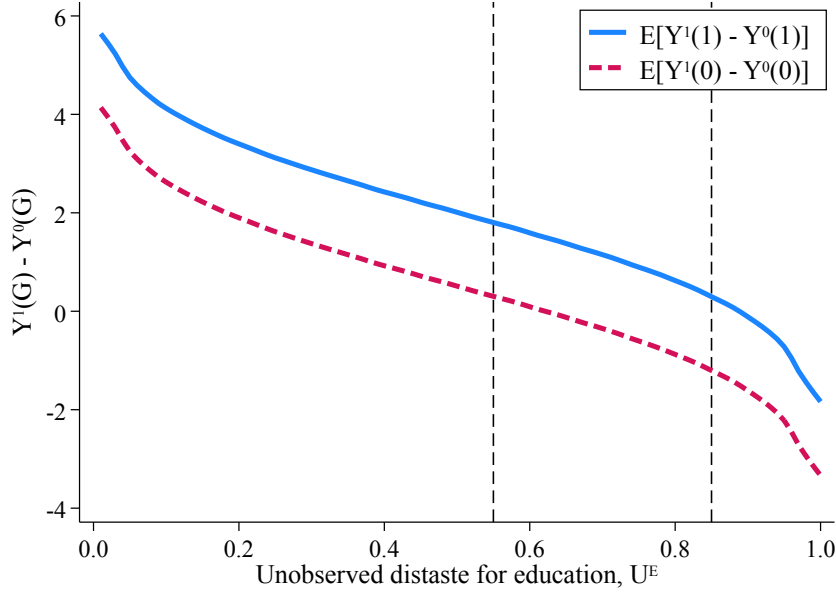


Figure 4: Marginal treatment effects of  $E_i$  on  $Y_i$  by  $G_i$  in the simulation model

*Notes:* This figure shows stylized marginal treatment effect curves in our simulation model. The differences in potential outcomes are defined by the data-generating process outlined in Appendix D using generated data.

Again, the interaction effect is the difference between the blue solid and the red dashed curve. In our simulation example, it is always 1.5, but in practice, the two curves do not need to be parallel. The interaction effect may differ along the x-axis, which, as before, represents the willingness to take education. Applying the MTE framework, this willingness is more precisely the unobserved distaste for education, which we call  $U_i^E$  (see Section 5.1).

There are several ways to translate the two curves into interaction effects reflecting the choice of subsamples for whom the effects are estimated: They can be evaluated at certain points on the x-axis or over intervals that represent the location and shares of different groups (always-takers, compliers, never-takers). For example, one possibility is to compute the difference at a specific value of  $U_i^E$ , say 0.4. The advantage of this method over 2SLS estimation is that unobserved heterogeneity is fixed. This allows us to make correct vertical comparisons of the two lines and yields a consistent, albeit local, estimate of the interaction at  $U_i^E = 0.4$ . MTEs can also be used to estimate all treatment parameters, depending on

how they are aggregated and how MTEs in different areas of the unit interval are weighted. In principle, it is possible to compute interaction effects using the MTE curves with 2SLS weights either for  $C(G_i = 1)$  or  $C(G_i = 0)$ . In our application below, we use a simpler solution. We will aggregate the MTE results to receive the average interaction effect for all individuals on the  $U_i^E$  interval between 0.55 and 0.85, visualized by the two vertical lines in Figure 4. We choose this interval since most of the compliers to the education instrument in our application are located in this area, see Section 5.<sup>9</sup> We do this mainly to maintain comparability to 2SLS/LATE estimates.

When does the problem with 2SLS, outlined in Section 4.1, not occur? First, whenever  $G_i$  does not affect the complier status of individuals such that groups with  $C(G_i = 1)$  and  $C(G_i = 0)$  do not differ on average. This is the case when  $G_i$  does not affect  $E_i$ . In this case – although there may be differences in potential outcomes between groups – there would be only one complier group, irrespective of  $G_i$  and 2SLS would estimate the interaction effect correctly according to Eq. (7). Second, without selection into gains (or losses), that is, when individuals do not self-select into education based on the unobserved gain (loss) from treatment. The simulated data leading to Figure 4 show the case of selection into gains. Those with the highest effects of education on cognition are those with the highest likelihood to take education. This holds both for the blue solid curve ( $G_i = 1$ ) and the red dashed curve ( $G_i = 0$ ). Without this type of selection, both curves would be horizontal and all red circles in Figure 3 would be on a horizontal line, as would all blue triangles. Then, the effects of  $E_i$  on  $Y_i$  would not differ by complier type. Even though 2SLS would still make the wrong comparison (filled blue triangle minus filled red circle), the resulting interaction effect would correspond in size to the correct group-specific interaction effect.

Nevertheless, we believe both conditions that lead to 2SLS estimates failing to represent true interaction effects are likely present in many real-world scenarios. Selection into gains has widely been documented in the context of education (Carneiro et al., 2011, Nybom, 2017, Kamhöfer et al., 2019, Westphal et al., 2022). Moreover, Barcellos et al. (2018) and Barcellos et al. (2021) show differences in first-stage responses to a compulsory schooling reform according to  $G_i$ . Such self-selection into environments according to genetic makeup has long been established in the  $G \times E$  literature as “active gene-environment correlation”, where the environment mediates the effect of genes on the outcome (Biroli et al., 2025; Plomin et al., 1977; Plomin, 2014). Nevertheless, it may not always be the case that an inappropriate 2SLS comparison leads to a significantly different or – as in our simulation above – even reversed sign of the estimate.

In the  $G \times E$  literature, many studies use an exogenous environment, such as a policy change, which we refer to as  $Z_i$ . Hence, these studies estimate reduced-form  $G \times Z$  interactions. Muslimova et al. (2025), for instance, use a firstborn indicator as the measure of

---

<sup>9</sup>We show robustness checks for other intervals in Table 5.

environment, which is not an individual decision, and hence, exogenous to the individual. In other studies, not every individual is affected by a change in  $Z_i$ . Examples include [Schmitz and Conley \(2017\)](#), where  $Z_i$  is the Vietnam draft lottery which provides incentives for education, the implicit  $E_i$ ), [van den Berg et al. \(2023a\)](#), where  $Z_i$  is a vaccination campaign, and the implicit  $E_i$  would be measles infections), [van den Berg et al. \(2023b\)](#), with  $Z_i$  constituting a sugar derationing policy, with the implicit  $E_i$  being the maternal sugar consumption), and [Ahlskog et al. \(2024\)](#), where  $Z_i$  is a compulsory schooling reform shifting education, as in our setting). As in our study,  $Z_i$  constitutes an incentive for the underlying individual decision (or behavior)  $E_i$  we are interested in. While the focus on  $G \times Z$  can be the policy-relevant effect (depending on the context), the focus on the reduced form does not solve problems with essential heterogeneity. As shown in Appendix E, the  $G \times Z$  interaction may be driven solely by a first-stage gradient, even if the (latent structural)  $G \times E$  interaction is absent.

The problem and possible solution we identify have implications that extend beyond the estimation of gene-environment interactions. In theory, they apply to any interaction effect of an endogenous, instrumented treatment with observable characteristics, provided there is essential heterogeneity and a first-stage gradient with respect to the interaction variable. The consequences of 2SLS failing to represent true interaction effects are difficult to generalize. Therefore, Appendix F covers several specific applications from the literature that go beyond the gene-environment setting in which these problems might occur and in which estimating marginal treatment effects could be warranted.

### 4.3 Going beyond a binary representation of $G$

The problem and its solution are not specific to cases where  $G_i$  is binary. On the one hand, our solution requires a discrete  $G_i$  because we will estimate separate curves by  $G_i$ . On the other hand, generating a binary indicator of genetic endowment from a continuous polygenic index entails a loss of information. Recall that in our application, we transform the continuous index into quintiles, i.e., a discrete and ordered measure that takes the values  $g \in \{1, 2, 3, 4, 5\}$ . Consequently, the number of potential outcomes we estimate increases from four to ten. In Table 3, we list these potential outcomes and how to calculate the effect of  $E_i$  on  $Y_i$  and the  $G \times E$  interaction by genetic type, i.e., quintile of the polygenic index. The reference group is the first (lowest) quintile. Accordingly, all interaction effects are calculated in comparison to this quintile. For example, the gene-environment interaction effect of the fifth (highest) quintile is the difference between the effect of  $E_i$  on  $Y_i$  for  $G = 5$  and the effect of  $E_i$  on  $Y_i$  for  $G = 1$ .

Extending the setting to a more complex (but still discrete) classification has advantages. We can make better use of the rich variation of the polygenic index and account for possible

Table 3: Potential outcomes and calculation of MTEs using quintiles of the polygenic index

		$E_i = j$		Individual treatment effects for	
		0	1	the effect of $E_i$ on $Y_i$	the gene-environment interaction
$G_i = g$	1	$Y_i^0(1)$	$Y_i^1(1)$	$Y_i^1(1) - Y_i^0(1)$	$(Y_i^1(1) - Y_i^0(1)) - (Y_i^1(1) - Y_i^0(1))$
	2	$Y_i^0(2)$	$Y_i^1(2)$	$Y_i^1(2) - Y_i^0(2)$	$(Y_i^1(2) - Y_i^0(2)) - (Y_i^1(1) - Y_i^0(1))$
	3	$Y_i^0(3)$	$Y_i^1(3)$	$Y_i^1(3) - Y_i^0(3)$	$(Y_i^1(3) - Y_i^0(3)) - (Y_i^1(1) - Y_i^0(1))$
	4	$Y_i^0(4)$	$Y_i^1(4)$	$Y_i^1(4) - Y_i^0(4)$	$(Y_i^1(4) - Y_i^0(4)) - (Y_i^1(1) - Y_i^0(1))$
	5	$Y_i^0(5)$	$Y_i^1(5)$	$Y_i^1(5) - Y_i^0(5)$	$(Y_i^1(5) - Y_i^0(5)) - (Y_i^1(1) - Y_i^0(1))$

Notes: This table lists all combinations of potential outcomes when  $G_i$  corresponds to quintiles of the PGI such that  $G \in \{1, 2, 3, 4, 5\}$  (left panel). The right panels show how to compute different individual treatment effects, including the interaction effects at every quintile we are after. We chose the first (the lowest) quintile as the reference. All effects are therefore calculated in relation to this group.

nonlinearities in the interaction effects between different sections of the distribution. Of course, the choice to use quintiles is arbitrary. [Barcellos et al. \(2018\)](#) and [Barcellos et al. \(2021\)](#) show differences in their results according to the terciles of the education polygenic index. This is already considerably less restrictive than using a binary representation. The use of quintiles offers a further improvement over terciles. While the general problem is present independent of the binning choice underlying the PGI, the more granular the binning, the more likely it is to detect it by finding a meaningful first-stage gradient or eventual interaction effects. At the same time, it allows us to estimate gene-environment interactions at more points across the polygenic index's distribution, which we can use to detect a possible non-linear evolution of interaction effects across the index. Lastly, using more bins of  $G_i$  increases the identifying variation when estimating MTEs with binary instruments.

## 5 MTE estimation of the $G \times E$ interaction

### 5.1 A brief introduction to MTEs

We start by briefly summarizing the classic MTE framework by [Heckman and Vytlacil \(2005\)](#), adjusted to our notation. See [Heckman and Vytlacil \(2005\)](#) for an extensive introduction to MTEs, their derivation, and traditional ways to estimate them with continuous instruments.

Assume that the potential outcomes of individual  $i$  are defined by the following functions:  $Y_i^j(G_i) = \mu^j(G_i, X_i) + \varepsilon_i^j(G_i)$ ,  $j \in \{0, 1\}$ ,  $G_i \in \{0, 1\}$ , where  $j$  denotes potential outcomes of educational status and  $G_i$  the attributes of genetic endowment.  $\mu^j(G_i, X_i)$  is a function of genetic endowment  $G_i$  and observable characteristics  $X_i$ , and  $\varepsilon_i^j(G_i)$  is an unobservable part.

We model the choice  $E_i$  in a generalized Roy framework (Roy, 1951), where individuals choose  $E_i$  if the (expected) returns to education exceed monetary and/or non-monetary costs  $C_i = \mu^C(G_i, X_i, Z_i) + U_i^C$ . Costs depend on  $G_i$ , the observable characteristics  $X_i$ , an instrumental variable  $Z_i$  and an unobservable term  $U_i^C$ . Note that  $Z_i$  does not directly affect  $Y_i^j(G_i)$ . The decision rule for  $E_i$  (depending on the realization of  $G_i = g$ ) reads:

$$\begin{aligned} E_i(G_i) = 1 &\Leftrightarrow Y_i^1(G_i) - Y_i^0(G_i) > C_i \\ &\Leftrightarrow \mu^1(G_i, X_i) - \mu^0(G_i, X_i) - \mu^C(G_i, X_i, Z_i) > -(\varepsilon_i^1(G_i) - \varepsilon_i^0(G_i) - U_i^C(G_i)) \\ &\Leftrightarrow \mu^E(G_i, X_i, Z_i) > V_i(G_i) \end{aligned}$$

While not necessary for any theoretical result,  $\mu^E(G_i, X_i, Z_i) = \mu^1(G_i, X_i) - \mu^0(G_i, X_i) - \mu^C(G_i, X_i, Z_i)$  can be represented as a linear index, such as:

$$\mu^E(G_i, X_i, Z_i) = \pi_0 + \pi_1 G_i + \pi_2 Z_i + \pi_3 Z_i \cdot G_i + \pi X_i + V_i(G_i)$$

where  $V_i(G_i) = -(\varepsilon_i^1(G_i) - \varepsilon_i^0(G_i) - U_i^C)$  is the unobservable term. The decision rule implies that  $E_i$  correlates with  $\varepsilon_i^1(G_i)$  and  $\varepsilon_i^0(G_i)$  and, thus,  $V_i(G_i)$ , which renders  $E_i$  endogenous. In the spirit of Heckman and Vytlacil (2005) we rewrite the choice equation as:

$$\begin{aligned} E(G_i) &= \mathbb{1}\{\mu^E(G_i, X_i, Z_i) \geq V_i(G_i)\} \\ &= \mathbb{1}\{F_V(\mu^E(G_i, X_i, Z_i)) \geq F_V(V_i(G_i))\} \\ &= \mathbb{1}\{\Pr(V_i(G_i) \leq \mu^E(G_i, X_i, Z_i)) \geq F_V(V_i(G_i))\} \\ &= \mathbb{1}\{\Pr(E(G_i = 1)|X_i, Z_i) \geq U_i^E\} \\ &= \mathbb{1}\{PS(G_i, X_i, Z_i) \geq U_i^E\} \end{aligned}$$

The second step applies a monotonic transformation  $F_V(\cdot)$  — which is the cumulative density of  $V_i(G)$  — to both sides of the inequality.  $F_V(\cdot)$  evaluated at the point  $\mu^E(G_i, X_i, Z_i)$  is defined as  $\Pr(V_i(G) \leq \mu^E(G_i, X_i, Z_i))$  and, referring to the choice equation, the same as  $\Pr(E(G_i = 1)|X_i, Z_i)$ . This choice probability based on observable characteristics is the propensity score, and we abbreviate it by  $PS(G_i, X_i, Z_i)$ . Irrespective of the underlying distribution of  $V_i(G)$ , the unobserved term  $U_i^E$  is uniformly distributed on the unit interval and comprises the unobserved heterogeneity correlating with the decision to take  $E_i$ . Low values of unobserved resistance to more education  $U_i^E$  increase  $PS(G_i, X_i, Z_i)$ , leading to  $E_i = 1$ . This corresponds to high unobserved preferences for  $E_i$ , whereas large values of  $U_i^E$  indicate a high distaste for  $E_i$ .

MTEs are estimates of the causal effect of education on the outcome  $Y_i$  at certain values of  $U_i^E = u$ . That is,  $\mathbb{E}[Y^1(G) - Y^0(G)|U_i^E = u]$ . The MTEs are identified by those individuals who, at  $U_i^E = u$ , are indifferent between choosing  $E_i = 0$  and  $E_i = 1$ . Referring to the choice equation, this is the group for whom the realization  $p$  of the propensity score  $PS(G_i, X_i, Z_i) = p = u$ . For our framework, the quantities  $\mathbb{E}[Y_i^1(G)|U_i^E = u]$  and  $\mathbb{E}[Y_i^0(G)|U_i^E = u]$  are essential (as their difference is the MTE). We follow the literature and call these quantities marginal treatment response curves (MTRs).

Note that manually conditioning on a more narrow  $PS$  range first and then estimating Eq. (6) in this more homogeneous sample only works if  $PS$  (and  $Z$ ) is continuous.<sup>10</sup> In our setting,  $PS$  is discrete at the cutoff. Our total variation in  $PS$  is equivalent to the instrument response types. Within  $G$ , only the binary  $Z$  affects  $PS$ , implying that between  $G$ , it is impossible to make the propensity scores more homogeneous. We have to estimate MTEs.

## 5.2 Options for MTE estimation

There are many different ways to estimate MTRs and MTEs, depending on the underlying data, setting (e.g., continuous or binary instrumental variables) and the assumptions the researcher wants to impose (e.g., functional form assumptions for the MTE, separability between observed and unobserved terms). In our case, with a binary instrument, there are three options.

1. Estimate different expected values of potential outcomes  $\mathbb{E}[Y_i^1(G)|AT]$ ,  $\mathbb{E}[Y_i^1(G)|C]$ ,  $\mathbb{E}[Y_i^0(G)|C]$ , and  $\mathbb{E}[Y_i^0(G)|NT]$  (see [Imbens and Rubin, 1997](#)) for each value of  $G_i$ . Plotted on the  $U_i^E$  unit interval and assuming linearity, we can fit lines through each pair of points, one for treated and one for untreated potential outcomes ([Brinch et al., 2017](#)), which provide the MTRs. The linearity ensures that the lines run through the respective (type-specific) midpoints on the  $U_i^E$  scale. The difference between the two MTR lines is the linear MTE by  $G_i$ , and the four differences between the five  $G_i$ -specific MTEs inform about the interaction effects.
2. Relax the linearity assumption but impose additive separability between controls  $X_i$  and error terms. That is, specify the potential outcomes  $Y_i^j(G_i) = \mu^j(G_i, X_i) + \varepsilon_i^j(G_i)$ , as we have already done above, instead of the more general form  $Y_i^j(G_i) = f^j(G_i, X_i, \varepsilon_i^j)$  with some arbitrary function  $f$ . This allows variation in  $X_i$  to parametrically or semi-parametrically identify the MTEs, since a binary instrument alone cannot provide this ([Brinch et al., 2017](#)).

---

<sup>10</sup>Observables need to be partialled out before a local linear regression can be applied. This is the traditional method of identifying the MTE semiparametrically. Conditioning on  $PS$  and estimating Eq. (6) by 2SLS does not work as one would likely drop observations based on realizations of  $Z$ , rendering the conditional sample endogenous.



3. Allow for a wide range of flexible polynomial shapes of the MTEs and subsequently restrict the shapes. This can be achieved by requiring the curves to reproduce observable sample analogs and imposing further reasonable assumptions derived from theory and the data. The target parameter the researcher aims to identify can be bounded by the two shapes that produce minimum and maximum values (Mogstad et al., 2018).

A linearity assumption is hard to justify a priori. Furthermore, although additive separability is commonly assumed across the entire literature that uses regression models, we do not benefit from it for a semi-parametric identification of the MTEs. This is because we only use a sparse set of control variables that do not add sufficient variation in the propensity score, which would help identify substantially more than the four points from the first approach. Overall, the third approach appears to be the most suitable for our setting. Nevertheless, we first estimate linear MTEs according to Brinch et al. (2017). They help illustrate the setting and show general trends. They are also informative about underlying shape restrictions. For our main approach, we relax this linearity restriction and allow for flexible polynomials.

We begin by estimating type-specific expected outcomes, that correspond to group means in the data:  $\mathbb{E}[Y_i^1(G)|AT]$ ,  $\mathbb{E}[Y_i^1(G)|C]$ ,  $\mathbb{E}[Y_i^0(G)|C]$ ,  $\mathbb{E}[Y_i^0(G)|NT]$ ; as well as the shares of AT, C, and NT for each quintile of the polygenic index. Appendix I presents the details on generating these 35 values by applying the Imbens and Rubin (1997) method. We visualize the 20 means (circles and diamonds, depending on treatment status) as well as the 15 type shares (horizontal lines at the bottom) in Figure 5. Again, we sort the three types according to their willingness to take education on the unit interval. Always-takers have the highest willingness and are located at the left. For example, the share of always-takers in the lowest PGI quintile (lowest horizontal line) is 22 percent. The share of compliers with  $G_i = 1$  is 68 percent. They are located between 0.22 and 0.9 on the  $U_i^E$  unit interval. The remaining 10 percent are never-takers. Following Kowalski (2023), we use the midpoints of the range where each type is located to place the potential outcomes (circles and diamonds) on the x-axis, while the y-axis measures the size of the estimated potential outcomes. The blue circles denote treated potential outcomes  $\mathbb{E}[Y_i^1(G)]$  while the red diamonds denote untreated potential outcomes  $\mathbb{E}[Y_i^0(G)]$ . The numbers next to the markers refer to the realization of  $G_i$ .

The lines through the points produce MTRs under a linearity assumption, which allows us to identify them by the two points. In principle, the lines can be extrapolated to the full unit interval and taking differences between  $\mathbb{E}[Y_i^1(G)|U_i^E = u]$  and  $\mathbb{E}[Y_i^0(G)|U_i^E = u]$  would yield the MTEs by  $G_i$ . The comparison of the resulting five MTEs provides insight into gene-environment interactions. However, as mentioned above, the linearity assumption, which will drive the final results, is hard to defend a priori. Nevertheless, this analysis has

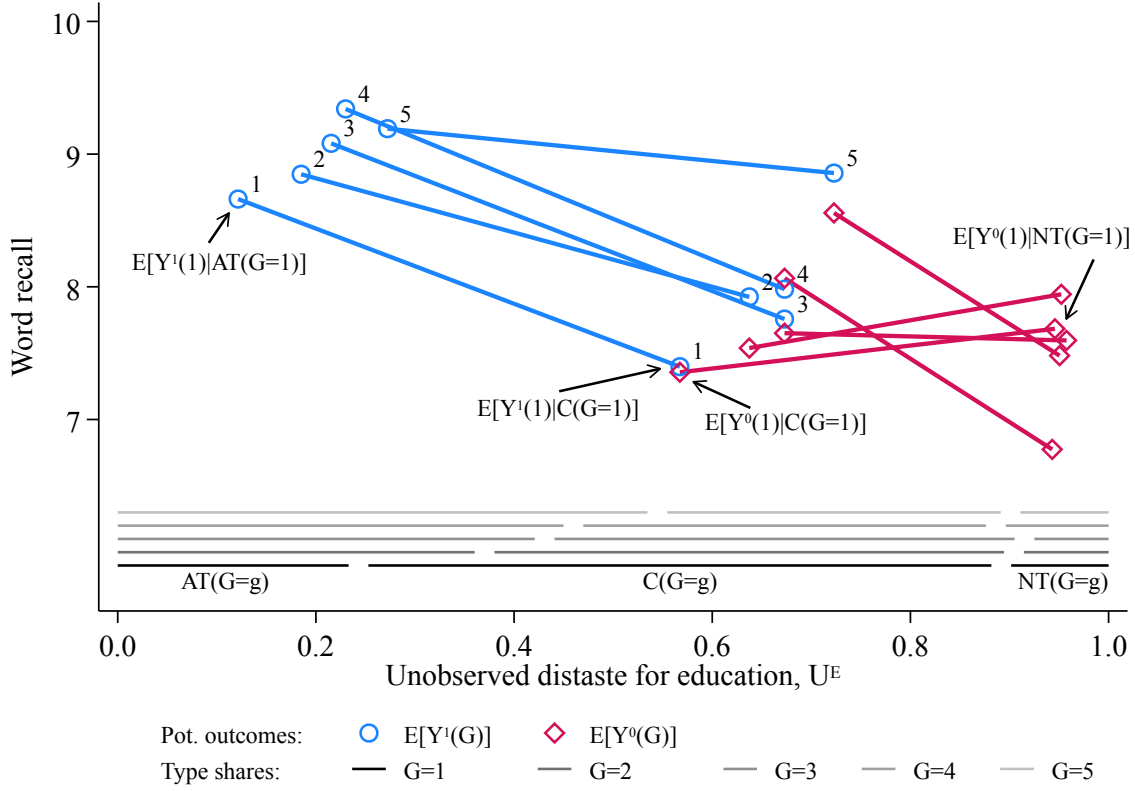


Figure 5: Linear potential outcome curves

*Notes:* This figure shows the 20 estimated potential outcomes  $E[Y_i^j(G)]$  using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. The lines through them represent linear MTRs. Red diamonds refers to potential outcomes for  $E_i = 0$ ; blue circles to  $E_i = 1$ . Thus, for example, the red line labeled “1” shows our estimate of the potential outcome curve of  $Y_i^0(1)$ ; the blue line labeled “3” shows the curve for  $Y_i^1(3)$ . Horizontal lines at the bottom show type shares by quintiles of the educational attainment PGI and their location on the unit interval in ascending order, starting with  $G_i = 1$  represented by the lowest (black) line and  $G_i = 5$  by the highest (lightest) line. We provide detailed descriptions of potential outcomes as text with arrows for  $G_i = 1$  (as an example and to maintain readability).

important implications for our bounding approach to compute our main results (Section 5.3). Treated potential outcomes (blue) are higher for always-takers than compliers, causing the  $E[Y_i^1(G)|U_i^E = u]$ -MTR curves to have a negative slope. Therefore, there seems to be a correlation between types and our dependent variable. Moreover, the  $E[Y_i^1(G)|U_i^E = u]$ -MTR curves are fairly parallel, with no substantial slope differences. They represent level shifts in treated outcomes by  $G_i$ . These results will be used to justify restrictions 3 and 5 in Section 5.3, see below.

The picture is less clear for the untreated outcomes (red diamonds). Here, we see that the outcomes of the untreated compliers are, on average, slightly smaller than those of the treated compliers, suggesting positive effects of education on recall. We can replicate the 2SLS finding of a zero effect for  $C(G_i = 1)$  and positive effects for  $C(G_i = 2)$  and  $C(G_i = 5)$ . However, the estimates for never-takers are less clear, as they are above those for the untreated compliers for the first and second quintiles, resulting in positive slopes of the two lowest  $E[Y_i^0(G)|U_i^E = u]$ -MTRs. Given this ambiguous result and the small share of never-takers in the data, we include a robustness check of our main result where we

estimate MTEs without relying on never-takers in Section 5.6. However, connecting to the argument in [Clark and Royer \(2013\)](#), never-takers are the youngest in their class and, when leaving school. Hence, as it contains valid information, we opted to use never-takers in our main specification.

### 5.3 MTE estimation following Mogstad, Santos, and Torgovitsky (2018)

As our main approach, we now sketch the partial estimation method suggested by [Mogstad et al. \(2018\)](#) and recently applied by [Rose and Shem-Tov \(2021\)](#). It allows a transparent and credible estimation of marginal treatment effects when the instrument is binary or when the variation in the instrument does not sufficiently identify marginal treatment effects over the whole support of the propensity score. This is the third approach mentioned in Section 5.2.

The approach assumes a parametric shape of the MTR, which, however, is extremely flexible. The parameters of the MTR are derived from a linear programming exercise as follows: (i) Among all theoretically possible MTRs, consider only those that fulfill certain restrictions. These restrictions are set by the researcher and we lay them out below. (ii) Among all MTRs that fulfill the restrictions find those that maximize the target parameter. (iii) Among all MTRs that fulfill the restrictions, find those that minimize the target parameter. As a result, we get a bound around the target parameter. Figure 6 presents a stylized and simplified representation of the procedure using simulated data.

Our target parameter is the average difference in the effect of education on recall between  $G_i = 5$  and  $G_i = 1$  on the interval  $U_i^E \in [0.55, 0.85]$ . That is, the interaction effect of  $E_i$  and  $G_i$  when turning from the lowest to the highest  $G$ -quintile. We choose the interval  $[0.55, 0.85]$  because this is the interval that is always covered with compliers from every quintile, see the horizontal lines in Figure 5. While this should approximate the underlying population that determines the 2SLS effects well, we will also assess the robustness of this choice for our main result. To get a sense for the width of this interval, note that by definition, it should cover 30 percent of the overall population. Additionally, the average  $U_i^E$  increases by approximately 0.07, on average, from one PGI quintile to the next. Both facts underscore that this interval encompasses a relevant share of the population.

The flexible parametric form of the MTRs is achieved by specifying them as Bernstein polynomials. The Bernstein polynomials are defined as

$$\mathbb{E}[Y_i^j(G)|U_i^E = u, G_i = g] = \sum_{v=0}^n \theta_v^{jg} \binom{n}{v} u^v (1-u)^{n-v},$$

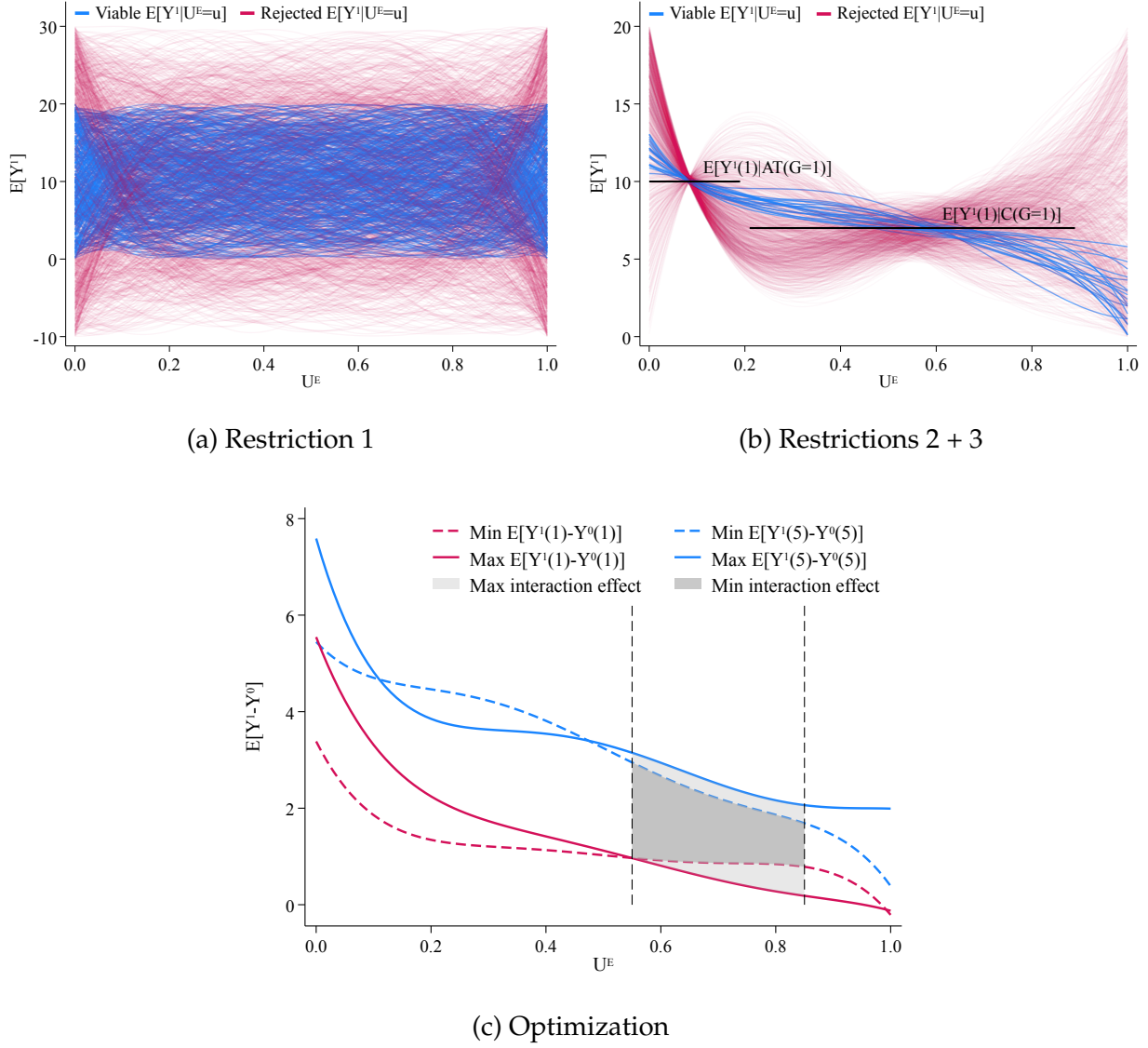
where  $u$  is a specific point on the unit interval,  $j$  refers to the treatment state,  $g$  is the PGI quintile, and  $n$  is the polynomial degree. We choose  $n = 5$ . Therefore, we have  $n + 1 = 6$  parameters  $(\theta_0^{jg}, \dots, \theta_n^{jg})$  that determine each MTR function.

For pedagogical reasons, the stylized representation in Figure 6 and our description of the approach are shown here as an iterative process, whereas in reality, it is a joint linear optimization. In Panels (a) and (b), red curves show those that have been discarded after applying the respective restrictions. Blue curves are the candidates that remain. In Panel (a) of Figure 6, we illustrate numerous different MTR curves for  $\mathbb{E}[Y_1^j(G)|U_i^E = u, G_i = 5]$  to demonstrate that, in principle, the choice of the  $\theta_v^{jg}$  for the Bernstein polynomial can yield virtually any conceivable shape of MTRs. In reality, many more shapes are possible, but we simply show an arbitrarily chosen selection here.

Next, we eliminate all MTRs that do not satisfy Restrictions 1 to 3:

- Restriction 1: All values of  $\mathbb{E}[Y_1^j(G)|U_i^E = u, G]$  and  $\mathbb{E}[Y_0^j(G)|U_i^E = u, G]$  are on the support of  $Y_i$ , that is between 0 and 20. MTRs that in part or completely fall outside that range are discarded.
- Restriction 2: Averaged over the type-specific  $U_i^E$  range, the resulting MTRs reproduce the type-specific outcome means  $\mathbb{E}[Y_i^1(G)|AT]$ ,  $\mathbb{E}[Y_i^1(G)|C]$ ,  $\mathbb{E}[Y_i^0(G)|C]$ , and  $\mathbb{E}[Y_i^0(G)|NT]$ . This means that the average  $y$  of each curve over the appropriate range of the  $x$ -axis must be equal to its respective potential outcome mean (the  $y$ -coordinates from Figure 5). This also implies that they reproduce the LATEs for each  $G_i$ .
- Restriction 3: Monotone treatment selection (see Manski, 1997):  $\mathbb{E}[Y_i^1(G)|U_i^E = u]$  does not increase in  $U_i^E$  along every  $G_i$  quintile. This assumption implies that, on average, individuals with a higher unobserved propensity to pursue more education (i.e., those with a lower  $U_i^E$ ) have the same or higher recall ability, given more education, than individuals with a lower propensity. This extrapolates our finding from Figure 5 that always-takers have higher recall ability than treated compliers to the whole unit interval of  $U_i^E$ . We do not impose a comparable restriction on  $\mathbb{E}[Y_0^j(G)|U_i^E = u, G]$  since the pattern in Figure 5 is not as clear here.

Panel (b) of Figure 6 shows how Restriction 3 is applied to MTRs that remain after implementing Restriction 1 ( $y$ -values between 0 and 20) and 2. The most compelling feature of Restriction 2 is that it introduces a completely data-driven selection of admissible MTRs by forcing them to reproduce outcome means over a specific  $U_i^E$  range. Restriction 3 then picks those that are monotonically decreasing over the entire unit interval (blue). We proceed by using all viable generated MTRs, that is,  $\mathbb{E}[Y_i^0(G)|U_i^E = u, G_i = 1]$  to  $\mathbb{E}[Y_i^1(G)|U_i^E = u, G_i = 5]$ , to compute the respective MTEs, and retain only those that satisfy Restrictions 4 and 5.



**Figure 6: Stylized representation of the constrained optimization of MTE curves**

*Notes:* This figure shows exemplarily the constrained optimization of the MTE curves to match reduced-form evidence and further restrictions. Panel (a) displays an arbitrarily large set of MTR curves (that have the same, highly flexible parametric structure: Bernstein polynomial curves of degree  $n = 5$ ) for  $Y^1$  without additional restrictions. Red lines indicate MTR curves incompatible with the support of our dependent variable (recall score, ranging from 0 to 20). Blue lines indicate remaining curves that lie completely on the support. Panel (b) adds two additional restrictions: Restriction 2 restricts the curve to reproduce the mean recall scores for always-takers, treated compliers, untreated compliers, and never-takers (we term this "reduced-form evidence"). We discard all curves whose average over the type-specific  $U_i^E$  interval does not match the reduced-form evidence. Restriction 3 discards all MTRs for  $Y^1$  with a positive slope at any point on the unit interval. Panel (c) shows the final MTE curves that minimize and maximize the interaction effect after Restrictions 4 and 5 are applied. Note that these are stylized graphs based on simulated data, not data from our empirical setting, in order to properly illustrate each incremental step.

- **Restriction 4: No selection into losses:** The MTE  $\mathbb{E}[Y_i^1(G)|U_i^E = u, G_i = g] - \mathbb{E}[Y_i^0(G)|U_i^E = u, G_i = g]$  is not allowed to increase in  $U_i^E$ . Suppose the treatment is a choice and the outcome is beneficial (or correlates with such a variable). This is likely in our setting with education as treatment and cognition as outcome. In that case, we may expect selection into gains (MTEs decrease in  $U_i^E$ ). The literature on the effect of education on earnings and cognitive skills documents overwhelming empirical evidence of selection into gains (Carneiro et al., 2011; Nybom, 2017; Kamhöfer et al., 2019; Westphal et al., 2022).<sup>11</sup> Note that we allow our MTEs to exhibit no essential heterogeneity (i.e., horizontal MTEs, a setting in which a 2SLS estimation of  $G \times E$  is non-problematic). In Appendix H, we provide suggestive evidence that selection into losses is unlikely in our setting.
- **Restriction 5: Additive Separability of  $G_i$  and the error term.** By specifying linear regression models such as the workhorse model in Eq. (1), this assumption—without justifying it explicitly—is made throughout almost all applied econometric regression analyses, also in the 2SLS estimations we ran before. In the MTE world, this implies that the slope of  $\mathbb{E}[Y_i^1(G)|U_i^E = u, G_i = g]$ ,  $\mathbb{E}[Y_i^0(G)|U_i^E = u, G_i = g]$ , and  $\mathbb{E}[Y_i^1(G)|U_i^E = u, G_i = g] - \mathbb{E}[Y_i^0(G)|U_i^E = u, G_i = g]$  does not depend on  $G_i$ , meaning that MTRs and MTEs for different values of  $G$  are parallel. While possibly a strong assumption, the ordered and parallel-running treated linear potential outcomes with  $E_i = 1$  of Figure 5 suggest that it can be reasonable. For the untreated outcomes with  $E_i = 0$ , in contrast, the picture is less clear. However, the points may still fit curves with the same slope between different values of  $G_i$  if we allow for nonlinearities along  $U_i^E$  (such that MTRs are increasing until  $U_i^E = 0.9$  and decreasing thereafter).

We consider all MTEs that result from Restrictions 1 to 5. Among these, we keep the two that maximize the area between the blue and the red line over the interval  $[0.55, 0.85]$ , and the two that minimize this area. This is visualized in Panel (c) of Figure 6. These are the resulting MTRs and MTEs, from which we derive the bounds on the interaction effect.

More formally, we maximize and minimize our target parameter

$$\begin{aligned} \beta_{G \times E}(0.55, 0.85) := & \\ & \frac{1}{5-1} \int_{0.55}^{0.85} \left[ \mathbb{E}[Y_i^1(G) | U_i^E = u, G_i = 5] - \mathbb{E}[Y_i^0(G) | U_i^E = u, G_i = 5] \right. \\ & \left. - \left( \mathbb{E}[Y_i^1(G) | U_i^E = u, G_i = 1] - \mathbb{E}[Y_i^0(G) | U_i^E = u, G_i = 1] \right) \right] du \end{aligned} \quad (8)$$

---

<sup>11</sup>This does not mean that all MTE applications find selection into gains. In the context of childcare, Cornelissen et al. (2018) find evidence of selection into losses.



over the choice of the parameters  $\theta_v^{jg}$  of the Bernstein polynomials and subject to Restrictions 1 to 5. In total, there are 60 parameters: 6 times 2 (treated and untreated cases) times 5 (different values of  $G_i$ ). Estimating the bounds (i.e., choosing the 60 parameters) involves solving a linear programming problem where constraints on the Bernstein polynomial shapes can be represented as constraints on the parameters  $\theta$  (Rose and Shem-Tov, 2021). The result is a linearized gene-environment interaction effect on the  $U_i^E$ -range that is always covered with compliers from every quintile. The denominator ensures a normalization of the effect to a one-unit increase in  $G_i$ . We optimize the interaction effect of the difference between the first and fifth quintile as the natural choice covering the entire PGI distribution. In Section 5.6, we report robustness checks to show this choice is not crucial.

We make the problem finite and evaluate  $u$  at 20 equidistant grid points (as Rose and Shem-Tov, 2021). Increasing the number of grid points does not affect our results significantly (but increases computation time substantially).

## 5.4 Results

Our main results are visualized in Figure 7. Each panel compares the bounded marginal treatment effects from the first PGI quintile (in red) to the remaining four (in blue). The MTE curves that produce the minimum possible interaction effect are the dashed curves, and the solid curves are MTEs that produce the maximum. Recall that we set up the linear programming approach to optimize the  $G \times E$  effect in the interval  $U_i^E \in [0.55, 0.85]$ . This is because the compliers from all quintiles are located in this range. In this optimization area, the bounds almost coincide, suggesting that the effects are practically point-identified. This tightness inside the optimization area indicates that the reduced-form evidence ( $G_i$ -specific averages for never-takers, always-takers, and compliers) in combination with a highly flexible polynomial and some additional structure (selection into gains, monotone treatment selection for  $E_i = 1$ , and additive separability) almost allow for a perfect interpolation of  $G$ -specific LATEs to the MTE. Note that the tight bounds outside the optimization area are instead a coincidence. The MTEs could look different in this region if the interaction effect were optimized over the whole unit interval. Hence, we only interpret MTE curves and their averages in this region. In Section 5.6, we show that our results are robust to variations of this range. Generally, the differences between the solid MTE curves for quintiles 2–4 and the reference category produce an estimate of the maximal gene-environment interaction effect. The difference between the blue and red dashed curves in each panel yields an estimate of the minimum interaction effect. For example, in the top panel, the area between the blue and red curves indicates how the effect of education on recall changes in the population when  $G_i$  “moves” from the first to the second quintile.

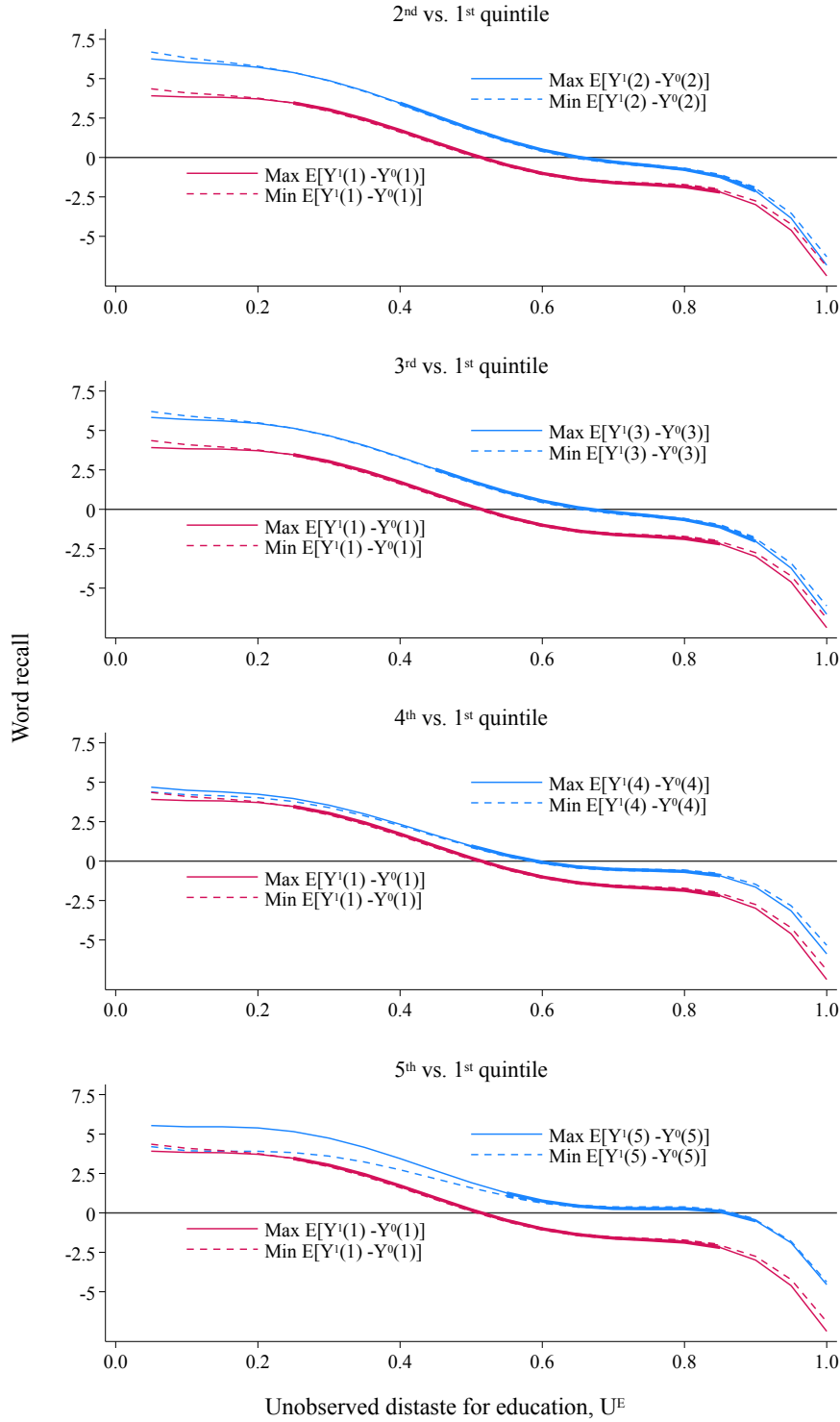
The results have the same sign as our 2SLS estimates. The interaction effect is positive for each quintile comparison. This suggests that individuals with a higher PGI for education benefit more from an additional year of education due to the compulsory schooling reform in terms of their cognition later in life. Our approach also allows us to capture possible nonlinearities in the interaction effect across the PGI. Indeed, the estimated magnitude of the interaction differs across comparisons. Not surprisingly, the highest quintile has the largest interaction effect. However, the size of the interaction for the second quintile is substantial. Those in the third and fourth quintiles have the smallest effects.

We present estimates of the nonlinear interaction effects in Panel A of Table 4. While the previously discussed 2SLS results from Table 2 are reported in column (1) as a benchmark, columns (2) and (3) present the bounds of the marginal treatment effects from Figure 7 aggregated over the  $U_i^E$  range from 0.55 to 0.85. As in Table 2, the effect on  $E_i$  in the first row indicates the baseline effect in the bottom quintile. The direct effects on  $G_i$  in the subsequent rows are not of immediate interest, but we present them for completeness. Our focus is on the interaction effects, which are again relative to the reference category, the bottom quintile. In addition, we present the linearized interaction effect from the quintile coefficients (Panel B, see Eq. 8) that is our main measure of the gene-environment interaction effect.<sup>12</sup> This measure is simply the slope of a line through the interaction effect estimates of the lowest and highest quintiles and can be thought of as the average interaction effect standardized to a one-quintile change.<sup>13</sup> This allows for a comparison of interaction effects from 2SLS and MTE in one number to infer whether unobserved effect heterogeneity and different proportions of compliers in  $G_i$  – which we fix by estimating marginal treatment effects – affected the 2SLS coefficients.

Overall, four features characterize our results. First, the MTE method yields informative and narrow upper and lower bounds of the interaction MTE, which almost point-identify the effect. Second, even the lower-bound MTE results indicate a relevant interaction effect that is substantially larger than 2SLS estimates. The linearized lower bound is about 4.7 times larger than the linearized 2SLS coefficient. While we could not detect significant gene-environment interaction effects with 2SLS, MTE estimation suggests statistically significant effects at the 5 percent level. However, note that the most important difference is the effect sizes and not the statistical significance. The significance level is stricter than necessary because it is based on a two-sided test, even though we are only interested in estimating the true MTE, not in exploring the full range of possible values that it could take. Imbens and Manski (2004) suggests that a one-sided test is sufficient and would lead all interaction coefficients except for  $(E_i \times G_i = 2)$  to shift one significance

<sup>12</sup>The linear slope is calculated as  $(\beta_{5,1}^f - \beta_{1,1}^f)/4$ . The interaction coefficient for the bottom quintile,  $\beta_{1,1}^f$ , is zero since this quintile serves as the reference category.

<sup>13</sup>Note that this measure differs from the linear interaction coefficients in Table 2, Panel C, where we present interactions with the standardized PGI as a continuous variable, which is conceptually different.



**Figure 7: Quintile comparisons of the interaction effect**

*Notes:* This figure shows the four comparisons of gene-environment interactions from our bounding approach using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. For every PGI quintile, we estimate bounds: maxima (solid lines) and minima (dashed lines) at which the interaction effect is maximized/minimized. The bounds for quintiles 2–5 (in blue) are compared to those of the bottom quintile (in red), our reference category, yielding four comparisons. The smallest possible gene-environment interaction is the difference between the blue and red dashed curves over  $U_i^E \in [0.55, 0.85]$ ; the largest possible interaction effect is calculated as the difference between blue and red solid curves over this interval. The thick part of the curves indicates the size of the complier share and its location on the  $U_i^E$  scale, both of which differ by PGI quintile.

Table 4: Estimates of the  $G \times E$  interaction

	Dependent variable – total recall score					
	2SLS (1)		MTE <sub>min</sub> (2)		MTE <sub>max</sub> (3)	
<i>Panel A: nonlinear <math>G \times E</math> effect with <math>G_i</math> as quintiles</i>						
$E_i$	−0.021	(0.449)	0.121	(0.450)	0.121	(0.450)
$G_i = 1$	reference category		reference category		reference category	
$G_i = 2$	0.239	(0.397)	−0.415	(0.463)	−0.415	(0.441)
$G_i = 3$	0.524	(0.451)	−0.576	(0.607)	−0.579	(0.579)
$G_i = 4$	0.784	(0.448)*	−0.249	(0.479)	−0.252	(0.547)
$G_i = 5$	0.871	(0.591)	0.095	(0.822)	0.078	(0.840)
$E_i \times (G_i = 1)$	reference category		reference category		reference category	
$E_i \times (G_i = 2)$	0.314	(0.497)	1.308	(0.582)**	1.342	(0.766)*
$E_i \times (G_i = 3)$	0.091	(0.552)	1.377	(0.666)**	1.418	(0.852)*
$E_i \times (G_i = 4)$	0.049	(0.577)	1.012	(0.637)	1.033	(0.771)
$E_i \times (G_i = 5)$	0.394	(0.698)	1.851	(0.810)**	1.883	(0.912)**
<i>Panel B: linearized <math>G_i \times E_i</math> effect from quintile coefficients</i>						
$E_i \times G_i$	0.098	(0.174)	0.463	(0.203)**	0.471	(0.228)**
Controls	Yes		Yes		Yes	
Observations	11,027		11,027		11,027	

*Notes:* This table presents 2SLS and MTE estimates of the effect of staying in school until at least age 15 ( $E_i$ ), an education PGI ( $G_i$ ), and their gene-environment interaction ( $G \times E$ ) on recall later in life using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. Panel A shows estimates for which we use quintiles of the PGI to estimate possible nonlinear effects across  $G_i$ . Estimates that include  $G_i$  are computed relative to the reference category, the bottom quintile. Panel B shows a linearized slope of a line through the coefficients for  $G = 1$  and  $G = 5$  from Panel A. 2SLS estimates from Table 2 are included for reference in Column (1). The MTE estimates in column (2) refer to the minimal effects where the underlying optimization minimizes the linearized interaction effect. Estimates in column (3) are the maximal effects estimated accordingly. The controls in each case include a linear cohort trend, its interaction with the instrument, gender, survey wave fixed effects, parental education, the first ten principal components of the genetic data, and their interactions with the instrument. Results in different panels are obtained from separate regressions. Standard errors clustered at the individual level are shown in parentheses. For MTE bounds, standard errors are bootstrapped with 100 repetitions. \* $p < 0.1$ , \*\* $p < 0.05$ , and \*\*\* $p < 0.01$ .

level (i.e., gain one star). The MTE standard errors are only slightly larger than those of the 2SLS. This is because our MTE estimation is only slightly more data demanding than a pure 2SLS estimation, our target parameter does not rely on extrapolation to non-complying groups, and we additionally impose the outlined restrictions that reduce sampling variation of the MTE relative to 2SLS.<sup>14</sup> Third, our estimates suggest that the gene-environment interaction is more substantial for individuals with higher PGI, while 2SLS estimates suggest a zero or small and statistically insignificant interaction effect. On average, “moving” to a higher PGI quintile leads to an additional increase of 0.46–0.47 words in the impact of compulsory education on recall due to the education reform. This finding reveals substantial heterogeneity and suggests a high complementarity between

<sup>14</sup>The method of Mogstad and Torgovitsky (2018) only requires estimating twice the number of parameters compared to 2SLS (see Eq. 14) and without restricting the sample. The method does not extrapolate to non-complying individuals (which could inflate the standard errors), because we report MTE-based estimates for an  $U_i^E$  interval that is only covered by compliers in every  $G$  quintile.

education and “nature” as measured by the PGI. Individuals with a higher PGI have higher returns to schooling in terms of cognitive abilities later in life. This result is independent of observable and unobservable factors, both of which we can fix by estimating marginal treatment effects. Fourth, the interaction effect size does not appear linear along the PGI, as indicated by the visual differences in the interaction effects between each panel in Figure 7. The MTE results suggest that individuals in the highest quintile experience a large additional increase in recall of between 1.85 and 1.88 words relative to those in the first quintile. The increases for individuals in the fourth quintile may not be statistically different from the interaction for individuals in the lowest quintile, although the point estimates are also positive and substantial.

To put our results into perspective, they suggest that individuals in the lowest quintile of the education PGI do not experience increased memory later in life. However, compared to them, the additional education increases memory for those in the highest quintile by about half a standard deviation (i.e., 50% of 3.37). Individuals in between experience lower benefits from schooling than those in the top quintile. The average linearized increase per quintile for quintiles 2–5 is 13.6 percent of a standard deviation. We also calculate MTE estimates of the total effect of  $E_i$  on  $Y_i$ , i.e., without interacting with  $G_i$ . The results are reported in Table B.4. The lower-bound MTE estimate is small, positive, and not statistically significant. The upper bound suggest a total effect of 1.64 words, or, again, half of a standard deviation. Conditional on the standard error of the linearized  $E_i \times G_i$  effect (0.203) and detecting a significant effect, we have a “minimal significant effect size” of  $1.96 \times 0.203 = 0.398$ .<sup>15</sup> This effect is 11.8 percent of the standard deviation of the unconditional recall error. This shows that our statistical power is too low to detect minor interaction effects. For example, it would be difficult to detect a significant 2SLS interaction given its small point estimate. Nonetheless, Appendix G documents that our power suffices to detect the average gradient along unobserved heterogeneity as found in the MTE literature.<sup>16</sup>

## 5.5 Discussion

How do our results relate to the previous literature? To begin with, we compare our results to studies that estimate the effects of schooling on cognition measures without considering gene-environment interactions. Related studies find relatively large effects of

<sup>15</sup>This minimal significant effect size differs from the minimal detectable effect size for power calculations because we condition on the test result. The minimal detectable effect size is typically defined as the minimal effect detected as significant at the 5 percent level in 80 percent of all cases. This measure, however, cannot be computed without further assumptions.

<sup>16</sup>A comparison between the interquintile heterogeneity along the PGI in our study and the quantiles of the unobserved (MTE) heterogeneity is appropriate because (i) this heterogeneity is supposed to correlate highly (as uncontrolled genes are supposed to be an important component of the unobserved heterogeneity) and (ii) it is easy to transfer the MTE heterogeneity (reported along quantiles) to a measure along quintiles.

schooling on recall ability. Our maximum MTE estimate (without interactions), as well as the estimated interaction effect comparing individuals at the lowest and highest ends of the education PGI distribution, are within the range of main effects reported in this literature (see the summary of our results above): Using the same setting as in this paper, [Banks and Mazzonna \(2012\)](#) find increases of about half a standard deviation in old-age memory from the additional year of schooling induced by the 1947 UK reform. [Gorman \(2023\)](#) finds increases of one-third to half a standard deviation in memory from the 1972 reform, and [Glymour et al. \(2008\)](#) report about a third of a standard deviation for U.S. compulsory schooling increases. [Carvalho \(2025\)](#), among other things, estimates the effect of the 1972 education reform on fluid intelligence, finding no effect. However, his outcome – answers to several reasoning questions – captures a different aspect of fluid intelligence related to problem-solving, while our outcome measures episodic memory.

Part of the effect of schooling on improved old-age cognitive abilities could arise through higher earnings. However, while [Harmon and Walker \(1995\)](#) and [Oreopoulos \(2006\)](#) find substantial effects of schooling on wages from the 1947 reform, [Devereux and Hart \(2010\)](#) uncover that the impact on earnings is considerably smaller. Taken together, findings in the literature imply that education could have a stronger influence on cognition than on wages. [Banks and Mazzonna \(2012\)](#) discuss channels that include and extend beyond income. They suggest that the effect of education on cognition could also come about via access to more cognitively demanding occupations, enabling greater engagement in cognitively stimulating activities; potentially increased social and cultural participation, or greater productive efficiency in maintaining cognitive health. They rule out effects via physical health improvements and mortality, and note that benefits likely emerged among lower-educated individuals due to diminishing marginal returns to school years and possibly the protective effect of delayed entry into the labor force.

To our knowledge, we are the first to estimate gene-environment interactions on old-age memory. Nevertheless, it is worthwhile to compare interactions of education and genetic markers on different outcomes found in related studies, especially income. Using the 1972 UK increase in the school-leaving age [Barcellos et al. \(2021\)](#) estimate average wage increases of 6-7% and interaction effects of 2% additional gains per standard deviation of an education PGI. When dividing the PGI into terciles, they find no wage effects for individuals with low PGI, and increases of 6-8% for those in the highest PGI tercile. [Ahlskog et al. \(2024\)](#) show that a Swedish schooling reform directly benefited earnings of women with lower education PGI, finding interaction effects of 12% of the average reform effect – however, only for women from wealthy families. Our gene-environment interaction estimates for cognition are larger than those for income, which fits the general notion that effects of education on cognition are larger than wage effects. In a recent addition, [Barcellos et al. \(2025\)](#) estimate a large  $G \times E$  interaction for an Alzheimer's diagnosis. They show that schooling is especially beneficial for people with higher genetic



risk, reducing their likelihood of an Alzheimer’s diagnosis by at least 40% of the pre-reform average.

## 5.6 Robustness

We perform several robustness checks and report our main measure, the linearized MTE bound estimates, in Table 5. As a baseline for reference, we provide our main result from Table 4. First, we estimate the interaction over alternative  $U_i^E$  ranges, in particular  $U_i^E \in [0.6, 0.8]$  and  $U_i^E \in [0.5, 0.9]$ . While the main range  $U_i^E \in [0.55, 0.85]$  covers most compliers from all quintiles well, we show that this choice is not critical to our main results (see Panel B). Both over a wider and a narrower  $U_i^E$  range, the distance between minimum and maximum bounds only marginally varies and remains statistically significant at the 5% level on a two-sided test. We do not go beyond  $U_i^E \in [0.5, 0.9]$  since the never-takers are predominantly located to the right of  $U_i^E = 0.9$ . Second, we show robustness checks for different sample compositions. Our dataset consists of repeated cross-sections (waves) of ELSA and we control for wave fixed effects. Nevertheless, some individuals are observed only once, while most are observed several times. We include a robustness check in which we use only the most recent observation of each individual. This reduces the number of observations (from 11,027 to 3,009) but not the number of individuals in the analysis. Both the upper and lower bound estimates are slightly larger, as are standard errors due to the lower number of observations. Both estimates are statistically significant at the 10% level on the two-sided test we are utilizing. To further demonstrate that the composition of our sample does not significantly alter our results, we include estimates when also considering individuals under the age of 65. The minimum MTE estimate is smaller than our main result, but the maximum is remains similar to our main specification. The choice of the appropriate polynomial to control for cohort trends is not obvious. We demonstrate how the results change when quadratic cohort trends are used instead of linear ones, or when cohort trends are allowed to vary linearly across quintiles of the PGI. Doing so increases flexibility. These changes only marginally affect the possible range of effects, as the estimates are similar to those in our main result. However, standard errors increase, especially when interacting cohort trends with  $G_i$ . With the latter, the minimum and maximum estimates both increase slightly. With squared trends, however, the minimum decreases and the maximum increases somewhat. Next, we show that controlling for principal components of the genetic data – while being sensible and an established norm in the literature – does not drive our main result. Removing them and their interactions with the instrument as control variables leads to slightly lower estimates for both bounds. Next, we remove never-takers from the analysis (see Section 5.2 for a discussion of never-takers). Their presence helps to tighten the bounds. However, even without them, the minimum and maximum MTEs are informative. The interaction effect’s lower and upper

bounds are still positive, although the lower bound may not be statistically different from zero and the upper bound is larger than in our main result. This is to be expected, since never-takers have lower expected outcomes. Not including them in the analysis means that the MTE bounds do not have to reproduce these lower means. As a result, the resulting MTE curves will look different. We visualize the quintile comparisons when computing interaction effects without never-takers in Figure I.3 in the Appendix. Next, we consider the delayed recall score as an alternative outcome variable. This measure counts the number of correctly recalled words (out of ten) five minutes after they are read to participants in studies like ELSA. Since delayed recall is a more difficult task than recalling the words immediately – the second component of our main outcome, total recall – the sample mean for delayed recall is 4.15, less than half the total recall mean. As expected, the estimated MTE effect sizes are much smaller, even though the lower bound is not statistically significant. Furthermore, we estimate our main result, but calculate standard errors with double the number of bootstrap repetitions (200). The standard errors barely change.

Table 5: Robustness

Linearized $G \times E$ effect	Dependent variable – total recall score	
	MTE <sub>min</sub> (1)	MTE <sub>max</sub> (2)
Baseline (main result, Table 4)	0.463 (0.203)**	0.471 (0.228)**
$U_i^E \in [0.6, 0.8]$	0.461 (0.203)**	0.466 (0.232)**
$U_i^E \in [0.5, 0.9]$	0.418 (0.201)**	0.507 (0.217)**
One observation per individual	0.480 (0.248)*	0.496 (0.277)*
Keeping individuals below age 65	0.355 (0.206)*	0.482 (0.201)**
Squared cohort trends	0.444 (0.210)**	0.512 (0.234)**
Interaction of $G_i$ and cohort trends	0.483 (0.281)*	0.514 (0.277)*
No principal components	0.390 (0.201)*	0.444 (0.238)*
No never-takers	0.203 (0.253)	0.861 (0.243)***
Different outcome: delayed recall	0.115 (0.116)	0.353 (0.118)***
200 bootstrap repetitions	0.463 (0.219)**	0.471 (0.224)**

Notes: This table presents robustness checks for the linearized gene-environment estimates (our main result) using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. For reference, we provide our main result from Table 4. Robustness checks include calculating our main estimate (Eq. 8) over larger ranges of  $U_i^E$ , using only the most recent panel observation of each individual, relaxing the age restriction by keeping individuals below age 65, adding squared cohort trends, interacting cohort trends with  $G_i$ , excluding principal components (and their interaction with the instrument) from the control variables, excluding never-takers, and using the delayed recall score (see Section 2.2) as an alternative outcome variable. Unless otherwise specified, standard errors are bootstrapped with 100 repetitions. \* $p < 0.1$ , \*\* $p < 0.05$ , and \*\*\* $p < 0.01$ .

Furthermore, we show robustness checks of our linearized 2SLS measure in Table B.5. Given the sample size, we are somewhat limited in how flexibly we can estimate the 2SLS model. Nevertheless, we show robustness checks in which we add flexibility. First, we additionally interact cohort trends and the interaction of cohort trends and the instrument with  $G_i$  to allow trends to vary by PGI quintile. Doing so increases the linearized 2SLS coefficient of the  $G \times E$  interaction marginally and increases the standard error. Second, we

estimate a fully interacted model in which the controls are the levels and all combinations of interactions between cohort trends, the instrument, PGI quintiles, and the remaining previous controls (gender, principal components of genetic data, parental education, and survey wave fixed effects). The linearized estimate of this flexible model is somewhat closer to zero and has larger standard errors, than our main result. Additionally, we estimate the 2SLS model with lower bandwidths of 8 and 5 years, respectively. The point estimate of the interaction changes marginally and standard errors increase as the bandwidth decreases.

Finally, we show the robustness of our estimation method for the underlying optimization of the interaction effect (see Eq. 8). In our main specification, we maximize/minimize the difference between the first and the fifth quintile. Here, we additionally estimate linearized effects of the final interaction effect when the optimization maximizes/minimizes the difference between the first and each of the other quintiles. We visualize the result in Figure J.1 in the Appendix. The results show that optimizing for different comparisons does not produce substantial changes in the final interaction effect, especially not in the crucial  $U_i^E$  range where we estimate our MTEs. The choice of quintile comparisons for the underlying optimization is not critical. The reason is that the MTRs must reproduce the data in form of group means, and given our shape constraints, the range of possible meaningful candidate MTRs that produce maximum and minimum MTEs is limited.

## 6 Conclusion

The growing gene-environment literature aims to estimate interactions between genetic endowments and environmental exposure (e.g., behavior or choice variables like education) in their effect on an outcome of interest. The goal is to assess whether the effect of the environment varies by genetic endowment (or vice versa) while all else is equal. Since environmental variables are often endogenous, a popular choice is using instruments or natural experiments as a source of exogenous variation. This usually involves estimating a two-stage least squares model. Estimating gene-environment interactions by two-stage least squares regression identifies gene-specific effects of the environment. However, they may not retain the desired interpretation as interaction effects if (1) the first stage is heterogeneous across different values of  $G_i$  and (2) the empirical setting entails essential heterogeneity in  $E_i$  (the unobserved heterogeneities for the outcome and treatment correlate). If both conditions hold, then two properties differ between gene-specific local average treatment effects: the genetic endowment and the unobserved effect heterogeneity. While the former is precisely what researchers want to isolate (the interaction), 2SLS cannot account for the latter. Thus, 2SLS estimates may not reflect complementarity between genes and the environment. We suggest solving this problem by estimating marginal

treatment effects. MTEs allow for the computation of  $G \times E$  estimates while accounting for unobserved heterogeneity.

While gene-environment interactions are a natural choice to illustrate this problem, since the central parameter is the instrumented interaction estimate, it theoretically applies to all interactions estimated by 2SLS. The two conditions that generate it, non-overlapping complier groups due to variations in the interaction variable and unobserved effect heterogeneity correlated with treatment propensity, could be present in other real-world scenarios involving choice variables. Nevertheless, there are likely also many settings where they are not present or the 2SLS comparisons are inconsequential. For example, [Barcellos et al. \(2021\)](#) find no differences between 2SLS and linear MTE estimates of their gene-education interaction. Moreover, in many applications, heterogeneous first stages by the interaction variable are unlikely and studies that estimate only reduced form (gene-environment) interaction effects avoid wrong 2SLS comparisons all together.

Our empirical study examines the long-term effects of education and genetic predisposition for education, as well as their interaction, on memory, our measure of cognition, using data from the English Longitudinal Study of Ageing. Word recall is frequently used as a measure of cognitive functioning and predicts cognitive decline and impairment. To identify the effect of education, we use a compulsory schooling reform from 1947 that increased the minimum school-leaving age in the UK to 15. Our baseline 2SLS estimates document a zero effect of education on recalled words for individuals in the lowest PGI quintile. Effects for higher quintiles are positive, but we lack the precision to estimate them precisely with 2SLS. We find evidence that both conditions for 2SLS to make the wrong comparisons apply in our setting. We see a strong gradient in the first stage across the quintiles of the education PGI and essential heterogeneity is present, more precisely, selection into gains. This is well documented for educational decisions. We estimate marginal treatment effects using the partial identification approach from [Mogstad et al. \(2018\)](#). Building on reduced-form evidence, we generate minimal and maximal  $G \times E$  effects consistent with the data. We add further benign restrictions (such as additive separability and negative MTE slopes that imply selection into gains) to gain precision and tighten the bounds. The resulting bounds almost point-identify the interaction effect.

Our main finding is that, holding unobserved heterogeneity across  $G_i$  fixed, even the lower bound  $G \times E$  effect is 4.7 times larger than the corresponding 2SLS estimate. In absolute terms, the gene-environment complementarity is substantial: on average, the effect of education on recalled words increases by 0.46–0.47 with each PGI quintile. This means that the MTE results imply higher returns to education for cognitive functioning later in life for those with a higher PGI. The complementarity between education and genetic predisposition that widens existing gaps in returns to education is larger than initially

estimated with two-stage least squares. Not accounting for essential heterogeneity limits the usefulness of the 2SLS estimates.

## References

- Agostinelli, F. and Wiswall, M. (2025). Estimating the Technology of Children's Skill Formation. *Journal of Political Economy*.
- Ahlskog, R., Beauchamp, J., Okbay, A., Oskarsson, S., and Thom, K. (2024). Testing for treatment effect heterogeneity: Educational reform, genetic endowments, and family background.
- Altonji, J. G. and Mansfield, R. K. (2018). Estimating Group Effects Using Averages of Observables to Control for Sorting on Unobservables: School and Neighborhood Effects. *American Economic Review*, 108(10):2902–2946.
- Anderson, E. L., Howe, L. D., Wade, K. H., Ben-Shlomo, Y., Hill, W. D., Deary, I. J., Sanderson, E. C., Zheng, J., Korologou-Linden, R., Stergiakouli, E., Davey Smith, G., Davies, N. M., and Hemani, G. (2020). Education, intelligence and Alzheimer's disease: evidence from a multivariable two-sample Mendelian randomization study. *International Journal of Epidemiology*, 49(4):1163–1172.
- Angrist, J. D. and Evans, W. N. (1998). Children and Their Parents' Labor Supply: Evidence from Exogenous Variation in Family Size. *The American Economic Review*, 88(3):450–477.
- Angrist, J. D. and Krueger, A. B. (1991). Does Compulsory School Attendance Affect Schooling and Earnings?\*. *The Quarterly Journal of Economics*, 106(4):979–1014.
- Apolinario, D., Lichtenthaler, D. G., Magaldi, R. M., Soares, A. T., Busse, A. L., das Gracas Amaral, J. R., Jacob-Filho, W., and Brucki, S. M. D. (2016). Using temporal orientation, category fluency, and word recall for detecting cognitive impairment: the 10-point cognitive screener (10-CS). *International Journal of Geriatric Psychiatry*, 31(1):4–12.
- Arold, B. W., Hufe, P., and Stoeckli, M. (2025). Genetic Endowments, Educational Outcomes and the Moderating Influence of School Investments. *Journal of Political Economy: Microeconomics*.
- Banks, J., Batty, G. D., Breedvelt, J., Coughlin, K., Crawford, R., Marmot, M., Nazroo, J., Oldfield, Z., Steel, N., Steptoe, A., Wood, M., and Zaninotto, P. (2023). English Longitudinal Study of Ageing: Waves 0-9, 1998-2019.
- Banks, J. and Mazzone, F. (2012). The Effect of Education on Old Age Cognitive Abilities: Evidence from a Regression Discontinuity Design. *The Economic Journal*, 122(560):418–448.
- Barcellos, S. H., Carvalho, L., Langa, K., Nimmagadda, S., and Turley, P. (2025). Education and Dementia Risk.
- Barcellos, S. H., Carvalho, L., and Turley, P. (2021). The Effect of Education on the Relationship between Genetics, Early-Life Disadvantages, and Later-Life SES. NBER Working Paper No. 28750, National Bureau of Economic Research.
- Barcellos, S. H., Carvalho, L. S., and Turley, P. (2018). Education can reduce health differences related to genetic risk of obesity. *Proceedings of the National Academy of Sciences*, 115(42).
- Barth, D., Papageorge, N. W., and Thom, K. (2020). Genetic Endowments and Wealth Inequality. *Journal of Political Economy*, 128(4):1474–1522.
- Barth, D., Papageorge, N. W., Thom, K., and Velásquez-Giraldo, M. (2022). Genetic Endowments, Income Dynamics, and Wealth Accumulation Over the Lifecycle.
- Behrman, J. R. and Taubman, P. (1989). Is Schooling "Mostly in the Genes"? Nature-Nurture Decomposition Using Data on Relatives. *Journal of Political Economy*, 97(6):1425–1446.

- Belsky, D. W., Domingue, B. W., Wedow, R., Arseneault, L., Boardman, J. D., Caspi, A., Conley, D., Fletcher, J. M., Freese, J., Herd, P., Moffitt, T. E., Poulton, R., Sicinski, K., Wertz, J., and Harris, K. M. (2018). Genetic analysis of social-class mobility in five longitudinal studies. *Proceedings of the National Academy of Sciences*, 115(31):E7275–E7284.
- Biroli, P., Galama, T. J., Hinke, S. v., Kippersluis, H. v., Rietveld, C. A., and Thom, K. (2025). The Economics and Econometrics of Gene-Environment Interplay. *The Review of Economic Studies*, page rdaf034.
- Björklund, A. and Salvanes, K. G. (2011). Education and Family Background. In *Handbook of the Economics of Education*, volume 3, pages 201–247. Elsevier.
- Blundell, R. and Powell, J. L. (2003). Endogeneity in Nonparametric and Semiparametric Regression Models. In Dewatripont, M., Hansen, L. P., and Turnovsky, S. J., editors, *Advances in Economics and Econometrics: Theory and Applications, Eighth World Congress*, Econometric Society Monographs, pages 312–357. Cambridge University Press.
- Bonsang, E., Adam, S., and Perelman, S. (2012). Does retirement affect cognitive functioning? *Journal of Health Economics*, 31(3):490–501.
- Brinch, C. N., Mogstad, M., and Wiswall, M. (2017). Beyond LATE with a Discrete Instrument. *Journal of Political Economy*, 125(4):985–1039.
- Brunello, G., Weber, G., and Weiss, C. T. (2017). Books are Forever: Early Life Conditions, Education and Lifetime Earnings in Europe. *The Economic Journal*, 127(600):271–296.
- Bruno, D., Reiss, P. T., Petkova, E., Sidtis, J. J., and Pomara, N. (2013). Decreased Recall of Primacy Words Predicts Cognitive Decline. *Archives of Clinical Neuropsychology*, 28(2):95–103.
- Cadar, D., Abell, J., Matthews, F. E., Brayne, C., Batty, G. D., Llewellyn, D. J., and Steptoe, A. (2020). Cohort Profile Update: The Harmonised Cognitive Assessment Protocol Sub-study of the English Longitudinal Study of Ageing (ELSA-HCAP). *International Journal of Epidemiology*, 50(3):725–726i.
- Card, D. (1993). Using geographic variation in college proximity to estimate the return to schooling.
- Carneiro, P., Heckman, J. J., and Vytlacil, E. J. (2011). Estimating Marginal Returns to Education. *American Economic Review*, 101(6):2754–2781.
- Carneiro, P. and Lee, S. (2009). Estimating distributions of potential outcomes using local instrumental variables with an application to changes in college enrollment and wage inequality. *Journal of Econometrics*, 149(2):191–208.
- Carvalho, L. S. (2025). Genetics and Socioeconomic Status: Some Preliminary Evidence on Mechanisms. *Journal of Political Economy Microeconomics*. Forthcoming. Available at: <https://www.journals.uchicago.edu/doi/abs/10.1086/732835>.
- Christelis, D., Jappelli, T., and Padula, M. (2010). Cognitive abilities and portfolio choice. *European Economic Review*, 54(1):18–38.
- Clark, D. (2023). School quality and the return to schooling in Britain: New evidence from a large-scale compulsory schooling reform. *Journal of Public Economics*, 223:104902.
- Clark, D. and Royer, H. (2013). The Effect of Education on Adult Mortality and Health: Evidence from Britain. *American Economic Review*, 103(6):2087–2120.
- Conti, G., Heckman, J., and Urzua, S. (2010). The Education-Health Gradient. *American Economic Review*, 100(2):234–238.
- Cornelissen, T., Dustmann, C., Raute, A., and Schönberg, U. (2018). Who benefits from universal child care? estimating marginal returns to early child care attendance. *Journal of Political Economy*, 126(6):2356–2409.
- Cunha, F. and Heckman, J. J. (2007). The Technology of Skill Formation. *American Economic Review*, 97(2):31–47.



- Currie, J. (2009). Healthy, Wealthy, and Wise: Socioeconomic Status, Poor Health in Childhood, and Human Capital Development. *Journal of Economic Literature*, 47(1):87–122.
- Devereux, P. J. and Hart, R. A. (2010). Forced to be Rich? Returns to Compulsory Schooling in Britain. *The Economic Journal*, 120(549):1345–1364.
- Ding, X., Barban, N., Tropf, F. C., and Mills, M. C. (2019). The relationship between cognitive decline and a genetic predictor of educational attainment. *Social Science & Medicine*, 239:112549.
- Gathmann, C., Jürges, H., and Reinhold, S. (2015). Compulsory schooling reforms, education and mortality in twentieth century Europe. *Social Science & Medicine*, 127:74–82.
- Glymour, M. M., Kawachi, I., Jencks, C. S., and Berkman, L. F. (2008). Does childhood schooling affect old age memory or mental status? Using state schooling laws as natural experiments. *Journal of Epidemiology & Community Health*, 62(6):532–537.
- Gorman, E. (2023). Does schooling have lasting effects on cognitive function? evidence from compulsory schooling laws. *Demography*, 60(4):1139–1161.
- Gottfredson, L. S. (1997). Mainstream science on intelligence: An editorial with 52 signatories, history, and bibliography. *Intelligence*, 24(1):13–23.
- Harmon, C. and Walker, I. (1995). Estimates of the Economic Return to Schooling for the United Kingdom. *The American Economic Review*, 85(5):1278–1286.
- Heckman, J. J., Urzua, S., and Vytlacil, E. (2006). Understanding Instrumental Variables in Models with Essential Heterogeneity. *The Review of Economics and Statistics*, 88(3):389–432.
- Heckman, J. J. and Vytlacil, E. (2005). Structural Equations, Treatment Effects, and Econometric Policy Evaluation1. *Econometrica*, 73(3):669–738.
- Houmark, M. A., Ronda, V., and Rosholm, M. (2024). The Nurture of Nature and the Nature of Nurture: How Genes and Investments Interact in the Formation of Skills. *American Economic Review*, 114(2):385–425.
- Imbens, G. W. and Angrist, J. D. (1994). Identification and Estimation of Local Average Treatment Effects. *Econometrica*, 62(2):467.
- Imbens, G. W. and Manski, C. F. (2004). Confidence Intervals for Partially Identified Parameters. *Econometrica*, 72(6):1845–1857.
- Imbens, G. W. and Newey, W. K. (2009). Identification and Estimation of Triangular Simultaneous Equations Models without Additivity. *Econometrica*, 77(5):1481–1512.
- Imbens, G. W. and Rubin, D. B. (1997). Estimating Outcome Distributions for Compliers in Instrumental Variables Models. *The Review of Economic Studies*, 64(4):555–574.
- Ito, K., Ida, T., and Tanaka, M. (2023). Selection on Welfare Gains: Experimental Evidence from Electricity Plan Choice. *American Economic Review*, 113(11):2937–2973.
- Jeong, Y., Papageorge, N. W., Skira, M., and Thom, K. (2024). Genetic Risk for Alzheimer’s Disease and Related Dementias: Cognition, Economic Behavior, and Actionable Information.
- Johnston, D. W., Lordan, G., Shields, M. A., and Suziedelyte, A. (2015). Education and health knowledge: Evidence from UK compulsory schooling reform. *Social Science & Medicine*, 127:92–100.
- Jürges, H., Kruk, E., and Reinhold, S. (2013). The effect of compulsory schooling on health—evidence from biomarkers. *Journal of Population Economics*, 26(2):645–672.
- Kamhöfer, D. A., Schmitz, H., and Westphal, M. (2019). Heterogeneity in Marginal Non-Monetary Returns to Higher Education. *Journal of the European Economic Association*, 17(1):205–244.
- Kline, P. and Walters, C. R. (2019). On Heckits, LATE, and Numerical Equivalence. *Econometrica*, 87(2):677–696.

- Kowalski, A. E. (2023). Reconciling Seemingly Contradictory Results from the Oregon Health Insurance Experiment and the Massachusetts Health Reform. *The Review of Economics and Statistics*, 105(3):646–664.
- Lee, J. J., Wedow, R., Okbay, A., Kong, E., Maghzian, O., Zacher, M., Nguyen-Viet, T. A., Bowers, P., Sidorenko, J., Karlsson Linnér, R., Fontana, M. A., Kundu, T., Lee, C., Li, H., Li, R., Royer, R., Timshel, P. N., Walters, R. K., Willoughby, E. A., Yengo, L., Alver, M., Bao, Y., Clark, D. W., Day, F. R., Furlotte, N. A., Joshi, P. K., Kemper, K. E., Kleinman, A., Langenberg, C., Mägi, R., Trampush, J. W., Verma, S. S., Wu, Y., Lam, M., Zhao, J. H., Zheng, Z., Boardman, J. D., Campbell, H., Freese, J., Harris, K. M., Hayward, C., Herd, P., Kumari, M., Lencz, T., Luan, J., Malhotra, A. K., Metspalu, A., Milani, L., Ong, K. K., Perry, J. R. B., Porteous, D. J., Ritchie, M. D., Smart, M. C., Smith, B. H., Tung, J. Y., Wareham, N. J., Wilson, J. F., Beauchamp, J. P., Conley, D. C., Esko, T., Lehrer, S. F., Magnusson, P. K. E., Oskarsson, S., Pers, T. H., Robinson, M. R., Thom, K., Watson, C., Chabris, C. F., Meyer, M. N., Laibson, D. I., Yang, J., Johannesson, M., Koellinger, P. D., Turley, P., Visscher, P. M., Benjamin, D. J., and Cesarini, D. (2018). Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals. *Nature Genetics*, 50(8):1112–1121.
- Maestas, N., Mullen, K. J., and Strand, A. (2013). Does Disability Insurance Receipt Discourage Work? Using Examiner Assignment to Estimate Causal Effects of SSDI Receipt. *American Economic Review*, 103(5):1797–1829.
- Manski, C. F. (1997). Monotone Treatment Response. *Econometrica*, 65(6):1311–1334.
- Mogstad, M., Santos, A., and Torgovitsky, A. (2018). Using Instrumental Variables for Inference About Policy Relevant Treatment Parameters. *Econometrica*, 86(5):1589–1619.
- Mogstad, M. and Torgovitsky, A. (2018). Identification and Extrapolation of Causal Effects with Instrumental Variables. *Annual Review of Economics*, 10(1):577–613.
- Mogstad, M. and Torgovitsky, A. (2024). Instrumental variables with unobserved heterogeneity in treatment effects. In *Handbook of Labor Economics*, volume 5, pages 1–114. Elsevier.
- Muslimova, D., Van Kippersluis, H., Rietveld, C. A., Von Hinke, S., and Meddens, S. F. W. (2025). Gene-Environment Complementarity in Educational Attainment. *Journal of Labor Economics*. Forthcoming. Available at: <https://www.journals.uchicago.edu/doi/10.1086/734087>.
- Nybom, M. (2017). The Distribution of Lifetime Earnings Returns to College. *Journal of Labor Economics*, 35(4):903–952.
- Oreopoulos, P. (2006). Estimating Average and Local Average Treatment Effects of Education when Compulsory Schooling Laws Really Matter. *American Economic Review*, 96(1):152–175.
- Papageorge, N. W. and Thom, K. (2020). Genes, Education, and Labor Market Outcomes: Evidence from the Health and Retirement Study. *Journal of the European Economic Association*, 18(3):1351–1399.
- Pereira, R. D., Rietveld, C. A., and Kippersluis, H. v. (2022). The Interplay between Maternal Smoking and Genes in Offspring Birth Weight. *Journal of Human Resources*, 58(6).
- Plomin, R. (2014). Genotype-Environment Correlation in the Era of DNA. *Behavior Genetics*, 44(6):629–638.
- Plomin, R., DeFries, J. C., and Loehlin, J. C. (1977). Genotype-environment interaction and correlation in the analysis of human behavior. *Psychological Bulletin*, 84(2):309–322.
- Plug, E. and Vijverberg, W. (2003). Schooling, Family Background, and Adoption: Is It Nature or Is It Nurture? *Journal of Political Economy*, 111(3):611–641.
- Powdthavee, N. (2010). Does Education Reduce the Risk of Hypertension? Estimating the Biomarker Effect of Compulsory Schooling in England. *Journal of Human Capital*,

4(2):173–202.

- Price, A. L., Patterson, N. J., Plenge, R. M., Weinblatt, M. E., Shadick, N. A., and Reich, D. (2006). Principal components analysis corrects for stratification in genome-wide association studies. *Nature Genetics*, 38(8):904–909.
- Rohwedder, S. and Willis, R. J. (2010). Mental Retirement. *Journal of Economic Perspectives*, 24(1):119–138.
- Rose, E. K. and Shem-Tov, Y. (2021). How Does Incarceration Affect Reoffending? Estimating the Dose-Response Function. *Journal of Political Economy*, 129(12):3302–3356.
- Roy, A. D. (1951). Some Thoughts in the Distribution of Earnings. *Oxford Economic Papers*, 3(2):135–146.
- Schiele, V. and Schmitz, H. (2023). Understanding cognitive decline in older ages: The role of health shocks. *European Economic Review*, 151:104320.
- Schmitz, H. and Westphal, M. (2025). Early- and later-life stimulation: how retirement shapes the effect of education on old-age cognitive abilities. Ruhr Economic Papers #1146, RWI - Leibniz Institute for Economic Research.
- Schmitz, L. L. and Conley, D. (2017). The effect of Vietnam-era conscription and genetic potential for educational attainment on schooling outcomes. *Economics of Education Review*, 61:85–97.
- Silles, M. A. (2009). The causal effect of education on health: Evidence from the United Kingdom. *Economics of Education Review*, 28(1):122–128.
- Steptoe, A., Breeze, E., Banks, J., and Nazroo, J. (2013). Cohort Profile: The English Longitudinal Study of Ageing. *International Journal of Epidemiology*, 42(6):1640–1648.
- Tsoi, K. K. F., Chan, J. Y. C., Hirai, H. W., Wong, A., Mok, V. C. T., Lam, L. C. W., Kwok, T. C. Y., and Wong, S. Y. S. (2017). Recall Tests Are Effective to Detect Mild Cognitive Impairment: A Systematic Review and Meta-analysis of 108 Diagnostic Studies. *Journal of the American Medical Directors Association*, 18(9):807.e17–807.e29.
- van den Berg, G. J., von Hinke, S., and Vitt, N. (2023a). Early life exposure to measles and later-life outcomes: Evidence from the introduction of a vaccine. arXiv:2301.10558.
- van den Berg, G. J., von Hinke, S., and Wang, R. A. H. (2023b). Prenatal Sugar Consumption and Late-Life Human Capital and Health: Analyses Based on Postwar Rationing and Polygenic Scores. arXiv:2301.09982.
- Westphal, M., Kamhöfer, D. A., and Schmitz, H. (2022). Marginal College Wage Premiums Under Selection Into Employment. *The Economic Journal*, 132(646):2231–2272.

## A Additional sample information

Table A.1: Descriptive statistics (extended)

	Main sample	By $E_i$		
	Mean (SD)	$E_i=1$	$E_i=0$	Difference (SE)
<i>Outcome <math>Y_i</math></i>				
Recall score	9.67 (3.37)	10.11	8.08	−2.03 (0.07)***
<i>Treatment <math>E_i</math></i>				
Left school $\geq 15$	0.78 (0.41)	1.00	0.00	−1.00 (0.00)
<i>Polygenic index <math>G_i</math></i>				
1st PGI quintile	0.20 (0.40)	0.18	0.25	0.07 (0.01)***
2nd PGI quintile	0.19 (0.40)	0.19	0.21	0.02 (0.01)**
3rd PGI quintile	0.20 (0.40)	0.21	0.19	−0.02 (0.01)**
4th PGI quintile	0.21 (0.41)	0.21	0.20	−0.01 (0.01)
5th PGI quintile	0.20 (0.40)	0.22	0.15	−0.07 (0.01)***
<i>Instrument <math>Z_i</math></i>				
Born 1933 or later	0.66 (0.47)	0.82	0.13	−0.69 (0.01)***
<i>Controls</i>				
Female	0.52 (0.50)	0.52	0.50	−0.02 (0.01)**
Birth year	1934.89 (5.00)	1936.29	1929.92	−6.37 (0.10)***
Parental education:				
Missing	0.25 (0.43)	0.20	0.41	0.21 (0.01)***
Both left school $\leq 14$	0.57 (0.49)	0.58	0.55	−0.03 (0.01)**
At least one left school $\geq 15$	0.18 (0.39)	0.22	0.04	−0.18 (0.01)***
Principal components (standardized):				
− 1 −	0.00 (1.00)	0.00	−0.01	−0.02 (0.02)
− 2 −	0.00 (1.00)	0.01	−0.02	−0.03 (0.02)
− 3 −	0.00 (1.00)	0.01	−0.04	−0.05 (0.02)**
− 4 −	0.00 (1.00)	−0.01	0.02	0.03 (0.02)
− 5 −	0.00 (1.00)	0.00	0.00	0.00 (0.02)
− 6 −	0.00 (1.00)	0.02	−0.07	−0.09 (0.02)***
− 7 −	0.00 (1.00)	0.01	−0.03	−0.04 (0.02)*
− 8 −	0.00 (1.00)	0.00	0.02	0.02 (0.02)
− 9 −	0.00 (1.00)	0.01	−0.02	−0.02 (0.02)
− 10 −	0.00 (1.00)	0.01	−0.02	−0.02 (0.02)
Wave:				
− 1 −	0.15 (0.36)	0.12	0.27	0.15 (0.01)***
− 2 −	0.18 (0.38)	0.16	0.26	0.10 (0.01)***
− 3 −	0.17 (0.37)	0.16	0.19	0.03 (0.01)***
− 4 −	0.19 (0.39)	0.20	0.15	−0.06 (0.01)***
− 5 −	0.17 (0.37)	0.19	0.09	−0.10 (0.01)***
− 6 −	0.14 (0.35)	0.17	0.04	−0.13 (0.01)***
Observations	11,027	8,590	2,437	

Notes: This table presents extended descriptive statistics using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. Here, we also report the standardized first 10 principal components of the genetic data and information on survey waves. The categories for parental education include: Missing information of at least one parent, both parents left full-time education at age 14 or before or have no education, and at least one parent stayed in school until age 15 or longer. We include the mean and standard deviation of the main sample as well as the means by  $E_i$ , the difference of means, and the standard errors of a t-test for equality of means. \* $p < 0.1$ , \*\* $p < 0.05$ , and \*\*\* $p < 0.01$ .

Table A.2: Descriptive statistics by availability of genetic information

	Full sample	By availability of genetic information		
	Mean (SD)	Yes	No	Difference (SE)
<i>Outcome <math>Y_i</math></i>				
Recall score	9.32 (3.50)	9.67	8.76	−0.91 (0.05)***
<i>Treatment <math>E_i</math></i>				
Left school $\geq 15$	0.76 (0.43)	0.78	0.72	−0.06 (0.01)***
<i>Instrument <math>Z_i</math></i>				
Born 1933 or later	0.65 (0.48)	0.66	0.62	−0.04 (0.01)***
<i>Controls</i>				
Female	0.52 (0.50)	0.52	0.52	0.00 (0.01)
Birth year	1934.67 (5.10)	1934.89	1934.32	−0.57 (0.08)***
Parental education:				
Missing	0.31 (0.46)	0.25	0.40	0.16 (0.01)***
Both left school $\leq 14$	0.53 (0.50)	0.57	0.45	−0.12 (0.01)***
At least one left school $\geq 15$	0.17 (0.37)	0.18	0.14	−0.04 (0.01)***
Wave:				
– 1 –	0.19 (0.39)	0.15	0.24	0.09 (0.01)***
– 2 –	0.17 (0.38)	0.18	0.15	−0.03 (0.01)***
– 3 –	0.15 (0.36)	0.17	0.13	−0.04 (0.01)***
– 4 –	0.19 (0.39)	0.19	0.18	−0.01 (0.01)
– 5 –	0.16 (0.37)	0.17	0.16	−0.01 (0.01)
– 6 –	0.14 (0.34)	0.14	0.13	−0.01 (0.01)*
Observations	17,884	11,027	6,857	

Notes: This table presents descriptive statistics by availability of genetic information in ELSA using data from waves 1–6 and different sample restrictions. “Full sample” includes all restrictions outlined in Chapter 2.2, except for the removal of individuals without genetic information. In the second part of this table, we split this sample based on the availability of genetic information. The sub-sample with genetic information corresponds to our main estimation sample. We display the means and difference of means, as well as the standard errors of a t-test for equality of means between the two groups.

\*  $p < 0.1$ , \*\*  $p < 0.05$ , and \*\*\*  $p < 0.01$ .

Table A.3: Descriptive statistics by PGI quintiles

	1st quintile	2nd quintile	3rd quintile	4th quintile	5th quintile
<i>Outcome <math>Y_i</math></i>					
Recall score	8.98	9.47	9.69	9.84	10.31
<i>Treatment <math>E_i</math></i>					
Left school $\geq 15$	0.72	0.76	0.80	0.78	0.84
<i>Instrument <math>Z_i</math></i>					
Born 1933 or later	0.67	0.65	0.66	0.67	0.67
<i>Controls</i>					
Female	0.54	0.51	0.53	0.50	0.50
Birth year	1934.87	1934.84	1935.05	1934.69	1934.98
Parental education:					
Missing	0.27	0.26	0.22	0.23	0.25
Both left school $\leq 14$	0.62	0.57	0.59	0.62	0.47
At least one left school $\geq 15$	0.11	0.17	0.19	0.15	0.29
Wave:					
– 1 –	0.16	0.15	0.14	0.16	0.15
– 2 –	0.19	0.18	0.17	0.19	0.17
– 3 –	0.17	0.17	0.17	0.17	0.17
– 4 –	0.19	0.19	0.20	0.18	0.20
– 5 –	0.16	0.17	0.17	0.16	0.17
– 6 –	0.13	0.14	0.15	0.14	0.15
Observations	2,152	2,145	2,216	2,284	2,230

Notes: This table presents sample means by quintiles of the educational attainment PGI using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2.



## B Additional regression results

Table B.1: The 1947 UK compulsory schooling reform and panel attrition

	DV: Dropped out of sample	
	Coefficient (1)	Standard error (2)
$Z_i$	0.001	(0.018)
Controls	Cohort trends only	
Observations	12,108	

*Notes:* This table presents estimates of the effect of the 1947 UK compulsory schooling reform ( $Z_i$ ) on a panel attrition indicator. The analysis uses data from ELSA waves 1–6, our main sample selection, as outlined in Chapter 2.2, but we fill up observations from the first wave an individual was observed until wave 6 to create the panel attrition indicator. Controls include a linear cohort trend and its interaction with the instrument. Standard errors clustered at the individual level shown are in parentheses. \* $p < 0.1$ , \*\* $p < 0.05$ , and \*\*\* $p < 0.01$ .

Table B.2: The 1947 UK compulsory schooling reform and providing genetic information to ELSA

	DV: Provided genetic information (1)	DV: Left school at 15 or later ( $E_i$ ) (2)
$Z_i$	-0.01 (0.030)	0.479 (0.030)***
Provided genetic information		0.035 (0.023)
Provided genetic information $\times Z_i$		-0.008 (0.025)
Controls	Yes	Yes
Observations	17,884	17,884

*Notes:* This table shows that our instrument, eligibility for the 1947 UK compulsory schooling reform, did not affect the probability of providing genetic information to ELSA (column 1) and that the first stage does not vary with the provision of genetic information to ELSA (column 2) using data from ELSA waves 1–6 and the sample selection outlined in Chapter 2.2, except for the removal of individuals without genetic information. Specifically, column 1 shows the estimates of a linear regression of the instrument  $Z_i$  on a dummy variable equal to one if a person provided genetic information to ELSA, and column 2 shows the estimates of a linear regression of the environment (leaving school at age 15 or later) on the instrument  $Z_i$ , being born in 1933 or later, the genetic information dummy, and the interaction between the two. The controls in each case include a linear cohort trend, its interaction with the instrument, gender, and survey wave fixed effects. Both regressions are estimated using a larger sample from Table A.2, where we apply all the sample restrictions outlined in Section 2.2, except for removing individuals without genetic information. Standard errors in both regressions are clustered at the individual level. \* $p < 0.1$ , \*\* $p < 0.05$ , and \*\*\* $p < 0.01$ .

Table B.3: Estimates of the first stages by PGI quintile

	DV: Left school at 15 or later ( $E_i$ )	
	Coefficient (1)	Standard error (2)
$Z_i \times (G_i = 1)$	0.649	(0.018)***
$Z_i \times (G_i = 2)$	0.534	(0.018)***
$Z_i \times (G_i = 3)$	0.484	(0.018)***
$Z_i \times (G_i = 4)$	0.426	(0.018)***
$Z_i \times (G_i = 5)$	0.357	(0.018)***
Controls	Yes	
Observations	11,027	

Notes: This table presents estimates of the effect of the 1947 UK compulsory schooling reform ( $Z_i$ ) on attending school until at least age 15 ( $E_i$ ) by quintiles of the education PGI using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. These effects are obtained from the coefficients  $\pi_{1,\Delta}^f$  to  $\pi_{5,\Delta}^f$  of Eq. (4), which correspond to the complier shares in the respective quintile. Standard errors clustered at the individual level shown are in parentheses. The controls include a linear cohort trend, its interaction with the instrument, gender, survey wave fixed effects, parental education, the first ten principal components of the genetic data as well as their interactions with the instrument. \* $p < 0.1$ , \*\* $p < 0.05$ , and \*\*\* $p < 0.01$ . The standard errors are seemingly the same here, but only due to rounding. They differ at the fourth digit.

Table B.4: MTE estimation without interaction terms

	DV: Total recall score	
	MTE <sub>min</sub> (1)	MTE <sub>max</sub> (2)
$E_i$	0.163 (0.464)	1.636 (0.620)***
Controls	Yes	Yes
Observations	11,027	11,027

*Notes:* This table shows MTE estimates of the total effect of  $E_i$  on recall, i.e., without interacting  $E_i$  with the genetic endowment. The analysis uses data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. We use our method as described in Chapter 5, but only construct two MTE curves, that maximize/minimize the total effect instead of two for each PGI quintile. Controls include a linear cohort trend, its interaction with the instrument, gender, and survey wave fixed effects. Standard errors in both regressions are bootstrapped with 100 repetitions. \* $p < 0.1$ , \*\* $p < 0.05$ , and \*\*\* $p < 0.01$ .

Table B.5: Robustness of the linearized 2SLS estimate

	DV: Total recall score	
	Linearized $G \times E$ coefficient (1)	Standard error (2)
Baseline (main result, Table 4)	0.098	(0.174)
Birth cohort interacted with $G_i$	0.121	(0.425)
Fully interacted	0.044	(0.410)
Bandwidth 8 years	0.128	(0.191)
Bandwidth 5 years	0.121	(0.245)
Controls	Yes	
Observations	11,027	

*Notes:* This table presents robustness checks of the 2SLS estimation using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. We only present our main 2SLS result, the linearized  $G \times E$  estimate, that represents a line through the  $G \times E$  coefficients of the lowest and highest PGI quintile. This is done to compare average effects easily across methods. For reference, we report the effect from our main results. For “Birth cohort interacted with  $G_i$ ”, we add interactions between  $G_i$  and cohort trend  $t$  as well as  $G_i \times t \times Z_i$  as controls on top of the baseline specification ( $Z_i$  being the instrument). For the fully interacted model, the controls include all baselines and interactions between  $t$ ,  $G_i$ ,  $Z_i$ , and  $X$ , where  $X$  is a vector of controls that includes gender, the principal components of the genetic data, parental education and survey wave fixed effects. Standard errors clustered at the individual level shown are in parentheses.  $*p < 0.1$ ,  $**p < 0.05$ , and  $***p < 0.01$ .

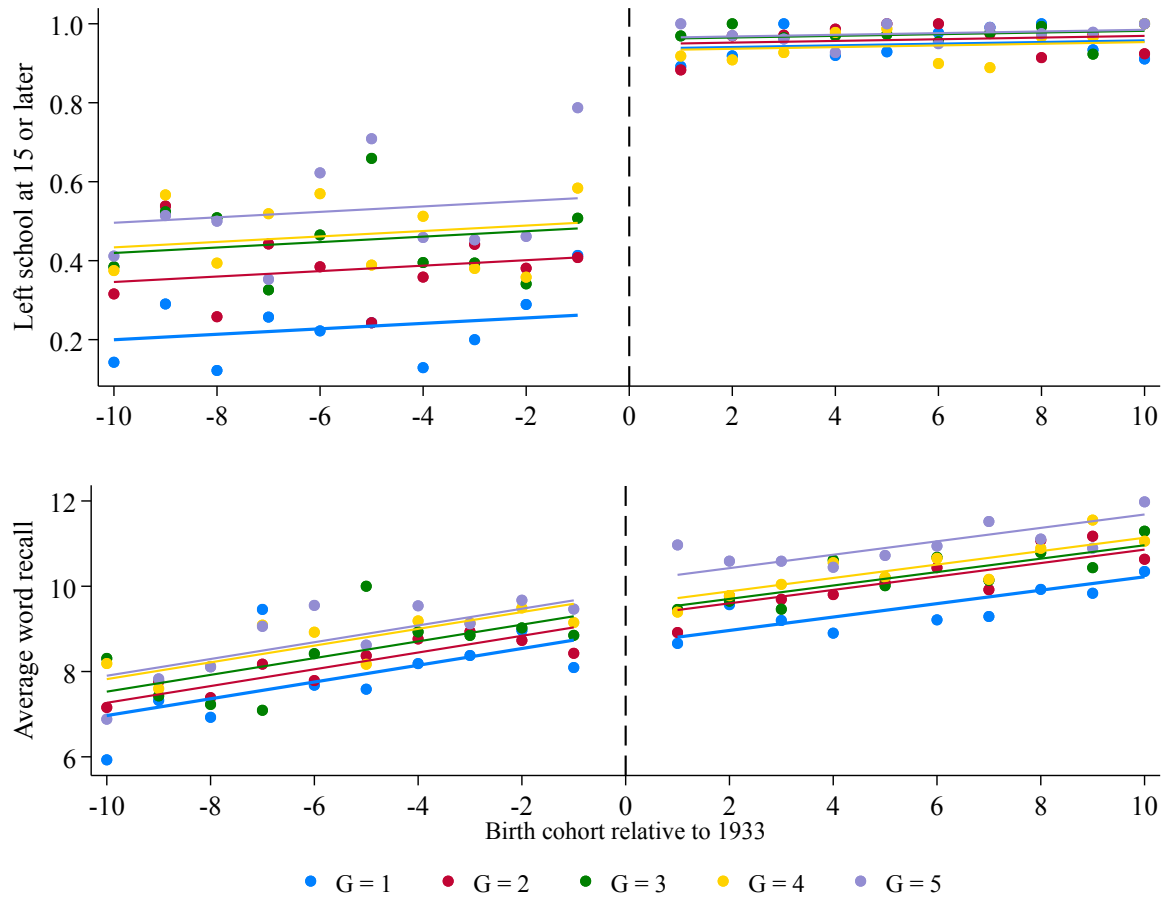


Figure B.1: First Stage and reduced form by  $G_i$

*Notes:* This figure shows a visualization of the first stage (upper panel) and the reduced form (lower panel) results of our regression discontinuity design by quintiles of the education PGI  $G_i$  using our main sample collapsed to cohort  $\times G_i$ -averages. Dots correspond to  $G_i$ -specific sample means, lines are linear fits that are allowed to vary by  $Z_i$ .



Table B.6: Estimates corresponding to Figure B.1

	Dependent variable:	
	Cohort average left school at 15 or later ( $E_i$ ) (1)	Cohort average word recall ( $Y_i$ ) (2)
$G_i = 1$	0.269 (0.032)***	8.933 (0.205)***
$G_i = 2$	0.415 (0.032)***	9.231 (0.205)***
$G_i = 3$	0.489 (0.032)***	9.493 (0.205)***
$G_i = 4$	0.503 (0.032)***	9.786 (0.205)***
$G_i = 5$	0.565 (0.032)***	9.866 (0.205)***
$Z_i \times (G_i = 1)$	0.667 (0.045)***	-0.284 (0.290)
$Z_i \times (G_i = 2)$	0.533 (0.045)***	0.054 (0.290)
$Z_i \times (G_i = 3)$	0.472 (0.045)***	-0.105 (0.290)
$Z_i \times (G_i = 4)$	0.430 (0.045)***	-0.222 (0.290)
$Z_i \times (G_i = 5)$	0.399 (0.045)***	0.244 (0.290)
$t$	0.007 (0.004)*	0.196 (0.024)***
$t \times Z_i$	-0.005 (0.005)	-0.039 (0.034)
Observations	100	100

Notes: This table presents estimates that correspond to averages and linear fits visualized in Figure B.1. This analysis uses our main sample collapsed to cohort  $\times$   $G_i$ -averages, yielding 100 observations (20 birth cohorts  $\times$  5 gene groups). Columns (1) and (2) are OLS regression of the average proportion of the sample who left school at 15 or above ( $E_i$ ) and average recall score ( $Y_i$ ), respectively, on  $G_i$ , interactions between  $G_i$  and the instrument  $Z_i$ , as well as a linear cohort trend  $t$  and the interaction of this trend with  $Z_i$ . We do not include a constant. Therefore, coefficients of  $G_i$  can be interpreted as sample means of the respective outcome variable in each quintile of  $G_i$  when  $Z_i = 0$ , i.e. on the left side of the cutoff. Interactions with  $Z$  are changes of these sample means for  $Z_i = 1$ , i.e. the right side. The coefficients of  $t$  and  $t \times Z_i$  can be interpreted as the slope of the linear cohort trend for  $Z_i = 0$  and its change for  $Z_i = 1$ . \* $p < 0.1$ , \*\* $p < 0.05$ , and \*\*\* $p < 0.01$ .

## C Polygenic indices

The human genome has about 3 billion base pairs, the pairs of nucleic acids that make up the DNA. However, any two people differ by only about 0.1 percent of the base pairs. Most of these genetic differences are substitutions of a single base (adenine, thymine, cytosine, or guanine) for another at a specific location in the genome, called "single nucleotide polymorphisms" (SNPs) that are common across the whole genome. These substitutions result in different genetic variants (alleles) that vary among parts of the population.<sup>17</sup> For example, at a specific SNP location, the DNA sequence might have an adenine base in some individuals, while others may have a thymine base at the same position. One is (arbitrarily) chosen as the reference variant. Then, each SNP can be represented as a count variable of occurrences of the reference variant at this location that can either be 0, 1 or 2, since there are two copies of each chromosome. Large research projects called genome-wide association studies (GWAS) correlate each of the SNPs with a disease or trait, e.g., diabetes, years of education, or smoking. This entails running one regression of type

$$Y_i = \beta_j S_{ij} + X_i' \delta + \zeta_i \quad (9)$$

for each of the SNPs, where  $Y_i$  is the outcome of interest (in our case educational attainment) of individual  $i$ ,  $\beta_j$  is the individual effect of each SNP  $j$ ,  $S_{ij}$  is the count variable of the reference variant of SNP  $j$  with  $S_{ij} \in \{0, 1, 2\}$ ,  $X_i$  is a vector of controls that typically include age, gender and principal components of the genetic data, which control for spurious correlations of genetic variants and outcomes that are due to population structure.<sup>18</sup> The PGI is then calculated as a weighted sum of all  $J$  SNPs that are relevant to the outcome<sup>19</sup>, where the weights correspond to the  $\beta_j$ 's obtained in the GWAS:

$$PGI_i = \sum_{j=1}^J \beta_j S_{ij} \quad (10)$$

Polygenic scores for various traits or behaviors (personality, mental and physical health, health behaviors, and more) have been calculated for the ELSA sample based on various GWAS and are readily available.

---

<sup>17</sup>The generally agreed-upon threshold for a substitution to be regarded a SNP is common occurrence in at least one percent of the population.

<sup>18</sup>For a more detailed description of principal components, see Chapter 2.2.

<sup>19</sup>The discovery study the PGI we are using is based on, [Lee et al. \(2018\)](#), finds 1,271 SNPs that are significantly correlated with educational attainment.

## D Simulation model

To visualize 2SLS being unable to disentangle interaction effects from shifts in complier groups, as discussed in Section 4, we set up a simple simulation model. Assume the following arbitrary parameterizations of the potential outcomes, where, for simplicity, we leave out observable variables  $X_i$ :

$$\begin{aligned} Y_i^1(1) &= 2.3 + \varepsilon_i^1, & Y_i^1(0) &= 0.5 + \varepsilon_i^1, & Y_i^0(1) &= 0.3 + \varepsilon_i^0, & Y_i^0(0) &= 0 + \varepsilon_i^0 \\ E_i &= \mathbb{1}\{0.23 + 2.5G_i - 4Z_i + 3Z_i \times G_i > -(\varepsilon_i^1 - \varepsilon_i^0)\} \\ Z_i, G_i &= \text{Bernoulli distributed with } p = 0.5, \end{aligned}$$

where  $\varepsilon_i^0 = \varepsilon_i^0(G_i)$  and  $\varepsilon_i^1 = \varepsilon_i^1(G_i)$  are error terms. Here, we make the simplifying assumption that  $\varepsilon^1(1) = \varepsilon_i^1(0) = \varepsilon_i^1$  and  $\varepsilon_i^0(1) = \varepsilon_i^0(0) = \varepsilon_i^0$ . In this simulation — but not in the application later — they are assumed to follow a bivariate normal distribution with the following parameters:

$$\begin{pmatrix} \varepsilon_i^1 \\ \varepsilon_i^0 \end{pmatrix} \sim N \left[ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 & 0.4 \\ 0.4 & 2 \end{pmatrix} \right].$$

Setting  $\varepsilon^1(1) = \varepsilon_i^1(0) = \varepsilon_i^1$  and  $\varepsilon_i^0(1) = \varepsilon_i^0(0) = \varepsilon_i^0$  does not affect the main line of argumentation and is merely for a simple exposition. It restricts the gene-environment effect to 1.5 for each individual and, thus, each complier type. Assuming four different error terms allows for a different gene-environment interaction effect by complier type. Our argument is not affected by that, and neither does our solution need this restriction, nor do we make this assumption in the application in Sections 2 to 5.

## E Interpreting $G \times Z$

Rearranging the Wald estimator demonstrates that the reduced-form effect is the product of the structural (i.e., the second stage) and the first stage effects. This expression holds at every conceivable value of  $G = g$  (which, for the sake of the argument, is assumed to be continuous):

$$\begin{aligned} \mathbb{E}(Y \mid Z = 1, G = g) - \mathbb{E}(Y \mid Z = 0, G = g) &= \mathbb{E}(Y^1(G = g) - Y^0(G = g) \mid C(G = g)) \\ &\quad \times \mathbb{E}(\mathbb{1}[C(G = g)]) \end{aligned} \quad (11)$$

Conceptually, you can think of this equation as the following form:  $w(g) := u(g, v(g)) \times v(g)$ , where we define the reduced form  $w(g)$  as a function of  $g$ . It is the product of  $u(\cdot)$  – the second stage – and  $v(g)$  – the first stage (if aggregated). The second stage varies in  $G = g$ , the PGI at which the effect is assessed, and the specific complier group determined by  $v(\cdot)$ , a complier-indicating function (if unaggregated).

Applying the product and chain rule to this simplified expression  $w(g)$  demonstrates that marginally changing  $g$  has three distinct effects (i.e., partial derivatives) on the reduced form: (i) the partial direct derivative of  $w(\cdot)$  with respect to the structural gene heterogeneity of the second stage (i.e.,  $\frac{\partial u(\cdot)}{\partial g}$ , holding  $v(g)$  fixed), (ii) the partial derivative of  $u(g)$  with respect to the ( $g$ -specific) complier groups (i.e.,  $\frac{\partial u(\cdot)}{\partial v(\cdot)}$  holding  $g$  fixed), and (iii) the partial derivative of  $g$  with respect to the first stage  $\frac{\partial v(\cdot)}{\partial g}$ .

Formally in the notation of Eq. (11), this reads:

$$\begin{aligned}
& \frac{\partial \mathbb{E}(Y \mid Z = 1, G = g) - \mathbb{E}(Y \mid Z = 0, G = g)}{\partial g} = \\
& \underbrace{\frac{\partial \mathbb{E}(Y^1(G = g) - Y^0(G = g) \mid C(G = g))}{\partial g}}_{\text{(i) Structural outcome interaction}} \times \mathbb{E}\left(\mathbb{1}[C(G = g)]\right) \\
& + \underbrace{\frac{\partial \mathbb{E}(Y^1(G = g) - Y^0(G = g) \mid C(G = g))}{\partial \mathbb{E}\left(\mathbb{1}[C(G = g)]\right)}}_{\text{(ii) Interference with essential heterogeneity}} \\
& \times \frac{\partial \mathbb{E}\left(\mathbb{1}[C(G = g)]\right)}{\partial g} \times \mathbb{E}(C(G = g)) \\
& + \underbrace{\mathbb{E}(Y^1(G = g) - Y^0(G = g) \mid C(G = g))}_{\text{(iii) First-stage gradient}} \times \overbrace{\frac{\partial \mathbb{E}\left(\mathbb{1}[C(G = g)]\right)}{\partial g}}^{<0 \text{ in our paper}}
\end{aligned}$$

The interaction between channels (i) and (ii) is the core of our paper. The focus on the reduced-form adds a third channel, which blurs the structural outcome interaction we want to identify. To see this clearly, assume there is no structural outcome interaction and also no essential heterogeneity – so channels (i) and (ii) are switched off. Yet, the first stage may exhibit a gradient. In this case, the reduced form interaction  $G \times Z$  may still differ from zero. It remains:

$$\begin{aligned}
& \frac{\partial \mathbb{E}(Y \mid Z = 1, G = g) - \mathbb{E}(Y \mid Z = 0, G = g)}{\partial g} = \mathbb{E}(Y^1(G = g) - Y^0(G = g) \mid C(G = g)) \\
& \times \frac{\partial \mathbb{E}\left(\mathbb{1}[C(G = g)]\right)}{\partial g}
\end{aligned}$$

This demonstrates that even without any structural outcome interaction and essential heterogeneity, the  $G \times Z$  interaction may differ from zero. This is because the first stage changes along  $G$ . While the  $G \times Z$  may be informative about whether a policy  $Z$  has heterogeneous effects along  $G$  (as measured by the sum of all three mechanisms), it is uninformative about the structural outcome interaction (mechanism (i)).

## F Implications for other empirical applications

Of course, the described problem and solution are not restricted to gene-environment interactions. Observed heterogeneity in treatment effects in general and the effect gradient along one particular variable, in particular, are the focus of many applied papers. Consider the following regression equation:

$$Y_i = \beta_0 + \beta_1 D_i + \beta_2 X_i + \beta_3 D_i \times X_i + \varepsilon_i, \quad (12)$$

where  $Y_i$  is the outcome,  $D_i$  is the (endogenous) treatment indicator, and  $X_i$  is an observable variable. Without claiming that the problems with 2SLS apply here — because the conditions for 2SLS to cause problems are not necessarily fulfilled — we briefly mention some studies in the literature that are interested in how the effect of treatment  $D$  varies by  $X$ .

The literature on effects of education  $D$  has been interested in how they vary by socio-economic status  $X$ . For instance, early research on local average treatment effects of education focused on identifying effect gradients by socio-economic status to provide additional support for the plausibility of the IV estimates (Angrist and Krueger, 1991; Card, 1993). Relatedly, Brunello et al. (2017) want to identify the observed heterogeneity in the returns to education explicitly. They distinguish between urban and rural areas, as well as between individuals who grow up in households with many books versus those with few. They interpret their heterogeneous returns to compulsory schooling both structurally and in terms of selection (via different marginal costs and returns of schooling for affected individuals). Nonetheless, applying our method in this context could help differentiate between the two explanations.

Relatedly,  $D$  could be the effect of children,  $Y$  labor supply,  $X$  household income: In their landmark study, Angrist and Evans (1998) use the same-sex instrument to estimate whether the effects of children on labor supply differ by household income (as economic theory would predict). To assess the validity of this theory, one would need to hold the complier population constant and compare the effects of children on labor supply between different groups formed by household income. Given that Angrist and Evans (1998) show that the first stage is increasing along the husband's income (the more affluent the household, the more the household can afford their same-sex preference), 2SLS might be problematic here.

Other examples relate to the study of origins of the socio-economic gradient in health (Currie, 2009) and the literature on skill formation (Cunha and Heckman, 2007). These are prominent examples of a research question where the interaction effect between the environment in which individuals grow up ( $X$ ) and a choice variable  $D$  is at the center of interest. For instance, Conti et al. (2010) estimate the interaction between education and measures ( $X$ ) of early-life cognitive and non-cognitive skills, as well as health endowments



(the dimensions of observed heterogeneity) on later-life health behavior using structural methods (such as latent factor models that allow proxying unobservables and measurement error). If one were to employ instrumental variables estimation for this research question, it would be essential to consider differential responses to the instruments for all observed heterogeneity dimensions.

To take a recent paper as another example of this literature, [Agostinelli and Wiswall \(2025\)](#), estimate a latent factor model and specify an empirical model, in which they explore the heterogeneity in the returns to parental investments concerning the endowment of the children. If endowments and investments are more directly observed without measurement error and if quasi-experimental instruments for parental investments exist, researchers would employ instrumental variable methods. However, the problems outlined in this paper would arise if there is a gradient in the first stage along children's endowment and if, for instance, parents have a higher propensity to invest more in children with high returns to this investment.

Related to our research question, [Houmark et al. \(2024\)](#) estimates the technology of skill formation using skills, genetic endowment of children and parents, and parental investments, documenting that all factors are interrelated.

Heterogeneity between unobserved heterogeneity and covariates also appears to interfere in the paper by [Maestas et al. \(2013\)](#). They estimate the effects of disability benefit receipt on employment and report heterogeneous first-stage and IV regressions that vary by observed characteristics. For the first stage, coefficients on the instrument and the intercept vary substantially by covariates, suggesting that the instrument affects complier groups with different unobserved characteristics. While interpreting the 2SLS coefficients as heterogeneous LATEs is (of course) appropriate (as they do), this suggests that generalizing these LATEs beyond the instrument-specific complying population, as covariate-specific average treatment effects or covariate-specific effects of a more lenient or stringent disability receipt allowance reform, is likely inappropriate. This is especially important because [Maestas et al. \(2013\)](#) detect essential heterogeneity: employment effects of disability receipt are less negative for individuals with a higher unobserved severity.

Lastly, if researchers are interested in the effects of education by gender to gain better insights into possible structural disadvantages for females, it would be preferable to compare, for example, average treatment effects across genders than local average treatment effects for two possibly different complier groups.

## G Effect sizes of essential heterogeneity in MTE applications

Table G.1: The gradient within essential heterogeneity of other studies

Study	Outcome	Treatment	Note (Model/sample)	(1)	(2)	(3)	(4)	(5)	(6)
				$MTE(U_D = p)$		Linearized effect per quintile ( $\frac{(1)-(2)}{5-1}$ )	Standard deviation of outcome	Standardized linearized effect size ( $\frac{(3)}{(\frac{3}{4})}$ )	% of our MSE size ( $\frac{(5)}{0.12}$ )
				$p = 0.1$	$p = 0.9$				
Nybom (2017)	Lifetime earnings	College ed.	Norm. sel. mod. Local IV	0.04	-0.02	0.015	0.4	0.03	37.5%
				0.08	0.1	0.005	0.4	0.0125	10.4%
Carneiro et al. (2011)	Wage in 1991	College ed.	Norm. sel. mod. Local IV	0.14	0	0.035	0.47	0.06	62%
				0.35	-0.15	0.1	0.47	0.21	175%
Kamhöfer et al. (2019)	Math literacy	College ed.	Local IV	2.00	0.5	0.375	1	0.375	321.5%
	Reading speed		Local IV	1.8	0.3	0.375	1	0.375	321.5%
	Reading competence		Local IV	2.5	0	0.625	1	0.625	520%
Cornelissen et al. (2018)	School readiness	Child care att.	Linear MTE	-0.1	0.22	0.08	0.082	0.975	813%
Kowalski (2023)	ER visits	HI coverage	Linear MTE	0.47	-0.58	0.2625	2.63	0.1	83.3%
Ito et al. (2023)	Electricity usage	Dynamic pricing	Local IV	-750	100	212.5	250	0.85	708.3%

Notes: Our (non-standardized) "minimal significant effect size" (MSE) with 95% (90%) confidence is 1.96 (1.64) times the standard error on the linearized interaction reported in Table 4, i.e., 0.203. Hence, the MSE yields 0.4 (0.33). Our dependent variable (the recall score) has a standard deviation of 3.37, yielding a standardized MSE of  $0.4/3.37 = 0.12$ . Note that this minimal significant effect size differs from the minimal detectable effect size for power calculations because we condition on the test result. The minimal detectable effect size is typically defined as the minimal effect detected as significant at the 5 percent level in 80 percent of all cases. This measure, however, cannot be computed without further assumptions. Our target parameter is the average difference in the effects of education on cognitive abilities between the fifth and the first quintile (i.e., four quintiles). The MTE maps the impact of a treatment for every quintile of the unobserved heterogeneity in the treatment choice. We approximate the average effect heterogeneity between quintiles of this unobserved heterogeneity by computing  $\frac{MTE(p=0.1) - MTE(0.9)}{5-1}$  to compare the genetic heterogeneity to those of MTE studies (a perfect approximation would require identification at infinity and computing the integral between 0 and 0.2 and 0.8 and 1 for the first and fifth quintile, respectively). The reported values for  $MTE(p = 0.1)$  and  $MTE(p = 0.9)$  are approximated through eyeballing the corresponding MTE graphs. Nybom (2017) interprets the degree of heterogeneity as low, justifying this finding with a low general heterogeneity in wages in Sweden. Further note that standard errors may be approximated, e.g., by averaging standard errors between treatment and control samples. The reported standard deviation of Ito et al. (2023) is approximated based on the average pre-intervention daily electricity consumption (the outcome measures hourly consumption during peak hours). Therefore, the reported SD may likely be a lower bound. See Mogstad and Torgovitsky (2024) for more applications of the MTE. Norm. sel. mod. = Normal selection model (maximum likelihood based functional form assumption (joint normality)). HI coverage = Health insurance coverage. ER visits = Emergency room visits, College ed. = College education.

## H Testing “no selection into losses” (non-positive MTE slopes)

A critical constraint we apply in our linear programming approach is “no selection-into-losses”, i.e., no MTEs that increase in  $U_E$ . To test this in our setting, we follow [Imbens and Rubin \(1997\)](#) and use the instrument to compute mean outcomes for always-takers, treated and untreated compliers, and never-takers. For simplicity, we test this condition globally and do not distinguish between cells of  $G_i$  (we show the complete  $G_i$ -specific means in Figure 5). We present the results in Table H.1. In Panel A, we focus on differences between always-takers and treated compliers (Column 3) and untreated compliers and never-takers (Column 6). The differences are informative about whether the treated outcome  $\mathbb{E}[Y_i^1|U_E = u]$  and the untreated outcome  $\mathbb{E}[Y_i^0|U_E = u]$  – the difference of which is the MTE – are heterogeneous in  $U_E$ .

Column (3) presents the mean recall difference between always-takers and treated compliers. It shows a substantial and statistically significant heterogeneity: Always-takers recall about 1.25 words more. Intuitively, this is unsurprising, as always-takers to a compulsory schooling reform will, on average, have more years of education, will be more likely to hold advanced degrees, or may be positively selected in terms of unobserved characteristics (if we have selection into gains, which we want to argue). Furthermore, this result shows that  $\mathbb{E}[Y_i^1|U_E = u]$  has a negative slope. Likewise, we do the same with untreated compliers and never-takers. Here, the heterogeneity is less pronounced and not statistically significant. If we conclude that both groups do not have different outcomes, we can stop as in this case, the difference in the first two groups proves that we have essential heterogeneity. If the insignificant difference is meaningful, things may change. The difference is also negative, contrasting the existing empirical evidence for the slope of the untreated outcome (see, e.g., [Carneiro and Lee, 2009](#); [Westphal et al., 2022](#)). However, it is essential to mention that never-takers should not exist with a compulsory schooling reform, where everyone should be forced to stay in school until age 15. If this group has never existed, this might be a measurement error. If these individuals had special exemptions from the rule change (and therefore existed), the difference between never-takers and untreated compliers may not inform about the global course of the curve. Assessing the multiple complier groups that we gain by stratifying by  $G_i$  (see Figure 5) indeed suggests that never-takers are different and  $\mathbb{E}[Y_i^0|U_E = u]$  indeed increases when  $U_E < 0.95$ .

Nonetheless, with only a binary instrument and without exploiting covariate heterogeneity together with the additive separability assumption (which we will do below), an additional linearity assumption is necessary (due to the never-takers) to point-identify a marginal treatment effect via the method introduced by [Brinch et al. \(2017\)](#). We document a formal

Table H.1: Mean outcomes by instrument response types and test for essential heterogeneity

	Unobserved heterogeneity					
	in the treated outcome			in the untreated outcome		
	(1) Always-takers	(2) Treated compliers	(3) Difference (2) – (1)	(4) Untreated compliers	(5) Never-takers	(6) Difference (5) – (4)
<i>Panel A:</i>						
Mean recall:	9.500 (0.215)	8.306 (0.332)	–1.245*** (0.454)	8.109 (0.215)	7.679 (0.340)	–0.353 (0.396)
Share:	0.456 (0.035)	0.489 (0.036)		0.489 (0.036)	0.055 (0.011)	
<i>Panel B:</i>						
Test for essential heterogeneity: (sufficient condition, may be uninformative if heterogeneity is nonlinear )						
Slope of $\mathbb{E}[Y_i^1 U_E = u]$			–2.631*** (0.961)			
Slope of MTE $\mathbb{E}[Y_i^1 - Y_i^0 U_E = u]$			–1.326 (1.423)			

*Notes:* This table presents estimates of mean outcomes for always-takers, treated and untreated compliers, and never-takers (panel A) as well as results of a test for essential heterogeneity (panel B) using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. We compute the type-specific shares using the specification of Eq. (2) without  $G_i$ . The complier share is the coefficient on  $Z_i$ , the always-taker share is the constant (as all variables are demeaned), and the never-taker share is the remainder. For the type-specific outcome means, we compute means by  $E_i$  and  $Z_i$  (and their interaction) using a reduced-form specification to regress recall on the same controls and full interactions of  $E_i$  and  $Z_i$ . As compliers never appear alone in these means, we use the formula provided in Imbens and Rubin (1997) and the type-specific shares. Standard errors are computed using 200 bootstrap replications and are shown in parentheses. \*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$  indicate significance levels for the differences.

test of the slope of  $\mathbb{E}[Y_i^1|U_E = u]$  and  $\mathbb{E}[Y_i^1 - Y_i^0|U_E = u]$  in Panel B.<sup>20</sup> It shows that the slope of the treated outcome is negative and statistically significant (as shown in Panel A). The slope of the linear MTE is also negative and still large in magnitude. However, likely due to the concerns about never-takers outlined above, it is not statistically significant, albeit with a negative sign. Again, evidence from the  $G_i$ -specific complier groups strongly suggests that the  $\mathbb{E}[Y_i^0|U_E = u]$  increases at least for a relevant range when  $U_E < 0.95$ . We conclude that we likely face essential heterogeneity in our setting. Combined differences in the first stage induced by  $G_i$ , the result may be that 2SLS cannot recover the true interaction parameter. We would need to make accurate statements about the interaction effect.

<sup>20</sup>The exact formula reads

$$\frac{\partial \mathbb{E}[Y_i^1|U_E = u]}{\partial U_E} = \frac{Y_i^{CT} - Y_i^{AT}}{\frac{\pi^C + \pi^{AT}}{2}}, \quad \frac{\partial \mathbb{E}[Y_i^1 - Y_i^0|U_E = u]}{\partial U_E} = \frac{Y_i^{CT} - Y_i^{AT}}{\frac{\pi^C + \pi^{AT}}{2}} - \frac{Y_i^{NT} - Y_i^{CU}}{\frac{\pi^C + \pi^{NT}}{2}},$$

where  $Y_i^{AT}$ ,  $Y_i^{CT}$ ,  $Y_i^{CU}$ , and  $Y_i^{NT}$  are means from Columns (1), (2), (4), and (5), respectively and  $\pi^{AT}$ ,  $\pi^C$ ,  $\pi^{NT}$  are the corresponding shares (compliers do not need to be differentiated).

# I Details on the MTE estimation

We run the following two regressions:

$$E_i = \sum_{g=1}^5 \sum_{k=0}^1 \left[ \pi_{g,k}^f \mathbb{1}[G_i = g] \times [Z_i = k] \right] + \text{controls} + \omega_i \quad (13)$$

$$Y_i = \sum_{g=1}^5 \sum_{j=0}^1 \sum_{k=0}^1 \left[ \delta_{g,j,k}^f \mathbb{1}[G_i = g] \times [E_i = j][Z_i = k] \right] + \text{controls} + \eta_i. \quad (14)$$

The first equation estimates  $G_i$ -specific first-stage from which the complier types can be inferred. The second equation estimates conditional means of  $Y_i$ , conditional on  $G_i$ ,  $Z_i$ , and  $Y_i$  when covariates are fixed. On these estimates, we apply the [Imbens and Rubin \(1997\)](#) formula to compute  $G_i$ -specific outcome means for always-takers (AT), never-takers (NT), and (treated on untreated) compliers (C) the are plotted in [Figure 5](#):

$$\begin{aligned} \mathbb{E}[Y_i^1(G)|C, G_i = g] &= \frac{\delta_{g,1,1}\pi_{g,1} - \delta_{g,1,0}\pi_{g,0}}{\pi_{g,1} - \pi_{g,0}} \\ \mathbb{E}[Y_i^0(G)|C, G_i = g] &= \frac{\delta_{g,0,0}\pi_{g,0} - \delta_{g,0,1}\pi_{g,1}}{\pi_{g,1} - \pi_{g,0}} \\ \mathbb{E}[Y_i^0(G)|NT, G_i = g] &= \delta_{g,0,1} \\ \mathbb{E}[Y_i^1(G)|AT, G_i = g] &= \delta_{g,1,0} \end{aligned}$$

These linear potential outcome curves could already solve the problems associated with 2SLS estimation of interactions while using richer variations of the PGI. Based on them, we can calculate the (interaction) effects according to [Table 3](#) in the interval  $0.55 \leq U_D \leq 0.85$ . Graphically, this would entail subtracting the blue from the red lines for each quintile. However, this would require extrapolating the lines with  $E_i = 0$  to the left or the lines with  $E_i = 1$  to the right, demonstrating the general extrapolation problem that we could solve here by a linearity restriction. If we are willing to make this extrapolation, it yields five MTE curves for the effect of  $E_i$  on  $Y_i$ , one for each quintile, which can then be used to calculate the interaction effects.

In the paper, we are unwilling to make such an assumption and apply the partial identification method by [Mogstad et al. \(2018\)](#). As one input, the method uses the conditional means that the coefficients ( $\delta_{g,j,k}^f$  and the corresponding  $\pi_{g,k}^f$ ) reflect. These are the "moments" for the linear programming method by [Mogstad et al. \(2018\)](#). [Figure I.1](#) plots the results of this approach, where the slightly transparent, horizontal lines are the "moments" ( $G_i$ -specific outcome means and their placement on the unit-interval, which we derive from the complier shares). The blue (for the treated outcome) and red (for the untreated) lines are the

output of this linear programming approach. They reflect the minimal (the dashed lines) and maximal (the solid lines) possible interaction effect (defined in the main text) that the MTR lines (Bernstein polynomials, see Figure I.2) produce while being consistent with the shape restrictions and matching the moments.

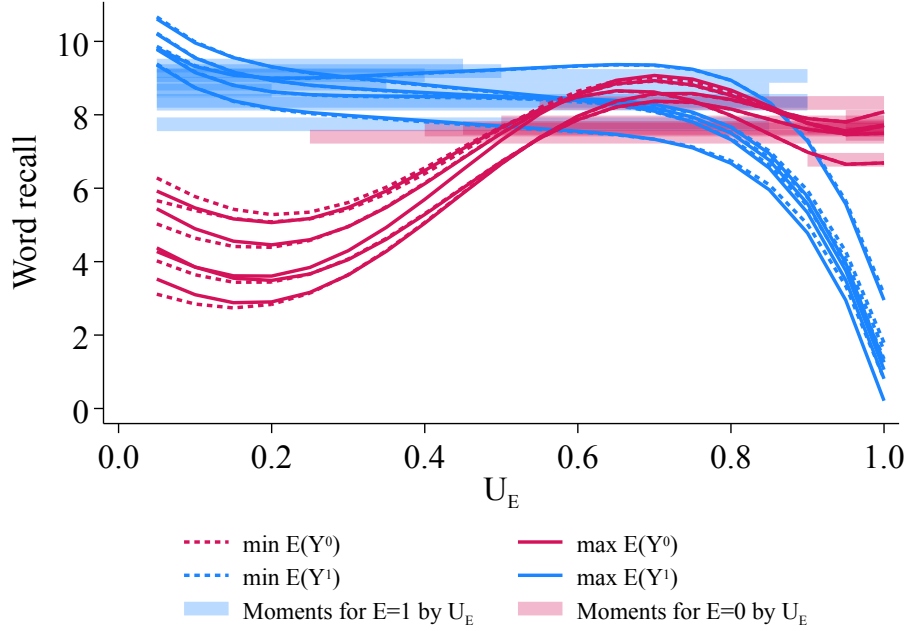


Figure I.1: Potential outcome curves estimated with Bernstein polynomials

*Notes:* This figure shows the minima and maxima of the ten potential outcome curves estimated via linear program with Bernstein polynomials using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. Blue indicates curves and moments for  $E_i = 1$ , and red indicates  $E_i = 0$ . Solid lines are maxima; dashed lines are minima of the potential outcome curves. There are five pairs of curves for  $E_i = 1$  and five for  $E_i = 0$ , one pair for every PGI quintile. Every pair consists of a minimum and a maximum that bound the potential outcome curve for its respective quintile. The vertical bars indicate the moments the curves must pass and the  $U_E$  ranges of individuals contributing to these means.

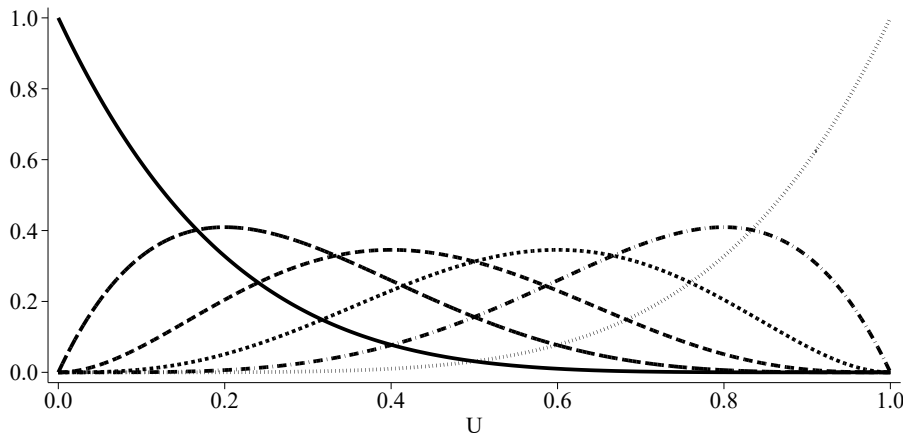


Figure I.2: Graphical representation of the Bernstein base functions

*Notes:* This figure depicts the six Bernstein base functions that compose a Bernstein polynomial of degree five in simulated data. The formula for each base function reads  $b_{v,n}(u) = \binom{n}{v} u^v (1-u)^{n-v}$ , where  $n = 5$  is the degree,  $v$  denotes the specific base function and  $u$  is a specific grid point on the unit interval. The formula that obtains the MTE by the sum of all base functions weighted by the corresponding parameter  $\theta_v^{jg}$  reads  $m^j(u, g) = \sum_{v=0}^n \theta_v^{jg} b_{v,n}(u)$ , where  $G_i$  is the genetic bin,  $j$  the treatment state (as defined above) in addition to the variables and parameters defined above.



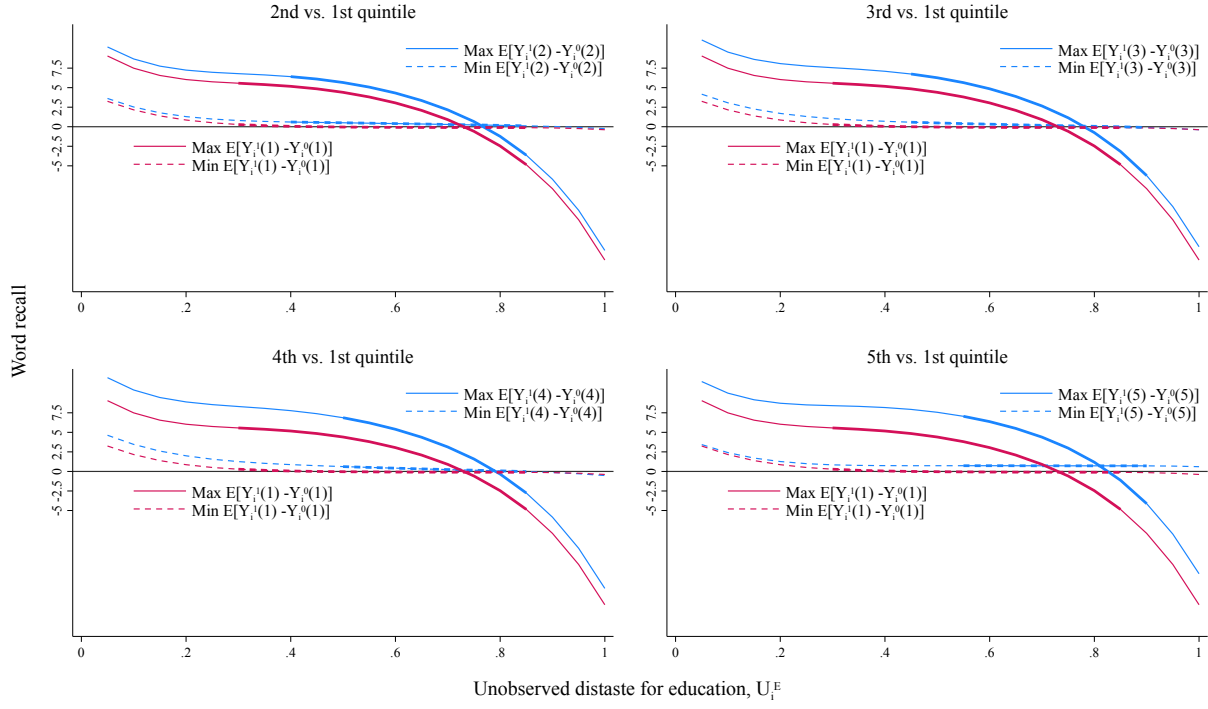


Figure I.3: Quintile comparisons of the interaction effect without never-takers

*Notes:* This figure shows the quintile comparisons of the interaction effect from Figure 7 when never-takers (their sample moments) are not used to construct the MTE bounds using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. For every PGI quintile, we estimate bounds: maxima (solid lines) and minima (dashed lines) at which the interaction effect is maximized/minimized. The bounds for quintiles 2–4 (in blue) are compared to those of the bottom quintile (in red), our reference category, yielding four comparisons. The gene-environment interaction is the difference between the blue and red curves at  $U_E \in [0.55, 0.85]$ . The thick part of the curves indicates the size of the complier share and its location on the  $U_E$  scale, both of which differ by PGI quintile.

## J Robustness checks for the linear programming approach

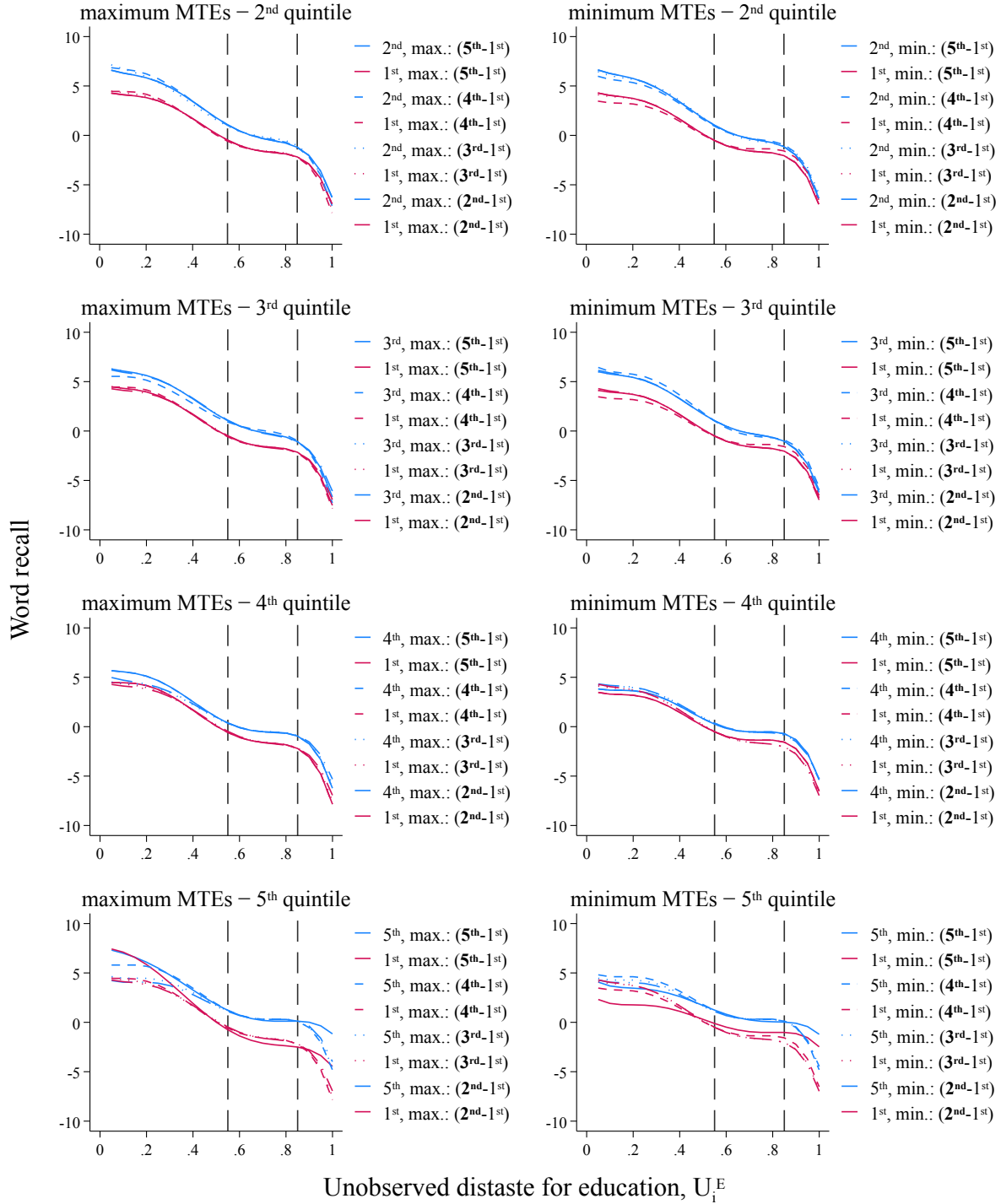


Figure J.1: MTEs when the target  $G \times E$  parameter is adjusted to specific quintiles

Notes: This figure shows robustness checks for our main result in Figure 7 using data from ELSA waves 1–6 and our main sample selection, as outlined in Chapter 2.2. Here, we optimize the interaction effect for different comparisons. Whereas our preferred specification optimizes the difference between the first and the fifth quintile (see Eq. 8), we generalize this approach and optimize differences between the first and any other quintile such that  $\beta_{G \times E}(0.55, 0.85, g) = \frac{1}{g-1} \int_{0.55}^{0.85} [m^1(u, g) - m^0(u, g)] - [m^1(u, 1) - m^0(u, 1)] du \quad \forall g \in \{2, 3, 4, 5\}$ . The solid lines correspond to optimizing  $g = 5$ , our main result. The dashed lines show the optimization for  $g = 4$ , the dotted for  $g = 3$ , and the dashed-dotted line for  $g = 2$ . The respective quintile  $G_i$  used for the target parameter  $\beta_{G \times E}(0.55, 0.85, g)$  is highlighted in bold. Maximized and minimized MTEs are shown separately, maximized MTEs in the left and minimized MTEs in the right column. The rows present pairwise comparisons between the first and another PGI quintile (the second quintile in the first row, the third in the second row, ...).