









Mendelian randomisation with coarsened exposures

Matthew J. Tudball^{1,2}  | Jack Bowden^{1,2,3}  | Rachael A. Hughes^{1,2}  |
Amanda Ly^{1,2}  | Marcus R. Munafò^{1,2,4}  | Kate Tilling^{1,2}  |
Qingyuan Zhao⁵  | George Davey Smith^{1,2} 

¹MRC Integrative Epidemiology Unit,
University of Bristol, Bristol, UK

²Population Health Sciences, Bristol
Medical School, University of Bristol,
Bristol, UK

³College of Medicine and Health,
University of Exeter, Exeter, UK

⁴School of Psychological Science,
University of Bristol, Bristol, UK

⁵Department of Pure Mathematics and
Mathematical Statistics, University of
Cambridge, Cambridge, UK

Correspondence

Matthew J. Tudball, MRC Integrative
Epidemiology Unit, University of Bristol,
Oakfield House, Oakfield Grove, Bristol,
BS8 2BN, UK.

Email: matt.tudball@bristol.ac.uk

Funding information

Wellcome Trust and the Royal Society,
Grant/Award Number: 220067/Z/20/Z,
215408/Z/19/Z, 220059/Z/19/Z;
Establishing Excellence in England (E3)
Medical Research Council Integrative
Epidemiology Unit at the University of
Bristol, Grant/Award Number:
MC_UU_00011/1, MC_UU_00011/7

Abstract

A key assumption in Mendelian randomisation is that the relationship between the genetic instruments and the outcome is fully mediated by the exposure, known as the exclusion restriction assumption. However, in epidemiological studies, the exposure is often a coarsened approximation to some latent continuous trait. For example, latent liability to schizophrenia can be thought of as underlying the binary diagnosis measure. Genetically driven variation in the outcome can exist within categories of the exposure measurement, thus violating this assumption. We propose a framework to clarify this violation, deriving a simple expression for the resulting bias and showing that it may inflate or deflate effect estimates but will not reverse their sign. We then characterise a set of assumptions and a straight-forward method for estimating the effect of *SD* increases in the latent exposure. Our method relies on a sensitivity parameter which can be interpreted as the genetic variance of the latent exposure. We show that this method can be applied in both the one-sample and two-sample settings. We conclude by demonstrating our method in an applied example and reanalysing two papers which are likely to suffer from this type of bias, allowing meaningful interpretation of their effect sizes.

KEYWORDS

biomarkers, latent variable modelling, Mendelian randomisation analysis, sensitivity analysis

1 | INTRODUCTION

Mendelian randomisation proposes to use genetic variants that alter, or mirror the biological effects of, modifiable exposures to study the causal effects of such exposures on downstream outcomes. The principle underlying Mendelian randomisation is that genetic variants are randomly passed from parents to offspring at

conception, resulting in a plausibly unconfounded source of variation in the exposures with which they are associated. For Mendelian randomisation estimates to inform policies or clinical practices, we must additionally assume that genetic and environmental modifiers of the exposure produce similar effects on the outcome (Davey Smith & Ebrahim, 2003). For example, Mendelian randomisation studies of pharmaceutical exposures typically

This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2021 The Authors. *Genetic Epidemiology* published by Wiley Periodicals LLC

use genetic variants that code for potential drug targets, assuming that similar effects would be observed if those targets were altered therapeutically (Plump & Davey Smith, 2019).

One of the crucial assumptions underlying the Mendelian randomisation approach is that the relationship between the genetic instruments and the outcome is fully mediated by the exposure, known as the exclusion restriction assumption. However, it is important to draw a distinction between the true exposure experienced by an individual and our attempt at measuring it. For practical purposes, we are often restricted to coarsened approximations which do not fully encapsulate the mechanism by which the true exposure of interest affects the outcome. Consistent with existing terminology, we define an exposure measurement as coarsened if it is a discrete measure approximating a continuous latent exposure (Marshall, 2016).

In the Mendelian randomisation context, coarsened exposures can violate the exclusion restriction assumption. If the genetic instruments are acting on a latent exposure, such as body mass index (BMI), but the measured exposure is a discretisation of it, such as obesity status, then there can exist genetically driven variation in the true exposure within categories of the measured exposure. We could imagine that counterfactually altering some BMI-raising single-nucleotide polymorphism (SNP) in an individual could result in a change in their BMI without necessarily changing their obesity status. This can be viewed a form of measurement error which opens up potential pathways from the genetic instruments to the outcome that do not pass through the exposure measure, thus violating the exclusion restriction assumption.

For example, Richardson et al. (2020) attempt to separate the effects of early and later life adiposity on disease risk. The adiposity variable is a three-category self-report measure (“thinner,” “plumper,” and “about average”). It is reasonable to conceptualise a continuous measure of body mass (e.g., BMI) underlying this coarsened categorical measure, such that genetic variation in this latent continuous measure could occur within categories of the self-report variable. We later reanalyse Richardson et al. (2020) in Box 2 using the approach proposed in this paper. Another example is Richmond et al. (2019), who apply Mendelian randomisation to investigate the effect of sleep traits (e.g., morning preference, sleep duration) on breast cancer risk, finding large causal effects of several traits. These traits are categorical measures, for example, morning preference is measured in six categories and sleep duration is split into several groups. It is reasonable to conceptualise the true exposures on which the genetic variants are acting as latent continuous sleep traits and preferences, for which the measured exposures are discrete markers.

An important class of latent exposures we consider in this paper is disease liabilities, for which binary disease diagnosis or case status is the typical exposure measurement. There are an increasing number of Mendelian randomisation studies investigating the effects of complex diseases such as asthma, schizophrenia and attention deficit hyperactivity disorder on various outcomes (Lawn et al., 2019; Martins-Silva et al., 2019; Pasman et al., 2018; Sun et al., 2019). Complex diseases which result from the interaction of environment and multiple genetic variants are likely to affect outcomes of interest through pathways other than diagnosis, for example, severity of subclinical symptoms. Since genetic instruments are, in turn, likely to influence the manifestation or severity of the underlying symptoms, rather than diagnosis alone, this represents a potential violation of the exclusion restriction.

This specific violation of the exclusion restriction assumption has been raised before in both the economics and political science literatures (Angrist & Imbens, 1995; Marshall, 2016). It has also been raised briefly in the Mendelian randomisation context in Burgess and Labrecque (2018), who discuss interpretation of estimates with binary exposures. The authors recommend that findings be framed in terms of this latent exposure but note that the estimates themselves have no meaningful causal interpretation. However, it remains to explore in more detail how this bias may distort estimates and clarify how to appropriately frame estimates in terms of the latent exposure, which will depend on the unobservable relationship between the latent exposure and its coarsened measurement.

We attempt to provide these clarifications in this paper. In particular, we derive an expression for the bias and introduce a clear set of identifying assumptions under which one can estimate the causal effect of the latent exposure. We hope to allow researchers to decide whether these assumptions are plausible in the context of their study. In Section 2, we outline our technical framework, which assumes a linear single threshold model for the relationship between the latent exposure and its measurement. That is, we assume that values of the coarsened exposure are determined by whether the latent exposure is above or below some threshold, which could be individual-specific. For example, an individual is classified as obese if their BMI is above 30 and not obese otherwise. This framework also contains the Falconer (1965) liability-threshold model, which assumes that a disease occurs in an individual, or is sufficiently pronounced to be diagnosed, if a build-up of underlying liability crosses some threshold. In this model, liability is assumed to capture all genetic, shared and nonshared environmental risk factors.

In Section 3.1, we derive an expression for the bias from the naive approach of using the coarsened measure as the exposure directly. Then, in Section 3.2, we show that, if the latent exposure is standardised to

have a SD of one, its causal effect can be identified if we have auxiliary information on the genetic variance of the latent exposure. This may be obtained from genome wide association study (GWAS) or treated as a sensitivity parameter and varied over a plausible range of values. In the context of disease liabilities, we may use the coefficient of determination developed by Lee et al. (2012).

Section 4 provides some generalisations to this framework, in particular, allowing two-sample estimation. Section 5 provides a real data example by creating artificially dichotomised variables from the continuous BMI measure in UK Biobank. Boxes 1 and 2 present re-analyses of two papers which could be interpreted within the framework proposed in this paper (Pasman et al., 2018; Richardson et al., 2020). In sections A and B of the appendix, we examine the bias that can emerge when the assumptions of our framework are violated.

2 | FRAMEWORK

We begin by outlining some key notation. Suppose there is a genetic instrument $Z \in \mathbb{R}$, other genetic variants (e.g., pleiotropic, weak) $X \in \mathbb{R}^K$ and an environmental risk factor $V \in \mathbb{R}$, where V is assumed to be continuously distributed with mean zero. We also assume that Z , X and V are mutually independent. We define $G = \mu + \alpha Z + \gamma'X$ as the genetic share of the latent exposure and define the latent exposure itself as

$$\begin{aligned} L &= G - V \\ &= \mu + \alpha Z + \gamma'X - V. \end{aligned} \quad (1)$$

It would be equally correct to define $L = G + V$, but the formulation in (1) simplifies some later expressions. In the Falconer framework described in Section 1, L would represent liability to some disease. We are able to observe a coarsened exposure characterised by a dichotomisation of the latent exposure.

$$D = \begin{cases} 1 & \text{if } L \geq 0 \\ 0 & \text{if } L < 0 \end{cases}. \quad (2)$$

If L is disease liability, then D would represent occurrence of the disease. In practice, we measure diagnosis of the disease, which does not necessarily correspond to occurrence due to under- or over-diagnosis. We will treat the two as equivalent throughout and discuss violations of this equivalence in Section 6.

Equation (2) is the crucial assumption underlying our approach; namely, that L is a linear index that relates to D according to a single threshold. Section A of the appendix elaborates on the importance of this structural assumption. Figure 3 illustrates our model within the

Box 1 Reanalysis of Pasman et al. (2018)

Pasman et al. (2018) performs a two-sample bidirectional Mendelian randomisation analysis of schizophrenia and cannabis use (Burgess et al., 2015). The gene-exposure associations for schizophrenia are pulled from a GWAS of cases and controls and are reported on the log-odds scale (Schizophrenia Working Group of the Psychiatric Genomics Consortium, 2015). While this avoids the problem of using the dichotomous diagnosis variable as the exposure (as discussed in Section 1), it means that the resulting estimates are interpreted as unit increases in the log-odds, which are scaled by the unobserved parameter σ_V . The authors report an odds ratio (OR) of 1.16 (95% confidence interval [95% CI] = 1.06–1.27) for the effect of genetic liability to schizophrenia. While we can infer the direction of the effect from this estimate, we cannot draw any conclusions about the magnitude.

We apply the two-sample generalisation of Section 4.4. One of the strengths of this generalisation is that we do not need to re-estimate the original inverse-variance weighted Mendelian randomisation estimates ourselves. In addition to the estimates reported in the original paper, we need only an estimate of σ_{G^*} , which can be computed from summary data from the schizophrenia GWAS, and some plausible choices for the sensitivity parameter θ^2 . The schizophrenia GWAS reports that their genome-wide significant loci explain roughly 3.4% of the variation in schizophrenia liability using the Lee et al. (2012) coefficient of determination. Using this estimate as a baseline, we select three choices for θ^2 : 0.02, 0.034, and 0.05.

Our findings are consistent with a modest positive effect of schizophrenia liability on the odds of cannabis use. As shown in Figure 1, a one SD increase in schizophrenia liability corresponds to a 1.15–1.26 increase in the odds of cannabis use, with 95% CI range of 1.10–1.44. It is important not to directly compare these estimates with the original estimates: the two are not on the same scale. We must interpret the estimates of Figure 1 in terms of SD increases in schizophrenia liability.

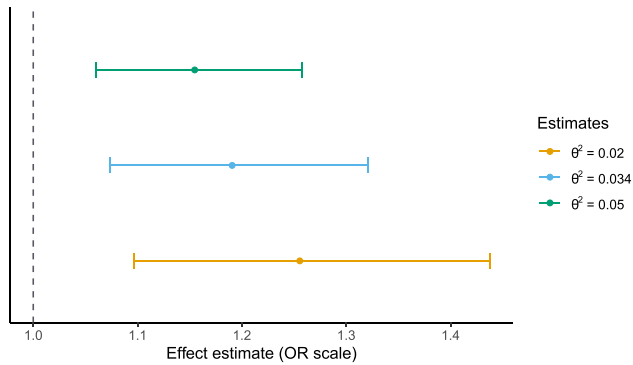


FIGURE 1 Effect of schizophrenia liability on risk of ever using cannabis for several choices of sensitivity parameter θ^2 . 95% confidence intervals are estimated as in section C of the appendix

Falconer framework. There is a distribution of disease liabilities and the disease occurs at the right tail of this distribution. The size of the grey region represents the prevalence of the disease in the population.

We also have an observed outcome $Y \in \mathbb{R}$. For ease of exposition, we restrict ourselves to a simple linear structural equation model

$$Y = \beta L + \varepsilon \quad (3)$$

which is implicitly conditional on covariates, where ε can be correlated with both V and X . However, this

framework can accommodate more general exposure-outcome relationships of the form $Y = f(L) + \varepsilon$, provided $E[f(L)|Z = z]$ is differentiable with respect to z . We make the standard instrumental variable assumptions, namely, that $\alpha \neq 0$ and Z is independent of ε conditional on covariates. The model (3) implicitly captures the assumption described in Section 1 that genetic and environmental modifiers of the exposure produce equivalent effects on the outcome. In this setting, the marginal effect (in absolute value) of both G and V is β . Figure 4 summarises this model in a directed acyclic graph. We can see that the exclusion restriction is violated since there exists a path from the latent exposure L to Y which does not pass through the measured exposure D . The structural Equation (3) assumes no effect of D itself. For a disease such as schizophrenia, liability could have a harmful effect on the outcome but being diagnosed will usually lead to receiving treatment and thus could have a protective effect. We cannot separately identify the two effects in this setting, although possibilities for doing so are discussed in Section 4.2. When D is believed to have a distinct effect on the outcome, we may instead identify the total effect of liability on the outcome; that is, the direct effect β and the indirect effect through D .

The structural assumptions made in this section can be summarised as follows:

Box 2 Reanalysis of Richardson et al. (2020)

Richardson et al. (2020) performs two-sample Mendelian randomisation analysis of child and adult BMI on risk of several diseases: coronary artery disease, type 2 diabetes, breast cancer and prostate cancer. The instrument-exposure relationship is estimated in the UK Biobank cohort. However, child BMI is not measured directly in UK Biobank, instead, there is a measure of self-reported adiposity in three discrete categories (“thinner,” “plumper,” or “about average”). In this context, the latent exposure is child BMI and the self-report measure is a coarsening of child BMI. Since the genetic instruments will act on child BMI directly, the exclusion restriction is likely to be violated.

Therefore, we apply the latent variable method of Section 3.2 to this data. We reanalyse the original univariable effect of child BMI on risk of type 2 diabetes (OR = 2.32, 95% CI = 1.76–3.05), coronary artery disease (1.49, 1.33–1.68), and breast cancer (0.59, 0.50–0.71).

We apply the two-sample generalisation of the inverse-variance weighted estimator of Section 4.4, estimating the instrument-exposure relationship in UK Biobank using an ordered probit model and the instrument-outcome relationships using the MR-Base platform (Hemani et al., 2018). We choose three values for θ^2 based on a large GWAS of adult BMI: 0.01, 0.02, and 0.05 (Locke et al., 2015). The genetic share of child BMI is estimated using an ordered probit model and standard errors are calculated using the formula in section C of the appendix.

Figure 2 shows our results for three of the diseases analysed in the paper. Our estimates are in the same direction as the original estimates, which is expected, however, the interpretation of the magnitudes is different. For example, the original paper estimates that a per-category increase in self-reported child adiposity corresponds to an increase in the odds of coronary artery disease of 1.49 (95% CI = 1.33–1.68), which could be inflated due to violation of the exclusion restriction. For $\theta^2 = 0.02$, we estimate that a one *SD* increase in child BMI corresponds to an increase in the odds of coronary artery disease of 1.13 (95% CI = 0.99–1.28). It is difficult to directly compare the two sets of estimates since the exposures are different, however, our estimate is suggestive of a modest effect of child BMI on the risk of coronary artery disease.

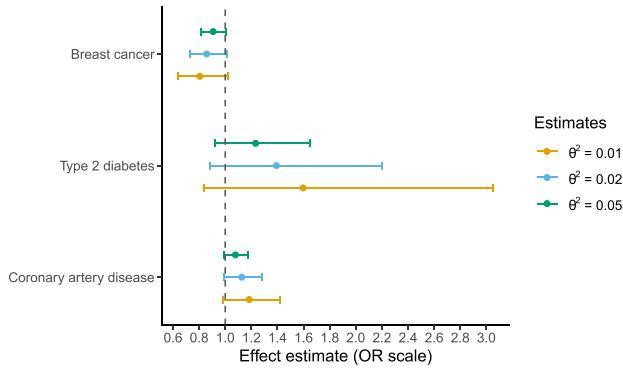


FIGURE 2 Effect of childhood body mass index on risk of several diseases for several choices of sensitivity parameter θ^2 . 95% confidence intervals are estimated as in section C of the appendix

Assumption 1 (Single threshold). The latent exposure L and its binary measurement D are related by a single threshold model of the form $D = I\{L \geq 0\}$.

$$\begin{aligned}
 \text{cov}(Z, D)/\text{var}(Z) &= \text{pr}(D = 1|Z = 1) - \text{pr}(D = 1|Z = 0) \\
 &= \text{pr}(L \geq 0|Z = 1) - \text{pr}(L \geq 0|Z = 0) && \text{(Assumption 1)} \\
 &= \text{pr}(V \leq \mu + \alpha Z|Z = 1) - \text{pr}(V \leq \mu + \alpha Z|Z = 0) && \text{(Assumptions 2, 3 and 4)} \\
 &= F_V(\mu + \alpha) - F_V(\mu) && \text{(Assumption 5)} \\
 &= \alpha f_V(\mu^*) && \text{(Assumption 4)}
 \end{aligned}$$

Assumption 2 (Additivity). $L = G - V$, where G and V are, respectively, the genetic and environmental shares of L .

Assumption 3 (Linearity). G is a linear function of the genetic instrument Z and other genetic variants X , such that $G = \mu + \alpha Z + \gamma'X$.

Assumption 4 (Environmental share). V has mean zero, SD σ_V and is in some family of continuous distributions, with cumulative distribution function given by $F(v/\sigma_V) = F_V(v)$ and density $f(v/\sigma_V) = f_V(v)$.

Assumption 5 (Risk factor independence). Z , X and V are mutually independent.

Assumption 6 (Gene–environment equivalence). The outcome model takes the form $Y = \beta L + \varepsilon$, where ε is a random disturbance and X and V may be correlated with ε .

Assumption 7 (Instrumental variable assumptions). Z is independent of ε and $\alpha \neq 0$.

3 | IDENTIFICATION

3.1 | Bias from the naive approach

The naive approach to Mendelian randomisation is to use the coarsened exposure D as the exposure directly. We show in this section that this results in a “multiplicative” bias which will scale the true effect β up or down, but not change its direction. When the distribution of L has a light tail (e.g., normal distribution), we will typically see inflation of effect estimates, with the degree of inflation increasing as the prevalence of D becomes smaller. If D is case status for a disease, for example, then effect estimates will be more inflated for rarer diseases. We see this pattern of inflation occurring in our real data examples in Section 5.

We call the naive Wald estimand $\beta_D = \text{cov}(Z, Y)/\text{cov}(Z, D)$. It is illustrative to derive a closed-form expression for β_D . Suppose Z is binary and $G = \mu + \alpha Z$ (i.e., there is no X). Begin by noting that

by the mean value theorem, where $\mu \leq \mu^* \leq \mu + \alpha$. Thus, the estimand can be written as

$$\beta_D = \text{cov}(Z, Y)/\text{cov}(Z, D) = \beta/f_V(\mu^*), \quad (4)$$

meaning that β_D is equal to the true latent exposure effect β divided by the density of V at the value μ^* . $f_V(\mu^*)$ is not identified since the distribution of V is unknown and μ^* is defined on the scale of the latent exposure.

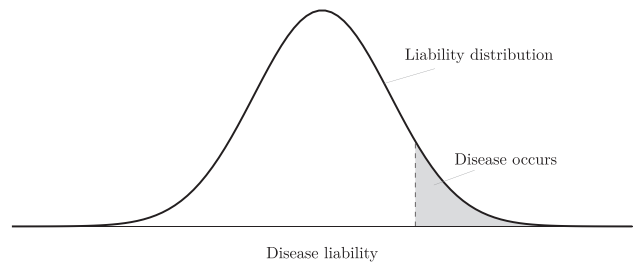


FIGURE 3 In the Falconer framework, liability to a disease is assumed to follow a smooth (often normal) distribution. The disease occurs at the tail of the distribution, with the grey region representing expected prevalence in the population

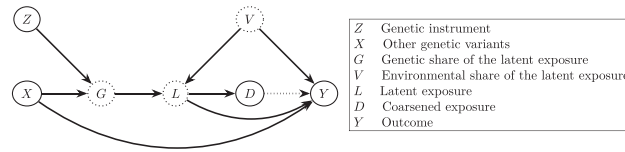


FIGURE 4 The framework proposed in Section 2 is summarised in a directed acyclic graph. Dotted circles represent latent variables and complete circles represent observed variables

3.2 | The latent variable approach

The bias formula (4) indicates that the nuisance term is $f_V(\cdot)$, which is the distribution of the environmental share V . Although D depends on this unobserved distribution, the genetic share G does not. Our latent variable approach therefore proceeds in four steps: (1) estimate the linear predictor of a generalised linear model of D on Z and X ; (2) normalise the linear predictor to have mean zero and variance one; (3) use this normalised linear predictor as the exposure in an instrumental variable model; and (4) scale the resulting effect estimate up by the genetic variance of the latent exposure. Step 4 is necessary to interpret effect estimates in terms of SD increases in the latent exposure, which is typically the desired scale.

To state this more precisely, define $\sigma_L = (\sigma_G^2 + \sigma_V^2)^{1/2}$ as the SD of L , where σ_G^2 and σ_V^2 are the variances of G and V , respectively. Within the framework described in Section 2, we claim that the four steps above allow us to identify $\beta_L = \sigma_L \beta$ from the observed data (Z, X, D, Y) .

The remainder of this section proves this claim given the assumptions outlined in Section 2 and discusses its implications. We begin by expressing the quantity $pr(D = 1|X = x, Z = z)$ within the framework of Section 2.

where $\tilde{\mu} = \mu/\sigma_V$, $\tilde{\alpha} = \alpha/\sigma_V$, and $\tilde{\gamma} = \gamma/\sigma_V$. F can be interpreted as the link function in a generalised linear model and $\tilde{\mu}$, $\tilde{\alpha}$, and $\tilde{\gamma}$ as parameters that can be identified from the observable data. In practice, we could specify F directly, for example, as a logistic or normal distribution (corresponding to logistic and probit regressions respectively). Alternatively, to avoid imposing potentially strong distributional assumptions, we could use semi-parametric estimation methods for generalised linear models, which only require some

smoothness conditions on F (Ichimura, 1993; Klein & Spady, 1993). Disease liabilities are often assumed to be the product of many small, independent traits. Therefore, by the central limit theorem, a normal distribution (i.e., probit model) is a natural choice of link function in this context (Curnow, 1972).

Step 1 is accomplished by constructing the predicted genetic share of the latent exposure

$$\tilde{G} = \tilde{\mu} + \tilde{\alpha}Z + \tilde{\gamma}'X = G/\sigma_V$$

using parameters estimated from the generalised linear model of D on Z and X . An immediate complication is that σ_V is unobserved. Treating σ_V as a sensitivity parameter is not tractable since its value is defined on the scale of the latent exposure, which is unknown. However, if we standardise \tilde{G} by its SD as in step 2, we can remove σ_V since

$$\tilde{G}/\sigma_{\tilde{G}} = (G/\sigma_V)/(\sigma_G/\sigma_V) = G/\sigma_G.$$

By using G/σ_G as our exposure, we can obtain effects in terms of SD increases in the genetic share of the latent exposure. The instrumental variable estimand of step 3 equals

$$\text{cov}(Z, Y)/\text{cov}(Z, G/\sigma_G) = \sigma_G \beta = \beta_G.$$

This estimand does not often have a natural interpretation. We would prefer to interpret our effects in terms of changes in the latent exposure itself.

Let $\theta^2 = \sigma_G^2/\sigma_L^2$ be defined as the genetic variance of the latent exposure. If we have a suitable choice of θ^2 , we can simply adjust our estimand as in step 4 such that

$$\beta_G/\theta = \sigma_G \beta / (\sigma_G/\sigma_L) = \sigma_L \beta = \beta_L$$

5

$$\begin{aligned} pr(D = 1|X = x, Z = z) &= pr(L \geq 0|X = x, Z = z) && \text{(Assumption 1)} \\ &= pr(V \leq \mu + \alpha z + \gamma'x|X = x, Z = z) && \text{(Assumptions 2, 3 and 4)} \\ &= F((\mu + \alpha z + \gamma'x)/\sigma_V) && \text{(Assumption 5)} \\ &= F(\tilde{\mu} + \tilde{\alpha}z + \tilde{\gamma}'x), \end{aligned}$$

(5)

which is our desired effect. The parameter θ^2 can be treated as a sensitivity parameter and varied over a plausible range of values or can, in some instances, be obtained from GWAS which report this measure.

For disease liabilities in particular, Lee et al. (2012) uses the Falconer liability-threshold model to develop a coefficient of determination for GWAS that is interpretable on the liability scale, which corresponds to θ^2 . Therefore, θ^2 can be estimated using this approach or selected from GWAS which report this coefficient. For ease of interpretation, liability is often assumed to have mean zero and variance one, in which case $\sigma_L = 1$ and β itself is identified on this scale (Lee et al., 2012).

4 | SOME GENERALISATIONS

4.1 | Individual-specific threshold

The formalisation of the relationship between disease and liability in Equation (2) and Figure 3 assumes a fixed threshold. That is, all individuals with liability above the threshold will develop or be diagnosed with the disease and all those below the threshold will not. In reality, we might imagine that diagnosis has a random component, driven, for example, by preferences of the diagnosing clinician or imprecision of the testing procedure. It might be more realistic to assume a model such that

$$D = \begin{cases} 1 & \text{if } L \geq R \\ 0 & \text{if } L < R \end{cases} \quad (6)$$

where R is a random individual-specific threshold. Provided R is independent of the instrument Z and other variants X , this random threshold will not affect identification of $\tilde{\mu}$, $\tilde{\alpha}$, and $\tilde{\gamma}$ of Equation (5) under correct model specification. However, the link function F of Equation (5) no longer corresponds to the distribution family of V ; instead, it corresponds to the distribution family of $V + R$. This could make correct specification of the link function more difficult and semiparametric approaches may be warranted.

4.2 | Identifying effects of the coarsened exposure

The structural model (3) assumes no direct effect of the binary exposure measure D on the outcome. As discussed in Section 3, when D is diagnosis of a disease, we might expect resulting treatment or therapy to have an effect on the outcome distinct from disease liability, suggesting a structural equation model of the form

$$Y = \beta L + \delta D + \varepsilon. \quad (7)$$

The exposure measure is downstream of the latent exposure and there are assumed to be no direct pathways from the genetic instruments to the exposure measure, as illustrated in Figure 4. Therefore, we cannot use our genetic instrument Z to estimate the independent effect of the exposure measure on the outcome; the genetic instruments induce no unique variation in the exposure measure independent of the latent exposure. However, consider the individual-specific threshold of Section 4.1. The variable R could represent preferences of the clinician for diagnosing the disease or a change in clinical practices affecting some individuals (Brookhart & Schneeweiss, 2007; Davies et al., 2013). If R is independent of each individual's liability, without directly affecting the outcome, then it is a potential instrument for disease diagnosis. The general rule for separately estimating the effects of the latent exposure and coarsened exposure is to have instruments which induce distinct variation in both.

4.3 | Multivalued discrete exposure

This method generalises easily to the multivalued discrete exposure setting. Suppose we observe a discretised variable characterised by

$$D = \begin{cases} 0 & \text{if } L \leq 0 \\ 1 & \text{if } 0 < L \leq d_1 \\ \vdots & \\ K & \text{if } d_{K-1} < L \end{cases} \quad (8)$$

where $0 < d_1 < \dots < d_{K-1}$ are latent thresholds. D could represent number of years in education and L could represent time in education as a continuous measure. Similar to how the dichotomous exposure can be formulated as a binary response model as in Equation (5), exposures of the form (8) can be formulated as an ordered response model and the parameters $\tilde{\mu}$, $\tilde{\alpha}$, and $\tilde{\gamma}$ are still identified, allowing the method to be applied as usual.

4.4 | Two-sample design with GWAS summary statistics

For rare diseases, it is not always possible to observe the coarsened exposure D and the outcome Y in the same sample. It is common practice in Mendelian randomisation studies to use summary statistics from separate GWAS of the exposure and outcome to obtain two-sample estimates (Burgess et al., 2015). This method also generalises to the two-sample setting using the popular inverse-variance weighted approach (Burgess et al., 2013).

Suppose there is a set $\mathcal{L}_J = \{Z_j : j = 1, \dots, J\}$ of SNPs from the exposure GWAS, of which a subset $\mathcal{L}_{J_0} = \{Z_j : j = 1, \dots, J_0\}$, $J_0 \leq J$, is selected as instruments from the outcome GWAS. Suppose we have estimates $\hat{\alpha}_j$ on the log-odds scale of the instrument-exposure relationship $\tilde{\alpha}_j$ for each instrument in \mathcal{L}_J and estimates of the instrument-outcome relationship $\hat{\Gamma}_j$ for each instrument in \mathcal{L}_{J_0} . Additionally, we need the variance $\sigma_{Z_j}^2$ for each instrument in \mathcal{L}_J , which can be obtained from reported allele frequencies. Lastly, we also need estimates for the inverse-variance weights $w_j = \hat{\alpha}_j^2 / \sigma_{\hat{\Gamma}_j}^2$, where $\sigma_{\hat{\Gamma}_j}$ is the standard error of $\hat{\Gamma}_j$. Under the assumption that the instruments in \mathcal{L}_J are mutually independent, the inverse-variance weighted estimator for $\beta_G = \text{cov}(Z, Y) / \text{cov}(Z, G / \sigma_G)$ can be obtained from the above summary statistics as

$$\left(\sum_{j=1}^J \hat{\alpha}_j^2 \sigma_{Z_j}^2 \right)^{1/2} \frac{\sum_{j=1}^{J_0} w_j \hat{\Gamma}_j / \hat{\alpha}_j}{\sum_{j=1}^{J_0} w_j} \quad (9)$$

which is derived in section C of the appendix. We can recover the effect in terms of σ_L (i.e., β_L) by rescaling by a suitable choice of θ^2 as described in Section 3. Conveniently, the second term in (9) is the standard form of the inverse-variance weighted estimator. This means that we can easily readjust existing Mendelian randomisation estimates of coarsened exposures using only the exposure GWAS and a choice for θ^2 . The large-sample distribution of the estimator (9) is derived in section C of the appendix.

5 | REAL DATA EXAMPLES

We can assess the performance of this method in a realistic setting by creating a dichotomised variable from an observed continuous measure, BMI. The idea is to dichotomise BMI at some threshold value and then treat only the dichotomisation as observed. We shall compare the true standardised effect of BMI on some outcome with our procedure described in Section 3 and with the naive approach of using the dichotomisation as the exposure.

Our example is based on the Mendelian randomisation analysis performed in Lyall et al. (2017), which estimates the effect of BMI on several cardiometabolic measures in the UK Biobank cohort. In particular, we look at the effect of BMI on systolic blood pressure. This is a convenient exposure-outcome relationship to estimate because we should not expect there to be threshold effects, that is, the dichotomisations of BMI should have no distinct effects on systolic blood pressure except through BMI itself.

Consistent with Lyall et al. (2017), we use as potential instruments the 93 genome-wide significant SNPs

reported in Locke et al. (2015) available in UK Biobank and we control for age, sex, assessment centre, alcohol intake, smoking status and Townsend deprivation index, along with genetic batch and the first 10 principal components of the genetic relatedness matrix. To avoid weak instrument bias, we prune these SNPs by including those which correlate with BMI with $|t| > 4$ (conditional on the other SNPs) as instruments. We estimate the “true” standardised effect of BMI on systolic blood pressure via two-stage least squares, finding that a one *SD* increase in BMI corresponds to an increase in systolic blood pressure of 1.53 mmHg (95% CI = 0.34–2.72). At each BMI threshold, we then generate a binary variable equal to 1 if an individual’s BMI is above the threshold and 0 otherwise. Treating only this binary measure as observed, we apply the latent variable approach of Section 3.2 using a probit link function.

The results of this example are summarised in Figure 5, which compares the estimated effects with the “true” effect of 1.53. The estimates using the dichotomised measure as the exposure are highly sensitive to the choice of threshold. Since we should not expect there to be distinct threshold effects in this setting, this demonstrates that the dichotomised exposure is not capturing the effect of the latent exposure, instead, it is picking up the shape of the distribution of the environmental risk factor for BMI, as discussed in Section 3.1. As predicted by the bias formula in Section 3.1, the estimates were inflated at the extreme thresholds where the distribution is flatter.

For the latent variable approach, we select a θ^2 of 0.0256 based on the R^2 of our first-stage regression of BMI on the genetic variants. The effect estimate from this approach is much less sensitive to the choice of threshold. Furthermore, the estimates appear to accurately recover the “true” effect of 1.53 regardless of the threshold

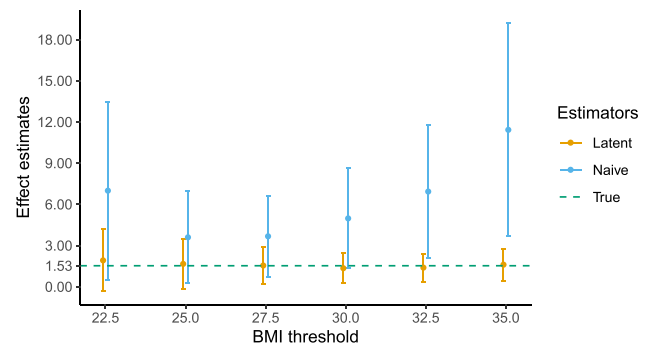


FIGURE 5 Comparison of estimated effect with “true” effect for various BMI thresholds. $N = 70,261$, $\theta^2 = 0.0256$, and 95% confidence intervals are generated over 1000 bootstrap resamples. “True” corresponds to the sample estimate using BMI as the exposure; “naive” corresponds to using the binary measure as the exposure β_D ; and “latent” corresponds to the latent variable estimator β_L of Section 3.2. BMI, body mass index

value, ranging from 1.35 at a BMI cut-off of 30 to 1.92 at a BMI cut-off of 22.5.

We can also investigate this approach in a more realistic setting by reanalysing two existing papers. Box 1 gives an example of how existing two-sample results which do not have interpretable effect sizes can be re-interpreted using this method. The original paper finds that schizophrenia liability increases one's likelihood of using cannabis, although the effect sizes are not interpretable (Pasman et al., 2018). Using our approach, we find that a one *SD* increase in liability corresponds to an OR in the range 1.15–1.26 (95% CI 1.10–1.44) for ever using cannabis. This approach allows us to infer the size of this effect which, in this instance, is very modest.

Box 2 gives an example of how this approach can correct exclusion restriction violations. In the original paper the exposure is self-reported adiposity which is measured on a three-point scale (“thinner,” “plumper,” and “about average”). Genetic instruments will be acting on the underlying measure of child adiposity (e.g., BMI) rather than the three-point scale, so the exclusion restriction is likely to be violated (Richardson et al., 2020). We use our latent variable approach to ameliorate this bias and to estimate the effect of child BMI directly, which is the exposure of interest.

6 | DISCUSSION

We propose a simple framework for estimation and interpretation of Mendelian randomisation for coarsened measurements of latent continuous exposures. We begin by demonstrating in Section 3.1 that using the coarsened measurement as the exposure results in a multiplicative bias which will inflate or deflate effect estimates without reversing their sign. However, under the assumptions of our framework, described in Section 2, we can recover the effect of the latent exposure in terms of *SD* increases. Section 4.4 shows that it is straight-forward to generalise this approach to the two-sample setting. The key sensitivity parameter in our approach is the genetic share of the variance of the latent exposure, which may be estimated or varied over a plausible range of values (Lee et al., 2012). Section 5 evaluates this approach by creating binary exposure measurements from the continuous BMI measure in UK Biobank. We show that we can accurately recover the effect of a *SD* increase in BMI on systolic blood pressure. We also demonstrate this approach in practice by re-analysing two papers which are likely to suffer from this type of exclusion restriction violation, allowing us to meaningfully interpret their effect sizes.

The approach proposed in this paper relies on a number of strong structural assumptions on the

relationship between the latent exposure and its corresponding measurement. The appropriateness of these assumptions must be assessed on a case-by-case basis. Exposure measurements which are defined by strict thresholds of the latent continuous exposure are easiest to conceptualise within this framework. In general, the assumption most difficult to justify is that the thresholds are independent of the genetic share of the latent exposure. One example where this assumption may be violated is self-report measures of mental health status, for example, feelings of depression on a 1–5 scale. Individuals who are genetically predisposed to depression may have different thresholds for reporting their mental wellbeing, either over- or under-reporting.

An additional complication occurs when this method is applied to disease exposures. We have assumed throughout that disease occurrence and disease diagnosis are equivalent; that is, everyone who develops the disease will receive a diagnosis. However, there are often barriers to seeking and accessing the healthcare services needed to receive a diagnosis. These might include stigma surrounding the disease, a lack of trust in healthcare providers or a lack of access to healthcare services due to cost, distance or institutional complexities (Cassim et al., 2019; Stangl et al., 2019). It is therefore possible that individuals with the disease will fail to be diagnosed. This can be viewed as a form of misclassification bias. Misclassification-robust methods for binary exposures could potentially be incorporated into this approach, which we leave for future work (Lewbel, 2000; Rekaya et al., 2016; Smith et al., 2013).

In studies where the assumptions in Section 2 are believed to be implausible, it is important for researchers to be transparent that the magnitude of their effect estimate will not be well-defined.

ACKNOWLEDGEMENTS

This study has been conducted using the UK Biobank Resource. UK Biobank received ethical approval from the Research Ethics Committee (REC reference for UK Biobank is 11/NW/0382). This study was approved as part of application 16391. Matthew J Tudball and Amanda Ly are funded Wellcome Trust studentships [220067/Z/20/Z, 220059/Z/19/Z]. Rachael A Hughes is supported by a Sir Henry Dale Fellowship jointly funded by the Wellcome Trust and the Royal Society [215408/Z/19/Z]. Jack Bowden is funded by an Establishing Excellence in England (E3) grant. George Davey Smith and Marcus Munafò are funded by the Medical Research Council Integrative Epidemiology Unit at the University of Bristol (MC_UU_00011/1, MC_UU_00011/7). For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission.

AUTHOR CONTRIBUTIONS

George Davey Smith and Matthew J. Tudball conceived the idea. Matthew J. Tudball designed the method and performed the analyses. George Davey Smith, Jack Bowden, Kate Tilling, Qingyuan Zhao and Rachael A. Hughes supervised the project. All authors contributed to the main ideas and the writing of the manuscript.

CONFLICT OF INTERESTS

The authors declare that there are no conflict of interests.


DATA AVAILABILITY STATEMENT

A replication kit for the analyses presented in this paper can be obtained from https://github.com/matt-tudball/mrlat_replication. Access to the full genetic and phenotype data from UK Biobank waves 1 and 2 is required to replicate Figures 3 and 5. UK Biobank is an open access resource available to bona fide scientists who are undertaking health-related research that is in the public good. Information regarding access to UK Biobank can be found at <https://www.ukbiobank.ac.uk>.

ORCID

Matthew J. Tudball  <https://orcid.org/0000-0002-7897-6180>

Jack Bowden  <https://orcid.org/0000-0003-2628-3304>

Rachael A. Hughes  <https://orcid.org/0000-0003-0766-1410>

Amanda Ly  <https://orcid.org/0000-0002-4348-5146>

Marcus R. Munafò  <https://orcid.org/0000-0002-4049-993X>

Kate Tilling  <https://orcid.org/0000-0002-1010-8926>

Qingyuan Zhao  <https://orcid.org/0000-0001-9902-2768>

George Davey Smith  <https://orcid.org/0000-0002-1407-8314>

REFERENCES

- Angrist, J. D., & Imbens, G. W. (1995). Two-stage least squares estimation of average causal effects in models with variable treatment intensity. *Journal of the American Statistical Association*, 90(430), 431.
- Brookhart, M. A., & Schneeweiss, S. (2007). Preference-based instrumental variable methods for the estimation of treatment effects: Assessing validity and interpreting results. *International Journal of Biostatistics*, 3(1), 1–19.
- Burgess, S., Butterworth, A., & Thompson, S. G. (2013). Mendelian randomization analysis with multiple genetic variants using summarized data. *Genetic Epidemiology*, 37(7), 658–665.
- Burgess, S., & Labrecque, J. A. (2018). Mendelian randomization with a binary exposure variable: interpretation and presentation of causal estimates. *European Journal of Epidemiology*, 33(10), 947–952.
- Burgess, S., Scott, R. A., Timpson, N. J., Smith, G. D., & Thompson, S. G. (2015). Using published data in Mendelian randomization: A blueprint for efficient identification of causal risk factors. *European Journal of Epidemiology*, 30(7), 543–552.
- Cassim, S., Chepulis, L., Keenan, R., Kidd, J., Firth, M., & Lawrenson, R. (2019). Patient and carer perceived barriers to early presentation and diagnosis of lung cancer: A systematic review 11 medical and health sciences 1117 public health and health services 11 medical and health sciences 1112 oncology and carcinogenesis. *BMC Cancer*, 19(1), 1–14.
- Curnow, R. (1972). The multifactorial model for the inheritance of liability to disease and its implications for relatives at risk. *Biometrics*, 28(4), 931–946.
- Davey Smith, G., & Ebrahim, S. (2003). 'Mendelian randomization': Can genetic epidemiology contribute to understanding environmental determinants of disease? *International Journal of Epidemiology*, 32(1), 1–22.
- Davies, N. M., Gunnell, D., Thomas, K. H., Metcalfe, C., Windmeijer, F., & Martin, R. M. (2013). Physicians' prescribing preferences were a potential instrument for patients' actual prescriptions of antidepressants. *Journal of Clinical Epidemiology*, 66(12), 1386–1396.
- Falconer, D. S. (1965). The inheritance of liability to certain diseases, estimated from the incidence among relatives. *Annals of Human Genetics*, 29(1), 51–76.
- Hemani, G., Zheng, J., Elsworth, B., Wade, K. H., Haberland, V., Baird, D., Laurin, C., Burgess, S., Bowden, J., Langdon, R., Tan, V. Y., Yarmolinsky, J., Shihab, H. A., Timpson, N. J., Evans, D. M., Relton, C., Martin, R. M., Smith, G. D., Gaunt, T. R., & Haycock, P. C. (2018). The MR-base platform supports systematic causal inference across the human phenome. *eLife*, 7, 1–29.
- Ichimura, H. (1993). Semiparametric least squares (SLS) and weighted SLS estimation of single-index models. *Journal of Econometrics*, 58(1–2), 71–120.
- Klein, R. W., & Spady, R. H. (1993). An efficient semiparametric estimator for binary response models. *Econometrica*, 61(2), 387–421.
- Lawn, R. B., Sallis, H. M., Taylor, A. E., Wootton, R. E., Smith, G. D., Davies, N. M., Hemani, G., Fraser, A., Penton-Voak, I. S., & Munafò, M. R. (2019). Schizophrenia risk and reproductive success: A Mendelian randomization study. *Royal Society Open Science*, 6(3).
- Lee, S. H., Goddard, M. E., Wray, N. R., & Visscher, P. M. (2012). A better coefficient of determination for genetic profile analysis. *Genetic Epidemiology*, 36(3), 214–224.
- Lewbel, A. (2000). Identification of the binary choice model with misclassification. *Econometric Theory*, 16(4), 603–609.
- Lewbel, A. (2019). The identification zoo: Meanings of identification in econometrics. *Journal of Economic Literature*, 57(4), 835–903.
- Locke, A. E., Kahali, B., Berndt, S. I., Justice, A. E., Pers, T. H., Day, F. R., Powell, C., Vedantam, S., Buchkovich, M. L., Yang, J., Croteau-Chonka, D. C., Esko, T., Fall, T., Ferreira, T., Gustafsson, S., Kutalik, Z., Luan, J., Mägi, R., Randall, J. C., ... Speliotes, E. K. (2015). Genetic studies of body mass index

- yield new insights for obesity biology. *Nature*, 518(7538), 197–206.
- Lyall, D. M., Celis-Morales, C., Ward, J., Iliodromiti, S., Anderson, J. J., Gill, J. M., Smith, D. J., Ntuk, U. E., Mackay, D. F., Holmes, M. V., Sattar, N., & Pell, J. P. (2017). Association of body mass index with cardiometabolic disease in the UK biobank: A Mendelian randomization study. *JAMA Cardiology*, 2(8), 882–889.
- Marshall, J. (2016). Coarsening bias: How coarse treatment measurement upwardly biases instrumental variable estimates. *Political Analysis*, 24(2), 157–171.
- Martins-Silva, T., Vaz, J. D. S., Hutz, M. H., Salatino-Oliveira, A., Genro, J. P., Hartwig, F. P., Moreira-Maia, C. R., Rohde, L. A., Borges, M. C., & Tovo-Rodrigues, L. (2019). Assessing causality in the association between attention-deficit/hyperactivity disorder and obesity: A Mendelian randomization study. *International Journal of Obesity*, 43, 2500–2508.
- Pasman, J. A., Verweij, K. J., Gerring, Z., Stringer, S., Sanchez-Roige, S., Treur, J. L., Abdellaoui, A., Nivard, M. G., Baselmans, B. M. L., Ong, J.-S., Ip, H. F., van der Zee, M. D., Bartels, M., Day, F. R., Fontanillas, P., Elson, S. L., the 23 and Me Research Team, de Wit, H., Davis, L. K., ... Vink, J. M. (2018). GWAS of lifetime cannabis use reveals new risk loci, genetic overlap with psychiatric traits, and a causal influence of schizophrenia liability. *Nature Neuroscience*, 21(9), 1161–1170.
- Plump, A., & DaveySmith, G. (2019). Identifying and validating new drug targets for stroke and beyond: Can Mendelian randomization help? *Circulation*, 140(10), 831–835.
- Rekaya, R., Smith, S., Hay, E. H., Farhat, N., & Aggrey, S. E. (2016). Analysis of binary responses with outcome-specific misclassification probability in genome-wide association studies. *Application of Clinical Genetics*, 9, 169–177.
- Richardson, T. G., Sanderson, E., Elsworth, B., Tilling, K., & DaveySmith, G. (2020). Use of genetic variation to separate the effects of early and later life adiposity on disease risk: Mendelian randomisation study. *The BMJ*, 369, 1–12.
- Richmond, R. C., Anderson, E. L., Dashti, H. S., Jones, S. E., Lane, J. M., Strand, L. B., Brumpton, B., Rutter, M. K., Wood, A. R., Straif, K., Relton, C. L., Munafò, M., Frayling, T. M., Martin, R. M., Saxena, R., Weedon, M. N., Lawlor, D. A., & Smith, G. D. (2019). Investigating causal relations between sleep traits and risk of breast cancer in women: Mendelian randomisation study. *British Medical Journal*, 365, 1–2.
- Schizophrenia Working Group of the Psychiatric Genomics Consortium. (2015). Biological insights from 108 schizophrenia-associated genetic loci. *Nature*, 511(7510), 421–427.
- Smith, S., Hay, E. H., Farhat, N., & Rekaya, R. (2013). Genome wide association studies in presence of misclassified binary responses. *BMC Genetics*, 14.
- Stangl, A. L., Earnshaw, V. A., Logie, C. H., Van Brakel, W., Simbayi, L. C., Barré, I., & Dovidio, J. F. (2019). The Health Stigma and Discrimination Framework: A global, crosscutting framework to inform research, intervention development, and policy on health-related stigmas. *BMC Medicine*, 17(1), 18–23.
- Sun, Y. Q., Brumpton, B. M., Langhammer, A., Chen, Y., Kvaløy, K., & Mai, X. M. (2019). Adiposity and asthma in adults: A bidirectional Mendelian randomisation analysis of the HUNT Study. *Thorax*, 75, 194–195.

How to cite this article: Tudball MJ, Bowden J, Hughes RA, et al. Mendelian randomisation with coarsened exposures. *Genetic Epidemiology*. 2021;45:338–350. <https://doi.org/10.1002/gepi.22376>

APPENDIX A: IMPORTANCE OF THE IDENTIFYING ASSUMPTIONS

Assumptions 1 and 3 require that the latent exposure and measurement are related by a linear single index model. This assumption imposes considerable structure on the relationship between the two. To see why this assumption is necessary for identification, consider the more general model $L = G - V = \nu(Z, X) - V$, where $\nu(\cdot)$ is some continuous function. D is invariant to any monotone transformation $t(\cdot)$ in the sense that

$$D = I\{\nu(Z, X) \geq V\} = I\{t(\nu(Z, X)) \geq t(V)\} \quad (\text{A1})$$

One such monotone transformation we can take is $t(\cdot) = F_V(\cdot)$, where $F_V(\cdot)$ is the cumulative distribution of V , such that

$$D = I\{F_V(\nu(Z, X)) \geq F_V(V)\} = I\{\text{pr}(D = 1 | Z, X) \geq U\} \quad (\text{A2})$$

where $U \sim \text{Unif}(0, 1)$. This means that the observable joint distribution of (Z, X, D, Y) is consistent with any monotone transformation of $\nu(Z, X)$, including $\text{pr}(D = 1 | Z, X)$ itself. By imposing the structural assumption that $G = \nu(Z, X) = \mu + \alpha Z + \gamma'X$, we reduce the class of models that the observed data distribution is consistent with to $\nu(Z, X)$ which are proportional to G . This allows us to separate G , which is linear in parameters, from the nonlinear link function. In the absence of this linear index assumption, this separation does not occur. This approach to identification is within the class of “identification by functional form” methods described in Lewbel (2019), which provides an overview of this class of methods and discusses their limitations. Section B provides some simulation results when other assumptions fail, namely, correct specification of the link function and independence between the threshold and Z and X , both of which can introduce considerable bias. Bias from misspecification of the link function can be ameliorated by using more flexible semiparametric

TABLE B1 Ratio of estimated to true β_L with link function misspecification

Choice of link function	Value of the skewness parameter a					
	0	1	2	3	4	5
Logistic	1.01	1.02	1.03	1.05	1.06	1.07
	[1.01, 1.02]	[1.01, 1.03]	[1.03, 1.04]	[1.04, 1.06]	[1.05, 1.07]	[1.06, 1.07]
Probit	1.01	1.02	1.03	1.05	1.06	1.07
	[1.00, 1.02]	[1.01, 1.03]	[1.02, 1.04]	[1.04, 1.06]	[1.05, 1.07]	[1.06, 1.08]
Semiparametric*	1.00	1.00	1.00	1.00	1.01	1.01
	[0.99, 1.00]	[0.99, 1.01]	[0.99, 1.01]	[1.00, 1.01]	[1.00, 1.01]	[1.00, 1.02]

*Klein and Spady estimator; mean over 1000 draws; $N = 2500$; $a = 0$; 95% Monte Carlo confidence.

binary outcome estimators (Ichimura, 1993; Klein & Spady, 1993). Independence of the threshold from Z and X is a reasonable assumption when these are genetic factors and D is disease diagnosis or when D is a deterministic categorisation of the latent exposure (i.e., splitting BMI into obesity status).

APPENDIX B: SIMULATING VIOLATIONS OF THE IDENTIFYING ASSUMPTIONS

In this section, we present some simulation results which violate the identifying assumptions stated in Section 2. Our data generating process is as follows:

$$\begin{aligned}
 Z &\stackrel{iid}{\sim} N(0, 1), \quad X \stackrel{iid}{\sim} \text{Exp}(1), \quad V \stackrel{iid}{\sim} SN(0, 1, a) \\
 G &= \alpha_Z Z + \alpha_X X \\
 L &= G - V, \\
 D &= I(L \geq bX) \\
 Y &= \beta_L L + \beta_X X + \beta_V V + \varepsilon, \quad \varepsilon \stackrel{iid}{\sim} N(0, 1)
 \end{aligned}$$

We set parameters $(\alpha_Z, \alpha_X, \beta_L)$ equal to 1, (β_X, β_V) equal to 0.2, normalise Z , X and V to have mean 0 and normalise the variances as follows: $\text{var}(Z) = 2\theta^2/5$, $\text{var}(X) = 3\theta^2/5$ and $\text{var}(V) = 1 - \theta^2$, where $\theta^2 = 0.1$, meaning that $\sigma_L = 1$.

$SN(0, 1, a)$ denotes the skew normal distribution with skewness parameter a . We vary the skewness parameter over a range of values in Table B1. When $a = 0$, this is equivalent to the standard normal distribution, meaning that the probit link will be correctly specified. As V becomes more skewed, the bias increases. This bias can be ameliorated with semi-parametric methods for binary outcomes (Ichimura, 1993; Klein & Spady, 1993).

Another assumption that can be violated is independence between the threshold of D and the observed variables Z and X . This dependence is captured by the parameter b . Since $\alpha_X = 1$, b can be interpreted as the relative contribution of X to the threshold compared to the latent exposure L (e.g., $b = 0.5$ means that X contributes half as much to the threshold as to the latent exposure). In Table B2, we vary the parameter b over a range of values and report the resulting bias. Despite the link function being correctly specified, there is significant bias from dependence in the threshold. Unlike misspecification of the link function, semiparametric techniques cannot correct this bias. When X determines the threshold value, we cannot separately identify G in this framework. This simulation also suggests that

TABLE B2 Ratio of estimated to true β_L with threshold dependence

Choice of link function	Value of the threshold dependence parameter b				
	0	0.1	0.25	0.5	1
Logistic	1.01	1.08	1.17	1.34	1.71
	[1.01, 1.02]	[1.07, 1.08]	[1.16, 1.18]	[1.33, 1.36]	[1.69, 1.72]
Probit	1.01	1.07	1.17	1.33	1.69
	[1, 1.02]	[1.06, 1.08]	[1.16, 1.18]	[1.32, 1.34]	[1.68, 1.70]
Semiparametric*	1.00	1.05	1.15	1.30	1.67
	[0.99, 1]	[1.05, 1.06]	[1.14, 1.16]	[1.29, 1.31]	[1.66, 1.69]

*Klein & Spady estimator; mean over 1000 draws; $N = 2500$; $a = 0$; 95% Monte Carlo confidence intervals.

threshold dependence may be bigger concern in this approach than misspecification of the link function.

APPENDIX C: TWO-SAMPLE ESTIMATOR AND VARIANCE DERIVATION

We begin by deriving Equation (9). For some instrument Z_{ki} in \mathcal{L}_{J_0} , the estimand β_G can be written as

$$\begin{aligned}\sigma_G \beta &= \text{cov}(Z_{ki}, Y) / \text{cov}(Z_{ki}, G / \sigma_G) \\ &= \sigma_G \sigma_{Z_k Y} / \text{cov}(Z_{ki}, G) \\ &= \sigma_G \sigma_{Z_k Y} / \text{cov}\left(Z_{ki}, \sum_{j=1}^J \alpha_j Z_{ji}\right) \\ &= \sigma_G \sigma_{Z_k Y} / \alpha_k \sigma_{Z_k}^2 \\ &= \sigma_{\tilde{G}} \Gamma_k / \tilde{\alpha}_k\end{aligned}\quad (\text{C1})$$

which we can estimate from GWAS summary data. We can use inverse-variance weighting to “meta-analyse” over these estimates for each Z_{ki} in \mathcal{L}_{J_0} , which recovers the estimator (9). Denote $\hat{\beta}$ as the inverse-variance weighted estimator for $\sigma_V \beta$, then our two-sample estimator can be written as

$$\hat{\sigma}_{\tilde{G}} \hat{\beta} = \left(\sum_{j=1}^J \hat{\alpha}_j^2 \sigma_{Z_j}^2 \right)^{1/2} \hat{\beta}. \quad (\text{C2})$$

If we make the common assumption that $\hat{\alpha}_j$ has negligible uncertainty (i.e., $\hat{\alpha}_j \approx \tilde{\alpha}_j$), then we can write an estimator for the variance of (C2) as

$$\left(\sum_{j=1}^J \tilde{\alpha}_j^2 \sigma_{Z_j}^2 \right) \sigma_{\hat{\beta}}^2. \quad (\text{C3})$$