

STATISTICAL INFERENCE IN TWO-SAMPLE SUMMARY-DATA MENDELIAN RANDOMIZATION USING ROBUST ADJUSTED PROFILE SCORE

BY QINGYUAN ZHAO¹, JINGSHU WANG², GIBRAN HEMANI³, JACK BOWDEN⁴ AND DYLAN S. SMALL⁵

¹*Statistical Laboratory, University of Cambridge, qyzhao@statslab.cam.ac.uk*

²*Department of Statistics, University of Chicago, jingshuw@galton.uchicago.edu*

³*MRC Integrative Epidemiology Unit, University of Bristol, g.hemani@bristol.ac.uk*

⁴*University of Exeter Medical School, j.bowden2@exeter.ac.uk*

⁵*Department of Statistics, The Wharton School, University of Pennsylvania, dsmall@wharton.upenn.edu*

Mendelian randomization (MR) is a method of exploiting genetic variation to unbiasedly estimate a causal effect in presence of unmeasured confounding. MR is being widely used in epidemiology and other related areas of population science. In this paper, we study statistical inference in the increasingly popular two-sample summary-data MR design. We show a linear model for the observed associations approximately holds in a wide variety of settings when all the genetic variants satisfy the exclusion restriction assumption, or in genetic terms, when there is no pleiotropy. In this scenario, we derive a maximum profile likelihood estimator with provable consistency and asymptotic normality. However, through analyzing real datasets, we find strong evidence of both systematic and idiosyncratic pleiotropy in MR, echoing the omnigenic model of complex traits that is recently proposed in genetics. We model the systematic pleiotropy by a random effects model, where no genetic variant satisfies the exclusion restriction condition exactly. In this case, we propose a consistent and asymptotically normal estimator by adjusting the profile score. We then tackle the idiosyncratic pleiotropy by robustifying the adjusted profile score. We demonstrate the robustness and efficiency of the proposed methods using several simulated and real datasets.

1. Introduction. A common goal in epidemiology is to understand the causal mechanisms of disease. If it was known that a risk factor causally influenced an adverse health outcome, effort could be focused to develop an intervention (e.g., a drug or public health intervention) to reduce the risk factor and improve the population's health. In settings where evidence from a randomized controlled trial is lacking, inferences about causality are made using observational data. The most common design of observational study is to control for confounding variables between the exposure and the outcome. However, this strategy can easily lead to biased estimates and false conclusions when one or several important confounding variables are overlooked.

Mendelian randomization (MR) is an alternative study design that leverages genetic variation to produce an unbiased estimate of the causal effect even when there is unmeasured confounding. MR is both old and new. It is a special case of the instrumental variable (IV) methods [21], which date back to the 1920s [54] and have a long and rich history in econometrics and statistics. The first MR design was proposed by Katan [33] over 3 decades ago and later popularized in genetic epidemiology by Davey Smith and Ebrahim [18]. As a public health study design, MR is rapidly gaining popularity from just 5 publications in 2003 to

Received February 2018; revised March 2019.

MSC2010 subject classifications. Primary 65J05; secondary 46N60, 62F35.

Key words and phrases. Causal inference, limited information maximum likelihood, weak instruments, errors in variables, path analysis, pleiotropy effects.

over 380 publications in the year 2016 [16]. However, due to the inherent complexity of genetics (the understanding of which is rapidly evolving) and the make-up of large international disease databases being utilized in the analysis, MR has many unique challenges compared to classical IV analyses in econometrics and health studies. Therefore, MR does not merely involve plugging genetic instruments in existing IV methods. In fact, the unique problem structure has sparked many recent methodological advancements [7, 8, 23, 32, 34, 50–52].

Much of the latest developments in Mendelian randomization has been propelled by the increasing availability and scale of genome-wide association studies (GWAS) and other high-throughput genomic data. A particularly attractive proposal is to automate the causal inference by using published GWAS data [14], and a large database and software platform is currently being developed [28]. Many existing IV and MR methods (e.g., [23, 40, 50]), though theoretically sound and robust to different kinds of biases, require having individual-level data. Unfortunately, due to privacy concerns, the access to individual-level genetic data is almost always restricted and usually only the GWAS summary statistics are publicly available. This data structure has sparked a number of new statistical methods anchored within the framework of meta-analysis (e.g., [7, 8, 26]). They are intuitively simple and can be conveniently used with GWAS summary data, thus are quickly gaining popularity in practice. However, the existing summary-data MR methods often make unrealistic simplifying assumptions and generally lack theoretical support such as statistical consistency and asymptotic sampling distribution results.

This paper aims to resolve this shortcoming by developing statistical methods that can be used with summary data, have good theoretical properties, and are robust to deviations of the usual IV assumptions. In the rest of the **Introduction**, we will introduce a statistical model for GWAS summary data and demonstrate the MR problem using a real data example. This example will be repeatedly used in subsequent sections to motivate and illustrate the statistical methods. We will conclude the **Introduction** by discussing the methodological challenges in MR and outlining our solution.

1.1. Two-sample MR with summary data. We are interested in estimating the causal effect of an exposure variable X on an outcome variable Y . The causal effect is confounded by unobserved variables, but we have p genetic variants (single nucleotide polymorphisms, SNPs), Z_1, Z_2, \dots, Z_p , that are approximately *valid* instrumental variables (validity of an IV is defined in Section 2.1). These IVs can help us to obtain unbiased estimate of the causal effect even when there is unmeasured confounding. The precise problem considered in this paper is two-sample Mendelian randomization with summary data, where we observe, for SNP $j = 1, \dots, p$, two associational effects: the SNP-exposure effect $\hat{\gamma}_j$ and the SNP-outcome effect $\hat{\Gamma}_j$. These estimated effects are usually computed from two different samples using a simple linear regression or logistic regression and are or are becoming available in public domain.

Throughout the paper, we assume the following.

ASSUMPTION 1. For every $j \in \{1, \dots, p\} := [p]$, $\hat{\gamma}_j \sim N(\gamma_j, \sigma_{X_j}^2)$, $\hat{\Gamma}_j \sim N(\Gamma_j, \sigma_{Y_j}^2)$, and the variances $(\sigma_{X_j}^2, \sigma_{Y_j}^2)_{j \in [p]}$ are known. Furthermore, the $2p$ random variables $(\hat{\gamma}_j)_{j \in [p]}$ and $(\hat{\Gamma}_j)_{j \in [p]}$ are mutually independent.

The first assumption is quite reasonable as typically there are hundreds of thousands of samples in modern GWAS, making the normal approximation very accurate. We assume the variances of the GWAS marginal coefficients are computed very accurately using the individual data (as they are typically based on tens of thousands of samples), but the methods

developed in this paper do not utilize individual data for statistical inference. The independence between $(\hat{\gamma}_j)_{j \in [p]}$ and $(\hat{\Gamma}_j)_{j \in [p]}$ is guaranteed because the effects are computed from independent samples. The independence across SNPs is reasonable if we only use uncorrelated SNPs by using a tool called linkage disequilibrium (LD) clumping [28, 43, 44]. See Section 2 for more justifications of the last assumption.

Our key modeling assumption for summary-data MR is the following.

MODEL FOR GWAS SUMMARY DATA. *There exists a real number β_0 such that*

$$(1.1) \quad \Gamma_j \approx \beta_0 \gamma_j \quad \text{for almost all } j \in [p].$$

In Section 2 and Appendix A of the Online Supplement [57], we explain why this model likely holds for a variety of situations and why the parameter β_0 may be interpreted as the causal effect of X on Y . However, by investigating a real data example, we will demonstrate in Section 3.5 that it is very likely that the strict equality $\Gamma_j = \beta_0 \gamma_j$ is not true for some if not most j . For now, we will proceed with the loose statement in (1.1), but it will be soon made precise in several ways.

Assumption 1 and model (1.1) suggest two different strategies of estimating β_0 :

1. Use the Wald ratio $\hat{\beta}_j = \hat{\Gamma}_j / \hat{\gamma}_j$ [53] as each SNP's individual estimate of β_0 , then aggregate the estimates using a robust meta-analysis method. Most existing methods for summary-data MR follow this line [7, 8, 26]; however, the Wald estimator $\hat{\beta}_j$ is heavily biased when γ_j is small, a phenomenon known as “weak instrument bias.” See Bound, Jaeger and Baker [6] and Section 1.3 below.

2. Treat equation (1.1) as an errors-in-variables regression problem [15], where we are regressing $\hat{\Gamma}_j$, whose expectation is Γ_j , on $\hat{\gamma}_j$, which can be regarded as a noisy observation of the actual regressor γ_j . Then we directly estimate β_0 in a robust way. This is the novel approach taken in this paper and will be described and tested in detail.

1.2. *A motivating example.* Next, we introduce a real data example that will be repeatedly used in the development of this paper. In this example, we are interested in estimating the causal effect of a person's Body Mass Index (BMI) on Systolic Blood Pressure (SBP). We obtained publicly available summary data from three GWAS with nonoverlapping samples:

BMI-FEM: BMI in females by the Genetic Investigation of ANthropometric Traits (GIANT) consortium [35] (sample size: 171977, unit: kg/m^2).

BMI-MAL: BMI in males in the same study by the GIANT consortium (sample size: 152893, unit: kg/m^2).

SBP-UKBB: SBP using the United Kingdom BioBank (UKBB) data (sample size: 317754, unit: mmHg).

Using the BMI-FEM dataset and LD clumping, we selected 25 SNPs that are genome-wide significant (p -value $\leq 5 \times 10^{-8}$) and uncorrelated (10,000 kilo base pairs apart and $R^2 \leq 0.001$). We then obtained the 25 SNP-exposure effects $(\hat{\gamma}_j)_{j=1}^{25}$ and the corresponding standard errors from BMI-MAL and the SNP-outcome effects $(\hat{\Gamma}_j)_{j=1}^{25}$ and the corresponding standard errors from SBP-UKBB. Later on in the paper we will consider an expanded set of 160 SNPs using the selection threshold p -value $\leq 10^{-4}$.

Figure 1 shows the scatter plot of the 25 pairs of genetic effects. Since they are measured with error, we added error bars of one standard error to every point on both sides. The goal of summary-data MR is to find a straight line through the origin that best fits these points. The statistical method should also be robust to violations of model (1.1) since not all SNPs satisfy the relation $\Gamma_j = \beta_0 \gamma_j$ exactly. We will come back to this example in Sections 3.5, 4.4 and 5.3 to illustrate our methods.

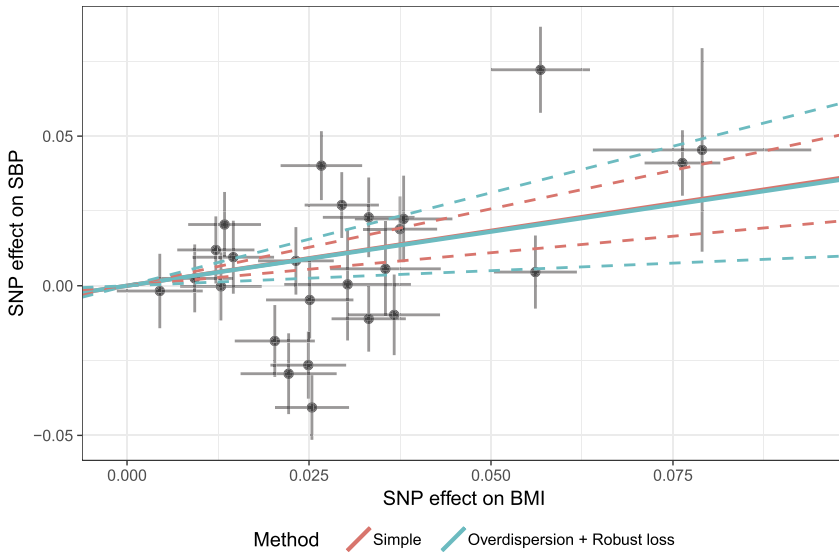


FIG. 1. Scatter plot of $\hat{\Gamma}_j$ versus $\hat{\gamma}_j$ in the BMI-SBP example. Each point is augmented by the standard error of $\hat{\Gamma}_j$ and $\hat{\gamma}_j$ on the vertical and horizontal sides. For presentation purposes only, we chose the allele codings so that all $\hat{\gamma}_j$ are positive. Solid lines are the regression slope fitted by two of our methods. Dashed lines are the 95% confidence interval of the slopes. The simple method using unadjusted profile score (PS, described in Section 3) has smaller standard error than the more robust method using robust adjusted profile score (RAPS, described in Section 5), because the simple method does not consider genetic pleiotropy. See also Section 3.5.

1.3. *Statistical challenges and organization of the paper.* Compared to classical IV analyses in econometrics and health studies, there are many unique challenges in two-sample MR with summary data:

1. **Measurement error:** Both the SNP-exposure and SNP-outcome effects are clearly measured with error, but most of the existing methods applicable to summary data assume that the sampling error of $\hat{\gamma}_j$ is negligible so a weighted linear regression can be directly used [13].

2. **Invalid instruments due to pleiotropy** (the phenomenon that one SNP can affect seemingly unrelated traits): A SNP Z_j may causally affect the outcome Y through other pathways not involving the exposure X . In this case, the approximate linear model $\Gamma_j \approx \beta_0 \gamma_j$ might be entirely wrong for some SNPs.

3. **Weak instruments:** Including a SNP j with very small γ_j can bias the causal effect estimates (especially when the meta-analysis strategy is used). It can also increase the variance of the estimator $\hat{\beta}$. See Section 3.4.2.

4. **Selection bias:** To avoid the weak instrument bias, the standard practice in MR is to only use the genome-wide significant SNPs as instruments (e.g., as implemented in the `TwoSampleMR` R package [28]). However, in many studies the same dataset is used for both selecting SNPs and estimating γ_j , resulting in substantial selection bias even if the selection threshold is very stringent.

Many previous works have considered one or some of these challenges. Bowden et al. [10] proposed a modified Cochran’s Q statistic to detect the heterogeneity due to pleiotropy instead of measurement error in $\hat{\gamma}_j$. Addressing the issue of bias due to pleiotropy has attracted lots of attention in the summary-data MR literature [7, 8, 26, 34, 51, 52], but no solid statistical underpinning has yet been given. Other methods with more rigorous statistical theory require individual-level data [23, 40, 50]. The weak instrument problem has been thoroughly studied in the econometrics literature (e.g., [6, 25, 49]), but all of this work operates in the

individual-level data setting. Finally, the selection bias has largely been overlooked in practice; common wisdom has been that the selection biases the causal effects toward the null (so it might be less serious) [27] and the bias is perhaps small when a stringent selection criterion is used (in Section 7 we show this is not necessarily the case).

In this paper, we develop a novel approach to overcome all the aforementioned challenges by adjusting the profile likelihood of the summary data. The measurement errors of $\hat{\gamma}_j$ and $\hat{\Gamma}_j$ (challenge 1) are naturally incorporated in computing the profile score. To tackle invalid IVs (challenge 2), we will consider three models for the GWAS summary data with increasing complexity:

MODEL 1 (No pleiotropy). *The linear model $\Gamma_j = \beta_0\gamma_j$ is true for every $j \in [p]$.*

MODEL 2 (Systematic pleiotropy). *Assume $\alpha_j = \Gamma_j - \beta_0\gamma_j \stackrel{i.i.d.}{\sim} N(0, \tau_0^2)$ for $j \in [p]$ and some small τ_0^2 .*

MODEL 3 (Systematic and idiosyncratic pleiotropy). *Assume $\alpha_j, j \in [p]$ are from a contaminated normal distribution: most α_j are distributed as $N(0, \tau_0^2)$ but some $|\alpha_j|$ may be much larger.*

The consideration of these three models is motivated by not only the theoretical models in Section 2 but also characteristics observed in real data (Sections 3.5, 4.4 and 5.3) and recent empirical evidence in genetics [12, 46].

The three models are considered in Sections 3 to 5, respectively. We will propose estimators that are provably consistent and asymptotically normal in Models 1 and 2. We will then derive an estimator that is robust to a small proportion of outliers in Model 3. We believe Model 3 best explains the real data and the corresponding Robust Adjusted Profile Score (RAPS) estimator is the clear winner in all the empirical examples.

Although weak IVs may bias the individual Wald's ratio estimator (challenge 3), we will show, both theoretically and empirically, that including additional weak IVs is usually helpful for our new estimators when there are already strong IVs or many weak IVs. Finally, the selection bias (challenge 4) is handled by requiring use of an independent dataset for IV selection as we have done in Section 1.2. This might not be possible in all practical problems, but failing to use a separate dataset for IV selection can lead to severe selection bias as illustrated by an empirical example in Section 7.

The rest of the paper is organized as follows. In Section 2, we give theoretical justifications of the model (1.1) for GWAS summary data. Then in Sections 3 to 5 we describe an adjusted profile score approach of statistical inference in Models 1 to 3, respectively. The paper is concluded with simulation examples in Section 6, another real data example in Section 7 and more discussion in Section 8.

2. Statistical model for MR. In this section, we explain why the approximate linear model (1.1) for GWAS summary data may hold in many MR problems. We will put structural assumptions on the original data and show that (1.1) holds in a variety of scenarios. Owing to this heuristic and the wide availability of GWAS summary datasets, we will focus on statistical inference for summary-data MR after Section 2.

2.1. Validity of instrumental variables. In order to study the origin of the linear model (1.1) for summary data and give a causal interpretation to the parameter β_0 , we must specify

how the original data (X, Y, Z_1, \dots, Z_p) are generated and how the summary statistics are computed. Consider the following structural equation model [42] for the random variables:

$$(2.1) \quad \begin{aligned} X &= g(Z_1, \dots, Z_p, U, E_X), \quad \text{and} \\ Y &= f(X, Z_1, \dots, Z_p, U, E_Y), \end{aligned}$$

where U is the unmeasured confounder, E_X and E_Y are independent random noises, $(E_X, E_Y) \perp (Z_1, \dots, Z_p, U)$ and $E_X \perp E_Y$. In two-sample MR, we observe n_X i.i.d. realizations of (X, Z_1, \dots, Z_p) and independently n_Y i.i.d. realizations of (Y, Z_1, \dots, Z_p) . We shall also assume that the SNPs Z_1, Z_2, \dots, Z_p are discrete random variables supported on $\{0, 1, 2\}$ and are mutually independent. To ensure the independence, in practice we only include SNPs with low pairwise LD score in our model by using standard genetics software like LD clumping [43].

A variable Z_j is called a *valid IV* if it satisfies the following three criteria:

1. Relevance: Z_j is associated with the exposure X . Notice that a SNP that is correlated (in genetics terminology, in LD) with the actual causal variant is also considered relevant and does not affect the statistical analysis below.
2. Effective random assignment: Z_j is independent of the unmeasured confounder U .
3. Exclusion restriction: Z_j only affects the outcome Y through the exposure X . In other words, the function f does not depend on Z_j .

The causal model and the IV conditions are illustrated by a directed acyclic graph (DAG) with a single instrument Z_1 in Figure 2. Readers who are unfamiliar with this language may find the tutorial by Baiocchi, Cheng and Small [3] helpful.

In Mendelian randomization, the first criterion—relevance—is easily satisfied by selecting SNPs that are significantly associated with X . Notice that the genetic instrument does not need to be a causal SNP for the exposure. The first criterion is considered satisfied if the SNP is correlated with the actual causal SNP [29]. For example, in Figure 2, Z_1 would be considered “relevant” even if it is not causal for X but it is correlated with \tilde{Z}_1 . Aside from the effects of population stratification, the second independence to unmeasured confounder assumption is usually easy to justify because most of the common confounders in epidemiology are postnatal, which are independent of genetic variants governed by Mendel’s second law of independent assortment [18, 20]. Empirically, there is generally a lack of confounding of genetic variants with factors that confound exposures in conventional observational epidemiological studies [19].

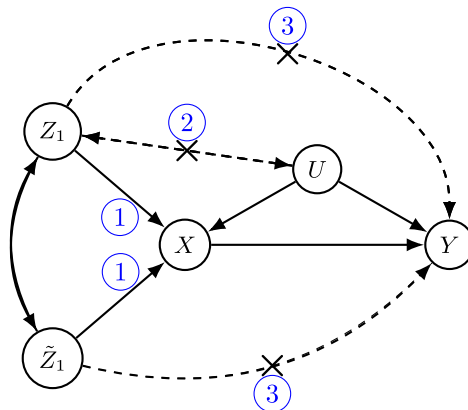


FIG. 2. Causal DAG and the three criteria for valid IV. The proposed IV Z_1 can either be a causal variant for X or correlated with a causal variant (\tilde{Z}_1 in the figure). Z_1 must be independent of any unmeasured confounder U and cannot have any direct effect on Y or be correlated with another variant that has direct effect on Y .

The main concern for Mendelian randomization is the possible violation of the third exclusion restriction criterion, due to a genetic phenomenon called pleiotropy [18, 47], a.k.a. the multifunction of genes. The exclusion restriction assumption does not hold if a SNP Z_j affects the outcome Y through multiple causal pathways and some do not involve the exposure X . It is also violated if Z_j is correlated with other variants (such as \tilde{Z}_1 in Figure 2) that affect Y through pathways that does not involve X . Pleiotropy is widely prevalent for complex traits [48]. In fact, a “universal pleiotropy hypothesis” developed by Fisher [22] and Wright [55] theorizes that every genetic mutation is capable of affecting essentially all traits. Recent genetics studies have found strong evidence that there is an extremely large number of causal variants with tiny effect sizes on many complex traits, which in part motivates our random effects Model 2.

Another important concept is the strength of an IV, defined as its association with the exposure X and usually measured by the F -statistic of an instrument-exposure regression. Since we assume all the genetic instruments are independent, the strength of SNP j can be assessed by comparing the statistic $\hat{\gamma}_j^2/\sigma_{\tilde{X}_j}^2$ with the quantiles of χ_1^2 (or equivalently $F_{1,\infty}$). When only a few weak instruments are available (e.g., F -statistic less than 10), the usual asymptotic inference is quite problematic [6]. In this paper, we primarily consider the setting where there is at least one strong IV or many weak IVs.

2.2. Linear structural model. We are now ready to derive the linear model (1.1) for GWAS summary data. Assuming all the IVs are valid, we start with the linear structural model where functions f and g in (2.1) are linear in their arguments (see also Bowden et al. [9]):

$$(2.2) \quad X = \sum_{j=1}^p \gamma_j Z_j + \eta_X U + E_X, \quad Y = \beta X + \eta_Y U + E_Y.$$

In this case, the GWAS summary statistics $(\hat{\gamma}_j)_{j \in [p]}$ and $(\hat{\Gamma}_j)_{j \in [p]}$ are usually computed from simple linear regressions:

$$\hat{\gamma}_j = \frac{\widehat{\text{Cov}}_{n_X}(X, Z_j)}{\widehat{\text{Cov}}_{n_X}(Z_j, Z_j)}, \quad \hat{\Gamma}_j = \frac{\widehat{\text{Cov}}_{n_Y}(Y, Z_j)}{\widehat{\text{Cov}}_{n_Y}(Z_j, Z_j)}.$$

Here, $\widehat{\text{Cov}}_n$ is the sample covariance operator with n i.i.d. samples. Using (2.2), it is easy to show that $\hat{\gamma}_j$ and $\hat{\Gamma}_j$ converge to normal distributions centered at γ_j and $\Gamma_j = \beta\gamma_j$.

However, $\hat{\gamma}_j$ and $\hat{\gamma}_k$ are not exactly uncorrelated when $j \neq k$ (same for $\hat{\Gamma}_j$ and $\hat{\Gamma}_k$), even if Z_j and Z_k are independent. After some simple algebra, one can show that

$$\text{Cor}^2(\hat{\gamma}_j, \hat{\gamma}_k) = 4 \cdot \frac{\gamma_j^2 \text{Var}(Z_j)}{\text{Var}(X) - \gamma_j^2 \text{Var}(Z_j)} \frac{\gamma_k^2 \text{Var}(Z_k)}{\text{Var}(X) - \gamma_k^2 \text{Var}(Z_k)}.$$

Notice that $\gamma_j^2 \text{Var}(Z_j)/\text{Var}(X)$ is the proportion of variance of X explained by Z_j . In the genetic context, a single SNP usually has very small predictability of a complex trait [12, 31, 41, 46]. Therefore, the correlation between $\hat{\gamma}_j$ and $\hat{\gamma}_k$ (similarly $\hat{\Gamma}_j$ and $\hat{\Gamma}_k$) is almost negligible. In conclusion, the linear model (1.1) is approximately true when the phenotypes are believed to be generated from a linear structural model.

To stick to the main statistical methodology, we postpone additional justifications of (1.1) in nonlinear structural models to Appendix A. In Appendix A.1, we will investigate the case where Y is binary and $\hat{\Gamma}_j$ is obtained via logistic regression, as is very often the case in applied MR investigations. In Appendix A.2, we will show the linearity between X and Z is also not necessary.

2.3. *Violations of exclusion restriction.* Equation (2.2) assumes that all the instruments are valid. In reality, the exclusion restriction assumption is likely violated for many if not most of the SNPs. To investigate its impact in the model for summary data, we consider the following modification of the linear structural model (2.2):

$$(2.3) \quad X = \sum_{j=1}^p \gamma_j Z_j + \eta_X U + E_X, \quad Y = \beta X + \sum_{j=1}^p \alpha_j Z_j + \eta_Y U + E_Y.$$

The difference between (2.2) and (2.3) is that the SNPs are now allowed to directly affect Y and the effect size of SNP Z_j is α_j . In this case, it is not difficult to see that the regression coefficient $\hat{\Gamma}_j$ estimates $\Gamma_j = \alpha_j + \gamma_j \beta$. This inspires our Models 2 and 3. In Model 2, we assume the direct effects α_j are normally distributed random effects. In Model 3, we further require the statistical procedure to be robust against any extraordinarily large direct effects α_j . See Section 8 for more discussion on the assumptions on the pleiotropy effects.

3. No pleiotropy: A profile likelihood approach. We now consider Model 1, the case with no pleiotropy effects.

3.1. *Derivation of the profile likelihood.* A good place to start is writing down the likelihood of GWAS summary data. Up to some additive constant, the log-likelihood function is given by

$$(3.1) \quad l(\beta, \gamma_1, \dots, \gamma_p) = -\frac{1}{2} \left[\sum_{j=1}^p \frac{(\hat{\gamma}_j - \gamma_j)^2}{\sigma_{X_j}^2} + \sum_{j=1}^p \frac{(\hat{\Gamma}_j - \gamma_j \beta)^2}{\sigma_{Y_j}^2} \right].$$

Since we are only interested in estimating β_0 , the other parameters, namely $\boldsymbol{\gamma} := (\gamma_1, \dots, \gamma_p)$, are considered nuisance parameters. There are two ways to proceed from here. One is to view $\boldsymbol{\gamma}$ as *incidental* parameters [39] and try to eliminate them from the likelihood. The other approach is to assume the sequence $\gamma_1, \gamma_2, \dots$ is generated from a fixed unknown distribution. When p is large, it is possible to estimate the distribution of $\boldsymbol{\gamma}$ to improve the efficiency using the second approach [38]. In this paper, we aim to develop a general method for summary-data MR that can be used regardless of the number of SNPs being used, so we will take the first approach.

The profile log-likelihood of β is given by profiling out $\boldsymbol{\gamma}$ in (3.1):

$$(3.2) \quad l(\beta) = \max_{\boldsymbol{\gamma}} l(\beta, \boldsymbol{\gamma}) = -\frac{1}{2} \sum_{j=1}^p \frac{(\hat{\Gamma}_j - \beta \hat{\gamma}_j)^2}{\sigma_{X_j}^2 \beta^2 + \sigma_{Y_j}^2}.$$

The maximum likelihood estimator of β is given by $\hat{\beta} = \arg \max_{\beta} l(\beta)$. It is also called a Limited Information Maximum Likelihood (LIML) estimator in the IV literature, a method due to Anderson and Rubin [2] with good consistency and efficiency properties. See also Pacini and Windmeijer [40].

Equation (3.2) can be interpreted as a linear regression of $\hat{\Gamma}$ on $\hat{\gamma}$, with the intercept of the regression fixed to zero and the variance of each observation equaling to $\sigma_{X_j}^2 \beta^2 + \sigma_{Y_j}^2$. There is another meta-analysis interpretation. Let $\hat{\beta}_j = \hat{\Gamma}_j / \hat{\gamma}_j$ be the individual Wald’s ratio, then (3.2) can be rewritten as

$$(3.3) \quad l(\beta) = -\frac{1}{2} \sum_{j=1}^p \frac{(\hat{\beta}_j - \beta)^2}{\sigma_{X_j}^2 \beta^2 / \hat{\gamma}_j^2 + \sigma_{Y_j}^2 / \hat{\gamma}_j^2}.$$

This expression is also derived by Bowden et al. [10] by defining a generalized version of Cochran’s Q statistic to test for the presence of pleiotropy that takes into account uncertainty in $\hat{\gamma}_j$.

3.2. *Consistency and asymptotic normality.* It is well known that the maximum likelihood estimator can be inconsistent when there are many nuisance parameters in the problem (e.g., [39]). Nevertheless, due to the connection with LIML, we expect and will prove below that $\hat{\beta}$ is consistent and asymptotically normal. However, we will also show that the profile likelihood (3.2) can be information biased [37], meaning the profile likelihood ratio test does not generally have a χ_1^2 limiting distribution under the null.

A major distinction between our asymptotic setting and the classical errors-in-variables regression setting is that our “predictors” $\hat{\gamma}_j$, $j \in [p]$ can be individually weak. This can be seen, for example, from the linear structural model (2.2) that

$$(3.4) \quad \text{Var}(X) = \sum_{j=1}^p \gamma_j^2 \text{Var}(Z_j) + \eta_X^2 \text{Var}(U) + \text{Var}(E_X).$$

Note that Z_j takes on the value 0, 1, 2 with probability $p_j^2, 2p_j(1 - p_j), (1 - p_j)^2$ where p_j is the allele frequency of SNP j . For simplicity, we assume p_j is bounded away from 0 and 1. In other words, only common genetic variants are used as IVs. Together with (3.4), this implies that, if $\text{Var}(X)$ exists, $\|\boldsymbol{\gamma}\|_2$ is bounded.

ASSUMPTION 2 (Collective IV strength is bounded). $\|\boldsymbol{\gamma}\|_2^2 = O(1)$.

As a consequence, the average effect size is decreasing to 0,

$$\frac{1}{p} \sum_{j=1}^p |\gamma_j| \leq \|\boldsymbol{\gamma}\|_2 / \sqrt{p} \rightarrow 0, \quad \text{when } p \rightarrow \infty.$$

This is clearly different from the usual linear regression setting where the “predictors” $\hat{\gamma}_j$ are viewed as random samples from a population. In the one-sample IV literature, this many weak IV setting ($p \rightarrow \infty$) has been considered by Bekker [5], Stock and Yogo [49], Hansen, Hausman and Newey [25] among many others in econometrics.

Another difference between our asymptotic setting and the errors-in-variables regression is that our measurement errors also converge to 0 as the sample size converges to infinity. Recall that n_X is the sample size of (X, Z_1, \dots, Z_p) and n_Y is the sample size of (Y, Z_1, \dots, Z_p) . We assume the following.

ASSUMPTION 3 (Variance of measurement error). Let $n = \min(n_X, n_Y)$. There exist constants c_σ, c'_σ such that $c_\sigma/n \leq \sigma_{Xj}^2 \leq c'_\sigma/n$ and $c_\sigma/n \leq \sigma_{Yj}^2 \leq c'_\sigma/n$ for all $j \in [p]$.

We write $a = O(b)$ if there exists a constant $c > 0$ such that $|a| \leq cb$, and $a = \Theta(b)$ if there exists $c > 0$ such that $c^{-1}b \leq |a| \leq cb$. In this notation, Assumption 3 assumes the known variances σ_{Xj}^2 and σ_{Yj}^2 are $\Theta(1/n)$.

In the linear structural model (2.2), $\text{Var}(\hat{\gamma}_j) \leq \text{Var}(X)/[\text{Var}(Z_j)/n_X]$. Thus Assumption 3 is satisfied when only common variants are used.

We are ready to state our first theoretical result.

THEOREM 3.1. *In Model 1 and under Assumptions 1 to 3, if $p/(n^2\|\boldsymbol{\gamma}\|_2^4) \rightarrow 0$, the maximum likelihood estimator $\hat{\beta}$ is statistically consistent, that is, $\hat{\beta} \xrightarrow{P} \beta_0$.*

A crucial quantity in Theorem 3.1 and the analysis below is the average strength of the IVs, defined as

$$\kappa = \frac{1}{p} \sum_{j=1}^p \frac{\gamma_j^2}{\sigma_{Xj}^2} = \Theta(n\|\boldsymbol{\gamma}\|_2^2/p).$$

An unbiased estimator of κ is the average F -statistic minus 1,

$$\hat{\kappa} = \frac{1}{P} \sum_{j=1}^p \frac{\hat{\gamma}_j^2}{\sigma_{\hat{X}_j}^2} - 1.$$

In practice, we require the average F -statistic to be large (say > 100) when p is small, or not too small (say > 3) when p is large. Thus the condition $p/(n^2 \|\boldsymbol{\gamma}\|_2^4) = \Theta(1/(p\kappa^2)) \rightarrow 0$ in Theorem 3.1 is usually quite reasonable. In particular, since this condition only depends on the average instrument strength κ , the estimator $\hat{\beta}$ remains consistent even if a substantial proportion of $\gamma_j = 0$ (e.g., if the selection step in Section 1.2 using BMI-FEM with less stringent p -value threshold finds many false positives).

Next, we study the asymptotic normality of $\hat{\beta}$. Define the *profile score* to be the derivative of the profile log-likelihood:

$$(3.5) \quad \psi(\beta) := -l'(\beta) = \sum_{j=1}^p \frac{(\hat{\Gamma}_j - \beta \hat{\gamma}_j)(\hat{\Gamma}_j \sigma_{\hat{X}_j}^2 \beta + \hat{\gamma}_j \sigma_{\hat{Y}_j}^2)}{(\sigma_{\hat{X}_j}^2 \beta^2 + \sigma_{\hat{Y}_j}^2)^2}.$$

The maximum likelihood estimator $\hat{\beta}$ solves the estimating equation $\psi(\hat{\beta}) = 0$, and we consider the Taylor expansion around the truth β_0 :

$$(3.6) \quad 0 = \psi(\hat{\beta}) = \psi(\beta_0) + \psi'(\beta_0)(\hat{\beta} - \beta_0) + \frac{1}{2} \psi''(\tilde{\beta})(\hat{\beta} - \beta_0)^2,$$

where $\tilde{\beta}$ is between $\hat{\beta}$ and β_0 . Since $\hat{\beta}$ is statistically consistent, the last term on the right-hand side of (3.6) can be proved to be negligible, and the asymptotic normality of $\hat{\beta}$ can be established by showing, for some appropriate V_1 and V_2 , $\psi(\beta_0) \xrightarrow{d} N(0, V_1)$ and $\psi'(\beta_0) \xrightarrow{P} -V_2$. When $V_1 = V_2$, the profile likelihood/score is called *information unbiased* [37].

THEOREM 3.2. *Under the assumptions in Theorem 3.1 and if at least one of the following two conditions are true: (1) $p \rightarrow \infty$ and $\|\boldsymbol{\gamma}\|_3/\|\boldsymbol{\gamma}\|_2 \rightarrow 0$; (2) $\kappa \rightarrow \infty$; then we have*

$$(3.7) \quad \frac{V_2}{\sqrt{V_1}}(\hat{\beta} - \beta_0) \xrightarrow{d} N(0, 1),$$

where

$$(3.8) \quad V_1 = \sum_{j=1}^p \frac{\gamma_j^2 \sigma_{\hat{Y}_j}^2 + \Gamma_j^2 \sigma_{\hat{X}_j}^2 + \sigma_{\hat{X}_j}^2 \sigma_{\hat{Y}_j}^2}{(\sigma_{\hat{X}_j}^2 \beta_0^2 + \sigma_{\hat{Y}_j}^2)^2}, \quad V_2 = \sum_{j=1}^p \frac{\gamma_j^2 \sigma_{\hat{Y}_j}^2 + \Gamma_j^2 \sigma_{\hat{X}_j}^2}{(\sigma_{\hat{X}_j}^2 \beta_0^2 + \sigma_{\hat{Y}_j}^2)^2}.$$

Notice that Theorem 3.2 is very general. It can be applied even in the extreme situation p is fixed and $\kappa \rightarrow \infty$ (a few strong IVs) or $p \rightarrow \infty$ and $\kappa \rightarrow 0$ (many very weak IVs). The assumption $\|\boldsymbol{\gamma}\|_3/\|\boldsymbol{\gamma}\|_2 \rightarrow 0$ is used to verify a Lyapunov’s condition for a central limit theorem. It essentially says the distribution of IV strengths is not too uneven and this assumption can be further relaxed.

Using our rate assumption for the variances (Assumption 3), $V_2 = \Theta(n \|\boldsymbol{\gamma}\|_2^2) = \Theta(p\kappa)$ and $V_1 = V_2 + \Theta(p)$. This suggests that the profile likelihood is information unbiased if and only if $\kappa \rightarrow \infty$. In general, the amount of information bias depends on the instrument strength κ . As an example, suppose $\beta_0 = 0$ and $\sigma_{\hat{Y}_j}^2 \equiv \sigma_{\hat{Y}_1}^2$. Then by (3.7) and (3.8), $\text{Var}(\hat{\beta}) \approx V_1/V_2^2 = (1 + \kappa^{-1})/V_2$. Alternatively, if we make the simplifying assumption that $\sigma_{\hat{Y}_j}^2/\sigma_{\hat{X}_j}^2$ does not depend on j , it is straightforward to show that

$$\text{Var}(\hat{\beta}) \propto \frac{1 + \kappa^{-1}}{p\kappa}.$$

This approximation can be used as a rule of thumb to select the optimal number of IVs.

In order to obtain standard error of $\hat{\beta}$, we must estimate V_1 and V_2 using the GWAS summary data. We propose to replace γ_j^2 and Γ_j^2 in (3.8) by their unbiased sample estimates, $\hat{\gamma}_j^2 - \sigma_{\hat{X}_j}^2$ and $\hat{\Gamma}_j^2 - \sigma_{\hat{Y}_j}^2$:

$$\hat{V}_1 = \sum_{j=1}^p \frac{(\hat{\gamma}_j^2 - \sigma_{\hat{X}_j}^2)\sigma_{\hat{Y}_j}^2 + (\hat{\Gamma}_j^2 - \sigma_{\hat{Y}_j}^2)\sigma_{\hat{X}_j}^2 + \sigma_{\hat{X}_j}^2\sigma_{\hat{Y}_j}^2}{(\sigma_{\hat{X}_j}^2\hat{\beta}^2 + \sigma_{\hat{Y}_j}^2)^2},$$

$$\hat{V}_2 = \sum_{j=1}^p \frac{(\hat{\gamma}_j^2 - \sigma_{\hat{X}_j}^2)\sigma_{\hat{Y}_j}^2 + (\hat{\Gamma}_j^2 - \sigma_{\hat{Y}_j}^2)\sigma_{\hat{X}_j}^2}{(\sigma_{\hat{X}_j}^2\hat{\beta}^2 + \sigma_{\hat{Y}_j}^2)^2}.$$

THEOREM 3.3. *Under the same assumptions in Theorem 3.2, we have $\hat{V}_1 = V_1(1 + o_p(1))$, $\hat{V}_2 = V_2(1 + o_p(1))$, and*

$$(3.9) \quad \frac{\hat{V}_2}{\sqrt{\hat{V}_1}}(\hat{\beta} - \beta_0) \xrightarrow{d} N(0, 1) \quad \text{as } n \rightarrow \infty.$$

3.3. Weak IV bias. As mentioned in Section 1.3, many existing statistical methods for summary-data MR ignore the measurement error in $\hat{\gamma}_j$. We briefly describe the amount of bias this may incur for the inverse variance weighted (IVW) estimator [13]. The IVW estimator is equivalent to the maximum likelihood estimator (3.2) assuming $\sigma_{\hat{X}_j}^2 = 0$, which has an explicit expression and can be approximated by

$$(3.10) \quad \begin{aligned} \hat{\beta}_{\text{IVW}} &= \frac{\sum_{j=1}^p \hat{\Gamma}_j \hat{\gamma}_j}{\sum_{j=1}^p \hat{\gamma}_j^2} \approx \frac{\mathbb{E}[\sum_{j=1}^p \hat{\Gamma}_j \hat{\gamma}_j]}{\mathbb{E}[\sum_{j=1}^p \hat{\gamma}_j^2]} \\ &= \frac{\beta \|\mathbf{y}\|^2}{\|\mathbf{y}\|^2 + \sum_{j=1}^p \sigma_{\hat{X}_j}^2} \approx \frac{\beta}{1 + (1/\kappa)}. \end{aligned}$$

Thus the amount of bias for the IVW estimator crucially depends on the average IV strength κ . In comparison, our consistency result (Theorem 3.1) only requires $\kappa \gg 1/\sqrt{p}$.

3.4. Practical issues. Next, we discuss several practical implications of the theoretical results above.

3.4.1. Influence of a single IV. Under the assumptions in Theorem 3.2, (3.6) and (3.5) lead to the following asymptotically linear form of $\hat{\beta}$:

$$\hat{\beta} = \frac{1 + o_p(1)}{V_2} \sum_{j=1}^p \frac{(\hat{\Gamma}_j - \beta_0 \hat{\gamma}_j)(\hat{\Gamma}_j \sigma_{\hat{X}_j}^2 \beta_0 + \hat{\gamma}_j \sigma_{\hat{Y}_j}^2)}{(\sigma_{\hat{X}_j}^2 \beta_0^2 + \sigma_{\hat{Y}_j}^2)^2}.$$

The above equation characterizes the influence of a single IV on the estimator $\hat{\beta}$ [24]. Intuitively, the IV Z_j has large influence if it is strong or it has large residual $\hat{\Gamma}_j - \beta_0 \hat{\gamma}_j$. Alternatively, we can measure the influence of a single IV by computing the leave-one-out estimator $\hat{\beta}_{-j}$ that maximizes the profile likelihood with all the SNPs except Z_j . In practice, it is desirable to limit the influence of each SNP to make the estimator robust against idiosyncratic pleiotropy (Model 3). This problem will be considered in Section 5.

3.4.2. *Selecting IVs.* The formulas (3.7) and (3.8) suggest that using extremely weak instruments may deteriorate the efficiency. Consider the following example in which we have a new instrument Z_{p+1} that is independent of X , so $\gamma_{p+1} = 0$. When adding Z_{p+1} to the analysis, V_1 increases but V_2 remains the same, thus the variance of $\hat{\beta}$ becomes larger. Generally, this suggests that we should screen out extremely weak IVs to improve efficiency. To avoid selection bias, we recommend to use two independent GWAS datasets in practice, one to screen out weak IVs and perform LD clumping and one to estimate the SNP-exposure effects γ_j unbiasedly.

3.4.3. *Residual quantile-quantile plot.* One way to check the modeling assumptions in Assumption 1 and Model 1 is the residual Quantile–Quantile (Q–Q) plot, which plots the quantiles of standardized residuals

$$\hat{t}_j = \frac{\hat{\Gamma}_j - \hat{\beta}\hat{\gamma}_j}{\sqrt{\hat{\beta}^2\sigma_{X_j}^2 + \sigma_{Y_j}^2}}$$

against the quantiles of the standard normal distribution. This is reasonable because when $\hat{\beta} = \beta_0, \hat{t}_j \sim N(0, 1)$ under Assumption 1 and Model 1. The Q–Q plot is helpful at identifying IVs that do not satisfy the linear relation $\Gamma_j = \beta_0\gamma_j$, most likely due to genetic pleiotropy.

Besides the residual Q–Q plot, other diagnostic tools can be found in related works. Bowden et al. [10] considered using each SNP’s contribution to the generalized Q statistic to assess whether it is an outlier. Bowden et al. [11] proposed a radial plot $\hat{\beta}_j\sqrt{w}_j$ versus \sqrt{w}_j , where w_j is the “weight” of the j th SNP in (3.3). Since these diagnostic methods are based on the Wald ratio estimates $\hat{\beta}_j$, they can suffer from the weak instrument bias.

3.5. *Example (continued).* We conclude this section by applying the profile likelihood or Profile Score (PS) estimator in the BMI-SBP example in Section 1.2. Here, we used 160 SNPs that have p -values $\leq 10^{-4}$ in the BMI–FEM dataset. The PS point estimate is 0.601 with standard error 0.054.

Figure 3 shows the Q–Q plot and the leave-one-out estimates discussed in Section 3.4. The Q–Q plot clearly indicates the linear model Model 1 is not appropriate to describe the summary data. Although the standardized residuals are roughly normally distributed, their standard deviations are apparently larger than 1. This motivates the random pleiotropy effects assumption in Model 2 which will be considered next.

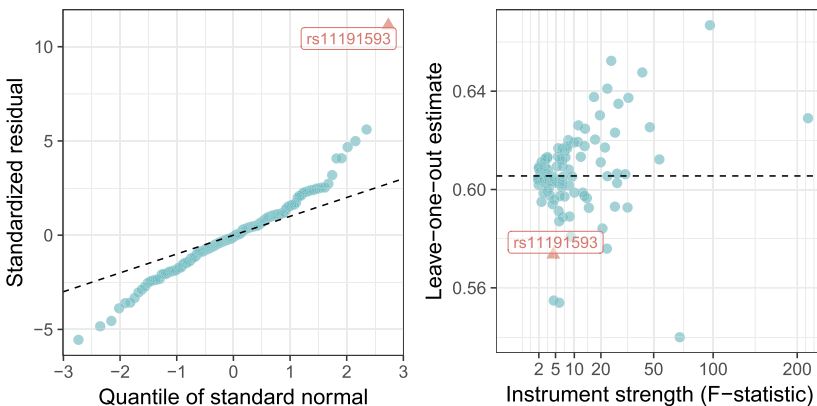


FIG. 3. Diagnostic plots of the Profile Score (PS) estimator. Left panel is a Q–Q plot of the standardized residuals against standard normal. Right panel is the leave-one-out estimates against instrument strength.

4. Systematic pleiotropy: Adjusted profile score.

4.1. *Failure of the profile likelihood.* Next, we consider Model 2, where the deviation from the linear relation $\Gamma_j = \beta_0\gamma_j$ is described by a random effects model $\alpha_j = \Gamma_j - \beta_0\gamma_j \sim N(0, \tau_0^2)$. The normality assumption is motivated by Figure 3 and does not appear to be very consequential in the simulation studies. In this model, the variance of $\hat{\Gamma}$ is essentially inflated by an unknown additive constant τ_0^2 :

$$\hat{\gamma}_j \sim N(\gamma_j, \sigma_{X_j}^2), \quad \hat{\Gamma}_j \sim N(\gamma_j\beta_0, \sigma_{Y_j}^2 + \tau_0^2), \quad j \in [p].$$

Similar to Section 3.1, the profile log-likelihood of (β, τ^2) is given by

$$l(\beta, \tau^2) = -\frac{1}{2} \sum_{j=1}^p \frac{(\hat{\Gamma}_j - \beta\hat{\gamma}_j)^2}{\sigma_{X_j}^2\beta^2 + \sigma_{Y_j}^2 + \tau^2} + \log(\sigma_{Y_j}^2 + \tau^2),$$

and the corresponding profile score equations are

$$\frac{\partial}{\partial \beta} l(\beta, \tau^2) = 0, \quad \frac{\partial}{\partial \tau^2} l(\beta, \tau^2) = 0.$$

It is straightforward to verify that the first estimating equation is unbiased, that is, it has expectation 0 at (β_0, τ_0^2) . However, the other profile score is

$$(4.1) \quad \frac{\partial}{\partial \tau^2} l(\beta, \tau^2) = \frac{1}{2} \sum_{j=1}^p \frac{(\hat{\Gamma}_j - \beta\hat{\gamma}_j)^2}{(\sigma_{X_j}^2\beta^2 + \sigma_{Y_j}^2 + \tau^2)^2} - \frac{1}{\sigma_{Y_j}^2 + \tau^2}.$$

It is easy to see that its expectation is not equal to 0 at the true value $(\beta, \tau^2) = (\beta_0, \tau_0^2)$. This means the profile score is biased in Model 2, thus the corresponding maximum likelihood estimator is not statistically consistent.

4.2. *Adjusted profile score.* The failure of maximizing the profile likelihood should not be surprising, because it is well known that the maximum likelihood estimator can be biased when there are many nuisance parameters [39]. There are many proposals to modify the profile likelihood; see, for example, Barndorff-Nielsen [4], Cox and Reid [17]. Here, we take the approach of McCullagh and Tibshirani [37] that directly modifies the profile score so it has mean 0 at the true value. The *Adjusted Profile Score (APS)* is given by $\psi(\beta, \tau^2) = (\psi_1(\beta, \tau^2), \psi_2(\beta, \tau^2))$, where

$$(4.2) \quad \begin{aligned} \psi_1(\beta, \tau^2) &= -\frac{\partial}{\partial \beta} l(\beta, \tau^2) \\ &= \sum_{j=1}^p \frac{(\hat{\Gamma}_j - \beta\hat{\gamma}_j)(\hat{\Gamma}_j\sigma_{X_j}^2\beta + \hat{\gamma}_j(\sigma_{Y_j}^2 + \tau^2))}{(\sigma_{X_j}^2\beta^2 + \sigma_{Y_j}^2 + \tau^2)^2}, \end{aligned}$$

$$(4.3) \quad \psi_2(\beta, \tau^2) = \sum_{j=1}^p \sigma_{X_j}^2 \frac{(\hat{\Gamma}_j - \beta\hat{\gamma}_j)^2 - (\sigma_{X_j}^2\beta^2 + \sigma_{Y_j}^2 + \tau^2)}{(\sigma_{X_j}^2\beta^2 + \sigma_{Y_j}^2 + \tau^2)^2}.$$

Compared to (4.1), we replaced $(\sigma_{Y_j}^2 + \tau^2)^{-1}$ by $(\sigma_{X_j}^2\beta^2 + \sigma_{Y_j}^2 + \tau^2)^{-1}$, so each summand in (4.3) has mean 0 at (β_0, τ_0^2) . We also weighted the IVs by $\sigma_{X_j}^2$ in (4.3), which is useful in the proof of statistical consistency.

Notice that both the denominators and numerators in ψ_1 and ψ_2 are polynomials of β and τ^2 . However, the denominators are of higher degrees. This implies that the APS estimating equations always have diverging solutions: $\psi(\beta, \tau^2) \rightarrow \mathbf{0}$ if $\beta \rightarrow \pm\infty$ or $\tau^2 \rightarrow \infty$. We define the APS estimator $(\hat{\beta}, \hat{\tau}^2)$ to be the nontrivial finite solution to $\psi(\beta, \tau^2) = \mathbf{0}$ if it exists.

4.3. *Consistency and asymptotic normality.* Because of the diverging solutions of the APS equations, we need to impose some compactness constraints on the parameter space to study the asymptotic property of $(\hat{\beta}, \hat{\tau}^2)$:

ASSUMPTION 4. $(\beta_0, p\tau_0^2)$ is in the interior of a bounded set $\mathcal{B} \subset \mathbb{R} \times \mathbb{R}^+$.

The overdispersion parameter τ_0^2 is scaled up in Assumption 4 by p . This is motivated by the linear structural model (2.3), where $\sum_{j=1}^2 \tau_0^2 \text{Var}(Z_j) = \Theta(p\tau_0^2)$ is the variance of Y explained by the direct effects of \mathbf{Z} . Thus it is reasonable to treat $p\tau_0^2$ as a constant.

We also assume, in addition to Assumption 2, that the variance of X explained by the IVs is nondiminishing.

ASSUMPTION 5. $\|\boldsymbol{\gamma}\|_2 = \Theta(1)$.

THEOREM 4.1. *In Model 2 and suppose Assumptions 1 and 3 to 5 hold, $p \rightarrow \infty$ and $p/n^2 \rightarrow 0$. Then with probability going to 1 there exists a solution of the APS equation such that $(\hat{\beta}, p\hat{\tau}^2)$ is in \mathcal{B} . Furthermore, all solutions in \mathcal{B} are statistically consistent, that is, $\hat{\beta} \xrightarrow{p} \beta_0$ and $p\hat{\tau}^2 - p\tau_0^2 \xrightarrow{p} 0$.*

Next, we consider the asymptotic distribution of the APS estimator.

THEOREM 4.2. *In Model 2 and under the assumptions in Theorem 4.1, if additionally $p = \Theta(n)$ and $\|\boldsymbol{\gamma}\|_3/\|\boldsymbol{\gamma}\|_2 \rightarrow 0$, then*

$$(4.4) \quad (\tilde{\mathbf{V}}_2^{-1} \tilde{\mathbf{V}}_1 \tilde{\mathbf{V}}_2^{-T})^{1/2} \begin{pmatrix} \hat{\beta} - \beta_0 \\ \hat{\tau}^2 - \tau_0^2 \end{pmatrix} \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{I}_2),$$

where

$$\tilde{\mathbf{V}}_1 = \sum_{j=1}^p \frac{1}{(\sigma_{X_j}^2 \beta_0^2 + \sigma_{Y_j}^2 + \tau_0^2)^2} \begin{pmatrix} (\gamma_j^2 + \sigma_{X_j}^2)(\sigma_{Y_j}^2 + \tau_0^2) + \Gamma_j^2 \sigma_{X_j}^2 & 0 \\ 0 & 2(\sigma_{X_j}^2)^2 \end{pmatrix},$$

$$\tilde{\mathbf{V}}_2 = \sum_{j=1}^p \frac{1}{(\sigma_{X_j}^2 \beta_0^2 + \sigma_{Y_j}^2 + \tau_0^2)^2} \begin{pmatrix} \gamma_j^2(\sigma_{Y_j}^2 + \tau_0^2) + \Gamma_j^2 \sigma_{X_j}^2 & \sigma_{X_j}^2 \beta_0 \\ 0 & \sigma_{X_j}^2 \end{pmatrix}.$$

Similar to Theorem 3.3, the information matrices $\tilde{\mathbf{V}}_1$ and $\tilde{\mathbf{V}}_2$ can be estimated by substituting γ_j^2 by $\hat{\gamma}_j^2 - \sigma_{X_j}^2$ and Γ_j^2 by $\hat{\Gamma}_j^2 - \sigma_{Y_j}^2 - \hat{\tau}^2$. We omit the details for brevity.

4.4. *Example (continued).* We apply the APS estimator to the BMI-SBP example. Using the same 160 SNPs in Section 3.5, the APS point estimate is $\hat{\beta} = 0.301$ (standard error 0.158) and $\hat{\tau}^2 = 9.2 \times 10^{-4}$ (standard error 1.7×10^{-4}). Notice that the APS point estimate of β is much smaller than the PS point estimate. One possible explanation of this phenomenon is that the PS estimator tends to use a larger β to compensate for the overdispersion in Model 2 (the variance of $\hat{\Gamma}_j - \beta \hat{\gamma}_j$ is $\beta^2 \sigma_{X_j}^2 + \sigma_{Y_j}^2$ in Model 1 and $\beta^2 \sigma_{X_j}^2 + \sigma_{Y_j}^2 + \tau_0^2$ in Model 2).

Figure 4 shows the diagnostic plots of the APS estimator. Compared to the PS estimator in Section 3.5, the overdispersion issue is much more benign. However, there is an outlier which corresponds to the SNP `rs11191593`. It heavily biases the APS estimate too: when excluding this SNP, the APS point estimate changes from 0.301 to almost 0.4 in the right panel of Figure 4. The outlier might also inflate $\hat{\tau}^2$ so the Q-Q plot looks a little underdispersed. These observations motivate the consideration of a robust modification of the APS in the next section.

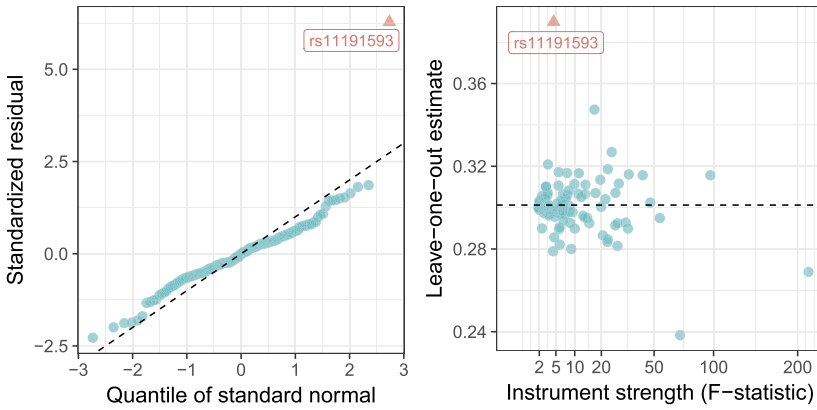


FIG. 4. Diagnostic plots of the Adjusted Profile Score (APS) estimator. Left panel is a Q-Q plot of the standardized residuals against standard normal. Right panel is the leave-one-out estimates against instrument strength.

5. Idiosyncratic pleiotropy: Robustness to outliers. Next, we consider Model 3 with idiosyncratic pleiotropy. As mentioned in Section 3.4.1, a single IV can have unbounded influence on the PS (and APS) estimators. When the IV Z_j has other strong causal pathways, its pleiotropy parameter α_j can be much larger than what is predicted by the random effects model $\alpha_j \sim N(0, \tau_0^2)$, leading to a biased estimate of the causal effect as illustrated in Section 4.4. In this section, we propose a general method to robustify the APS to limit the influence of outliers such as SNP rs11191593 in the example.

5.1. *Robustify the adjusted profile score.* Our approach is an application of the robust regression techniques pioneered by Huber [30]. As mentioned in Section 3.1, the profile likelihood (3.2) can be viewed as a linear regression of $\hat{\Gamma}_j$ on $\hat{\gamma}_j$ using the l_2 -loss. To limit the influence of a single IV, we consider changing the l_2 -loss to a robust loss function. Two celebrated examples are the Huber loss

$$\rho_{\text{huber}}(r; k) = \begin{cases} r^2/2, & \text{if } |r| \leq k, \\ k(|r| - k/2), & \text{otherwise,} \end{cases}$$

and Tukey’s biweight loss

$$\rho_{\text{tukey}}(r; k) = \begin{cases} 1 - (1 - (r/k)^2)^3, & \text{if } |r| \leq k, \\ 1, & \text{otherwise.} \end{cases}$$

This heuristic motivates the following modification of the profile log-likelihood when $\tau_0^2 = 0$:

$$(5.1) \quad l_\rho(\beta) := - \sum_{j=1}^p \rho \left(\frac{\hat{\Gamma}_j - \beta \hat{\gamma}_j}{\sqrt{\sigma_{X_j}^2 \beta^2 + \sigma_{Y_j}^2}} \right).$$

It is easy to see that $l_\rho(\beta)$ reduces to the regular profile log-likelihood (3.2) if $\rho(r) = r^2/2$.

When $\tau_0^2 > 0$, we cannot directly use the profile score $(\partial/\partial\tau^2)l(\beta, \tau^2)$ as discussed in Section 4.1. This issue can be resolved using the APS approach in Section 4.2 by using ψ_2 in (4.3). However, a single IV can still have unbounded influence in ψ_2 . We must further robustify ψ_2 , which is analogous to estimating a scale parameter robustly.

Next, we briefly review the robust M-estimation of scale parameter. Consider repeated measurements of a scale family with density $f_0(r/\sigma)/\sigma$. Then a general way of robust estimation of σ is to solve the following estimating equation [36], Section 2.5:

$$\hat{\mathbb{E}}[(R/\sigma) \cdot \rho'(R/\sigma)] = \delta,$$

where $\hat{\mathbb{E}}$ stands for the empirical average and $\delta = \mathbb{E}[R \cdot \rho'(R)]$ for $R \sim f_0$.

Based on the above discussion, we propose the following *Robust Adjusted Profile Score (RAPS)* estimator of β . Denote

$$t_j(\beta, \tau^2) = \frac{\hat{\Gamma}_j - \beta \hat{\gamma}_j}{\sqrt{\sigma_{X_j}^2 \beta^2 + \sigma_{Y_j}^2 + \tau^2}}.$$

Then the RAPS $\psi^{(\rho)} = (\psi_1^{(\rho)}, \psi_2^{(\rho)})$ is given by

$$(5.2) \quad \psi_1^{(\rho)}(\beta, \tau^2) = \sum_{j=1}^p \rho'(t_j(\beta, \tau^2)) u_j(\beta, \tau^2),$$

$$(5.3) \quad \psi_2^{(\rho)}(\beta, \tau^2) = \sum_{j=1}^p \sigma_{X_j}^2 \frac{t_j(\beta, \tau^2) \cdot \rho'(t_j(\beta, \tau^2)) - \delta}{\sigma_{X_j}^2 \beta^2 + \sigma_{Y_j}^2 + \tau^2},$$

where $\rho'(\cdot)$ is the derivative of $\rho(\cdot)$, $u_j(\beta, \tau^2) = -(\partial/\partial\beta)t_j(\beta, \tau^2)$ and $\delta = \mathbb{E}[R \cdot \rho'(R)]$ for $R \sim N(0, 1)$. Notice that $\psi^{(\rho)}$ reduces to the nonrobust APS ψ in (4.2) and (4.3) when $\rho(r) = r^2/2$ is the squared error loss. Finally, the RAPS estimator $(\hat{\beta}, \hat{\tau}^2)$ is given by the nontrivial finite solution of $\psi^{(\rho)}(\beta, \tau^2) = \mathbf{0}$.

5.2. *Asymptotics.* Because the RAPS estimator is the solution of a system of nonlinear equations, its asymptotic behavior is very difficult to analyze. For instance, it is difficult to establish statistical consistency because there could be multiple roots for the RAPS equations in the population level. Thus β might not be globally identified. We can, nevertheless, verify the local identifiability [45].

THEOREM 5.1 (Local identification of RAPS). *In Model 2, $\mathbb{E}[\psi^{(\rho)}(\beta_0, \tau_0^2)] = \mathbf{0}$ and $\mathbb{E}[\nabla\psi^{(\rho)}]$ has full rank.*

In practice, we find that the RAPS estimating equation usually only has one finite solution. To study the asymptotic normality of the RAPS estimator, we will assume $(\hat{\beta}, p\hat{\tau}^2)$ is consistent under Model 2. We further impose the following smoothness condition on the robust loss function ρ .

ASSUMPTION 6. The first three derivatives of $\rho(\cdot)$ exist and are bounded.

THEOREM 5.2. *In Model 2 and under the assumptions in Theorem 4.2, if additionally we assume*

1. *the RAPS estimator is consistent: $\hat{\beta} - \beta_0 \xrightarrow{p} 0, p(\hat{\tau}^2 - \tau_0^2) \xrightarrow{p} 0,$*
2. *Assumption 6 holds, and*
3. $\|\gamma\|_3^3 / \|\gamma\|_2^3 = O(p^{-1/2}),$

then

$$(5.4) \quad ((\tilde{V}_2^{(\rho)})^{-1} \tilde{V}_1^{(\rho)} (\tilde{V}_2^{(\rho)})^{-T})^{1/2} \begin{pmatrix} \hat{\beta} - \beta_0 \\ \hat{\tau}^2 - \tau_0^2 \end{pmatrix} \xrightarrow{d} N(\mathbf{0}, I_2),$$

where

$$\tilde{V}_1^{(\rho)} = \begin{pmatrix} c_1(\tilde{V}_1)_{11} & 0 \\ 0 & c_2(\tilde{V}_1)_{22} \end{pmatrix},$$

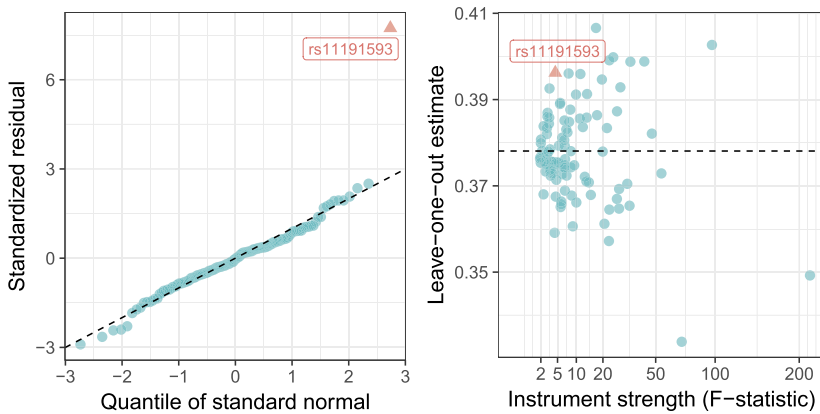
$$\tilde{V}_2^{(\rho)} = \begin{pmatrix} \delta(\tilde{V}_2)_{11} & \delta(\tilde{V}_2)_{12} \\ 0 & [(\delta + c_3)/2](\tilde{V}_2)_{22} \end{pmatrix},$$

and the constants are: for $R \sim N(0, 1)$, $c_1 = \mathbb{E}[\rho'(R)^2]$, $c_2 = \text{Var}(R\rho'(R))/2$, $c_3 = \mathbb{E}[R^2\rho''(R)]$.

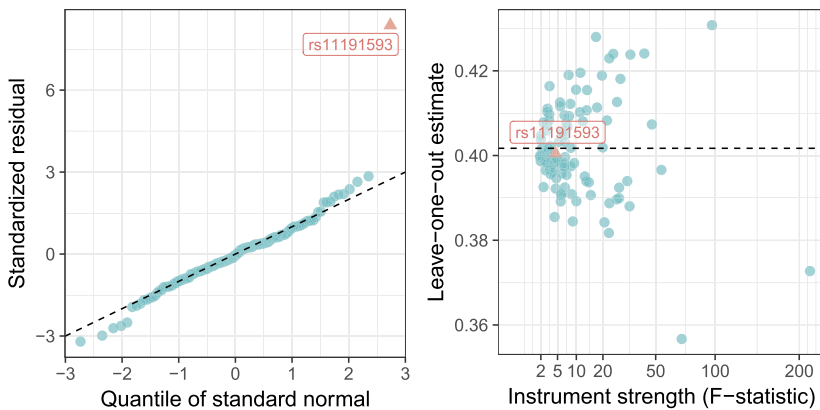
It is easy to verify that when $\rho(r) = r^2/2$, $\delta = c_1 = c_2 = c_3 = 1$, so $\tilde{V}_1^{(\rho)}$ and $\tilde{V}_2^{(\rho)}$ reduce to \tilde{V}_1 and \tilde{V}_2 . In other words, the asymptotic variance formula in Theorem 5.2 is consistent with the one in Theorem 4.2. However, additional technical assumptions are needed in Theorem 5.2 to bound the higher-order terms in the Taylor expansion.

5.3. *Example (continued)*. As before, we illustrate the RAPS estimator using the BMI-SBP example. Using the Huber loss with $k = 1.345$ (corresponding to 95% asymptotic efficiency in the simple location problem), the point estimate is $\hat{\beta} = 0.378$ (standard error 0.121), $\hat{\tau}^2 = 4.7 \times 10^{-4}$ (standard error 1.0×10^{-4}). Using the Tukey loss with $k = 4.685$ (also corresponding to 95% asymptotic efficiency in the simple location problem), the point estimate is $\hat{\beta} = 0.402$ (standard error 0.106), $\hat{\tau}^2 = 3.4 \times 10^{-4}$ (standard error 7.8×10^{-5}).

Figure 5 shows the diagnostic plots of the two RAPS estimators. Compared to Figure 4, the robust loss functions limit the influence of the outlier (SNP rs11191593), and the resulting $\hat{\beta}$ becomes larger. In Figure 5b, the outlier’s influence is essentially zero because the Tukey loss function is redescending. This shows the robustness of our RAPS estimator to the idiosyncratic pleiotropy.



(a) RAPS using the Huber loss.



(b) RAPS using the Tukey loss.

FIG. 5. Diagnostic plots of the Robust Adjusted Profile Score (RAPS) estimator: Left panels are $Q-Q$ plots of the standardized residuals against standard normal. Right panels are the leave-one-out estimates against instrument strength.

6. Simulation. Throughout the paper, all of our theoretical results are asymptotic. We usually require both the sample size n and the number of IVs p to go to infinity (except for Theorem 3.2 where finite p is allowed). We now assess if the asymptotic approximations are reasonably accurate in practical situations, where p may range from tens to hundreds.

6.1. *Simulating summary data directly from Assumption 1.* To this end, we first created simulated summary-data MR datasets that mimic the BMI-SBP example in Section 1.2. In particular, we considered two scenarios: $p = 25$, which corresponds to using the selection threshold 5×10^{-8} as described in Section 1.2, and $p = 160$, which corresponds to using the threshold 1×10^{-4} as in Sections 3.5, 4.4 and 5.3. The model parameters are chosen as follows: the variances of the measurement error, $\{\sigma_{X_j}^2, \sigma_{Y_j}^2\}_{j \in [p]}$, are the same as those in the BMI-SBP dataset. The true marginal IV-exposure effects, $\{\gamma_j\}_{j \in [p]}$, are chosen to be the observed effects in the BMI-SBP dataset, and $\hat{\gamma}_j$ is generated according to Assumption 1 by $\hat{\gamma}_j \stackrel{ind.}{\sim} N(\gamma_j, \sigma_{X_j}^2)$. The true marginal IV-outcome effects, $\{\Gamma_j\}_{j \in [p]}$, are generated in six different ways with $\beta_0 = 0.4$:

1. $\Gamma_j = \gamma_j \beta_0$;
2. $\Gamma_j = \gamma_j \beta_0 + \alpha_j, \alpha_j \stackrel{i.i.d.}{\sim} N(0, \tau_0^2)$, where $\tau_0 = 2 \cdot (1/p) \sum_{j=1}^p \sigma_{Y_j}$;
3. Γ_j is generated according to setup 2 above, except that α_1 has mean $5 \cdot \tau_0$ (the IVs are sorted so that the first IV has the largest $|\gamma_j|/\sigma_{X_j}$).
4. $\Gamma_j = \gamma_j \beta_0 + \alpha_j, \alpha_j \stackrel{i.i.d.}{\sim} \tau_0 \cdot \text{Lap}(1)$, where Lap(1) is the Laplace (double exponential) distribution with rate 1.
5. $\Gamma_j = \gamma_j \beta_0 + \alpha_j, \alpha_j = |\gamma_j|/(p^{-1} \sum_{j=1}^p |\gamma_j|) \cdot N(0, \tau_0^2)$.
6. Γ_j is generated according to setup 2 above, except that for 10% randomly selected IVs, their direct effects α_j have mean $5 \cdot \tau_0$.

The first three setups correspond to Models 1 to 3, respectively, and the last three setups violate our modeling assumptions and are used to assess the robustness of the procedures.

Finally, $\hat{\Gamma}_j$ is generated according to Assumption 1 by $\hat{\Gamma}_j \stackrel{ind.}{\sim} N(\Gamma_j, \sigma_{Y_j}^2)$.

We applied six methods to the simulated data (10,000 replications in each setting). The first three are existing methods to benchmark our performance: the inverse variance weighting (IVW) estimator [13], MR-Egger regression [7] and the weighted median estimator [8]. The next three methods are proposed in this paper: the profile score (PS) estimator in Section 3, the adjusted profile score (APS) estimator in Section 4, and the robust adjusted profile score (RAPS) estimator in Section 5 with Tukey’s loss function ($k = 4.685$).

The simulation results are reported in Table 1 for $p = 25$ and Table 2 for $p = 160$. Here is a summary of the results:

1. In setup 1, the PS estimator has the smallest root-mean square error (RMSE) and the shortest confidence interval (CI) with the desired coverage rate. The IVW estimator performs very well when $p = 25$ but has considerable bias and less than nominal coverage when $p = 160$. The APS and RAPS estimators have slightly longer CI than PS. The MR-Egger and weighted median estimators are less accurate than the other methods.
2. In setup 2, the PS estimator, as well as the weighted median, have substantial bias and perform poorly. The APS estimator is overall the best with very small bias and desired coverage, followed very closely by RAPS. The IVW and MR-Egger estimators also perform quite well, though their relative biases are more than 10% when $p = 160$.
3. In setup 3, all estimators besides RAPS have very large bias and poor CI coverage. The RMSE of the RAPS estimator is slightly larger than the RMSE in Model 2, and the coverage of RAPS is slightly below the nominal rate.

TABLE 1

Simulation results for $p = 25$. The summary statistics reported are: bias divided by β_0 , root-median-square error (RMSE) divided by β_0 , length of the confidence interval (CI) divided by β_0 and the coverage rate of the CI (nominal rate is 95%), all in %

Setup	Method	Bias %	RMSE %	CI Len. %	Cover. %
1	IVW	-2.9	12.7	73.8	95.4
	Egger	-7.4	24.4	142.3	95.3
	W. Median	-5.2	17.0	105.5	96.5
	PS	-0.1	12.7	74.9	95.1
	APS	-0.4	12.7	76.8	96.0
	RAPS	-0.4	13.0	79.0	96.1
2	IVW	-3.0	29.3	167.9	93.3
	Egger	-8.2	59.7	319.2	92.1
	W. Median	-12.8	39.9	121.4	70.6
	PS	14.7	36.1	71.4	49.2
	APS	-0.2	28.8	165.4	93.4
	RAPS	-0.1	30.1	170.2	93.1
3	IVW	-115.5	115.2	225.6	48.1
	Egger	-264.2	262.8	409.1	25.5
	W. Median	-80.7	79.5	151.4	47.3
	PS	-122.3	121.3	66.1	6.9
	APS	-86.2	85.6	207.0	65.0
	RAPS	-11.6	40.6	168.7	84.3
4	IVW	-5.1	25.1	159.5	96.0
	Egger	-54.5	58.8	300.9	90.0
	W. Median	-22.5	26.0	113.2	83.8
	PS	13.4	31.2	71.7	55.9
	APS	4.0	25.6	158.4	96.1
	RAPS	2.6	20.3	117.5	93.3
5	IVW	-2.4	48.2	169.7	76.3
	Egger	-8.2	98.0	321.0	72.9
	W. Median	-24.4	60.4	136.7	56.0
	PS	15.8	57.2	71.6	33.0
	APS	0.9	46.8	183.0	81.1
	RAPS	1.5	44.9	169.0	78.3
6	IVW	-8.1	64.2	382.8	94.8
	Egger	-102.2	134.8	723.7	90.7
	W. Median	-30.8	50.3	130.6	63.1
	PS	200.2	309.6	82.1	4.1
	APS	13.7	62.1	327.1	92.8
	RAPS	12.3	50.3	298.2	85.4

4. In setup 4, the direct effects α_j are distributed as Laplace instead of normal. The RAPS estimator has the smallest bias and RMSE, though the coverage is slightly below the nominal level.

5. In setup 5, the variance of α_j is proportional to $|\gamma_j|$. In this case, APS and RAPS are approximately unbiased but the coverage is significantly lower than 95%.

6. In setup 6, 10% of the IVs have very large but roughly balanced pleiotropy effects α_j . All estimators are biased in this case. The RAPS estimator has the smallest RMSE but the CI coverage is slightly below 95%. The IVW and APS estimators have slightly larger RMSE and the CI has the desired coverage rate.

TABLE 2

Simulation results for $p = 160$. The summary statistics reported are: bias divided by β_0 , root-median-square error (RMSE) divided by β_0 , length of the confidence interval (CI) divided by β_0 and the coverage rate of the CI (nominal rate is 95%), all in %

Setup	Method	Bias %	RMSE %	CI Len. %	Cover. %
1	IVW	-11.1	12.2	51.0	87.0
	Egger	-10.1	15.2	79.9	92.6
	W. Median	-12.6	15.6	84.3	93.9
	PS	0.1	9.6	57.0	95.2
	APS	-0.4	9.5	58.3	95.8
	RAPS	-0.5	9.8	59.9	95.8
2	IVW	-11.6	23.2	122.5	92.6
	Egger	-10.8	34.9	191.5	93.6
	W. Median	-25.7	34.3	105.5	68.9
	PS	119.2	119.8	51.0	6.2
	APS	-0.4	23.0	134.8	95.1
	RAPS	-0.4	23.8	138.7	95.1
3	IVW	-70.1	69.9	131.3	44.7
	Egger	-125.5	125.6	203.8	32.3
	W. Median	-65.0	65.0	111.5	41.5
	PS	4.1	77.9	44.6	15.5
	APS	-47.9	48.3	139.3	73.2
	RAPS	-3.9	27.4	137.9	90.6
4	IVW	-11.9	20.5	121.5	94.7
	Egger	-13.6	31.5	189.5	94.7
	W. Median	-24.1	24.8	93.9	80.2
	PS	134.7	114.3	51.4	7.1
	APS	4.8	20.8	133.6	96.5
	RAPS	4.3	16.1	91.3	93.6
5	IVW	-11.0	53.9	139.7	62.2
	Egger	-9.8	92.5	217.7	56.9
	W. Median	-56.0	63.7	125.2	49.3
	PS	-819.8	244.0	57.8	4.7
	APS	-0.3	55.3	170.7	71.6
	RAPS	1.5	48.6	120.4	59.8
6	IVW	-12.7	47.2	278.8	95.0
	Egger	-16.4	74.2	435.3	94.9
	W. Median	-34.9	43.6	115.2	63.1
	PS	>999.9	>999.9	>999.9	12.8
	APS	13.6	50.2	291.2	95.2
	RAPS	10.8	42.7	258.4	91.2

Finally, we briefly remark on the bias of IVW and other existing estimators. In Section 3.3, we have derived that the IVW estimator is biased toward 0 and the relative bias is approximately $1/\kappa$. The average instrument strength κ in the two settings are $\kappa = 33.1$ ($p = 25$) and $\kappa = 9.1$ ($p = 160$). The simulation results for setup 1 in Tables 1 and 2 almost exactly match the prediction from our approximation formula (3.10).

Overall, the RAPS estimator is the clear winner in this simulation: when there is no idiosyncratic outlier (setups 1 and 2), it behaves almost as well as the best performer; when there is an idiosyncratic outlier (setup 3), it still has very small bias and close-to-nominal coverage; when our modeling assumptions are not satisfied (setups 4, 5, 6), it still has the smallest bias and RMSE, though the CI may fail to cover β_0 at the nominal rate.

6.2. *Simulating from real genotypes.* As pointed out by an anonymous reviewer, the marginal GWAS coefficients might not perfectly follow the distributional assumptions in Assumption 1. In fact, in Section 2.2 we already showed that even in linear structural models the marginal coefficients have small but nonzero covariances. As a proof of concept, we perform another simulation study using real genotypes from the 1000 Genomes Project [1].

In total, the 1000 Genomes Project phase 1 dataset contains the genotypes of 1092 individuals. We simulated the exposure X and outcome Y according to the linear structural equation model (2.3) using the entire 10th chromosome as Z (containing 1,882,663 genetic variants). 100 random entries of $\boldsymbol{\gamma}$ are set to be nonzero and follow the Laplace distribution with rate 1. The unmeasured confounder U is simulated from the standard normal distribution and the parameters were set to $\eta_X = 3$, $\eta_Y = 5$. The noise variables were simulated from $E_X \sim N(0, 3^2)$ and $E_Y \sim N(0, 5^2)$. The direct effects $\boldsymbol{\alpha}$ had p_α random nonzero entries that were simulated from the Laplace distribution with rate r_α . In total, we considered five settings:

1. $\beta = 0$, $p_\alpha = 0$;
2. $\beta = 0$, $p_\alpha = 200$, $r_\alpha = 0.5$;
3. $\beta = 1$, $p_\alpha = 0$;
4. $\beta = 1$, $p_\alpha = 200$, $r_\alpha = 0.5$;
5. $\beta = 1$, $p_\alpha = 200$, $r_\alpha = 1.5$.

In this dataset, 368,977 variants have minor allele frequency greater 5% and are considered as potential instrumental variables. We used 292, 400 and 400 individuals (random partition) as the selection, exposure and outcome data and obtained GWAS summary data by running marginal linear regressions. We simulated Y using one of the five settings described above. After LD clumping (p -value $\leq 5 \times 10^{-3}$), 121 independent variants were selected as IVs, and we applied existing and our methods to these 121 SNPs. To provide a more comprehensive comparison, we also applied two classical IV estimator, two-stage least squares (2SLS) and limited information maximum likelihood (LIML), to the outcome sample of 400 individuals. For the LIML estimator, we computed the standard error using the “many weak IV asymptotics” [25]. Note that 2SLS and LIML cannot be computed using just the GWAS summary data and they assume all the IVs are valid.

We used 5000 replications to obtain the same performance metrics in Section 6.1, which are reported in Table 3. Overall, our estimators (in particular, APS and RAPS) are unbiased and maintain the nominal CI coverage rate in all 5 settings. The three existing estimators—IVW, MR-Egger and weighted median—are heavily biased toward 0 when $\beta \neq 0$. Also, notice that their RMSE and CI length are (abnormally) smaller than the RMSE and CI length of the “oracle” LIML estimator that uses individual genotypes. The 2SLS estimator is also heavily biased by weak instruments.

Although the simulation results in Table 3 are encouraging, we want to point out that the sample size and simulation parameters we used might be quite different from actual MR studies. The pleiotropy models (parametrized by p_α and r_α) being tested here are also limited. Nonetheless, this simulation shows that using the statistical framework developed in this paper, it is possible to obtain summary-data MR estimators that perform almost as well as the “oracle” LIML estimator that uses individual data.

7. Comparison in real data examples.

7.1. *In the BMI-SBP example.* Table 4 briefly summarize the results using different estimators in this and previous papers for the BMI-SBP example introduced in Section 1.2. Since the ground truth is unknown, we do not know which estimate is closer to the truth. Nevertheless, we can still make three remarks. First, the point estimates varied considerably between

TABLE 3

Results for the numerical simulation using real genotypes. The performance metrics reported are: bias (median $\hat{\beta}$ minus β), root-median-square error (RMSE), median length of the confidence interval (CI) and the coverage rate of the CI (nominal rate is 95%)

Setup	Method	Bias	RMSE	CI Len.	Coverage %
1	IVW	0.00	0.08	0.42	93.1
	Egger	0.00	0.11	0.62	95.1
	W. Median	0.00	0.12	0.74	96.8
	PS	0.01	0.26	1.42	92.9
	APS	0.01	0.23	1.61	98.9
	RAPS	0.00	0.23	1.76	98.2
	2SLS	-0.46	0.46	0.41	0.9
	LIML	0.00	0.26	1.40	94.5
2	IVW	-0.02	0.08	0.45	94.0
	Egger	-0.02	0.11	0.65	95.4
	W. Median	-0.04	0.12	0.78	97.2
	PS	-0.06	0.29	1.42	89.2
	APS	-0.05	0.25	1.67	98.6
	RAPS	-0.05	0.25	1.82	97.5
	2SLS	-0.47	0.47	0.43	1.1
	LIML	0.02	0.28	1.56	95.4
3	IVW	-0.63	0.63	0.43	0.1
	Egger	-0.45	0.45	0.61	21.1
	W. Median	-0.64	0.64	0.76	8.7
	PS	0.08	0.22	1.35	96.9
	APS	0.02	0.22	1.78	97.6
	RAPS	0.01	0.22	1.87	93.1
	2SLS	-0.46	0.46	0.41	1.2
	LIML	-0.01	0.26	1.41	94.8
4	IVW	-0.65	0.65	0.46	0.2
	Egger	-0.47	0.47	0.65	22.4
	W. Median	-0.61	0.61	0.79	13.6
	PS	0.13	0.26	1.39	95.1
	APS	0.01	0.25	1.86	96.6
	RAPS	-0.01	0.24	1.95	92.2
	2SLS	-0.46	0.46	0.43	1.4
	LIML	0.03	0.28	1.57	95.4
5	IVW	-0.68	0.68	0.62	0.9
	Egger	-0.50	0.50	0.90	40.4
	W. Median	-0.44	0.44	0.97	57.0
	PS	0.41	0.49	1.72	87.3
	APS	0.01	0.37	2.40	96.8
	RAPS	-0.04	0.33	2.48	94.8
	2SLS	-0.47	0.47	0.57	10.4
	LIML	0.23	0.51	2.93	97.9

the methods, so the choice of estimator may make a difference in practice. Second, the PS, IVW and MR-Egger point estimates changed substantially when all 160 SNPs were used instead of just the 25 strongest ones, whereas the RAPS estimators and the weighted median were more stable. Finally, all the standard errors are computed under the modeling assumptions each method makes, thus they may not be directly comparable. For example, in theory the RAPS estimators are less efficient than the APS estimator in Model 2, but their standard

TABLE 4
Comparison of results in the BMI-SBP example

Method	$p = 25$		$p = 160$	
	$\hat{\beta}$	SE	$\hat{\beta}$	SE
PS	0.367	0.075	0.601	0.054
APS	0.364	0.133	0.301	0.158
RAPS (Huber)	0.354	0.131	0.378	0.121
RAPS (Tukey)	0.361	0.133	0.402	0.106
IVW	0.332	0.140	0.514	0.102
MR-Egger	0.647	0.283	0.472	0.176
Weighted median	0.516	0.125	0.514	0.102

errors are indeed smaller in Table 4. This should not be too surprising because Model 2 does not hold in this example as illustrated in Section 4.4.

7.2. An illustration of weak IV bias and selection bias. Finally, we consider another real data validation example, which shall be referred to as the BMI-BMI example. In this example, both the “exposure” and the “outcome” are BMI. Although there is no “causal” effect of BMI on itself, Model 1 for GWAS summary data should technically hold with $\beta_0 = 1$. Therefore, this is a rare scenario where we know the truth in real data. Since there are many SNPs that are only weakly associated with BMI, this example also offers a good opportunity to probe the issue of weak instrument bias and the efficiency gain by including many weak IVs. The downside is that this example does not test the methods’ robustness to pleiotropy because the exposure and outcome are the same trait.

We obtained three GWAS datasets for this example:

BMI-GIANT: full dataset from the GIANT consortium [35] (i.e., combining BMI-FEM and BMI-MAL), used to select SNPs.

BMI-UKBB-1: half of the UKBB data, used as the “exposure.”

BMI-UKBB-2: another half of UKBB data, used as the “outcome.”

We applied in total six methods. Four have been previously developed: besides the three estimators considered in Section 6, we also included the weighted mode estimator of Hartwig, Davey Smith and Bowden [26]. We use the implementation in the `TwoSampleMR` software package [28] for the existing methods. The last two methods were the PS and RAPS estimators developed in this paper (APS performs similar to PS and RAPS and is omitted).

The results are reported in Table 5. Overall, the PS and RAPS estimators provided very accurate estimate of $\beta_0 = 1$. PS has the smallest standard error because there is no pleiotropy at all in this example. When there is pleiotropy (as expected in most real studies), PS can perform poorly as demonstrated in Section 6. All the existing methods are biased especially when there are many weak IVs.

In Table 6, we illustrate the danger of selection bias. In this example, we discard the BMI-GIANT dataset and use BMI-UKBB-1 for both selection and inference (estimating γ_j). The estimators are biased toward 0 in almost all cases, even if we only use the genome-wide significant p -value threshold 10^{-9} or 10^{-8} . This is because the assumption $\hat{\gamma}_j \sim N(\gamma_j, \sigma_{X_j}^2)$ is violated. In fact, due to selection bias, the selected $\hat{\gamma}_j$ are stochastically larger than their mean γ_j (if $\gamma_j > 0$). Compared with other methods, the MR-Egger regression seems to be less affected by the selection bias.

TABLE 5

Results of the BMI-BMI example. The true β_0 should be 1. We considered 8 selection thresholds p_{sel} from 1×10^{-9} to 1×10^{-2} . The mean and median of the F -statistics $\hat{\gamma}_j^2/\sigma_{\hat{X}_j}^2$ are reported. In each setting, we report the point estimate and the standard error of all the methods

p_{sel}	# SNPs	Mean F	IVW	W. Median	W. Mode
1e-9	48	78.6	0.983 (0.026)	0.945 (0.039)	0.941 (0.042)
1e-8	58	69.2	0.983 (0.024)	0.945 (0.039)	0.939 (0.044)
1e-7	84	55.0	0.988 (0.024)	0.945 (0.036)	0.933 (0.041)
1e-6	126	44.1	0.986 (0.022)	0.944 (0.034)	0.931 (0.038)
1e-5	186	34.3	0.986 (0.019)	0.943 (0.033)	0.928 (0.039)
1e-4	287	26.1	0.981 (0.017)	0.941 (0.031)	0.929 (0.035)
1e-3	474	18.8	0.955 (0.015)	0.903 (0.027)	0.917 (0.231)
1e-2	812	12.7	0.928 (0.014)	0.879 (0.023)	0.739 (7.130)
p_{sel}	# SNPs	Median F	Egger	PS	RAPS
1e-9	48	51.8	0.926 (0.055)	0.999 (0.023)	0.998 (0.026)
1e-8	58	42.0	0.928 (0.050)	0.999 (0.023)	0.998 (0.025)
1e-7	84	32.1	0.905 (0.048)	1.012 (0.021)	1.004 (0.025)
1e-6	126	27.4	0.881 (0.043)	1.017 (0.019)	1.009 (0.023)
1e-5	186	21.0	0.874 (0.036)	1.020 (0.018)	1.013 (0.020)
1e-4	287	15.8	0.921 (0.031)	1.023 (0.017)	1.018 (0.018)
1e-3	474	10.8	0.913 (0.027)	1.010 (0.016)	1.006 (0.016)
1e-2	812	5.6	0.909 (0.022)	1.010 (0.015)	1.005 (0.015)

8. Discussion. In this paper, we have proposed a systematic approach for two-sample summary-data Mendelian randomization based on modifying the profile score function. By considering increasingly more complex models, we arrived at the Robust Adjusted Profile

TABLE 6

Illustration of selection bias. The same BMI-UKBB-1 dataset is used for both selecting SNPs and estimating the SNP-exposure effects γ_j . All estimators are biased (true $\beta_0 = 1$) due to not accounting for selection bias

p_{sel}	# SNPs	Mean F	IVW	W. Median	W. Mode
1e-9	110	68.63	0.851 (0.02)	0.83 (0.025)	0.896 (0.046)
1e-8	168	57.00	0.823 (0.017)	0.8 (0.022)	0.885 (0.053)
1e-7	228	50.08	0.799 (0.016)	0.768 (0.019)	0.886 (0.058)
1e-6	305	43.92	0.761 (0.015)	0.736 (0.019)	0.865 (0.079)
1e-5	443	36.98	0.721 (0.013)	0.667 (0.016)	0.824 (0.12)
1e-4	652	30.68	0.678 (0.012)	0.616 (0.015)	0.593 (0.122)
1e-3	929	25.36	0.629 (0.011)	0.57 (0.014)	0.576 (0.096)
1e-2	1289	20.70	0.592 (0.01)	0.528 (0.013)	0.554 (0.093)
p_{sel}	# SNPs	Median F	Egger	PS	RAPS
1e-9	110	49.20	1.071 (0.051)	0.871 (0.015)	0.862 (0.021)
1e-8	168	41.12	1.018 (0.046)	0.848 (0.014)	0.831 (0.018)
1e-7	228	37.12	1.016 (0.043)	0.824 (0.012)	0.803 (0.016)
1e-6	305	33.68	1.006 (0.041)	0.793 (0.011)	0.763 (0.016)
1e-5	443	28.74	0.957 (0.037)	0.762 (0.01)	0.716 (0.015)
1e-4	652	23.23	0.89 (0.033)	0.724 (0.009)	0.66 (0.014)
1e-3	929	19.12	0.823 (0.03)	0.687 (0.008)	0.594 (0.013)
1e-2	1289	15.26	0.749 (0.025)	0.657 (0.008)	0.541 (0.012)

Score (RAPS) estimator which is robust to both systematic and idiosyncratic pleiotropy and performed excellently in all the numerical examples. Thus we recommend to routinely use the RAPS estimator in practice, especially if the exposure and the outcome are both complex traits.

Our theoretical and empirical results advocate for a new design of two-sample MR. Instead of using just a few strong SNPs (those with large $|\hat{\gamma}_j|/\sigma_{X_j}$), we find that adding many (potentially hundreds of) weak SNPs usually substantially decreases the variance of the estimator. This is not feasible with existing methods for MR because they usually require the instruments to be strong. An additional advantage of using many weak instruments is that outliers in the sense of Model 3 are more easily detected, so the results are generally more robust to pleiotropy. There is one caveat: selection bias is more significant for weaker instruments, so a sample-splitting design (such as the one in Section 1.2) should be used.

In Models 2 and 3, we have assumed that the pleiotropy effects are completely independent and normally or nearly normally distributed. We view this assumption as an approximate modeling assumption rather than the precise data generating mechanism. It is motivated by the real data (Section 3.5) and seems to fit the data very well (Section 5.3). It is a special instance of the INstrument Strength Independent of Direct Effect (INSIDE) assumption [9] that is common in the summary-data MR literature. Apart from normality, two other implicit but key assumptions we made are:

1. The pleiotropy effects α_j are additive rather than multiplicative (the variance of α_j is proportional to σ_{Y_j}) [7]. Multiplicative random effects model are easier to fit especially if the measurement error in $\hat{\gamma}_j$ is ignored; however, it is quite unrealistic because α_j is a population quantity and thus is unlikely to be dependent on a sample quantity (for example, σ_{Y_j} may vary due to missing data). In contrast, the additive model is well motivated by the linear structural model in 2.3.

2. The pleiotropy effects α_j have mean 0. In comparison, the MR-Egger regression [7] assumes α_j has an unknown mean μ and refers to the case $\mu \neq 0$ as “directional pleiotropy.” We have not seen strong evidence of “directional pleiotropy” in real datasets, and, more importantly, assuming $\mu \neq 0$ implies that there is a “special” allele coding so that $\alpha_j \sim N(\mu, \tau^2)$. It is thus impossible to obtain estimators of β that are invariant to allele recoding without completely reformulating the MR-Egger model. For further details, see Bowden et al. [11].

There are many technical challenges in the development of this paper. Due to the nature of the many weak IV problem, the asymptotics we considered are quite different from the classical measurement error literature. In Section 3, we showed the profile likelihood is information biased when there are many weak IVs, and in Section 4.1 we showed the profile likelihood is biased when there is overdispersion caused by systematic pleiotropy. This issue is solved by adjusting the profile score, but the proof of the consistency of the APS estimator is nontrivial. Consistency of the the RAPS estimator is even more challenging and still open because the estimating equations may have multiple roots, although we found its practical performance is usually quite benign. A possible solution is to initialize by another robust and consistent estimator (similar to the MM-estimation in robust regression, see Yohai [56]). However, we are not aware of any other provably robust and consistent estimator in our setting, and deriving such estimator is beyond the scope of this paper.

Software and reproducibility. R code for the methods proposed in this paper can be found in the package `mr.raps` that is publicly available at <https://github.com/qingyuanzha0/mr.raps> and can be directly called from `TwoSampleMR`. Numerical examples can be reproduced by running examples in the R package.

SUPPLEMENTARY MATERIAL

Supplement to “Statistical inference in two-sample summary-data mendelian randomization using robust adjusted profile score” (DOI: [10.1214/19-AOS1866SUPP](https://doi.org/10.1214/19-AOS1866SUPP); .pdf). In this supplement, we provide additional justifications of the linear model for GWAS summary data and detailed proof for the theoretical results.

REFERENCES

- [1] 1000 GENOMES PROJECT CONSORTIUM, AUTON, A., BROOKS, L. D., DURBIN, R. M., GARRISON, E. P., KANG, H. M., KORBEL, J. O., MARCHINI, J. L., MCCARTHY, S. et al. (2015). A global reference for human genetic variation. *Nature* **526** 68–74. <https://doi.org/10.1038/nature15393>
- [2] ANDERSON, T. W. and RUBIN, H. (1949). Estimation of the parameters of a single equation in a complete system of stochastic equations. *Ann. Math. Stat.* **20** 46–63. MR0028546 <https://doi.org/10.1214/aoms/1177730090>
- [3] BAIOCCHI, M., CHENG, J. and SMALL, D. S. (2014). Instrumental variable methods for causal inference. *Stat. Med.* **33** 2297–2340. MR3257582 <https://doi.org/10.1002/sim.6128>
- [4] BARNDORFF-NIELSEN, O. (1983). On a formula for the distribution of the maximum likelihood estimator. *Biometrika* **70** 343–365. MR0712023 <https://doi.org/10.1093/biomet/70.2.343>
- [5] BEKKER, P. A. (1994). Alternative approximations to the distributions of instrumental variable estimators. *Econometrica* **62** 657–681. MR1281697 <https://doi.org/10.2307/2951662>
- [6] BOUND, J., JAEGER, D. A. and BAKER, R. M. (1995). Problems with instrumental variables estimation when the correlation between the instruments and the endogenous explanatory variable is weak. *J. Amer. Statist. Assoc.* **90** 443–450.
- [7] BOWDEN, J., DAVEY SMITH, G. and BURGESS, S. (2015). Mendelian randomization with invalid instruments: Effect estimation and bias detection through Egger regression. *Int. J. Epidemiol.* **44** 512–525.
- [8] BOWDEN, J., DAVEY SMITH, G., HAYCOCK, P. C. and BURGESS, S. (2016). Consistent estimation in Mendelian randomization with some invalid instruments using a weighted median estimator. *Genet. Epidemiol.* **40** 304–314.
- [9] BOWDEN, J., DEL GRECO M., F., MINELLI, C., DAVEY SMITH, G., SHEEHAN, N. and THOMPSON, J. (2017). A framework for the investigation of pleiotropy in two-sample summary data mendelian randomization. *Stat. Med.* **36** 1783–1802. MR3648622 <https://doi.org/10.1002/sim.7221>
- [10] BOWDEN, J., FABIOLA DEL GRECO, M., MINELLI, C., LAWLOR, D., SHEEHAN, N., THOMPSON, J. and SMITH, G. D. (2019). Improving the accuracy of two-sample summary-data Mendelian randomization: Moving beyond the NOME assumption. *Int. J. Epidemiol.* **48** 728–742. <https://doi.org/10.1093/ije/dyy258>
- [11] BOWDEN, J., SPILLER, W., DEL-GRECO, F., SHEEHAN, N., THOMPSON, J., MINELLI, C. and SMITH, G. D. (2018). Improving the visualization, interpretation and analysis of two-sample summary data Mendelian randomization via the Radial plot and Radial regression. *Int. J. Epidemiol.* **47** 1264–1278. <https://doi.org/10.1093/ije/dyy101>
- [12] BOYLE, E. A., LI, Y. I. and PRITCHARD, J. K. (2017). An expanded view of complex traits: From polygenic to omnigenic. *Cell* **169** 1177–1186.
- [13] BURGESS, S., BUTTERWORTH, A. and THOMPSON, S. G. (2013). Mendelian randomization analysis with multiple genetic variants using summarized data. *Genet. Epidemiol.* **37** 658–665.
- [14] BURGESS, S., SCOTT, R. A., TIMPSON, N. J., SMITH, G. D., THOMPSON, S. G. and CONSORTIUM, E.-I. (2015). Using published data in Mendelian randomization: A blueprint for efficient identification of causal risk factors. *Eur. J. Epidemiol.* **30** 543–552.
- [15] CARROLL, R. J., RUPPERT, D., STEFANSKI, L. A. and CRAINICEANU, C. M. (2006). *Measurement Error in Nonlinear Models: A Modern Perspective*, 2nd ed. *Monographs on Statistics and Applied Probability* **105**. CRC Press/CRC, Boca Raton, FL. MR2243417 <https://doi.org/10.1201/9781420010138>
- [16] CLARIVATE ANALYTICS (2017). Web of Science Topic: Mendelian Randomization. Available at <http://www.webofknowledge.com>.
- [17] COX, D. R. and REID, N. (1987). Parameter orthogonality and approximate conditional inference. *J. Roy. Statist. Soc. Ser. B* **49** 1–39. MR0893334
- [18] DAVEY SMITH, G. and EBRAHIM, S. (2003). “Mendelian randomization”: Can genetic epidemiology contribute to understanding environmental determinants of disease? *Int. J. Epidemiol.* **32** 1–22.
- [19] DAVEY SMITH, G. and HEMANI, G. (2014). Mendelian randomization: Genetic anchors for causal inference in epidemiological studies. *Hum. Mol. Genet.* **23** R89–R98.

- [20] DAVEY SMITH, G., LAWLOR, D. A., HARBORD, R., TIMPSON, N., DAY, I. and EBRAHIM, S. (2007). Clustered environments and randomized genes: A fundamental distinction between conventional and genetic epidemiology. *PLoS Med.* **4** e352.
- [21] DIDELEZ, V. and SHEEHAN, N. (2007). Mendelian randomization as an instrumental variable approach to causal inference. *Stat. Methods Med. Res.* **16** 309–330. MR2395652 <https://doi.org/10.1177/0962280206077743>
- [22] FISHER, R. A. (1930). *The Genetical Theory of Natural Selection*. Oxford Univ. Press, Oxford.
- [23] GUO, Z., KANG, H., CAI, T. T. and SMALL, D. S. (2018). Confidence intervals for causal effects with invalid instruments by using two-stage hard thresholding with voting. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **80** 793–815. MR3849344 <https://doi.org/10.1111/rssb.12275>
- [24] HAMPEL, F. R. (1974). The influence curve and its role in robust estimation. *J. Amer. Statist. Assoc.* **69** 383–393. MR0362657
- [25] HANSEN, C., HAUSMAN, J. and NEWEY, W. (2008). Estimation with many instrumental variables. *J. Bus. Econom. Statist.* **26** 398–422. MR2459342 <https://doi.org/10.1198/073500108000000024>
- [26] HARTWIG, F. P., DAVEY SMITH, G. and BOWDEN, J. (2017). Robust inference in summary data Mendelian randomization via the zero modal pleiotropy assumption. *Int. J. Epidemiol.* **46** 1985–1998. <https://doi.org/10.1093/ije/dyx102>
- [27] HAYCOCK, P. C., BURGESS, S., WADE, K. H., BOWDEN, J., RELTON, C. and SMITH, G. D. (2016). Best (but oft-forgotten) practices: The design, analysis, and interpretation of Mendelian randomization studies. *Am. J. Clin. Nutr.* **103** 965–978.
- [28] HEMANI, G., ZHENG, J., ELSWORTH, B., WADE, K. H., HABERLAND, V., BAIRD, D., LAURIN, C., BURGESS, S., BOWDEN, J. et al. (2018). The MR-Base platform supports systematic causal inference across the human phenome. *eLife* **7** e34408. <https://doi.org/10.7554/eLife.34408>
- [29] HERNÁN, M. A. and ROBINS, J. M. (2006). Instruments for causal inference: An epidemiologist’s dream? *Epidemiology* **17** 360–372.
- [30] HUBER, P. J. (1964). Robust estimation of a location parameter. *Ann. Math. Stat.* **35** 73–101. MR0161415 <https://doi.org/10.1214/aoms/1177703732>
- [31] IOANNIDIS, J. P., TRIKALINOS, T. A. and KHOURY, M. J. (2006). Implications of small effect sizes of individual genetic variants on the design and interpretation of genetic association studies of complex diseases. *Am. J. Epidemiol.* **164** 609–614.
- [32] KANG, H., ZHANG, A., CAI, T. T. and SMALL, D. S. (2016). Instrumental variables estimation with some invalid instruments and its application to mendelian randomization. *J. Amer. Statist. Assoc.* **111** 132–144. MR3494648 <https://doi.org/10.1080/01621459.2014.994705>
- [33] KATAN, M. (1986). Apoprotein E isoforms, serum cholesterol, and cancer. *Lancet* **327** 507–508.
- [34] LI, S. (2017). Mendelian randomization when many instruments are invalid: Hierarchical empirical Bayes estimation. Available at [arXiv:1706.01389](https://arxiv.org/abs/1706.01389).
- [35] LOCKE, A. E., KAHALI, B., BERNDT, S. I., JUSTICE, A. E., PERS, T. H., DAY, F. R., POWELL, C., VEDANTAM, S., BUCHKOVICH, M. L. et al. (2015). Genetic studies of body mass index yield new insights for obesity biology. *Nature* **518** 197–206.
- [36] MARONNA, R. A., MARTIN, R. D. and YOHAI, V. J. (2006). *Robust Statistics: Theory and Methods*. Wiley Series in Probability and Statistics. Wiley, Chichester. MR2238141 <https://doi.org/10.1002/0470010940>
- [37] MCCULLAGH, P. and TIBSHIRANI, R. (1990). A simple method for the adjustment of profile likelihoods. *J. Roy. Statist. Soc. Ser. B* **52** 325–344. MR1064420
- [38] MURPHY, S. A. and VAN DER VAART, A. W. (1996). Likelihood inference in the errors-in-variables model. *J. Multivariate Anal.* **59** 81–108. MR1424904 <https://doi.org/10.1006/jmva.1996.0055>
- [39] NEYMAN, J. and SCOTT, E. L. (1948). Consistent estimates based on partially consistent observations. *Econometrica* **16** 1–32. MR0025113 <https://doi.org/10.2307/1914288>
- [40] PACINI, D. and WINDMEIJER, F. (2016). Robust inference for the two-sample 2SLS estimator. *Econom. Lett.* **146** 50–54. MR3542584 <https://doi.org/10.1016/j.econlet.2016.06.033>
- [41] PARK, J.-H., WACHOLDER, S., GAIL, M. H., PETERS, U., JACOBS, K. B., CHANOCK, S. J. and CHATTERJEE, N. (2010). Estimation of effect size distribution from genome-wide association studies and implications for future discoveries. *Nat. Genet.* **42** 570–575.
- [42] PEARL, J. (2009). *Causality: Models, Reasoning, and Inference*, 2nd ed. Cambridge Univ. Press, Cambridge. MR2548166 <https://doi.org/10.1017/CBO9780511803161>
- [43] PURCELL, S. PLINK (software V1.07). <http://pngu.mgh.harvard.edu/purcell/plink/>.
- [44] PURCELL, S., NEALE, B., TODD-BROWN, K., THOMAS, L., FERREIRA, M. A., BENDER, D., MALLER, J., SKLAR, P., DE BAKKER, P. I. et al. (2007). PLINK: A tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81** 559–575.

- [45] ROTHENBERG, T. J. (1971). Identification in parametric models. *Econometrica* **39** 577–591. [MR0436944 https://doi.org/10.2307/1913267](https://doi.org/10.2307/1913267)
- [46] SHI, H., KICHAEV, G. and PASANIUC, B. (2016). Contrasting the genetic architecture of 30 complex traits from summary association data. *Am. J. Hum. Genet.* **99** 139–153.
- [47] SOLOVIEFF, N., COTSAPAS, C., LEE, P. H., PURCELL, S. M. and SMOLLER, J. W. (2013). Pleiotropy in complex traits: Challenges and strategies. *Nat. Rev. Genet.* **14** 483–495.
- [48] STEARNS, F. W. (2010). One hundred years of pleiotropy: A retrospective. *Genetics* **186** 767–773.
- [49] STOCK, J. H. and YOGO, M. (2005). Asymptotic distributions of instrumental variables statistics with many instruments. In *Identification and Inference for Econometric Models* 109–120. Cambridge Univ. Press, Cambridge. [MR2232141 https://doi.org/10.1017/CBO9780511614491.007](https://doi.org/10.1017/CBO9780511614491.007)
- [50] TCHETGEN TCHETGEN, E. J., SUN, B. and WALTER, S. (2017). The GENIUS approach to robust Mendelian randomization inference. Available at [arXiv:1709.07779](https://arxiv.org/abs/1709.07779).
- [51] VAN KIPPERSLUIS, H. and RIETVELD, C. A. (2017). Pleiotropy-robust Mendelian randomization. *Int. J. Epidemiol.* In press.
- [52] VERBANCK, M., CHEN, C.-Y., NEALE, B. and DO, R. (2018). Detection of widespread horizontal pleiotropy in causal relationships inferred from Mendelian randomization between complex traits and diseases. *Nat. Genet.* **50** 693.
- [53] WALD, A. (1940). The fitting of straight lines if both variables are subject to error. *Ann. Math. Stat.* **11** 285–300. [MR0002739 https://doi.org/10.1214/aoms/1177731868](https://doi.org/10.1214/aoms/1177731868)
- [54] WRIGHT, P. G. (1928). *Tariff on Animal and Vegetable Oils*. MacMillan, New York.
- [55] WRIGHT, S. (1968). *Evolution and the Genetics of Populations, Volume 1: Genetic and Biometric Foundations* **1**. Univ. Chicago Press, Chicago.
- [56] YOHAI, V. J. (1987). High breakdown-point and high efficiency robust estimates for regression. *Ann. Statist.* **15** 642–656. [MR0888431 https://doi.org/10.1214/aos/1176350366](https://doi.org/10.1214/aos/1176350366)
- [57] ZHAO, Q., WANG, J., HEMANI, G., BOWDEN, J. and SMALL, D. S. (2020). Supplement to “Statistical inference in two-sample summary-data Mendelian randomization using robust adjusted profile score.” <https://doi.org/10.1214/19-AOS1866SUPP>.