



北京工商大學
BEIJING TECHNOLOGY AND BUSINESS UNIVERSITY



Identification and Estimation of the Bi-Directional MR with Some Invalid Instruments

Feng Xie, Zhen Yao, Lin Xie, Yan Zeng, Zhi Geng

School of Mathematics and Statistics

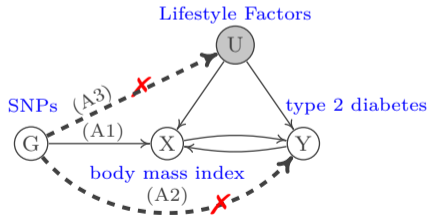
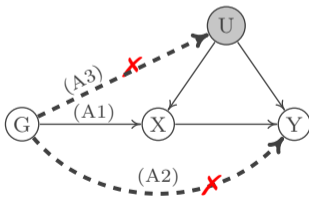
December 14, 2024



Background



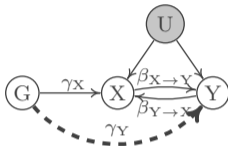
- Mendelian Randomization (MR)
 - A powerful method to **estimate causal effects**.
 - Widely applied in fields of clinics, socioeconomics, drug targets, etc.
 - Implemented as a form of instrumental variables analysis.
- Instrumental Variables (IVs)
 - (A1) Relevance
 - (A2) Exclusion Restriction
 - (A3) Randomness



- Most existing methods focus on the one-directional MR. However, bi-directional relationships are ubiquitous in real-life scenarios.

How to identify valid IV sets for bi-directional MR models from observational data?





Two Stage Least Square(TSLS) IV Estimator:

$$\hat{\beta}_{X \rightarrow Y} = [X^T P X]^{-1} X^T P Y = \beta_{X \rightarrow Y},$$

where $P = (G_V^{X \rightarrow Y})^T [G_V^{X \rightarrow Y} (G_V^{X \rightarrow Y})^T]^{-1} G_V^{X \rightarrow Y}$.

Definition (Bi-directional MR Causal Models)

$$X = Y\beta_{Y \rightarrow X} + G^T \gamma_X + \varepsilon_X,$$

$$Y = X\beta_{X \rightarrow Y} + G^T \gamma_Y + \varepsilon_Y.$$

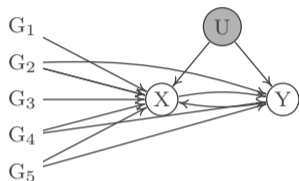
$$X = (G^T \gamma_X + G^T \gamma_Y \beta_{Y \rightarrow X} + \varepsilon_X + \varepsilon_Y \beta_{Y \rightarrow X}) \Delta,$$

$$Y = (G^T \gamma_X \beta_{X \rightarrow Y} + G^T \gamma_Y + \varepsilon_X \beta_{X \rightarrow Y} + \varepsilon_Y) \Delta.$$

Goals

1. Identify valid IVs and invalid IVs.
2. Estimate the causal effects $\beta_{X \rightarrow Y}$ and $\beta_{Y \rightarrow X}$.

Motivating Example



$$X = Y\beta_{Y \rightarrow X} + G_1\gamma_{X,1} + G_2\gamma_{X,2} + G_3\gamma_{X,3} + G_4\gamma_{X,4} + G_5\gamma_{X,5} + \varepsilon_X,$$

$$Y = X\beta_{X \rightarrow Y} + G_2\gamma_{Y,2} + G_4\gamma_{Y,4} + G_5\gamma_{Y,5} + \varepsilon_Y.$$

$$\begin{cases} \text{corr}(Y - X\omega_{\{G_3\}}, G_1) = 0, \\ \text{corr}(Y - X\omega_{\{G_1\}}, G_3) = 0. \end{cases}$$

$$\begin{cases} \text{corr}(Y - X\omega_{\{G_4, G_5\}}, G_2) \neq 0, \\ \text{corr}(Y - X\omega_{\{G_2, G_5\}}, G_4) \neq 0, \\ \text{corr}(Y - X\omega_{\{G_2, G_4\}}, G_5) \neq 0. \end{cases}$$

where,

$\text{corr}(\cdot)$ denotes the Pearson's correlation coefficient between two random variables,

$\omega_{\{G_i\}} = \text{TSLS}(X, Y, \{G_i\})$ with $i \in \{1, 3\}$, and

$\omega_{\{G_i, G_j\}} = \text{TSLS}(X, Y, \{G_i, G_j\})$ with $i \neq j$ and $i, j \in \{2, 4, 5\}$.

Identifying Invalid IV Sets



Definition (Pseudo-Residual)

Let \mathbb{G} be a subset of candidate genetic variants. A pseudo-residual of $\{X, Y\}$ relative to \mathbb{G} is defined as $\mathcal{PR}_{(X, Y | \mathbb{G})} := Y - X\omega_{\mathbb{G}}$, where $\omega_{\mathbb{G}} = \text{TSLS}(X, Y, \mathbb{G})$.

Assumption 1. (Valid IV Set)

For a given causal relationship (if the relationship exists), there exists a valid IV set that consists of at least two valid IVs. For example, for the causal relationship $X \rightarrow Y$, $|\mathbb{G}_{\mathcal{V}}^{X \rightarrow Y}| \geq 2$.

Proposition (Identifying Invalid IV Sets)

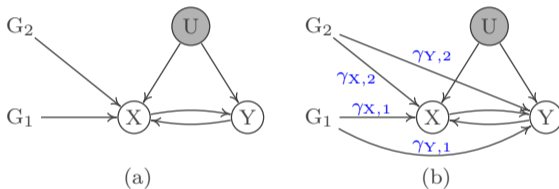
Let $\mathbb{G} = \{G_1, \dots, G_p\}$, $p \geq 2$ be a subset of candidate genetic variants \mathbb{G} . Suppose that Assumption 1 and the size of participants $n \rightarrow \infty$ hold. If there exists a $G_j \in \mathbb{G}$ such that, $\text{corr}(\mathcal{PR}_{(X, Y | \mathbb{G} \setminus G_j)}, G_j) \neq 0$, then \mathbb{G} is an invalid IV set for any one of the causal relationships.

Counter Example



Question: Does the above proposition identify all invalid IV sets?

⇒ Assumption 1 is not sufficient. See the following example.



$$\begin{cases} \text{corr}(\mathcal{PR}_{(X,Y|\{G_2\})}, G_1) = 0 \\ \text{corr}(\mathcal{PR}_{(X,Y|\{G_1\})}, G_2) = 0 \end{cases}$$

$$\begin{cases} \text{corr}(\mathcal{PR}_{(X,Y|\{G_2\})}, G_1) = \frac{\gamma_{Y,1}\gamma_{X,2} - \gamma_{Y,2}\gamma_{X,1}}{\beta_{Y \rightarrow X}\gamma_{Y,2} + \gamma_{X,2}} \\ \text{corr}(\mathcal{PR}_{(X,Y|\{G_1\})}, G_2) = \frac{\gamma_{Y,2}\gamma_{X,1} - \gamma_{Y,1}\gamma_{X,2}}{\beta_{Y \rightarrow X}\gamma_{Y,1} + \gamma_{X,1}} \end{cases}$$

Assumption 2. (Generic Identifiability)

For a given MR causal model, parameters γ and β live in a set of Lebesgue measure non-zero.

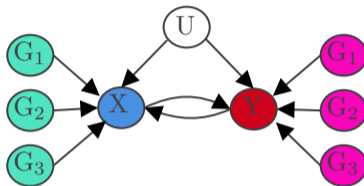
Proposition (Identifying Valid IV Sets)

Let $\mathbb{G} = \{G_1, \dots, G_p\}$, $p \geq 2$ be a subset of candidate genetic variants G . Suppose that Assumption 1 and Assumption 2, and the size of participants $n \rightarrow \infty$ hold. If each $G_j \in \mathbb{G}$ satisfies

$$\text{corr}(\mathcal{PR}_{(X,Y|\mathbb{G}\setminus G_j)}, G_j) = 0, \quad (1)$$

then \mathbb{G} is a proper valid IV set for one of the causal relationships, i.e., $\mathbb{G} \subseteq G_V^{X \rightarrow Y}$ or $\mathbb{G} \subseteq G_V^{Y \rightarrow X}$.

Identifying the Causal Direction



Assumption 3.

In a bi-directional MR model, for a given causal relationship $X \rightarrow Y$, $\beta_{X \rightarrow Y}^2 \cdot \text{Var}(X) < \text{Var}(Y)$. Similarly, for a given causal relationship $Y \rightarrow X$, $\beta_{Y \rightarrow X}^2 \cdot \text{Var}(Y) \leq \text{Var}(X)$ [Xue and Pan, 2020].

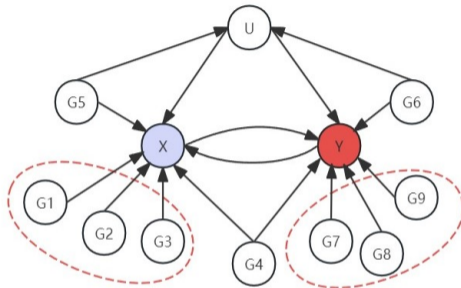
Proposition (Identifying the Direction for the Identified IV set)

Suppose Assumption 3 holds. Given a valid IV set denoted as \mathbb{G} , if each $G_j \in \mathbb{G}$ satisfies the condition $|\text{corr}(G_j, Y)/\text{corr}(G_j, X)| < 1$, then \mathbb{G} is a valid IV set for the causal relationship $X \rightarrow Y$, i.e., $\mathbb{G} \subseteq G_V^{X \rightarrow Y}$. Conversely, if this condition is not satisfied, \mathbb{G} is a valid IV set for the causal relationship $Y \rightarrow X$, i.e., $\mathbb{G} \subseteq G_V^{Y \rightarrow X}$.

Identifiability of Bi-directional MR model

Suppose that Assumptions 1 - 3, and the size of participants $n \rightarrow \infty$ hold, the Bi-directional MR model is fully identifiable, which includes the identification of valid IV sets, the causal relationships the identified sets relate to, the causal directions between X and Y, and the causal effects of X on Y as well as of Y on X.

- Find valid IV sets.
- Infer directions and estimate effects.



Algorithm 1 PReBiM

Input: A dataset of measured genetic variants $\mathbf{G} = (G_1, \dots, G_g)^\top$, two phenotypes X and Y , significance level α , and parameter W , maximum number of IVs to consider.

1: Valid IV sets $\mathcal{V} \leftarrow \text{FindValidIVSets}(\mathbf{G}, X, Y, \alpha, W)$;

▷ *Step I*

2: $(\hat{\beta}_{X \rightarrow Y}, \hat{\beta}_{Y \rightarrow X}) \leftarrow \text{InferCausalDirectionEffects}(X, Y, \mathcal{V})$;

▷ *Step II*

Output: $\hat{\beta}_{X \rightarrow Y}$ and $\hat{\beta}_{Y \rightarrow X}$, the causal effects of X on Y and Y on X , respectively.

Experiments on Synthetic Data



Table 1: Performance comparison of NAIVE, MR-Egger, sisVIVE, IV-TETRAD, TSHT, and PRe-BiM in estimating bi-directional MR models across various sample sizes and three scenarios.

Size	Algorithm	$S(2, 2, 6)$				$S(3, 3, 8)$				$S(4, 4, 10)$			
		X \rightarrow Y		Y \rightarrow X		X \rightarrow Y		Y \rightarrow X		X \rightarrow Y		Y \rightarrow X	
		CSR \uparrow	MSE \downarrow	CSR \uparrow	MSE \downarrow	CSR \uparrow	MSE \downarrow	CSR \uparrow	MSE \downarrow	CSR \uparrow	MSE \downarrow	CSR \uparrow	MSE \downarrow
2k	NAIVE	-	0.354	-	0.349	-	0.351	-	0.320	-	0.335	-	0.304
	MR-Egger	-	0.800	-	0.734	-	0.571	-	0.569	-	0.457	-	0.494
	sisVIVE	0.30	0.398	0.25	0.447	0.34	0.379	0.32	0.379	0.37	0.346	0.38	0.330
	IV-TETRAD	0.60	0.859	0.30	0.867	0.60	0.773	0.28	0.782	0.67	0.688	0.23	0.622
	TSHT	0.07	0.606	0.06	0.660	0.07	0.549	0.07	0.613	0.12	0.653	0.11	0.689
	PReBiM	0.85	0.046	0.85	0.070	0.89	0.083	0.87	0.075	0.88	0.039	0.89	0.057
5k	NAIVE	-	0.350	-	0.339	-	0.347	-	0.319	-	0.357	-	0.306
	MR-Egger	-	0.737	-	0.836	-	0.604	-	0.508	-	0.515	-	0.520
	sisVIVE	0.32	0.421	0.30	0.407	0.33	0.378	0.35	0.384	0.39	0.439	0.38	0.365
	IV-TETRAD	0.62	0.804	0.31	0.844	0.63	0.836	0.29	0.659	0.67	0.686	0.25	0.640
	TSHT	0.04	0.544	0.06	0.565	0.08	0.528	0.06	0.570	0.10	0.549	0.07	0.541
	PReBiM	0.93	0.020	0.91	0.027	0.90	0.019	0.92	0.020	0.90	0.009	0.91	0.010
10k	NAIVE	-	0.366	-	0.348	-	0.319	-	0.342	-	0.349	-	0.300
	MR-Egger	-	0.800	-	0.763	-	0.524	-	0.620	-	0.487	-	0.424
	sisVIVE	0.31	0.457	0.29	0.503	0.37	0.383	0.31	0.394	0.41	0.399	0.41	0.328
	IV-TETRAD	0.64	0.811	0.31	0.804	0.64	0.748	0.29	0.696	0.71	0.635	0.24	0.577
	TSHT	0.04	0.449	0.02	0.475	0.05	0.525	0.05	0.466	0.05	0.479	0.07	0.498
	PReBiM	0.95	0.044	0.93	0.018	0.93	0.026	0.93	0.039	0.93	0.027	0.94	0.011

- Major IV information for the bi-directional dataset. (Here X is obesity and Y denotes Vitamin D Status)
 - Valid IVs ($X \rightarrow Y$):
 - rs9939609 (FTO)
 - rs2867125 (TMEM18)
 - rs4074134 (BDNF)
 - rs7647305 (ETV5)
 - rs7138803 (FAIM2)
 - rs3101336 (NEGR1)
 - rs10938397 (GNPDA2)
 - Valid IVs ($Y \rightarrow X$):
 - rs12785878 (DHCR7)
 - rs10741657 (CYP2R1)
 - rs2282679 (GC)
 - rs6013897 (CYP24A1)
 - **Our Results:** FTO and FAIM2 ($X \rightarrow Y$); DHCR7 and CYP24A1 ($Y \rightarrow X$)



- We provided the identifiability conditions for the bi-directional MR model, enabling both valid IV sets for each direction and the causal effects of interests to be correctly identified.
- We proposed a practical and effective cluster fusion-like algorithm for unbiased estimation of valid IV sets.
- Experimental results on synthetic and real-world data verified the effectiveness of our algorithm.

Thank you!

