Using Mendelian randomization to discover biological mechanisms

Qingyuan Zhao

Statistical Laboratory, University of Cambridge

August 10, 2023 @ JSM

Example: Causal effect of the LDL-cholesterol



Basic idea: People who inherited certain alleles of *rs17238484* and *rs12916* have **naturally** higher concentration of LDL cholesterol.

Example: Causal effect of the LDL-cholesterol



Basic idea: People who inherited certain alleles of *rs17238484* and *rs12916* have **naturally** higher concentration of LDL cholesterol.

Different perspectives

- Statistical: much of the previous methodological work focuses on violations of exclusion restriction due to pleiotropy.
- Biological: the fundamental assumption of MR is gene-environment equivalence.
- Drug development: the most useful is to discover relevant pathways.

This talk

How can we use MR to discover biological mechanisms?

MR-PATH Daniel long, Qingyuan Zhao, and Yang Chen. "A Latent Mixture Model for Heterogeneous Causal Mechanisms in Mendelian Randomization". In: (2020). arXiv: 2007.06476 [stat.AP].

PathGPS Zijun Gao, Trevor Hastie, and Qingyuan Zhao. *PathGPS: Discover* shared genetic architecture using biobank data. 2022. bioRxiv: 10.1101/2022.05.01.490230.



Mechanistic Heterogeneity

MR-PATH

PathGPS

Conclusion

◆□ > ◆□ > ◆ Ξ > ◆ Ξ > → Ξ → のへで

Outline

Mechanistic Heterogeneity

Reivew Conceptual illustration

MR-PATH

PathGPS

Conclusion

Review: Linear structural equation model for MR



For exposure X, outcome Y, unobserved confounding variables U, and SNPs Z_1, \ldots, Z_p , the commonly assumed linear structural equation model is given by

$$X = \sum_{i=1}^{p} \theta_{X_i} Z_i + \eta_X U + E_X, \quad Y = \beta X + \sum_{i=1}^{p} \alpha_i Z_i + \eta_Y U + E_Y.$$

Review: Linear structural equation model for MR

$$X = \sum_{i=1}^{p} \theta_{X_i} Z_i + \eta_X U + E_X,$$

$$Y = \beta X + \sum_{i=1}^{p} \alpha_i Z_i + \eta_Y U + E_Y.$$

$$\Rightarrow Y = \sum_{i=1}^{p} (\theta_{X_i} \beta + \alpha_i) Z_i + (\beta \eta_X + \eta_Y) U + E_Y.$$

▶ If Z_i is a valid instrument, $\theta_{X_i} \neq 0$, $Z_i \perp \{U, E_X, E_Y\}$, and $\alpha_i = 0$.

- ▶ However, often $\alpha_i \neq 0$ due to pleiotropy and mechanistic heterogeneity.
- If α_i ≠ 0 for some SNPs, then the causal effect β cannot be estimated consistently without further assumptions on α_i.

• e.g. $\alpha_i \sim N(0, \tau^2)$ for most SNPs.

Mechanistic heterogeneity: Two scenarios



(a) Scenario 1: Multiple pathways of correlated pleiotropy.



(b) Scenario 2: Multiple mechanisms for the exposure *X*.

Mechanistic heterogeneity: Two scenarios

If we interpret the diagrams in the previous slide as linear structural equations as before, we can derive the Wald estimands for each pathway.

Instruments Z	Pathway M	Effect of M on X	Effect of M on Y	Wald estimand
Scenario 1				
$Z_{1,1}, \ldots, Z_{1,p_1}$	M_1	θ_1	$\theta_1 eta$	β
$Z_{2,1}, \ldots, Z_{2,p_2}$	M_2	θ_2	$\theta_2\beta + \alpha_2$	$\beta + \alpha_2/\theta_2$
$Z_{3,1},\ldots,Z_{3,p_3}$	M ₃	θ_3	$\theta_3\beta + \alpha_3$	$\beta + \alpha_3/\theta_3$
Scenario 2				
$Z_{1,1}, \ldots, Z_{1,p_1}$	M_1	θ_1	$\theta_1 \beta_1$	β_1
$Z_{2,1}, \ldots, Z_{2,p_2}$	M_2	θ_2	$\theta_2 \beta_2$	β_2
$Z_{3,1},\ldots,Z_{3,p_3}$	M_3	θ_3	$ heta_3eta_3$	eta_{3}

- SNPs on the same pathway have the same Wald estimand, while SNPs across different pathways generally have different estimands.
- Mechanistic heterogeneity can arise even when all SNPs are valid instruments (Scenario 2).

Nonparametric perspective: Mechanism-specific causal effect

The same clustering phenomenon also occurs in nonlinear models.

- It is well known that assuming monotonicity, IV nonparametrically estimates the complier average treatment effect (Angrist, Imbens, Rubin, JASA, 1996).
- By assuming monotonicity and Pearl's nonparametric structural equation model with independent errors (NPSEM-IE), our paper showed that (if X, Z, M are all binary variables)

$$\mathbb{E}[Y(X=1)-Y(X=0) \mid X(Z_{kj}=1) > X(Z_{kj}=0)] \ = \mathbb{E}[Y(X=1)-Y(X=0) \mid X(M_k=1) > X(M_k=0)],$$

where k indexes the mechanism and j indexes the gene within.

Outline

Mechanistic Heterogeneity

MR-PATH

Model Assumptions Statistical method Results

PathGPS

Conclusion

MR-PATH: Model Assumptions

We make the following assumptions on GWAS summary data.

Assumption (Error-in-variables regression)

The observed (marginal) SNP-exposure and SNP-outcome associations are distributed as $\begin{pmatrix} \hat{n} \\ \hat{n} \end{pmatrix} = \begin{pmatrix} -2 \\ -2 \\ -2 \end{pmatrix} = \begin{pmatrix} -2 \\ -2 \\ -2 \end{pmatrix}$

$$\begin{pmatrix} \hat{\theta}_{X_i} \\ \hat{\theta}_{Y_i} \end{pmatrix} \stackrel{indep.}{\sim} \mathsf{N}\Big(\begin{pmatrix} \theta_{X_i} \\ \beta_i \theta_{X_i} \end{pmatrix}, \begin{pmatrix} \sigma_{X_i}^2 & 0 \\ 0 & \sigma_{Y_i}^2 \end{pmatrix} \Big), \quad i = 1, \dots, p,$$

where σ_{X_i} , σ_{Y_i} are (fixed) measurement errors.

Assumption (Mixture model for mechanistic heterogeneity)

$$Z_i \sim \textit{Multinomial}(\pi_1, \dots, \pi_K),$$

 $eta_i | Z_i = k \sim \textit{N}(\mu_k, \sigma_k^2), \quad k = 1, \dots, K.$

・ロト・西ト・山田・山田・山下

MR-PATH: Statistical Inference

- 1. Monte-Carlo EM algorithm for obtaining model parameter estimates
- 2. Approximate confidence intervals for quantifying uncertainty of the estimates
- 3. Modified Bayesian Information criterion (BIC) for selecting number of clusters
- ▶ We perform simulation studies to verify the efficacy of these inference procedures.

See paper for implementation details.

Example: HDL-CHD



◆□ > ◆昼 > ◆臣 > ◆臣 > ○ ● ○ ●

Example: HDL-CHD



Results of MR-PATH (http://danieliong.me/mr-path/.)

・ロット (四)・ (目)・ (日)・ (日)

Example: HDL-CHD (individual β_i)



Example: HDL-CHD (heatmap against metabolites)



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Example: BMI-T2D



◆□ > ◆□ > ◆ Ξ > ◆ Ξ > → Ξ = の < @

Example: BMI-T2D



Results of MR-PATH.

◆□ → ◆□ → ◆三 → ◆三 → ◆○ ◆

Example: BMI-T2D (individual β_i)



▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ 三臣 - のへ⊙

Example: BMI-T2D (possible correlated pleiotropy: insulin)



Outline

Mechanistic Heterogeneity

MR-PATH

PathGPS Model Assumptions

Statistical method Results

Conclusion



PathGPS: PATHway discovery through Genome-Phenome Summary data

Goal: Discover SNP-trait clusters corresponding to the same genetic pathway.



Genetic and environmental pathways.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00

PathGPS: PATHway discovery through Genome-Phenome Summary data

Goal: Discover SNP-trait clusters corresponding to the same genetic pathway.



Two SNP-trait clusters: (1) X_1 , X_2 , Y_1 , Y_2 ; (2) X_2 , X_3 , Y_3 .

Statistical model: Linear SEM

• Observed data: **X** genotypes/SNPs $(n \times p)$, **Y** phenotypes/traits $(n \times q)$;

▶ Latent mediators/confounders: **M** genetic $(n \times r)$, **m** environmental $(n \times s)$.



Statistical model: Linear SEM

• Observed data: **X** genotypes/SNPs $(n \times p)$, **Y** phenotypes/traits $(n \times q)$;

▶ Latent mediators/confounders: **M** genetic $(n \times r)$, **m** environmental $(n \times s)$.

$$\underbrace{\mathbf{M}}_{n \times r} = \underbrace{\mathbf{X}}_{n \times p} \underbrace{\mathbf{U}}_{p \times r} + \underbrace{\mathbf{\varepsilon}}_{n \times r}, \quad \underbrace{\mathbf{Y}}_{n \times q} = \underbrace{\mathbf{M}}_{n \times r} \underbrace{\mathbf{V}}_{r \times q}^{\top} + \underbrace{\mathbf{m}}_{n \times s} \underbrace{\mathbf{W}}_{s \times q}^{\top} + \underbrace{\mathbf{\varepsilon}}_{n \times q}.$$
$$\Longrightarrow \mathbf{Y} = \mathbf{X} \mathbf{U} \mathbf{V}^{\top} + \mathbf{m} \mathbf{W}^{\top} + \varepsilon, \quad \varepsilon = \varepsilon_{M} \mathbf{V}^{\top} + \varepsilon_{Y}.$$

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

▶ Non-zero elements in $U_{\cdot k}$, $V_{\cdot k} \rightarrow k$ -th SNP-trait cluster.

Statistical model: Linear SEM

• Observed data: **X** genotypes/SNPs $(n \times p)$, **Y** phenotypes/traits $(n \times q)$;

▶ Latent mediators/confounders: **M** genetic $(n \times r)$, **m** environmental $(n \times s)$.

$$\underbrace{\mathbf{M}}_{n \times r} = \underbrace{\mathbf{X}}_{n \times p} \underbrace{\mathbf{U}}_{p \times r} + \underbrace{\mathbf{\varepsilon}}_{n \times r}, \quad \underbrace{\mathbf{Y}}_{n \times q} = \underbrace{\mathbf{M}}_{n \times r} \underbrace{\mathbf{V}}_{r \times q}^{\top} + \underbrace{\mathbf{m}}_{n \times s} \underbrace{\mathbf{W}}_{s \times q}^{\top} + \underbrace{\mathbf{\varepsilon}}_{n \times q}.$$
$$\Longrightarrow \mathbf{Y} = \mathbf{X} \mathbf{U} \mathbf{V}^{\top} + \mathbf{m} \mathbf{W}^{\top} + \varepsilon, \quad \varepsilon = \varepsilon_{M} \mathbf{V}^{\top} + \varepsilon_{Y}.$$

▶ Non-zero elements in $U_{\cdot k}$, $V_{\cdot k} \rightarrow k$ -th SNP-trait cluster.

Assumptions

- ► Low-rank UV^{\top} : #genetic mediators \ll #SNPs, #traits.
- Sparse **U**, **V**: a few influential SNPs/influenced traits.

Singular value decomposition (SVD) for $\hat{\beta}$.

• Suppose $\operatorname{Var}(\boldsymbol{m}) = \sigma_m^2 \boldsymbol{I}_s$, $\operatorname{Var}(\boldsymbol{\varepsilon}_M) = \sigma_M^2 \boldsymbol{I}_r$, $\operatorname{Var}(\boldsymbol{\varepsilon}_Y) = \sigma_Y^2 \boldsymbol{I}_q$,

$$\hat{\boldsymbol{\beta}}^{\top} \hat{\boldsymbol{\beta}} \approx \underbrace{\boldsymbol{\mathcal{V}} \boldsymbol{\mathcal{U}}^{\top} \boldsymbol{\mathcal{U}} \boldsymbol{\mathcal{V}}^{\top}}_{\text{gene}} + \frac{p}{n} \left(\underbrace{\boldsymbol{\sigma}_{m}^{2} \boldsymbol{\mathcal{W}} \boldsymbol{\mathcal{W}}^{\top}}_{\text{environment}} + \boldsymbol{\sigma}_{M}^{2} \boldsymbol{\mathcal{V}} \boldsymbol{\mathcal{V}}^{\top} + \boldsymbol{\sigma}_{Y}^{2} \boldsymbol{\boldsymbol{I}}_{q} \right).$$

▶ Right singular vectors: genetic effect (**V**) & environmental effect (**W**).

Singular value decomposition (SVD) for $\hat{\beta}$.

• Suppose $\operatorname{Var}(\boldsymbol{m}) = \sigma_m^2 \boldsymbol{I}_s$, $\operatorname{Var}(\boldsymbol{\varepsilon}_M) = \sigma_M^2 \boldsymbol{I}_r$, $\operatorname{Var}(\boldsymbol{\varepsilon}_Y) = \sigma_Y^2 \boldsymbol{I}_q$,

$$\hat{\boldsymbol{\beta}}^{\top} \hat{\boldsymbol{\beta}} \approx \underbrace{\boldsymbol{\mathcal{V}} \boldsymbol{\mathcal{U}}^{\top} \boldsymbol{\mathcal{U}} \boldsymbol{\mathcal{V}}^{\top}}_{\text{gene}} + \frac{p}{n} \left(\underbrace{\sigma_m^2 \boldsymbol{\mathcal{W}} \boldsymbol{\mathcal{W}}^{\top}}_{\text{environment}} + \sigma_M^2 \boldsymbol{\mathcal{V}} \boldsymbol{\mathcal{V}}^{\top} + \sigma_Y^2 \boldsymbol{\boldsymbol{I}}_q \right).$$

▶ Right singular vectors: genetic effect (**V**) & environmental effect (**W**).

Challenge 1: Genetics vs. Environment

Solution: Use negative control SNPs ($\boldsymbol{U}_{nc}=0$) and apply SVD to the contrast.

$$\hat{\boldsymbol{\beta}}_{nc}^{\top}\hat{\boldsymbol{\beta}}_{nc} \approx \underbrace{\mathbf{0}}_{\text{gene}} + \frac{p_{nc}}{n} \left(\underbrace{\sigma_m^2 \boldsymbol{W} \boldsymbol{W}^{\top}}_{\text{environment}} + \sigma_M^2 \boldsymbol{V} \boldsymbol{V}^{\top} + \sigma_Y^2 \boldsymbol{I}_q \right)$$

▶ Next step: Cluster SNPs and traits by non-zero entries of $\hat{\boldsymbol{V}}$ and $\hat{\boldsymbol{U}} = \hat{\boldsymbol{\beta}} \hat{\boldsymbol{V}}^{\top}$.

Challenge 2: Instability of unsupervised learning

Many hyper-parameters: number of principal components and clusters, rotation and thresholding, clustering algorithm.

Next step: Cluster SNPs and traits by non-zero entries of $\hat{\boldsymbol{V}}$ and $\hat{\boldsymbol{U}} = \hat{\boldsymbol{\beta}} \hat{\boldsymbol{V}}^{\top}$.

Challenge 2: Instability of unsupervised learning

Many hyper-parameters: number of principal components and clusters, rotation and thresholding, clustering algorithm.

Solution: Bootstrap aggregating (bagging)



◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶

Example: Metabolomics dataset



Example: UK Biobank





Mechanistic Heterogeneity

MR-PATH

PathGPS

Conclusion

Concluding remarks

Genetic variation can be used as instrumental variables for biological mechanisms.

- ▶ MR-PATH: many SNPs, 2 traits (exposure and outcome).
- PathGPS: many SNPs, many traits.
- Biobanks provide a rich data source for MR and pathway discovery.
- Many open problems: unified methodology? theoretical guarantees? more sophisticated/realistic models? integration with drug discovery?