

观察性研究

① $Z \perp\!\!\!\perp (Y(1), Y(0)) \mid X$ — 技术性

② $Z \not\perp\!\!\!\perp (Y(1), Y(0)) \mid X$ — 概念性

Peng Ding* and Luke W. Miratrix

To Adjust or Not to Adjust? Sensitivity Analysis of *M*-Bias and Butterfly-Bias

Pearl
启发

Abstract: “*M*-Bias,” as it is called in the epidemiologic literature, is the bias introduced by conditioning on a pretreatment covariate due to a particular “*M*-Structure” between two latent factors, an observed treatment, an outcome, and a “collider.” This potential source of bias, which can occur even when the treatment and the outcome are not confounded, has been a source of considerable controversy. We here present formulae for identifying under which circumstances biases are inflated or reduced. In particular, we show that the magnitude of *M*-Bias in linear structural equation models tends to be relatively small compared to confounding bias, suggesting that it is generally not a serious concern in many applied settings. These theoretical results are consistent with recent empirical findings from simulation studies. We also generalize the *M*-Bias setting (1) to allow for the correlation between the latent factors to be nonzero and (2) to allow for the collider to be a confounder between the treatment and the outcome. These results demonstrate that mild deviations from the *M*-Structure tend to increase confounding bias more rapidly than *M*-Bias, suggesting that choosing to condition on any given covariate is generally the superior choice. As an application, we re-examine a controversial example between Professors Donald Rubin and Judea Pearl.



Keywords: causality, collider, confounding, controversy, covariate

参考 Rosenbaum
& Rubin

$$e(x) = \Pr(Z=1 \mid x)$$

Instrumental variables as bias amplifiers with general outcome and confounding

By P. DING

Department of Statistics, University of California, 425 Evans Hall, Berkeley,
California 94720, U.S.A.
pengdingpku@berkeley.edu

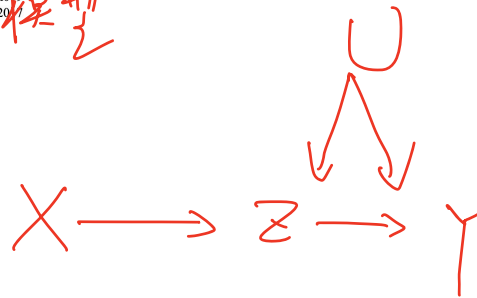
T. J. VANDERWEELE AND J. M. ROBINS

Department of Epidemiology, Harvard T. H. Chan School of Public Health, 677 Huntington
Avenue, Boston, Massachusetts 02115, U.S.A.
tvanderw@hsph.harvard.edu robins@hsph.harvard.edu

SUMMARY

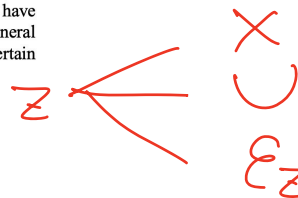
Drawing causal inference with observational studies is the central pillar of many disciplines. One sufficient condition for identifying the causal effect is that the treatment-outcome relationship is unconfounded conditional on the observed covariates. It is often believed that the more covariates we condition on, the more plausible this unconfoundedness assumption is. This belief has had a huge impact on practical causal inference, suggesting that we should adjust for all pretreatment covariates. However, when there is unmeasured confounding between the treatment and outcome, estimators adjusting for some pretreatment covariate might have greater bias than estimators that do not adjust for this covariate. This kind of covariate is called a bias amplifier, and includes instrumental variables that are independent of the confounder and affect the outcome only through the treatment. Previously, theoretical results for this phenomenon have been established only for linear models. We fill this gap in the literature by providing a general theory, showing that this phenomenon happens under a wide class of models satisfying certain monotonicity assumptions.

Some key words: Causal inference; Directed acyclic graph; Interaction; Monotonicity; Potential outcome.



Instrumental Variable (IV)

调整 X 增加偏差



The treatment assignment is a function of the instrumental variable, the unmeasured confounder and some other independent random error, which are the three sources of variation of the treatment. If we adjust for the instrumental variable, the treatment variation is driven more by the unmeasured confounder, which could result in increased bias due to this confounder. Seemingly paradoxically, without adjusting for the instrumental variable, the observational study is more like a randomized experiment, and the bias due to confounding is smaller. Although applied researchers (Myers et al., 2011; Walker, 2013; Brooks & Ohsfeldt, 2013; Ali et al., 2014) have confirmed through extensive simulation studies that this bias amplification phenomenon exists in a wide range of reasonable models, definite theoretical results have been established only for linear models. We fill this gap in the literature by showing that adjusting for an instrumental variable amplifies bias for estimating causal effects under a wide class of models satisfying certain monotonicity assumptions. However, we also show that there exist data-generating processes under which an instrumental variable is not a bias amplifier.

流行病学子集

Generalized Cornfield conditions for the risk difference

BY PENG DING

Department of Statistics, Harvard University, One Oxford Street, Cambridge,
Massachusetts 02138, U.S.A.
pengding@fas.harvard.edu

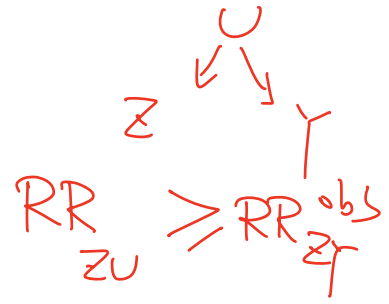
AND TYLER J. VANDERWEELE

Department of Epidemiology, Harvard School of Public Health, Boston, Massachusetts 02115, U.S.A.
tvanderw@hsph.harvard.edu

SUMMARY

A central question in causal inference with observational studies is the sensitivity of conclusions to unmeasured confounding. The classical Cornfield condition allows us to assess whether an unmeasured binary confounder can explain away the observed relative risk of the exposure on the outcome. It states that for an unmeasured confounder to explain away an observed relative risk, the association between the unmeasured confounder and the exposure and the association between the unmeasured confounder and the outcome must both be larger than the observed relative risk. In this paper, we extend the classical Cornfield condition in three directions. First, we consider analogous conditions for the risk difference and allow for a categorical, not just a binary, unmeasured confounder. Second, we provide more stringent thresholds that the maximum of the above-mentioned associations must satisfy, rather than weaker conditions that both must satisfy. Third, we show that all the earlier results on Cornfield conditions hold under weaker assumptions than previously used. We illustrate the potential applications by real examples, where our new conditions give more information than the classical ones.

Some key words: Causal inference; Confounding; Observational study; Sensitivity analysis.



RD 有类似结果
用处不大因为
结果依赖于
U 取值

流行病学子集

Lee (2011) obtained the above results (3) and (4) under Assumption 3, which can in fact be weakened to Assumption 2. Furthermore, in the Supplementary Material, we show that under Assumption 1, the following conditions must hold:

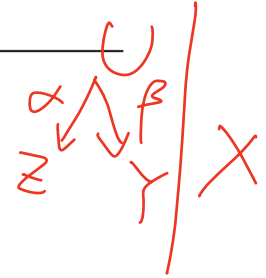
$$\min(U_E, U'_D) \geq RR_{ED}, \quad \max(U_E, U'_D) \geq \{RR_{ED}^{1/2} + (RR_{ED} - 1)^{1/2}\}^2,$$

where $U'_D = \max(U_D, U_D^*)$ replaces U_D in conditions (3) and (4).

Reviewer 反馈 3

E-value

Sensitivity Analysis Without Assumptions

Peng Ding^a and Tyler J. VanderWeele^b

Theorem 17.1 $\sum_n Z \perp\!\!\!\perp Y \mid X, U$

$$RR_{ZY|X}^{obs} \leq \frac{\alpha \beta}{\alpha + \beta - 1}$$

$$\begin{aligned} \frac{H}{c} \cdot \alpha &= RR_{ZU|X} \\ \beta &= RR_{UY|X} \end{aligned}$$

\Downarrow
 $\hat{\alpha}, \hat{\beta}$
 E-value

The claim that our technique is “without assumptions” requires some clarification. As we will see below, we will, without any assumptions, be able to make statements of the form: “For an observed association to be due solely to unmeasured confounding, two sensitivity analysis parameters must satisfy [a specific inequality].” We will also, without assumptions, be able to make statements of the form: “For unmeasured confounding alone to be able to reduce an observed association [to a given level], two sensitivity analysis parameters must satisfy [another specific inequality].” We believe the ability to make statements of this form without imposing any specific structure on the nature of the unmeasured confounder or confounders constitutes a major advance in the literature.

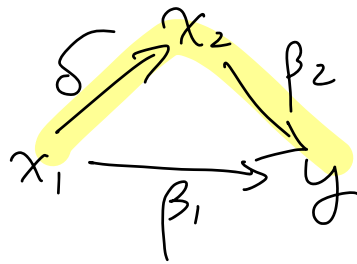
Chapter 16 HW

給分 Cochran's formula

if $\frac{1}{\delta}$ omitted-variable bias formula

$$\begin{cases} y = \beta_1^T x_1 + \beta_2^T x_2 + \varepsilon \\ y = \tilde{\beta}_1^T x_1 + e \end{cases}$$

$$\Rightarrow \tilde{\beta}_1 - \beta_1 = \delta \beta_2$$



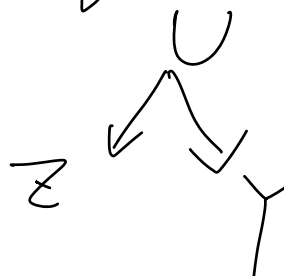
$$x_2 = \delta^T x_1 + v$$

信息论也常见

Lihua Lei
Bin Yu

Data processing inequality

$$\text{若 } Z \perp\!\!\!\perp Y \mid U$$



$$\begin{aligned} \text{则 } I(Z, Y) &\leq I(Z, U) \\ I(Z, Y) &\leq I(U, Y) \end{aligned}$$

~~定义~~ $I(A, B) = \iint p(a, b) \log_2 \frac{p(a, b)}{p(a)p(b)} da db$
mutual information

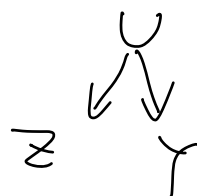
定理 17.1 的初步证明:

Z, Y, U 二值

X 给定: 省略

RR_{ZY}^{obs}

$$\text{定义} = \frac{\Pr(Y=1 | Z=1)}{\Pr(Y=1 | Z=0)}$$



全概率公式

$$= \frac{\Pr(Y=1 | Z=1, U=1) \Pr(U=1 | Z=1) + \Pr(Y=1 | Z=1, U=0) \Pr(U=0 | Z=1)}{\Pr(Y=1 | Z=0, U=1) \Pr(U=1 | Z=0) + \Pr(Y=1 | Z=0, U=0) \Pr(U=0 | Z=0)}$$

$Z \perp\!\!\!\perp Y | U$

$$= \frac{\Pr(Y=1 | U=1) \Pr(U=1 | Z=1) + \Pr(Y=1 | U=0) \Pr(U=0 | Z=1)}{\Pr(Y=1 | U=0) \Pr(U=0 | Z=0) + \Pr(Y=1 | U=0) \Pr(U=0 | Z=0)}$$

f_0 f_1 $1-f_1$
 除此项 $1-f_0$

$$\text{依然正确} \quad = \frac{RR_{UY} f_1 + 1 - f_1}{RR_{UY} f_0 + 1 - f_0}$$

$$= \frac{(RR_{UY} - 1) f_1 + 1}{\frac{(RR_{UY} - 1) f_1 + 1}{RR_{ZU}}}$$

参数化

关于 f_1 \rightarrow

若 $RR_{UY} > 1$

$RR_{ZU} > 1$

$f_1 \rightarrow 1$

$$\frac{\cancel{RR_{UY} - 1} + \cancel{1}}{\frac{RR_{UY} - 1}{RR_{ZU}} + 1}$$

$$= \frac{RR_{ZU} \cdot RR_{UY}}{RR_{ZU} + RR_{UY} - 1}$$

□

$$\Rightarrow RR_{ZU} \geq RR_{ZY}^{obs}$$

$$RR_{UY} \geq RR_{ZY}^{obs}$$

$$\max(RR_{ZU}, RR_{UY}) \geq E\text{-value}$$

$$= RR_{ZY}^{obs} + \sqrt{RR_{ZY}^{obs} (RR_{ZY}^{obs} - 1)}$$

流行病学. 病例-对照研究
 Case-control study
 (case-referent study)
 A2.6.3

$$Pr(S=1 | x, y) = Pr(S=1 | y)$$

↓ ↓
 协变量 结果
 干预

$$(x_i, y_i, S_i=1)_{i=1}^n : \text{not}$$

Cornfield

Ross Prentice : logistic regression

结论: logistic 对 $\begin{cases} \text{iid} \\ \text{case-control} \end{cases}$

都 对

$$P(Y=1 | Z, X) = \frac{e^{\beta_0 + \beta_1 Z + \beta_2^T X}}{1 + e^{\beta_0 + \beta_1 Z + \beta_2^T X}}$$

$\Rightarrow \beta_1 = \text{conditional OR}_{ZY/X}$

rare disease

\approx conditional $RR_{ZY/X}$

关于 τ 的敏感性分析

Sensitivity analysis
(Rosenbaum & Rubin
(1983 JRSSB))

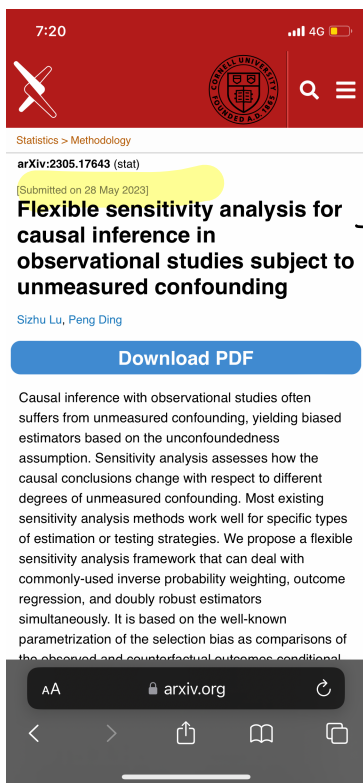
比较参数化与做没

Chapter 18 关于 τ

非常新的一章

疫情第一年教学

目的: 一个极其简单的办法
或视角



理论基础

Chapter 18 核心结果

ignorability | 找共同点 | bounds

$E(Y_{(1)} | Z=1, X) = E(Y_{(1)} | Z=0, X)$

$E(Y_{(0)} | Z=1, X) = E(Y_{(0)} | Z=0, X)$

Manski

想 / 3

上下界

Manski $\tau = E(Y_{(1)} - Y_{(0)})$

$$= \left[E(Y_{(1)} | Z=1) Pr(Z=1) + E(Y_{(1)} | Z=0) Pr(Z=0) \right] \text{ 易算 } X$$

$$- \left[E(Y_{(0)} | Z=1) Pr(Z=1) + E(Y_{(0)} | Z=0) Pr(Z=0) \right]$$

反事实 容易

若 $Y \in [l, u]$, 则 τ 有上下界

理想状态:

$$E(Y_{(1)} | Z=0) = E \left\{ E(Y | Z=1, X) \mid Z=0 \right\}$$

$$E(Y_{(0)} | Z=1) = E \left\{ E(Y | Z=0, X) \mid Z=1 \right\}$$

定义 敏感系数

$$\frac{E(Y(1) | Z=1, X) = \mu_1(X)}{E(Y(1) | Z=0, X)} = E_1(X)$$

$$\frac{E(Y(0) | Z=1, X)}{E(Y(0) | Z=0, X) = \mu_0(X)} = E_0(X)$$

\Rightarrow 推导出一个结果

比 $Z \perp (Y(1), Y(0)) | X$ 下结果更好

当 $E_1(X) = E_0(X) = 1$ 时, 所有结果

和 Part IV 一样

定理 18.1

图 18.1 $E(x)$, $e(x)$.

$$\tau = E \left\{ Z \mu_1(x) + (1-Z) \mu_0(x) / e(x) \right\} \\ - E \left\{ Z \mu_0(x) e(x) + (1-Z) \mu_0(x) \right\}$$

$$\text{其中 } \mu_1(x) = E(Y | Z=1, x) \\ \mu_0(x) = E(Y | Z=0, x)$$

$$= E \left(w_1(x) \frac{ZY}{e(x)} \right) - E \left(w_0(x) \frac{(1-Z)Y}{1-e(x)} \right)$$

$$\text{其中 } w_1(x) = e(x) + \frac{1-e(x)}{e_1(x)}$$

$$w_0(x) = e(x) e_0(x) + 1 - e(x)$$

= 双稳态定理的推广也见
R. Luk Dmg (2023)

$\frac{1}{L} dr$ 公式

$$\frac{1}{L} ht = \frac{1}{n} \sum_{i=1}^n \frac{\left(\hat{e}(x_i) \varepsilon_i(x_i) + 1 - \hat{e}(x_i) \right) z_i y_i}{\varepsilon_i(x_i) \hat{e}(x_i)}$$

$$- \frac{1}{n} \sum_{i=1}^n \frac{\left(\hat{e}(x_i) \varepsilon_i(x_i) + 1 - \hat{e}(x_i) \right) (1 - z_i) y_i}{1 - \hat{e}(x_i)}$$

$$\frac{1}{L} dr = \frac{1}{L} ht - \frac{1}{n} \sum_{i=1}^n \left(z_i - \hat{e}(x_i) \right) \left(\frac{\hat{\mu}_1(x_i)}{\hat{e}(x_i) \varepsilon_i(x_i)} + \frac{\hat{\mu}_0(x_i) \varepsilon_i(x_i)}{1 - \hat{e}(x_i)} \right)$$

关于 ATT : 平行验证

Chapter 19 Rosenbaum is 准确性 分析

观察性研究

1-1 匹配

精确匹配 (有白区!!)

① 无混杂 \Rightarrow MPE

\Rightarrow FRT

② 有混杂? \nRightarrow MPE

$$X_{i1} = X_{i2}$$

$$U_{i1} \neq U_{i2}$$

$$e_{ij} = \Pr(Z_{ij} = 1 \mid X_i, \underbrace{Y_{ij}(1), Y_{ij}(0)}_{\text{潜在结果本身}})$$

i 对 $j=1,2$ 共同 也是共同点

$$\pi_{i1} = \Pr(Z_{i1} = 1 \mid X_i, S_i, Z_{i1} + Z_{i2} = 1)$$

其中 $S_i = \{Y_{i1}(1), Y_{i1}(0), Y_{i2}(1), Y_{i2}(0)\}$

$$= \frac{\Pr(Z_{i1} = 1, Z_{i2} = 0 \mid X_i, S_i)}{\Pr(Z_{i1} + Z_{i2} = 1 \mid X_i, S_i)}$$

→ 分子有两项

$$= \frac{e_{i1} (1 - e_{i2})}{e_{i1} (1 - e_{i2}) + (1 - e_{i1}) e_{i2}}$$

$$\hat{O}_{ij} = \frac{O_{i1}}{O_{i1} + O_{i2}}$$

$$O_{ij} = \frac{e_{ij}}{1 - e_{ij}}$$

当 $O_{i1} = O_{i2}$

$$\Rightarrow = \frac{1}{2}$$

假设 $\pi \leq \frac{O_{i1}}{O_{i2}} \leq \pi$

与 logistic
回归函数有关

$$\Rightarrow \frac{1}{1 + \pi} \leq \pi_{i1} \leq \frac{\pi}{1 + \pi}$$

同: $P(\pi) =$ 最差情况下的 p 值

$$\pi = \sum_{i=1}^n S_i Z_i$$

$$1\left(\frac{1}{c_i} > 0\right) \quad \text{或} \quad \left|\frac{1}{c_1}\right| \dots \left|\frac{1}{c_n}\right|$$

是正数

最后4个说: $S_i \sim \text{iid Bern}\left(\frac{T}{1+P}\right)$

