

# Causal Inference for Health Data

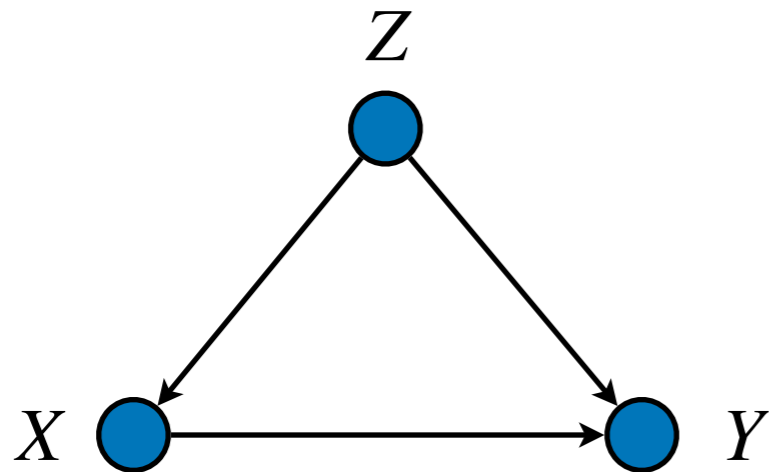
(STATS C160/C260 – Winter 2026)

Lecture 18:  
Unobserved Confounding – Part I

Drago Plečko

# How well do causal estimates from observational data translate to RCTs?

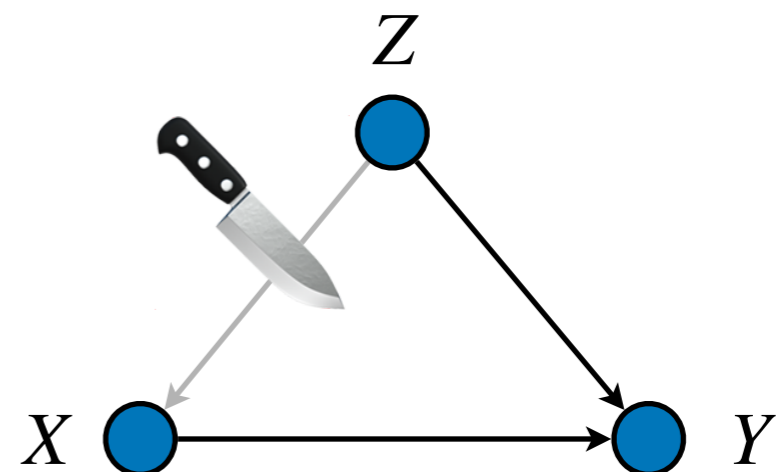
## Observational Setting



- $Z$  is back-door for  $(X, Y)$
- $P(y | do(x))$  identified through adjustment on  $Z$

$$\sum_z P(y | x, z)P(z)$$

## Interventional Setting (RCT)

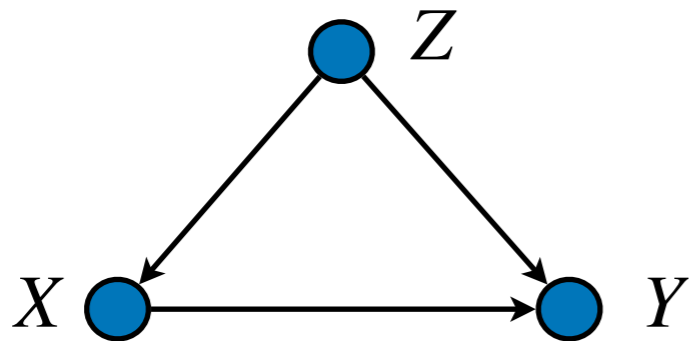


- We can obtain  $P(y | do(x))$  directly,
- However, we already know the answer, and the RCT is not really needed!

# What could go wrong?

- Most common culprit: **unobserved confounding**

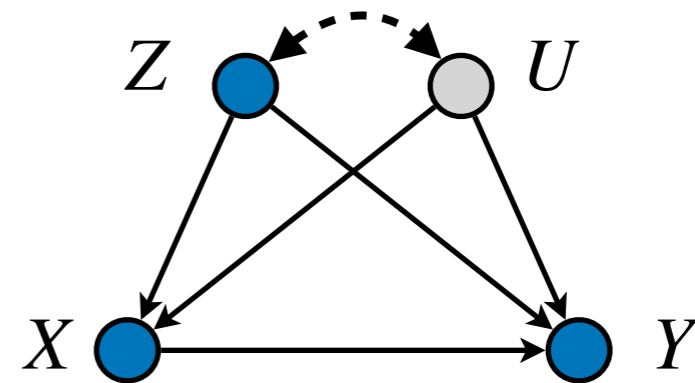
Our Analysis



- $Z$  is back-door for  $(X, Y)$
- $P(y | do(x))$  identified through adjustment on  $Z$

$$\sum_z P(y | x, z)P(z)$$

Reality



- $Z$  is **not** back-door for  $(X, Y)$
- $P(y | do(x))$  **not** identified through adjustment on  $Z$

$$P(y | do(x)) \neq \sum_z P(y | x, z)P(z)$$

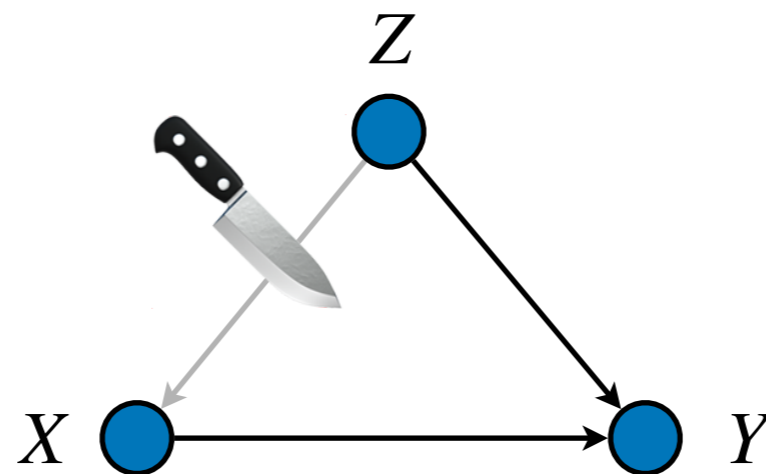
# What could go wrong?

- Most common culprit: **unobserved confounding**

Our Analysis

Reality

In this case, our findings will not translate to the RCT setting



# Does this happen in practice?

## I: Menopausal HRT

- Post-menopause, hormone replacement therapy (HRT) may have benefits, such as vasodilation, lower inflammatory activation, lower LDL cholesterol

> [N Engl J Med](#). 1991 Sep 12;325(11):756-62. doi: 10.1056/NEJM199109123251102.

### Postmenopausal estrogen therapy and cardiovascular disease. Ten-year follow-up from the nurses' health study

[M J Stampfer](#)<sup>1</sup>, [G A Colditz](#), [W C Willett](#), [J E Manson](#), [B Rosner](#), [F E Speizer](#), [C H Hennekens](#)

Affiliations + expand

PMID: 1870648 DOI: [10.1056/NEJM199109123251102](#)

[Free article](#)

*“Current estrogen use is associated with a reduction in the incidence of coronary heart disease as well as in mortality from cardiovascular disease”*

# Does this happen in practice?

## I: Menopausal HRT

- Post-menopause, hormone replacement therapy (HRT) may have benefits, such as vasodilation, lower inflammatory activation, lower LDL cholesterol



**HRT therapy seems like a good idea**

Clinical Trial > [JAMA](#). 2002 Jul 17;288(3):321-33. doi: 10.1001/jama.288.3.321.

### **Risks and benefits of estrogen plus progestin in healthy postmenopausal women: principal results From the Women's Health Initiative randomized controlled trial**

Jacques E Rossouw<sup>1</sup>, Garnet L Anderson, Ross L Prentice, Andrea Z LaCroix, Charles Kooperberg, Marcia L Stefanick, Rebecca D Jackson, Shirley A A Beresford, Barbara V Howard, Karen C Johnson, Jane Morley Kotchen, Judith Ockene; Writing Group for the Women's Health Initiative Investigators

Affiliations + expand

PMID: 12117397 DOI: [10.1001/jama.288.3.321](#)

# Does this happen in practice?

## I: Menopausal HRT

- Post-menopause, hormone replacement therapy (HRT) may have benefits, such as vasodilation, lower inflammatory activity, and improved bone density.

*“Overall health risks exceeded benefits from use of combined estrogen plus progestin for an average 5.2-year follow-up among healthy postmenopausal US women”*

*Issue: in observational data, patients on HRT were (i) healthier; (ii) had better healthcare utilization; (iii) had better socio-economic status.*

HRT t  
like a

1.  
n in  
results  
nized

# II: Vitamin E and coronary heart disease (CHD)

- Hypothesized benefits of vitamin E included reduced atherosclerosis, reduced inflammation, and antithrombotic effects.

> [N Engl J Med](#). 1993 May 20;328(20):1444-9. doi: 10.1056/NEJM199305203282003.

## Vitamin E consumption and the risk of coronary disease in women

[M J Stampfer](#)<sup>1</sup>, [C H Hennekens](#), [J E Manson](#), [G A Colditz](#), [B Rosner](#), [W C Willett](#)

Affiliations + expand

PMID: 8479463 DOI: [10.1056/NEJM199305203282003](#)

*“Among middle-aged women the use of vitamin E supplements is associated with a reduced risk of coronary heart disease”* RCT makes sense?

# II: Vitamin E and coronary heart disease (CHD)

- Hypothesized benefits of vitamin E included reduced atherosclerosis, reduced inflammation, and antithrombotic effects.

Clinical Trial > N Engl J Med. 2000 Jan 20;342(3):154-60.

doi: 10.1056/NEJM200001203420302.

## Vitamin E supplementation and cardiovascular events in high-risk patients

*“In patients at high risk for cardiovascular events, treatment with vitamin E for a mean of 4.5 years had no apparent effect on cardiovascular outcomes.”*

Supplement users exercised more, smoked less, and had healthier diets, and better SES.

# III: Beta-carotene -> Lung Cancer

- Hypothesized benefits of beta-carotene included antioxidant protection against DNA damage, reduced lipid peroxidation, and enhanced immune surveillance, thereby lowering the risk of lung cancer.

> [Br J Cancer](#). 1988 Apr;57(4):428-33. doi: 10.1038/bjc.1988.97.

## Serum beta-carotene and subsequent risk of cancer: results from the BUPA Study

[N J Wald](#)<sup>1</sup>, [S G Thompson](#), [J W Densem](#), [J Boreham](#), [A Bailey](#)

Affiliations + expand

PMID: 3390380 PMCID: [PMC2246576](#) DOI: [10.1038/bjc.1988.97](#)

*“Men in the top two quintiles of serum beta-carotene had only about 60% of the risk of developing cancer compared with men in the bottom quintile.”*

# III: Beta-carotene -> Lung Cancer

- Hypothesized benefits of beta-carotene included antioxidant protection against DNA damage, reduced lipid peroxidation, and enhanced immune surveillance, thereby lowering the risk of lung cancer.

Clinical Trial > N Engl J Med. 1994 Apr 14;330(15):1029-35.

doi: 10.1056/NEJM199404143301501.

**The effect of vitamin E and beta carotene on the incidence of lung cancer and other cancers in male smokers**

*“We found no reduction in the incidence of lung cancer among male smokers after five to eight years of dietary supplementation with alpha-tocopherol or beta carotene.”*

er:

# III: Beta-carotene -> Lung Cancer

- Hypothesized benefits of beta-carotene included antioxidant protection against DNA damage, reduced lipid peroxidation, and enhanced immune surveillance, thereby lowering the risk of lung cancer.

Clinical Trial > N Engl J Med. 1994 Apr 14;330(15):1029-35.

doi: 10.1056/NEJM199404143301501.

Supplement users and individuals with higher serum beta-carotene levels smoked less, consumed more fruits and vegetables, had overall healthier dietary patterns, and had higher SES

*cancer among male smokers after five to eight years of dietary supplementation with alpha-tocopherol or beta carotene.”*

# Unobserved Confounding

- There are two components of protecting our findings against unobserved confounding:

Knowing what we  
didn't measure

- In all studies, it is important to **report** factors that may be relevant, but have not been included in the analysis,
- These factors can be measured in prospective studies.

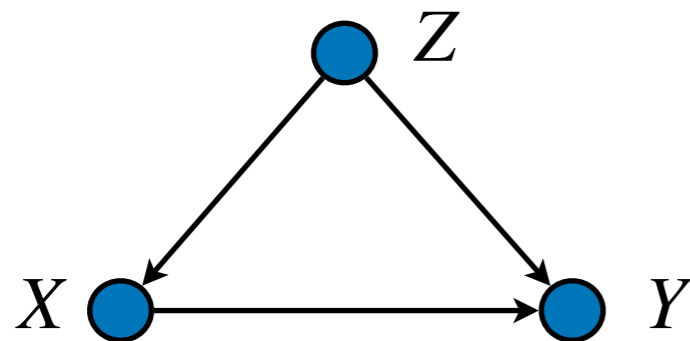
Knowing what we  
don't know

- This is a lot more difficult — what if there is an unobserved confounder previously not hypothesized (discussion: what is treatment allocation),
- Performing **sensitivity analysis**.



# Sensitivity Analysis for Unobserved Confounding

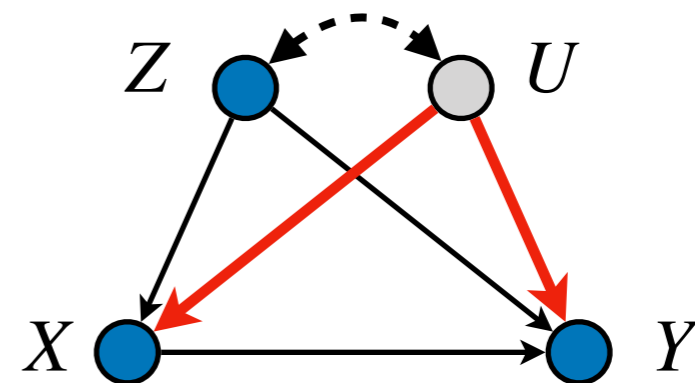
Run the Analysis



- $Z$  is back-door for  $(X, Y)$
- $P(y | do(x))$  identified and compute as

$$\sum_z P(y | x, z)P(z)$$

Key Question



- How strong would an unobserved confounder have to be to change our findings (red arrows)?
- How do we quantify this?

# Quantifying UC Strength

---

- In the general, non-parametric setting, quantifying the strength of UCs is hard,
- Finding nice ways to do this is an open research area,
- Some useful tools for parametric settings exist:

Binary Outcomes and  
Treatment

Continuous Outcomes with  
Linearity

# Risk Ratio Scale:

## Total Risk Ratio (TRR) vs. Total Effect (TE)

---

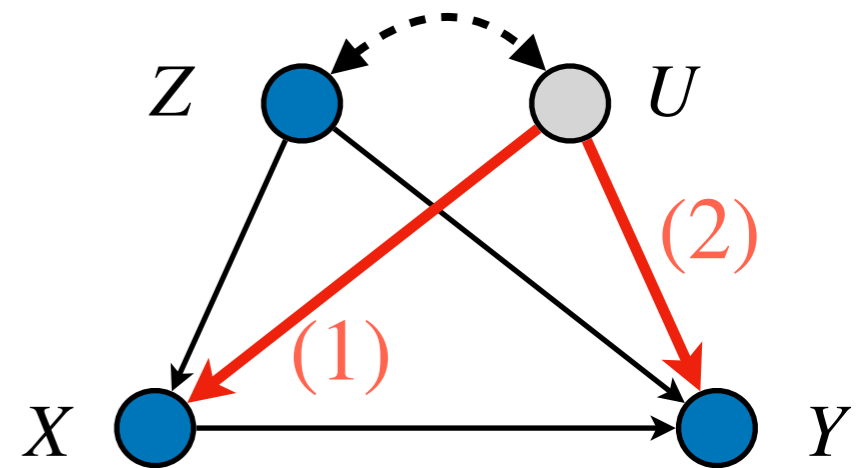
- So far, we considered the total effect (TE), given by  $E[Y | do(x_1)] - E[Y | do(x_0)]$ ,
- For binary outcomes, another quantification is useful, through the total risk ratio (TRR)

$$TRR_{x_0, x_1}(y) = \frac{P(Y = 1 | do(x_1))}{P(Y = 1 | do(x_0))}$$

Note:  $TE > 0 \iff TRR > 0$

# UC Sensitivity for TRR

- Let's assume the UC  $U$  is **binary**,
- Our analysis is  **$z$ -specific**,
- We start with the definitions of our **sensitivity parameters**:



$$\mathbf{RR}_{XU}^{(1)} = \max_u \frac{P(u \mid x_1, z)}{P(u \mid x_0, z)}$$

Captures the strength of  $X \rightarrow U$  relationship: maximal risk-ratio from flipping  $x_0 \rightarrow x_1$

$$\mathbf{RR}_{UY}^{(2)} = \max_x \frac{\max_u P(y \mid x, u, z)}{\min_u P(y \mid x, u, z)}$$

Captures the strength of  $U \rightarrow Y$  relationship: maximal  $u$ -variation induced in  $P(y \mid \cdot)$

# UC Sensitivity for TRR

**Theorem.** Let  $RR^{obs}(y | z)$  be the risk ratio  $\frac{P(y | x_1, z)}{P(y | x_0, z)}$ ,

and let  $RR^{true}(y | z)$  be the true risk ratio obtained after adjusting for  $U$ ,

$$\frac{\sum_u P(y | x_1, z, u)P(u)}{\sum_u P(y | x_0, z, u)P(u)}$$

Then, we have that

$$\frac{RR^{obs}(y | z)}{RR^{true}(y | z)} \leq \frac{RR_{XU}RR_{UY}}{RR_{XU} + RR_{UY} - 1}$$

# Proof: Part I

---

We are assuming that  $RR^{obs}(y | z) > 1$ , WLOG.

We first note that

homework

$$RR^{true}(y | z) = \lambda RR^{true}(y | z, x_1) + (1 - \lambda) RR^{true}(y | z, x_0), \lambda \in [0, 1]$$

meaning that the true RR is a convex combination of  $x$ -specific RRs.

Next, note that we have

whiteboard

$$\frac{RR^{obs}(y | z)}{RR^{true}(y | z, x_1)} = \frac{\sum_u P(y | x_0, z, u) P(u | z, x_1)}{\sum_u P(y | x_0, z, u) P(u | z, x_0)}$$

# Proof: Part II

Our next step is to re-write both the numerator and the denominator as a convex combination of min/max over  $P(y | x_0, z, u)$ :

$$\frac{\text{RR}^{obs}(y | z)}{\text{RR}^{true}(y | z, x_1)} = \frac{w_1 \max_u P(y | x_0, z, u) + (1 - w_1) \min_u P(y | x_0, z, u)}{w_0 \max_u P(y | x_0, z, u) + (1 - w_0) \min_u P(y | x_0, z, u)}.$$

with

$$w_x = \frac{\sum_u [P(y | x_0, z, u) - \min_{u'} P(y | x_0, z, u')] P(u | x, z)}{\max_u P(y | x_0, z, u) - \min_u P(y | x_0, z, u)}.$$

# Proof: Part III

Dividing through by  $\min_u P(y | x_0, z, u)$ , we get

$$\frac{\text{RR}^{obs}(y | z)}{\text{RR}^{true}(y | z, x_1)} = \frac{w_1 \left( \frac{\max_u P(y | x_0, z, u)}{\min_u P(y | x_0, z, u)} - 1 \right) + 1}{\left(\frac{w_1}{w_0}\right)^{-1} w_1 \left( \frac{\max_u P(y | x_0, z, u)}{\min_u P(y | x_0, z, u)} - 1 \right) + 1}$$

Note that we have

$\text{RR}_{UY|x_0}$

can be exchanged using  $\text{RR}_{XU}$

$$\begin{aligned} \Gamma &\triangleq \frac{w_1}{w_0} = \frac{\sum_u [P(y | x_0, z, u) - \min'_u P(y | x_0, z, u')] P(u | z, x_1)}{\sum_u [P(y | x_0, z, u) - \min'_u P(y | x_0, z, u')] P(u | z, x_0)} \\ &\leq \frac{\sum_u [P(y | x_0, z, u) - \min'_u P(y | x_0, z, u')] \text{RR}_{XU} P(u | z, x_0)}{\sum_u [P(y | x_0, z, u) - \min'_u P(y | x_0, z, u')] P(u | z, x_0)} = \text{RR}_{XU}. \end{aligned}$$

# Proof: Part IV

Putting together, we can bound the ratio as

$$\begin{aligned}
 \frac{\mathbf{RR}^{obs}(y | z)}{\mathbf{RR}^{true}(y | z, x_1)} &= \frac{w_1 \left( \frac{\max_u P(y | x_0, z, u)}{\min_u P(y | x_0, z, u)} - 1 \right) + 1}{\Gamma^{-1} w_1 \left( \frac{\max_u P(y | x_0, z, u)}{\min_u P(y | x_0, z, u)} - 1 \right) + 1} && f(x) = \frac{x + 1}{\Gamma^{-1}x + 1} \\
 &&& \text{increasing iff } \Gamma > 1 \\
 &\leq \frac{(\mathbf{RR}_{UY|x_0} - 1) + 1}{\Gamma^{-1}(\mathbf{RR}_{UY|x_0} - 1) + 1} \\
 &= \frac{\Gamma \mathbf{RR}_{UY|x_0}}{\mathbf{RR}_{UY|x_0} + \Gamma - 1} && f(x, y) = \frac{xy}{x + y + 1} \\
 &\leq \frac{\mathbf{RR}_{XU} \mathbf{RR}_{UY}}{\mathbf{RR}_{UY} + \mathbf{RR}_{XU} - 1} && \text{increasing in } x, y
 \end{aligned}$$

# Proof: Part V

---

Finally, we can combine the bounds treated and untreated groups:

$$\begin{aligned} \text{RR}^{true}(y | z) &= \lambda \text{RR}^{true}(y | z, x_1) + (1 - \lambda) \text{RR}^{true}(y | z, x_0), \lambda \in [0, 1] \\ &\leq \lambda \text{BRR}^{obs}(y | z) + (1 - \lambda) \text{BRR}^{obs}(y | z) \\ &= B \cdot \text{RR}^{obs}(y | z) \end{aligned}$$

where  $B \triangleq \frac{\text{RR}_{XU} \text{RR}_{UY}}{\text{RR}_{UY} + \text{RR}_{XU} - 1}$ , completing the proof.

# z-TRR to TRR

- Our obtained bound is for the *z-specific TRR*,
- For bounds on the population-level TRR, we have

$$\begin{aligned}
 \mathbf{RR}^{true}(y) &= \frac{\sum_z P(y_{x_1} | z)P(z)}{\sum_z P(y_{x_0} | z)P(z)} \\
 &= \frac{\sum_z \mathbf{RR}_{x_0,x_1}^{true}(y | z)P(y_{x_0} | z)P(z)}{\sum_z P(y_{x_0} | z)P(z)} \\
 &\geq \frac{\sum_z \mathbf{RR}_{x_0,x_1}^{obs}(y | z) \frac{\mathbf{RR}_{UY} + \mathbf{RR}_{XU} - 1}{\mathbf{RR}_{XU} \mathbf{RR}_{UY}} P(y_{x_0} | z)P(z)}{\sum_z P(y_{x_0} | z)P(z)} \\
 &\geq \frac{\mathbf{RR}_{UY} + \mathbf{RR}_{XU} - 1}{\mathbf{RR}_{XU} \mathbf{RR}_{UY}} \min_z \mathbf{RR}_{x_0,x_1}^{obs}(y | z).
 \end{aligned}$$

**Note: bound is based on the minimal z-specific TRR**

# E-value

---

- Now, in the very last step, we can further simplify the setting, and search for a minimum joint strength  $RR_{XU}$ ,  $RR_{UY}$  that explain away the observed effect,
- Here, it is convenient to assume that  $RR_{XU} = RR_{UY} = r$  so that we can say a confounder (i) for which  $x_0 \rightarrow x_1$  causes an  $r$ -fold increase in risk-ratio for  $u_1$  or  $u_0$ , and (ii) which causes an  $r$ -fold increase in risk-ratio for  $Y = 1$  is capable of explaining away the effect under study.
- When  $RR_{XU} = RR_{UY} = r$ , we can find the minimal  $r$  capable of explaining away the effect.

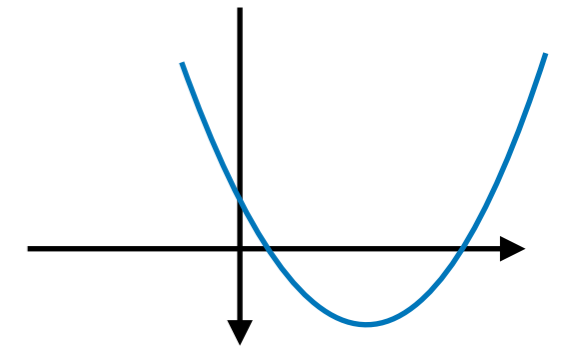
# E-value

- Note that we can write

$$\text{RR}^{true}(y | z) = \frac{\text{RR}_{UY} + \text{RR}_{XU} - 1}{\text{RR}_{XU}\text{RR}_{UY}} \text{RR}^{obs}(y | z)$$

$$= \frac{2r - 1}{r^2} \text{RR}^{obs}(y | z) > 1$$

$$\iff r^2 - (2r - 1)\text{RR}^{obs}(y | z) < 0$$



searching for  $r_{\min}$  so that this holds always  $\implies$  larger quadratic solution

$$r_{\min} = \text{RR}^{obs}(y | z) + \sqrt{\text{RR}^{obs}(y | z)(\text{RR}^{obs}(y | z) - 1)}$$

# E-value

- Note that we can write

$$RR^{true}(y | z) = \frac{RR_{UY} + RR_{XU} - 1}{RR} RR^{obs}(y | z)$$

This is known as the **E-value**:  
minimal joint strength of the confounder that can explain  
away the estimated risk ratio!

searching for  $r_{\min}$  so that this holds always  $\implies$  larger quadratic solution

$$r_{\min} = RR^{obs}(y | z) + \sqrt{RR^{obs}(y | z)(RR^{obs}(y | z) - 1)}$$

# Heuristics in Practice: Assume Away Heterogeneity

---

- In practice, often the  $\min_z \text{RR}^{obs}(y | z)$  is not taken for obtaining bounds on  $\text{RR}^{true}(y)$ ,
- Instead, practitioners are willing to assume away heterogeneity (i.e., that the risk ratio is constant across  $z$ ),
- Upon doing so, a heuristic to understand how large (or small) the obtained E-value is, proceeds as follows:  
leave-out each of the observed confounders one by one, and compute how large the  $\text{RR}_{XU}$ ,  $\text{RR}_{UY}$  would be for them!