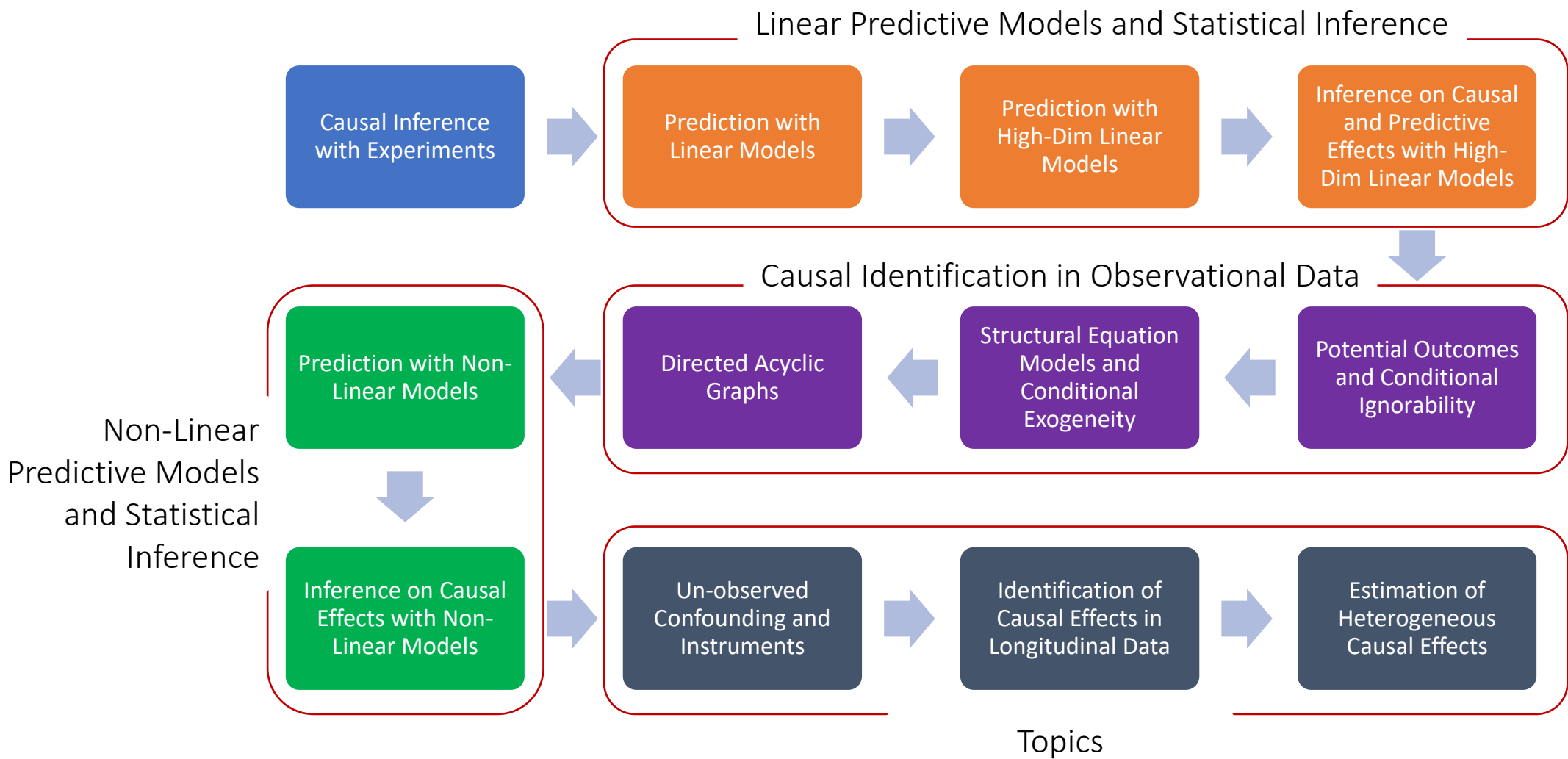
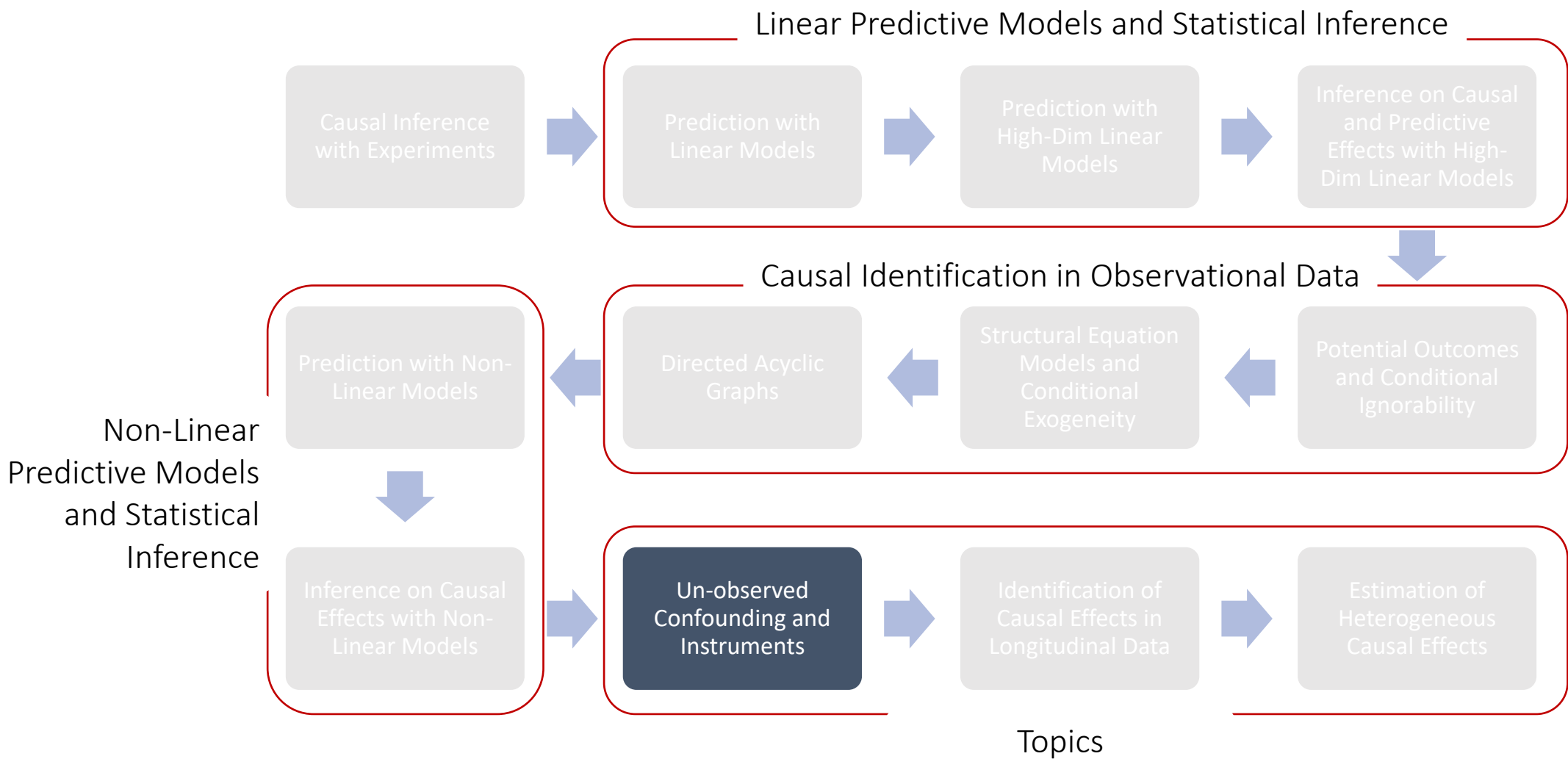


MS&E 228: Unobserved Confounding and Instruments

Vasilis Syrgkanis

MS&E, Stanford



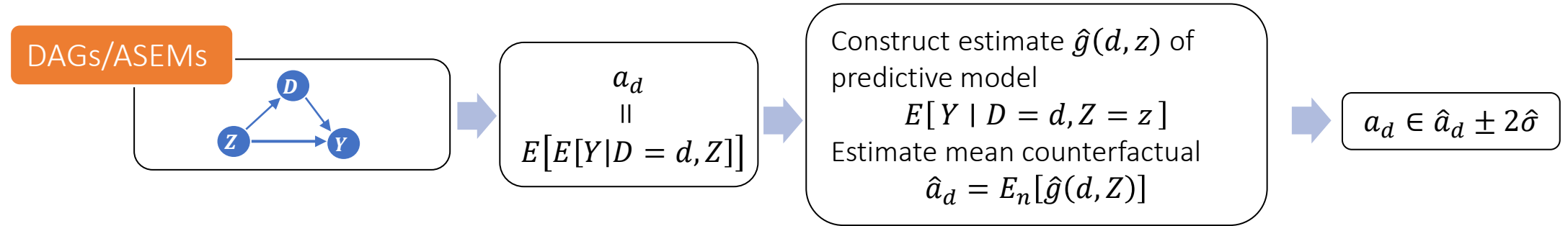


Goals for Today

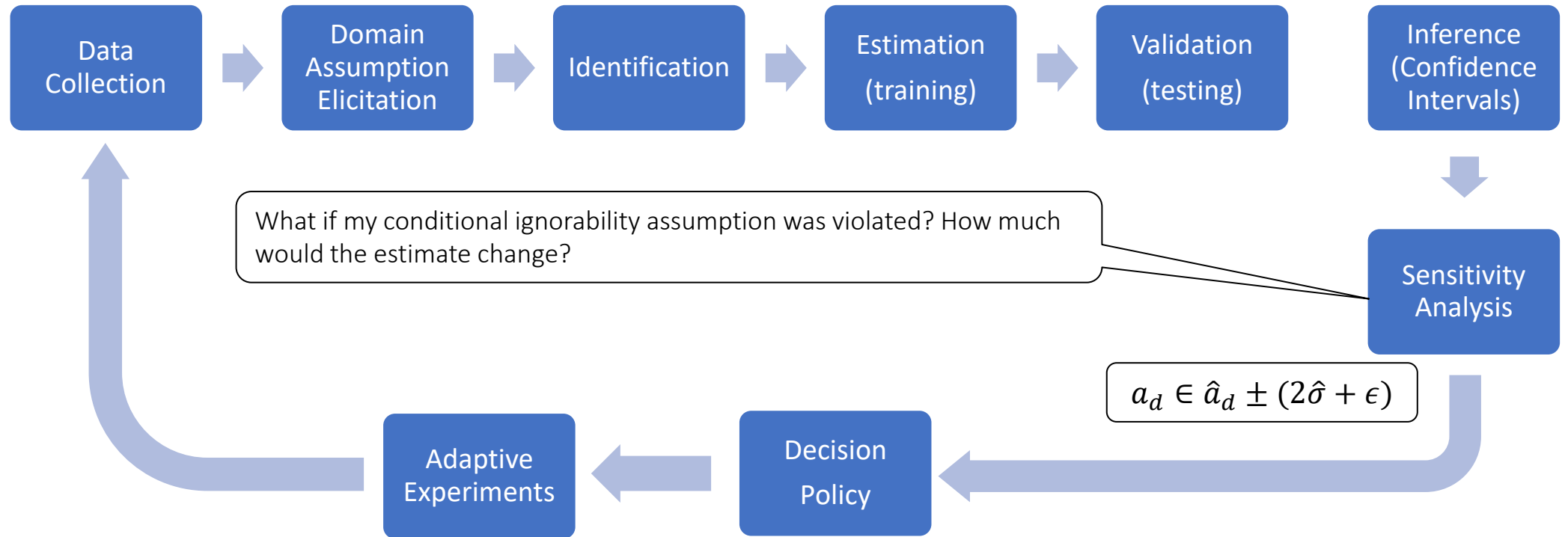
- What can we do when we have un-observed confounding
- Omitted variable bias bounds
- Introduction to “Instruments”

Causal Inference Pipeline

Theory



Practice



Bias Bounds

Reduction in unexplained variance of Y when adding A in the model that predicts Y from treatment and controls

- The analyst provides bounds on the partial R^2

$$R_{Y \sim A|D,X}^2 \leq C_Y^2,$$

$$R_{D \sim A|X}^2 \leq C_D^2$$

Reduction in unexplained variance of D when adding A in the model that predicts D from controls

- Based on these bounds we can conclude that

$$\theta_0 \in \theta_s \pm \sqrt{C_Y^2 \frac{C_D^2}{1 - C_D} \frac{E[(\tilde{Y} - \theta_s \tilde{D})^2]}{E[\tilde{D}^2]}}$$

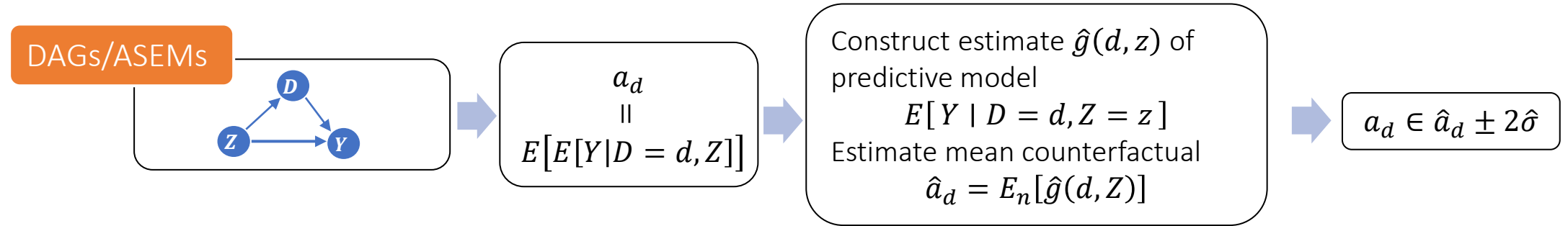
Measurable from the data

For more details:
[Making Sense of Sensitivity:
Extending Omitted Variable Bias](#)

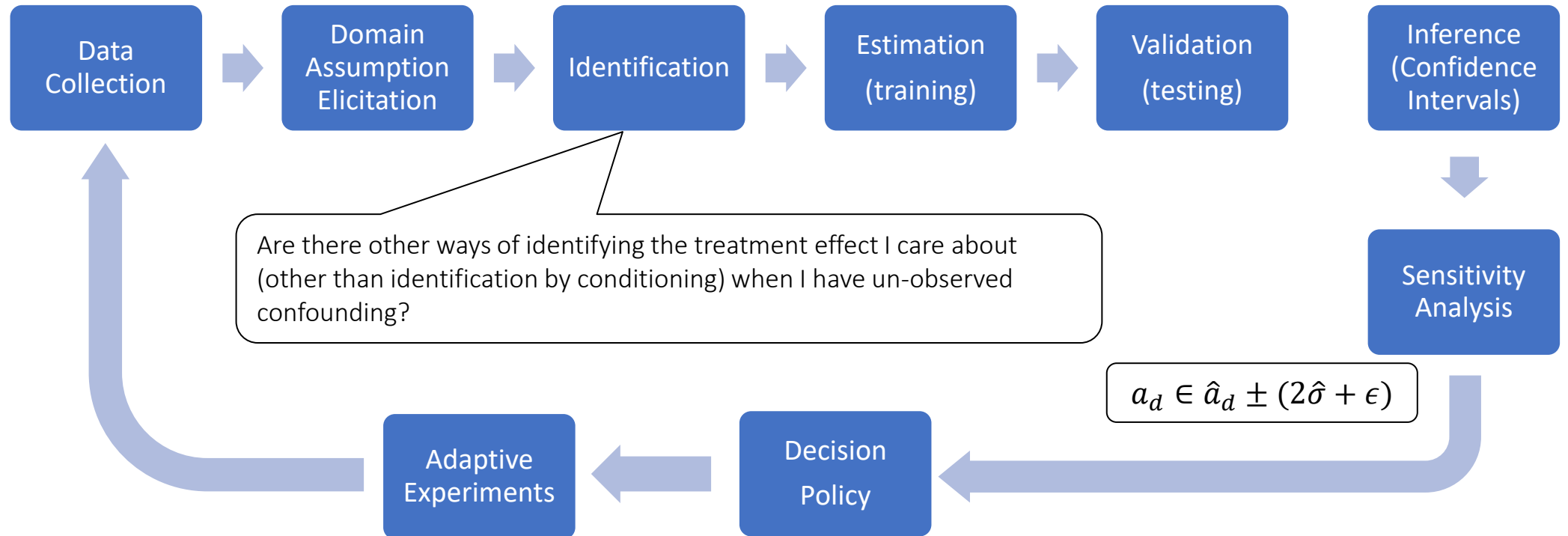
For more general analysis see:
[Long Story Short: Omitted Variable
Bias in Causal Machine Learning](#)

Causal Inference Pipeline

Theory

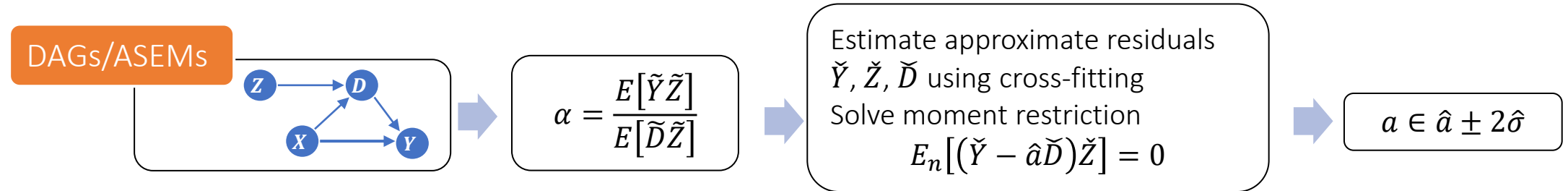


Practice

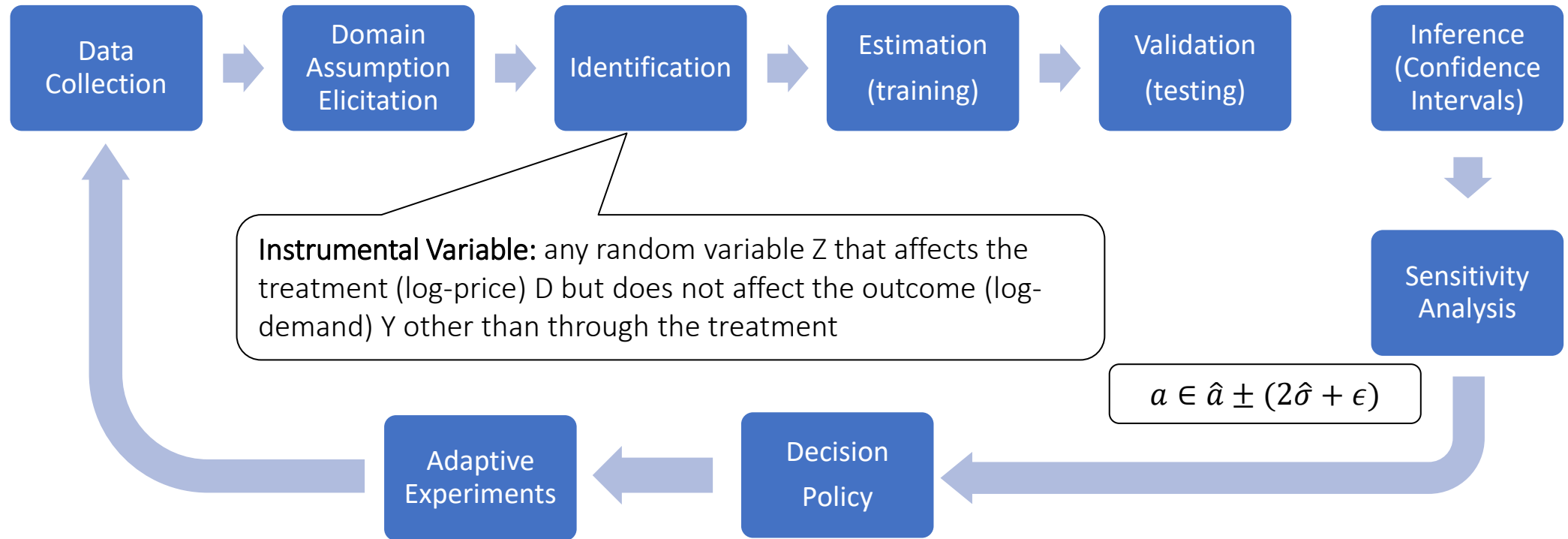


Causal Inference Pipeline

Theory

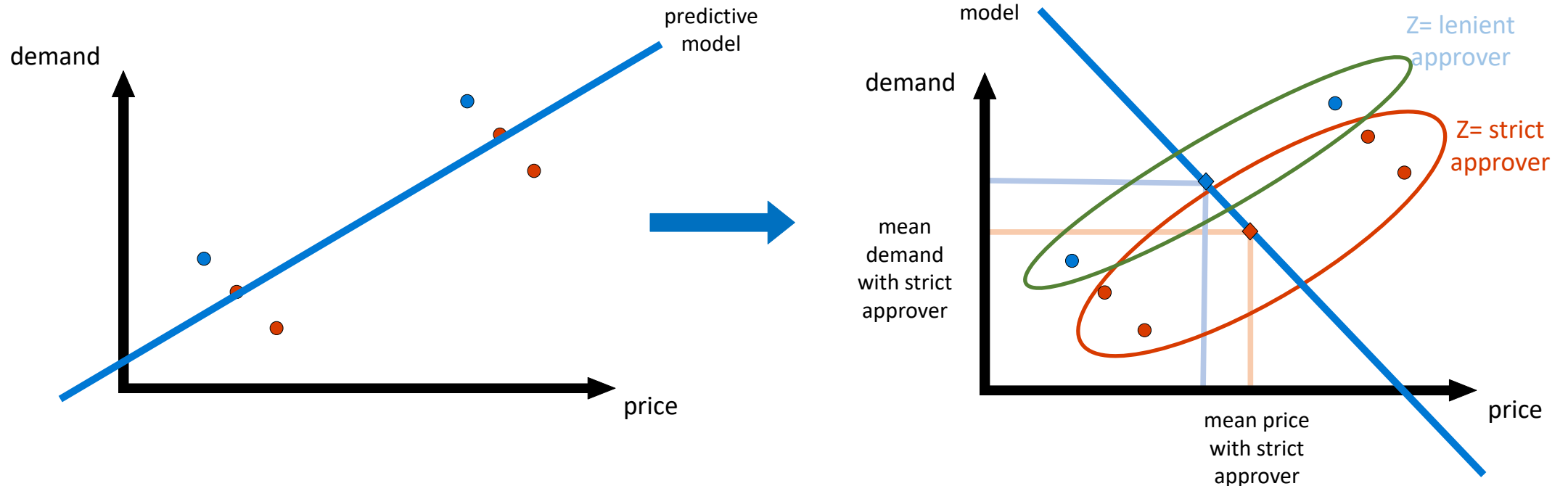
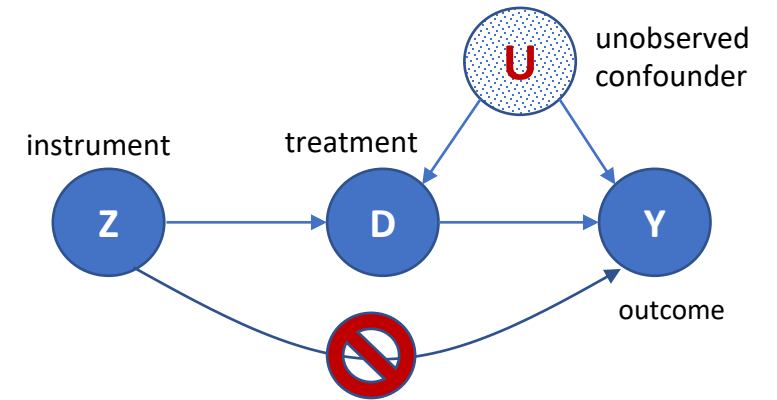


Practice



Instrumental Variables and 2SLS

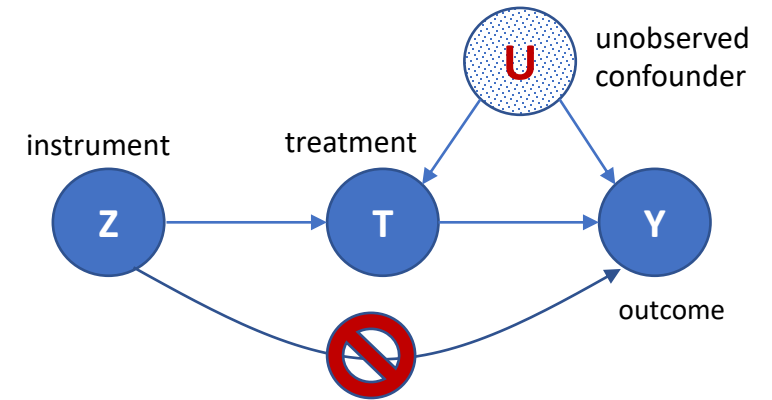
Instrumental Variable: any random variable Z that affects the treatment (log-price) D but does not affect the outcome (log-demand) Y other than through the treatment [Wright'28, Bowden-Turkington'90, Angrist-Krueger'91, Imbens-Angrist'94]



Instrumental Variables and 2SLS

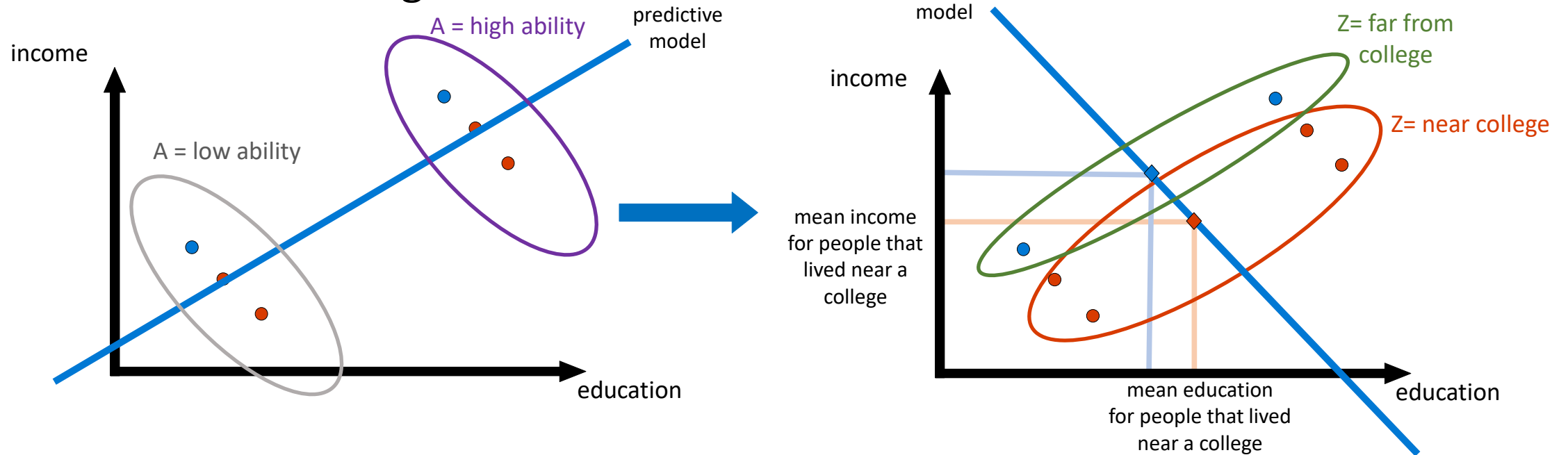
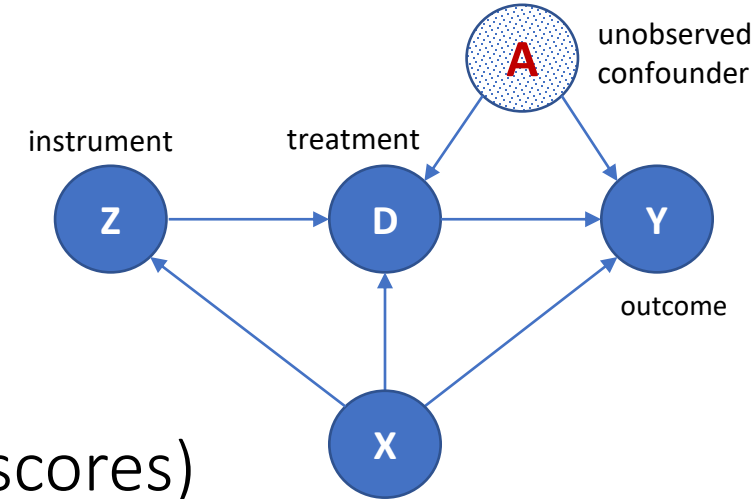
Instruments are widely used

- In the discount example (see also [Kling AER06] for effects of incarceration)
 - Discounts are sent to an approver desk
 - Approver assignment is random and different approvers are more or less “lenient”
 - Approver leniency is an instrument
- In healthcare [Doyle et al., JPE15]
 - Random assignment to ambulance companies of nearby patients is an instrument for measuring hospital quality
- In Tech [S., NeurIPS19]
 - Recommendation A/B tests as instruments for the effects of downstream actions



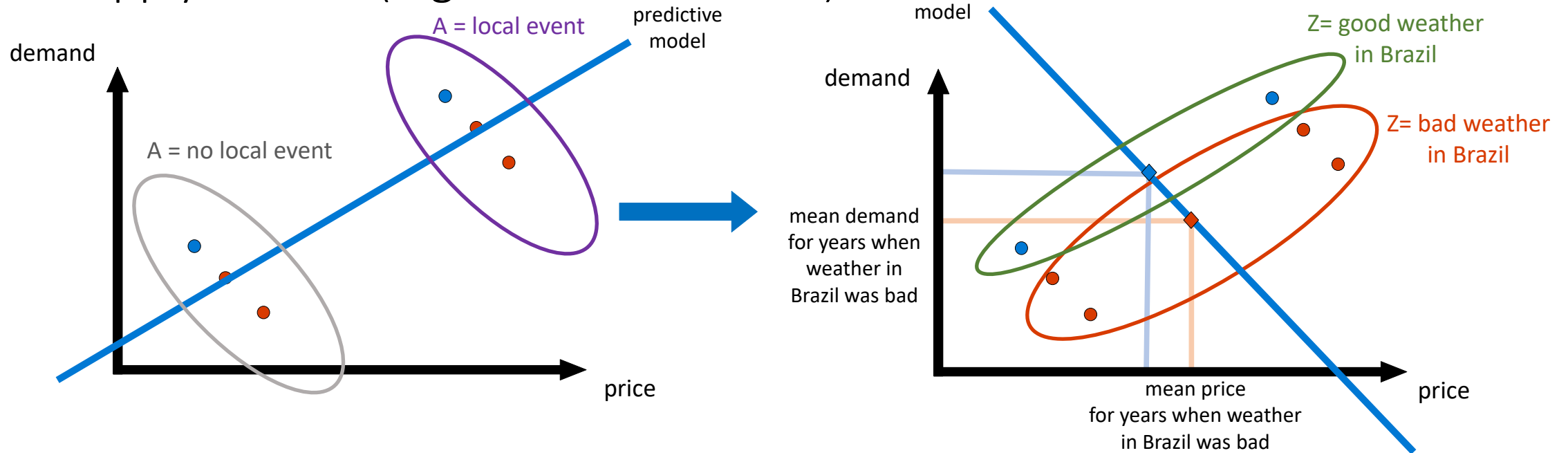
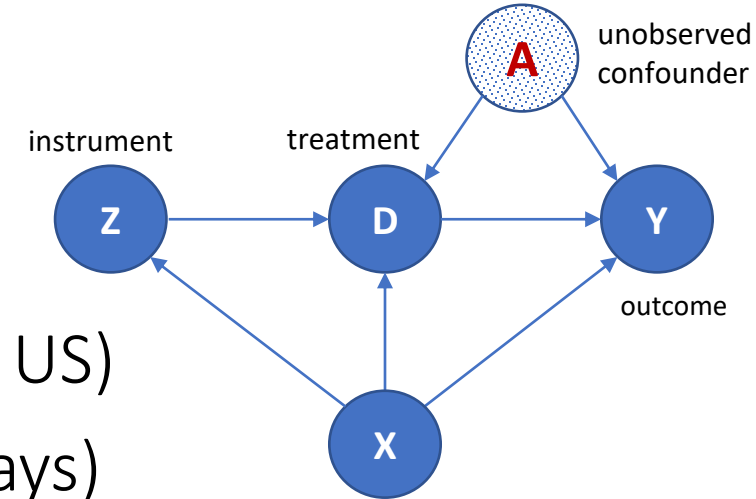
Returns to Education

- D: years of college, Y: income
- X: observable characteristics of a student (e.g. test scores)
- A: unobserved “ability”
- Z: distance to college



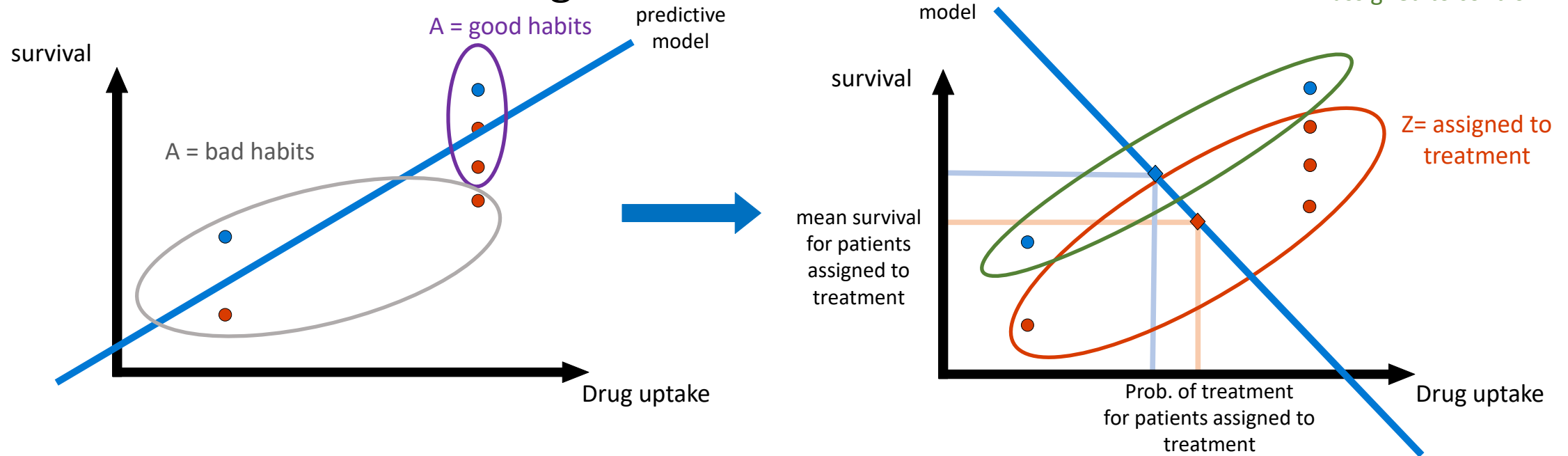
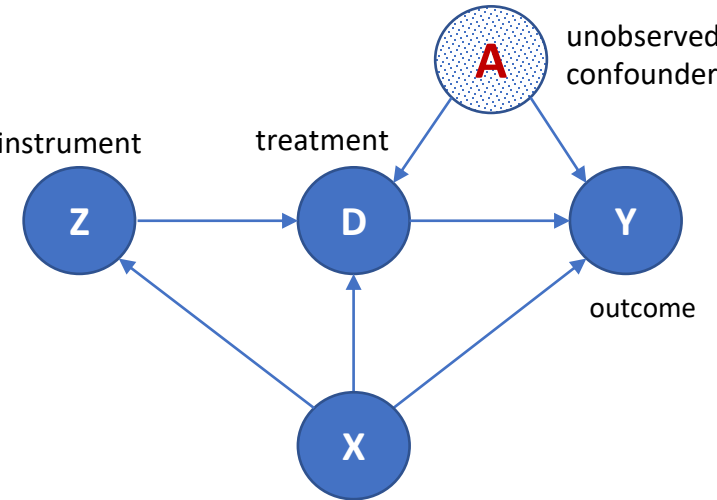
Demand Estimation

- D: price (e.g. of coffee), Y: demand (e.g. of coffee in US)
- X: observable characteristics of a market (e.g. holidays)
- A: unobserved “demand shocks” (e.g. local event)
- Z: supply shifters (e.g. weather in brazil)



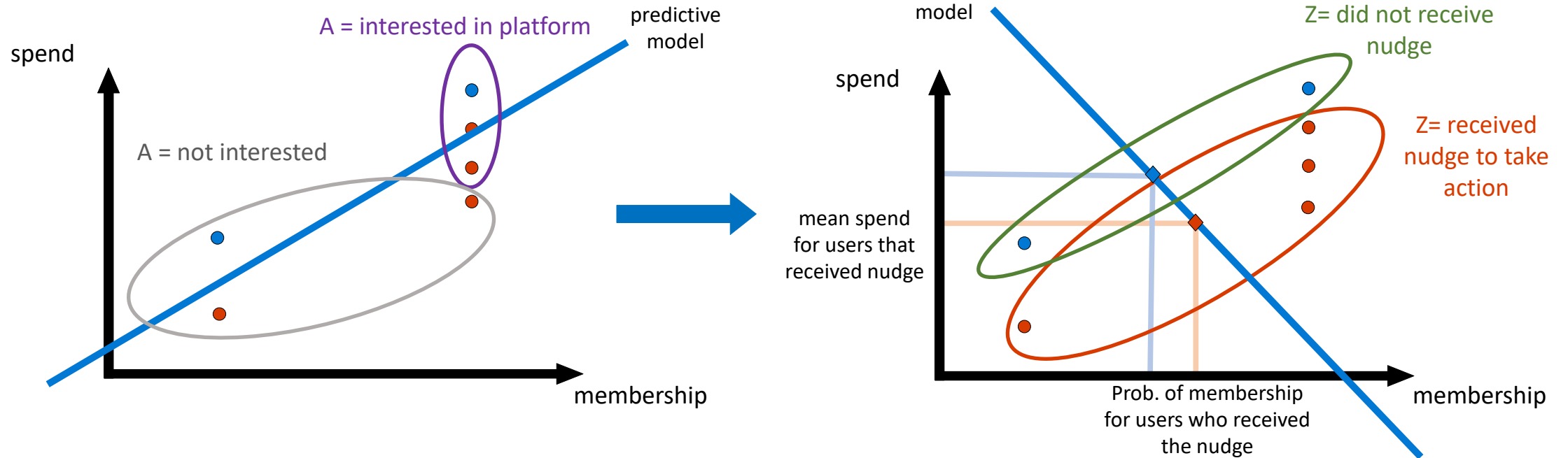
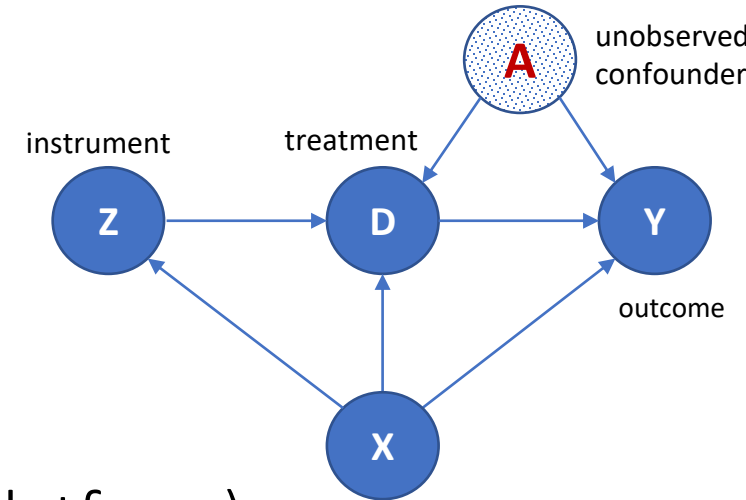
Clinical Trials with Non-Compliance

- D: drug treatment, Y: survival
- X: observable characteristics of a patient
- A: unobserved “compliance factors” (e.g. health habits)
- Z: randomized cohort assignment



Digital Recommendation A/B tests

- D: action taken by user (e.g. membership), Y: spend
- X: observable characteristics of a user
- A: unobserved confounding factors (e.g. interest in platform)
- Z: randomized nudge to take action (e.g. one-click sign-up pop-up)



Identification of Causal Effects via Instruments

Phillip Wright's idea (1928): the first causal path diagram analysis

- ◆ We can estimate effect of Z on y via a regression

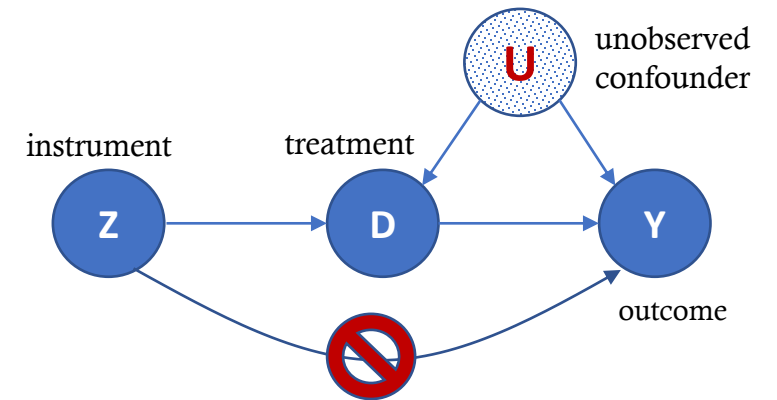
$$\gamma = \frac{\mathbb{E}[\tilde{Z}\tilde{y}]}{\mathbb{E}[\tilde{Z}^2]}$$

- ◆ We can estimate the effect of Z on D via a regression

$$\delta = \frac{\mathbb{E}[\tilde{Z}\tilde{D}]}{\mathbb{E}[\tilde{Z}^2]}$$

- ◆ The effect of Z on Y (γ) is the product of the effect of Z on T (δ) multiplied by the effect of T on y (θ)

$$\theta = \frac{\gamma}{\delta} = \frac{\mathbb{E}[\tilde{Z}\tilde{y}]}{\mathbb{E}[\tilde{Z}\tilde{D}]}$$



Partially Linear Instrumental Variable Model

- Typically for continuous treatment/instrument a partially linear structural equation assumed

$$\begin{aligned}Y &:= \theta_0 D + f_Y(X) + \delta A + \epsilon_Y \\D &:= \beta Z + f_D(X) + \gamma A + \epsilon_D \\Z &:= f_Z(X) + \epsilon_Z \\A &:= f_A(X) + \epsilon_A\end{aligned}$$

All errors are exogenous and un-correlated

Partially Linear Instrumental Variable Model

- After partialling out the observed controls X

$$\begin{aligned}\tilde{Y} &:= \theta_0 \tilde{D} + \delta \tilde{A} + \epsilon_Y \\ \tilde{D} &:= \beta \tilde{Z} + \gamma \tilde{A} + \epsilon_D \\ \tilde{Z} &:= \epsilon_Z \\ \tilde{A} &:= \epsilon_A\end{aligned}$$

- We see immediately that:

$$\tilde{Y} := \theta_0 \tilde{D} + U, \quad U := \delta \tilde{A} + \epsilon_Y \perp \tilde{Z}$$

- Since $\epsilon_A, \epsilon_Y, \epsilon_Z$ are un-correlated: $E[(\delta \tilde{A} + \epsilon_Y)\tilde{Z}] = 0$
- Thus we have the moment restriction: $E[(\tilde{Y} - \theta_0 \tilde{D})\tilde{Z}] = 0$

Partially Linear Instrumental Variable Model

- After partialling out the observed controls X

$$\begin{aligned}\tilde{Y} &:= \theta_0 \tilde{D} + \delta \tilde{A} + \epsilon_Y \\ \tilde{D} &:= \beta \tilde{Z} + \gamma \tilde{A} + \epsilon_D \\ \tilde{Z} &:= \epsilon_Z \\ \tilde{A} &:= \epsilon_A\end{aligned}$$

- Thus we have the moment restriction: $E[(\tilde{Y} - \theta_0 \tilde{D})\tilde{Z}] = 0$
- We re-derive a generalization of Wright's formula

$$\theta_0 = \frac{E[\tilde{Y}\tilde{Z}]}{E[\tilde{D}\tilde{Z}]}$$

Partially Linear Instrumental Variable Model

- After partialling out the observed controls X

$$\begin{aligned}\tilde{Y} &:= \theta_0 \tilde{D} + \tilde{A} + \epsilon_Y \\ \tilde{D} &:= \beta \tilde{Z} + \gamma \tilde{A} + \epsilon_D \\ \tilde{Z} &:= \epsilon_Z \\ \tilde{A} &:= \epsilon_A\end{aligned}$$

- Setting falls into the general moment estimation framework

$$M(\theta, h, p, m) = E \left[\left(Y - h(X) - \theta (D - p(X)) \right) (Z - m(X)) \right] = 0$$

- Where $h(X) = E[Y|X]$, $p(X) = E[D|X]$, $m(Z) = E[Z|X]$

Orthogonal Method: Double ML for IV

Double ML. Split samples in half

- Regress $Y \sim X$ with ML on first half, to get estimate $\hat{h}(S)$ of $E[Y|X]$
- Regress $D \sim X$ with ML on first half, to get estimate $\hat{p}(S)$ of $E[D|X]$
- Regress $Z \sim X$ with ML on first half, to get estimate $\hat{m}(S)$ of $E[Z|X]$
- Construct residuals on other half, $\hat{Z} = Z - \hat{m}(X)$, $\hat{D} := D - \hat{p}(X)$ and $\hat{Y} := Y - \hat{h}(X)$
- Solve moment condition:

$$E_n[(\hat{Y} - \theta \hat{D})\hat{Z}] = 0$$

```
from econml.iv.dml import OrthoIV
orthoiv = OrthoIV()
orthoiv.fit(y, D, Z, W=X).effect_inference()
```

Limits of Identification via Instruments

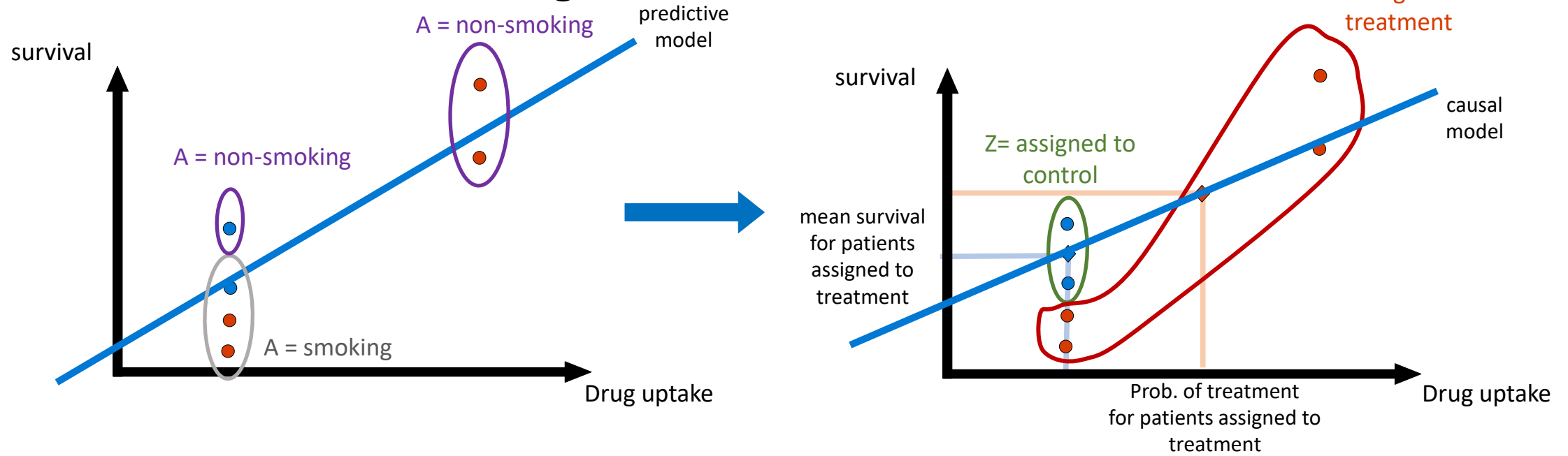
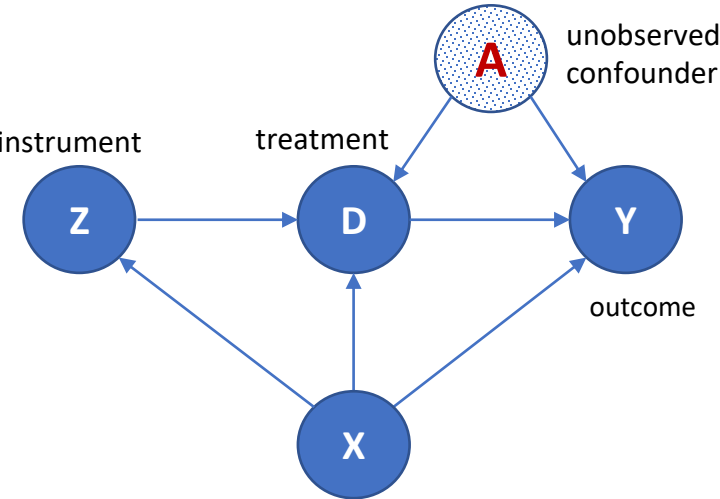
- ATE identification via Instruments not based solely on DAG restrictions
- Requires further restrictions on structural equation models

Example

- Binary treatment D (drug) and binary instrument Z (drug recommendation)
- Consider an unobserved confounder $A = \text{"smoking"}$
- Suppose that smokers ($A=1$) never take the drug (never comply) and non-smokers ($A=0$) always follow the recommendation (comply)
- Suppose that drug has positive effects for non-smokers but has severe side-effects for smokers

Clinical Trials with Non-Compliance

- D: drug treatment, Y: survival
- X: observable characteristics of a patient
- A: unobserved “compliance factors” (e.g. health habits)
- Z: randomized cohort assignment



Limits of Identification via Instruments

- ATE identification via Instruments not based solely on DAG restrictions
- Requires further restrictions on structural equation models

Example

- IV regression will never be able to uncover the side effects of drug treatment on smokers
- Nothing in the data is informative of that
- Effect will be biased as compared to average effect in whole population

What do we need for ATE

- Either the compliance behavior (effect of instrument on treatment) does not vary with A (or X)
- Or the treatment effect (effect of treatment on outcome) does not vary with A (or X)

$$\begin{aligned} Y &:= g_Y(\epsilon_Y) D + f_Y(X, A, \epsilon_Y) \\ D &:= f_D(Z, X, A, \epsilon_D) \\ Z &= f_Z(X, \epsilon_Z) \\ A &:= f_A(X, \epsilon_A) \end{aligned}$$

$$\begin{aligned} Y &:= g_Y(X, A, \epsilon_Y) D + f(X, A, \epsilon_Y) \\ D &:= g_D(\epsilon_D) Z + f_D(X, A, \epsilon_D) \\ Z &= f_Z(X) + \epsilon_Z \\ A &:= f_A(X, \epsilon_A) \end{aligned}$$

Joint Variation on Observables

- If joint variation is captured through observables then ATE is feasible

$$\begin{aligned} Y &:= g_Y(X, \epsilon_Y) D + f_Y(X, A, \epsilon_Y) \\ D &:= f_D(Z, X, A, \epsilon_D) \\ Z &= f_Z(X, \epsilon_Z) \\ A &:= f_A(X, \epsilon_A) \end{aligned}$$

$$\begin{aligned} Y &:= g_Y(X, A, \epsilon_Y) D + f_Y(X, A, \epsilon_Y) \\ D &:= g_D(X, \epsilon_D) Z + f_D(X, A, \epsilon_D) \\ Z &= f_Z(X, \epsilon_Z) \\ A &:= f_A(X, \epsilon_A) \end{aligned}$$

- We just need to do our identification analysis conditional on X and then average

$$\beta(X) = \frac{E[\tilde{Y}\tilde{Z} \mid X]}{E[\tilde{D}\tilde{Z} \mid X]}, \quad a = E[\beta(X)]$$

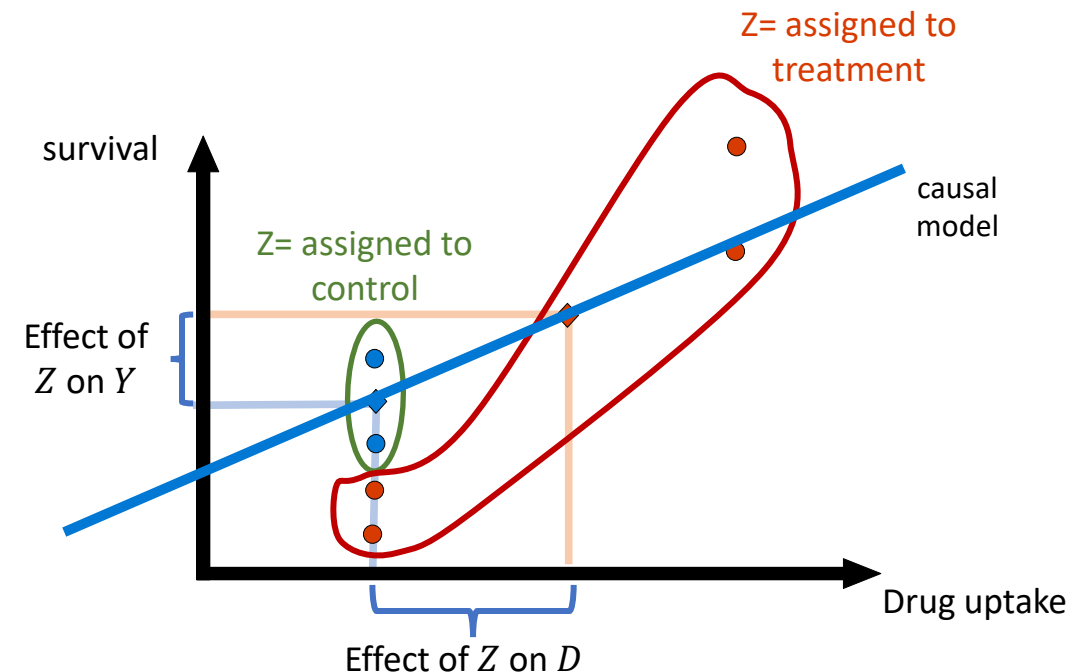
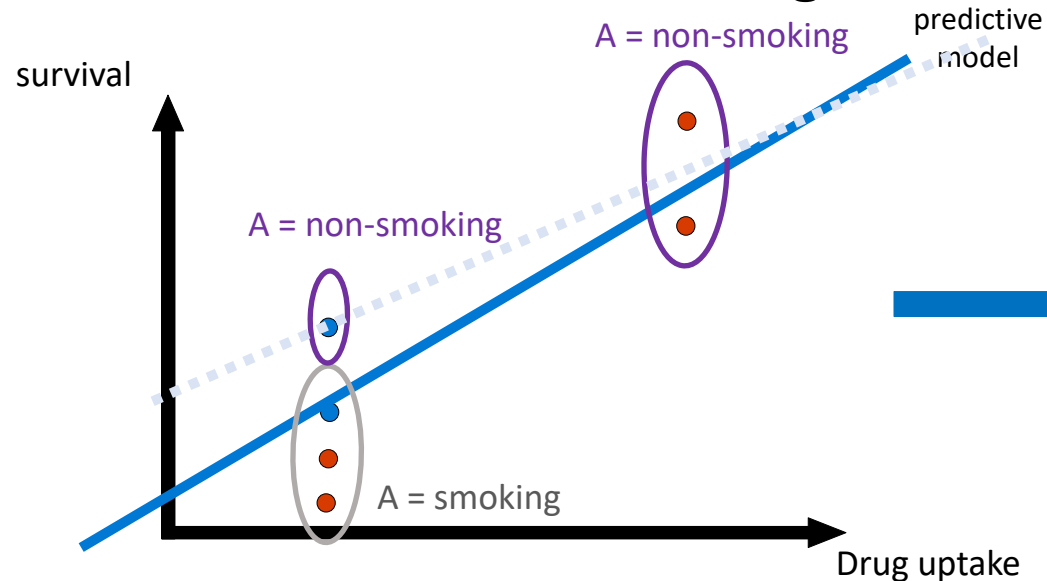
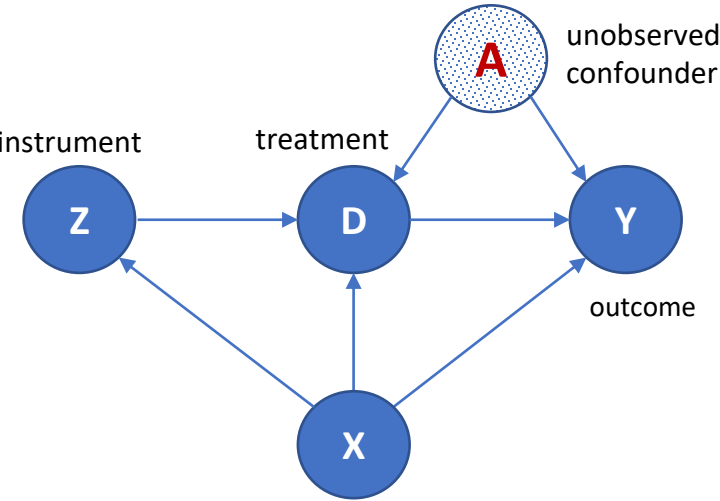
- Roughly: reweighting data based on compliance level $E[\tilde{D}\tilde{Z} \mid X]$

What if joint variation happens through unobservables?

$$\begin{aligned} Y &:= f_Y(D, X, A, \epsilon_Y) \\ D &:= f_D(Z, X, A, \epsilon_D) \\ Z &= f_Z(X, \epsilon_Z) \\ A &:= f_A(X, \epsilon_A) \end{aligned}$$

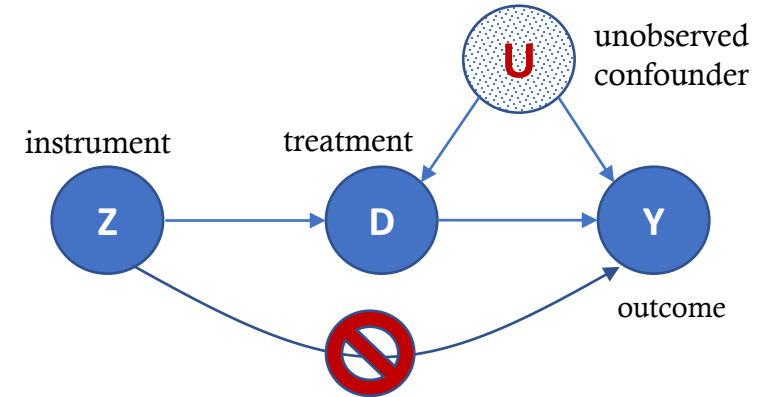
Clinical Trials with Non-Compliance

- D: drug treatment, Y: survival
- X: observable characteristics of a patient
- A: unobserved “compliance factors” (e.g. health habits)
- Z: randomized cohort assignment



Does the IV estimate coincide with the average effect for some sub-population?

The Binary Case



Imbens-Angrist (1994): core contribution of Nobel 2022 award

- Instrument/Treatment are binary (instrument=recommended treatment)
- Assume monotonicity: $D^{(1)} \geq D^{(0)}$
- Recommended treatment cannot reverse taken treatment
- Object of interest: Local Average Treatment Effect (ATE among compliers)

$$\theta_0 = E[Y^{(1)} - Y^{(0)} | D^{(1)} > D^{(0)}]$$

- Proof [Angrist-Imbens'94]:

$$\theta_0 = \frac{E[(Y^{(1)} - Y^{(0)})1\{D^{(1)} > D^{(0)}\}]}{E[1\{D^{(1)} > D^{(0)}\}]} = \frac{E[Y^{(D^{(1)})} - Y^{(D^{(0)})}]}{E[D^{(1)} - D^{(0)}]} = \frac{ATE(Z \rightarrow Y)}{ATE(Z \rightarrow D)}$$

$\gamma = \frac{E[\tilde{Z}\tilde{y}]}{E[\tilde{Z}^2]}$

$\delta = \frac{E[\tilde{Z}\tilde{D}]}{E[\tilde{Z}^2]}$

The Binary Case

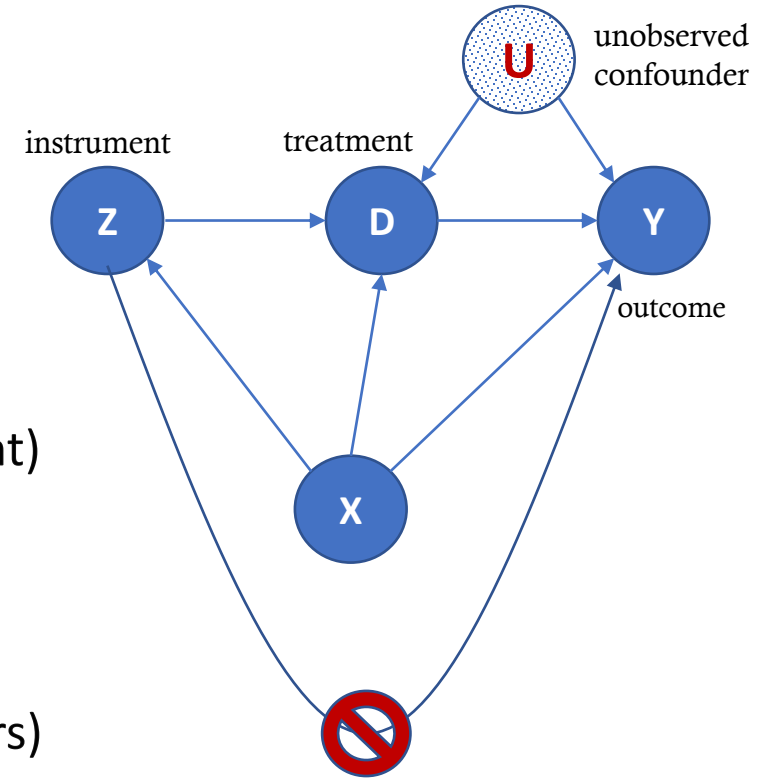
Imbens-Angrist (1994): core contribution of Nobel 2022 award

- Instrument/Treatment are binary (instrument=treatment)
- Assume monotonicity: $D^{(1)} \geq D^{(0)}$
- Recommended treatment cannot reverse taken treatment
- Object of interest: Local Average Treatment Effect (ATE among compliers)

$$\theta_0 = E[Y^{(1)} - Y^{(0)} | D^{(1)} > D^{(0)}]$$

- Proof [Angrist-Imbens'94]:

$$\theta_0 = \frac{E[(Y^{(1)} - Y^{(0)})1\{D^{(1)} > D^{(0)}\}]}{E[1\{D^{(1)} > D^{(0)}\}]} = \frac{E[Y^{(D(1))} - Y^{(D(0))}]}{E[D^{(1)} - D^{(0)}]} = \frac{ATE(Z \rightarrow Y)}{ATE(Z \rightarrow D)}$$



$$E[E[Y|Z = 1, X] - E[Y|Z = 0, X]]$$

ATE(Z → Y)

ATE(Z → D)

$$E[E[D|Z = 1, X] - E[D|Z = 0, X]]$$

LATE in the Binary Case

- Under monotonicity

$$\theta_0 = \frac{E[E[Y | Z = 1, X] - E[Y | Z = 0, X]]}{E[E[D | Z = 1, X] - E[D | Z = 0, X]]}$$

- Moment formulation

$$E[E[Y | Z = 1, X] - E[Y | Z = 0, X] - \theta_0(E[D | Z = 1, X] - E[D | Z = 0, X])] = 0$$

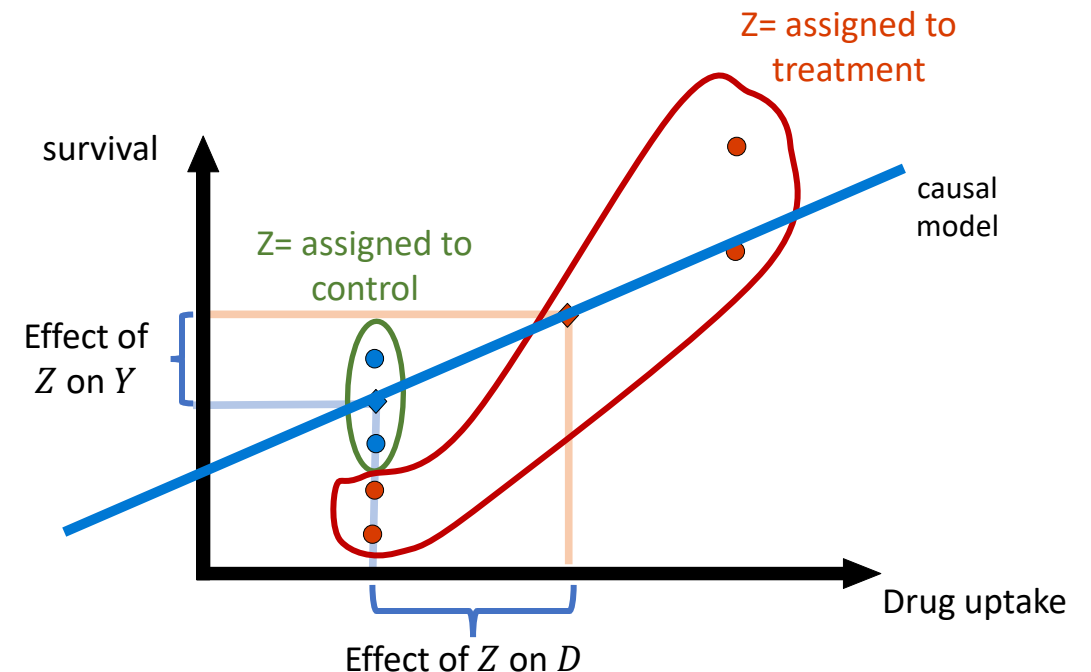
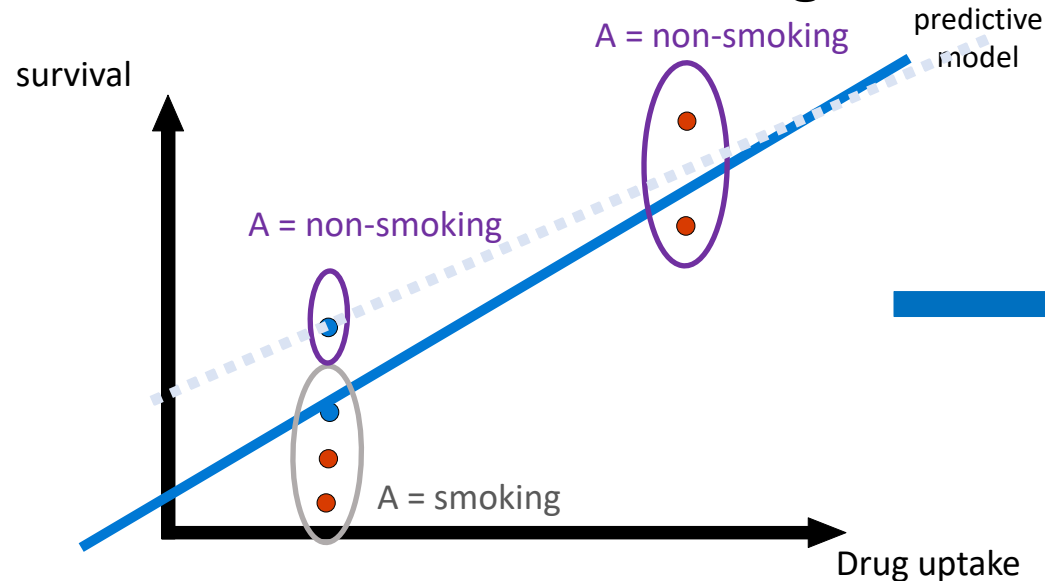
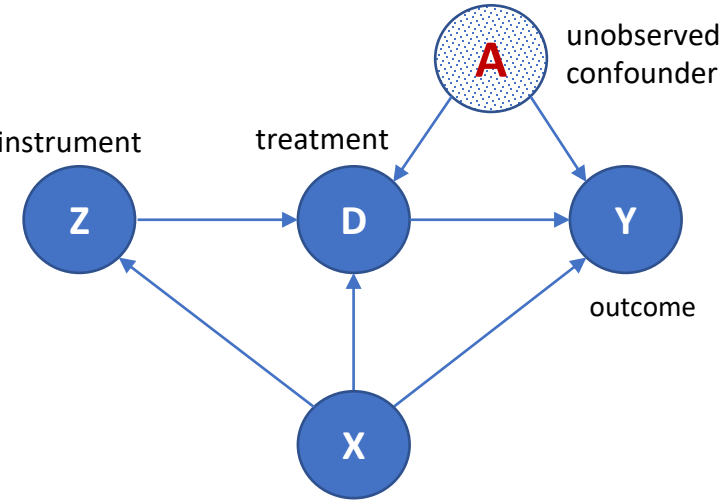
$$\begin{aligned} &+ \\ &\mathbf{H(Z, X)}(\mathbf{Y} - \mathbf{E[Y|Z, X]}) \qquad \qquad \qquad + \\ &\mathbf{H(Z, X)}(\mathbf{D} - \mathbf{E[D|Z, X]}) \end{aligned}$$

$$H(Z, X) = \frac{Z}{P(Z = 1|X)} - \frac{1 - Z}{1 - P(Z = 1|X)}$$

- Orthogonal moment formulation: apply ATE debiasing twice

Clinical Trials with Non-Compliance

- D: drug treatment, Y: survival
- X: observable characteristics of a patient
- A: unobserved “compliance factors” (e.g. health habits)
- Z: randomized cohort assignment



Confidence Intervals

Partially Linear Instrumental Variable Model

- Suppose that we can identify parameter of interest via moment

$$E[(\tilde{Y} - \theta_0 \tilde{D})\tilde{Z}] = 0$$

- Setting falls into the general moment estimation framework

$$M(\theta, h, p, m) = E \left[\left(Y - h(X) - \theta (D - p(X)) \right) (Z - m(X)) \right] = 0$$

Where $h(X) = E[Y|X]$, $p(X) = E[D|X]$, $m(Z) = E[Z|X]$

Inference with DML in PLIV Setting

- The estimate can be written as:

$$\hat{\theta} = \frac{E_n[\hat{Y}\hat{Z}]}{E_n[\hat{Z}\hat{D}]}$$

- If RMSE of propensity models and outcome model goes down at rate $n^{1/4}$, plus regularity conditions

$$\sqrt{n}(\hat{\theta} - \theta_0) = \sqrt{n}\left(\frac{E_n[\hat{Y}\hat{D}]}{E_n[\hat{Z}\hat{D}]} - \theta_0\right) \approx \sqrt{n}\left(\frac{E_n[\tilde{Y}\tilde{D}]}{E_n[\tilde{D}\tilde{Z}]} - \theta_0\right) = \sqrt{n}\left(\frac{E_n[\tilde{Y}\tilde{D}]}{E_n[\tilde{D}\tilde{Z}]} - \frac{E_n[\tilde{D}\tilde{Z}]}{E_n[\tilde{D}\tilde{Z}]} \theta_0\right) = \sqrt{n}\left(\frac{E_n[(\tilde{Y} - \theta_0\tilde{D})\tilde{Z}]}{E_n[\tilde{D}\tilde{Z}]}\right) \approx \sqrt{n}\left(\frac{E_n[(\tilde{Y} - \theta_0\tilde{D})\tilde{Z}]}{E[\tilde{D}\tilde{Z}]}\right)$$

- Consequently, it is *asymptotically normal*

$$\sqrt{n}(\hat{\theta} - \theta_0) \sim_a N(0, V), \quad V := \frac{E[(\tilde{Y} - \theta_0\tilde{D})^2 \tilde{Z}^2]}{E[\tilde{D}\tilde{Z}]^2}$$

- Confidence intervals* for any projection based on estimate of variance are asymptotically valid

$$\ell'\theta \in \left[\ell'\hat{\theta} \pm c \sqrt{\frac{\ell'\hat{V}\ell}{n}}\right], \quad \hat{V} = \frac{E_n[(\hat{Y} - \hat{\theta}\hat{D})^2 \hat{Z}^2]}{E_n[\hat{Z}\hat{D}]^2}$$

LATE in the Binary Case

- Under monotonicity

$$\theta_0 = \frac{E[E[Y | Z = 1, X] - E[Y | Z = 0, X]]}{E[E[D | Z = 1, X] - E[D | Z = 0, X]]}$$

- Moment formulation

$$E[E[Y | Z = 1, X] - E[Y | Z = 0, X] - \theta_0(E[D | Z = 1, X] - E[D | Z = 0, X])] = 0$$

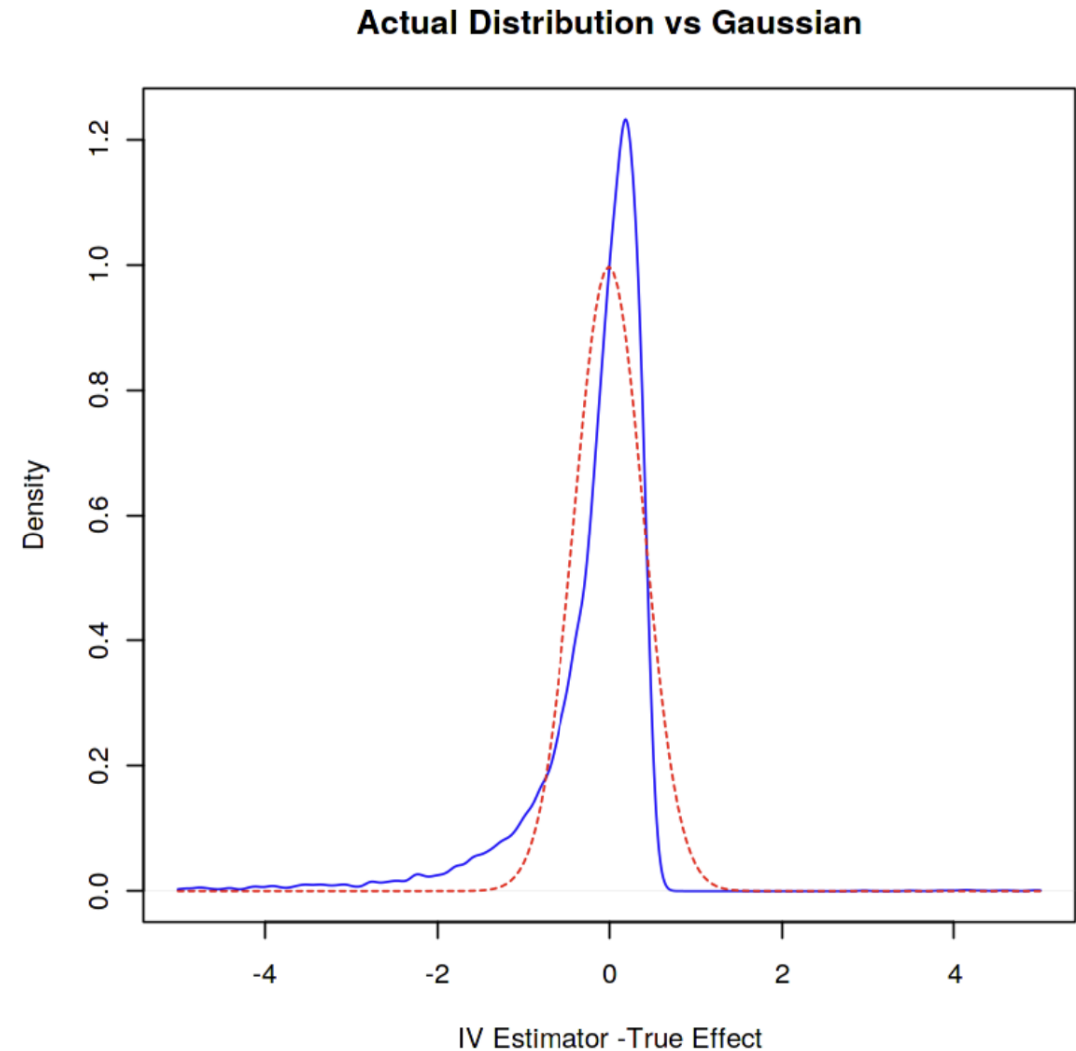
$$\begin{aligned} &+ \\ &\mathbf{H(Z, X)}(\mathbf{Y} - \mathbf{E[Y|Z, X]}) \qquad \qquad \qquad + \\ &\mathbf{H(Z, X)}(\mathbf{D} - \mathbf{E[D|Z, X]}) \end{aligned}$$

$$H(Z, X) = \frac{Z}{P(Z = 1|X)} - \frac{1 - Z}{1 - P(Z = 1|X)}$$

- Orthogonal moment formulation: apply ATE debiasing twice

Weak Identification

- If $E[\tilde{D}\tilde{Z}]$ is small and comparable with the sample size, then approximation $E_n[\tilde{D}\tilde{Z}]^{-1} \approx E[\tilde{D}\tilde{Z}]^{-1}$
- Can be inaccurate in finite samples and normal based approximation will yield in-correct confidence intervals



A More Robust Inference Approach

- Even in the weak regime the moment constraint is still well-behaved

$$E[(\tilde{Y} - \theta \tilde{D})\tilde{Z}]$$

- At the true parameter θ_0 we know that:

$$C(\theta) := \frac{(\sqrt{n} E_n[(\tilde{Y} - \theta \tilde{D})\tilde{Z}])^2}{Var_n((\tilde{Y} - \theta \tilde{D})\tilde{Z})} \sim_a (N(0,1))^2 = \chi^2(1)$$

- This statistic does not hinge on inversion of $E[\tilde{D}\tilde{Z}]$; approximation remains valid even with cross-fitted approximate residuals due to Neyman orthogonality
- We can perform a grid search over candidate parameters θ and for every such parameter test whether (for confidence interval with confidence α)

$$C(\theta) \leq (1 - \alpha) \text{ quantile of } \chi^2(1)$$

- Then by construction: $\Pr(\theta_0 \in C(\theta)) \approx 1 - \alpha$

General Moments and Weak Identification

- For a general Neyman orthogonal moment

$$E[m(Z; \theta_0, g_0)] = 0$$

- We can construct a statistic that is robust to weak identification (i.e. Jacobian $\partial_\theta E[m(Z; \theta_0, g_0)]$ very small)

$$C(\theta) = \frac{(\sqrt{n}E_n[m(Z; \theta, \hat{g})])^2}{Var_n(m(Z; \theta, \hat{g}))} \sim_a \chi^2(1)$$

- Construct a α -confidence region by including all parameter values θ s.t.

$$C(\theta) \leq (1 - \alpha) \text{ quantile of } \chi^2(1)$$

- Then by construction: $\Pr(\theta_0 \in C(\theta)) \approx 1 - \alpha$

Main Theorem (expanded) Define RMSE: $\|h\|_{L^2} = \sqrt{E[h(X)^2]}$

- If moment is Neyman orthogonal and RMSE of \hat{g} goes down at rate $n^{1/4}$, plus regularity conditions

$$n^{1/4} \|\hat{g} - g_0\|_{L^2} \approx 0$$

- Then the estimate $\hat{\theta}$ is *asymptotically linear*

$$\sqrt{n}(\hat{\theta} - \theta_0) \approx \sqrt{n} E_n[\phi_0(Z)], \quad \phi_0(Z) = -J_0^{-1} m(Z; \theta_0, g_0), \quad J_0 := \partial_{\theta} E[m(Z; \theta_0, g_0)]$$

- Consequently, it is *asymptotically normal*

$$\sqrt{n}(\hat{\theta} - \theta_0) \sim_a N(0, V), \quad V := E[\phi_0(Z)\phi_0(Z)']$$

- *Confidence intervals* for any projection based on estimate of variance are asymptotically valid

$$\ell' \theta \in \left[\ell' \hat{\theta} \pm c \sqrt{\frac{\ell' \hat{V} \ell}{n}} \right], \quad \hat{V} = \text{Var}_n(\hat{\phi}(Z)), \quad \hat{\phi}(Z) := -\hat{J}^{-1} m(Z; \hat{\theta}, \hat{g}), \quad \hat{J} = \partial_{\theta} E_n[m(Z; \hat{\theta}, \hat{g})]$$