

Slide 12 (Introduction)

- n experimental units (for measurement)
 - Unit i :
Covariates $X_i \in \mathcal{X}$ measured before treatment.
Outcome $Y_i \in \mathcal{Y} \subseteq \mathbb{R}$ measured after treatment.
 - Treatment $Z \in \mathcal{Z}$. Often $Z = (Z_1, \dots, Z_n)$, but not required.
 - Exposure: $A_i = A_i(Z) \in \mathcal{A}$. Usually $A=Z$ (identity mapping)
 $\mathcal{A} = \{0, 1\}$ (binary exposure)
 - $X = (X_1, \dots, X_n)$, same for A, Y .
-

Slide 13 (Design of Experiments)

- Bernoulli trial: $\pi(Z|X) = \prod_{i=1}^n (\pi(X_i))^{Z_i} (1 - \pi(X_i))^{1-Z_i}$
[Simple/sampling w. replacement: $\pi(X_i) = \pi$]
- Sampling w/o replacement: $\pi(Z|X) = \begin{cases} \binom{n}{n_i}^{-1}, & \text{if } \sum_{i=1}^n Z_i = n, \\ 0, & \text{otherwise.} \end{cases}$
- Randomized complete block design:
 k treatment levels: $0, 1, \dots, k-1$. m blocks.
 $Z = (Z_{ij})_{i \in [m], j \in [k]}$
 $\pi(Z) = \begin{cases} (k!)^{-m}, & \text{if } (Z_{ij})_{j \in [k]} \text{ is permutation of } (0, 1, \dots, k-1) \text{ for all } i. \\ 0, & \text{otherwise.} \end{cases}$

Slide 14 (Towards a formal theory)

- Potential outcomes: $Y(z) = (Y_1(z), \dots, Y_n(z))$, $z \in \mathcal{Z}$.
- Causal effect is understood as the contrast between $Y(z)$ for different z .

Assumption (Consistency of p.o.) $Y = Y(Z)$
↑ realized treatment.

- Fundamental problem: only 1 p.o. is observed for every unit.

Assumption

$A: \mathcal{Z} \rightarrow \mathcal{A}^n$ is a valid exposure mapping
 $z \mapsto (A_1(z), \dots, A_n(z))$

in the sense that $Y_i(z) = Y_i(z')$ for all i, z, z' s.t. $A_i(z) = A_i(z')$.

[In this case, we often write p.o. as $Y_i(a)$ for $a \in \mathcal{A}$.]

Definition The Neyman - Rubin causal model assumes

1. Consistency of p.o.
2. Identity exposure mapping is valid.

Usually $\mathcal{A} = \{0, 1\}$, and $Y_i(1) - Y_i(0)$ is called individual treatment effect.

Slide 15 (Example: Interference)

- Network: $G = (V = [n], E \subseteq [n] \times [n])$
 \swarrow vertex set \swarrow edge set
- $A_i(z) = (z_i, \sum_{(i', i) \in E} z_{i'}) \in \{0, 1\} \times \{0, \dots, n-1\}$.

Slide 16 (Example: Lady tasting tea)

- Data and potential outcome schedule.

[Coding: 0 means milk first, 1 means tea first.]

i	$z_i \in A_i$	Y_i	$Y_i(0)$	$Y_i(1)$	$Y_i(1) - Y_i(0)$
1	1	1	?	1	?
2	1	1	?	1	?
3	0	0	0	?	?
4	0	1	1	?	?
5	1	1	?	1	?
6	0	0	0	?	?
7	1	0	?	0	?
8	0	0	0	?	?

- Data can be summarized by the following 2×2 contingency table.

	$Y=0$	$Y=1$
$A=0$	3	1
$A=1$	1	3 ← test stat.

- Fisher's exact test: $\binom{8}{4} = 70$ possibilities.

Test statistic: how many correct tea-first guesses.

$$P(T=4) = \frac{\binom{4}{4} \binom{4}{4}}{70} = 1/70$$

$$P(T=3) = \frac{\binom{4}{3} \binom{4}{1}}{70} = 16/70.$$

p-value is $17/70 \approx 24.3\%$.

Slide 17 (Randomization inference: General tests)

Assumption (Exogeneity of randomization)

$$\textcircled{1} \quad Z \perp\!\!\!\perp (Y(z))_{z \in \mathcal{Z}} \mid X$$

$$[Z = W]$$

$$\textcircled{2} \quad P(Z=z \mid X=x) = \pi(z \mid x) \text{ is known}$$

- Fisher's sharp null: $H_0: Y(z)$ doesn't depend on z .
- In N-R model: $H_0: Y_i(z_0) = Y_i(z_1), \forall i \in [n]$.

This allows us to impute the p.o. schedule

i	$\sum_i^{\{A_i\}} Z_i$	Y_i	$Y_i(z_0)$	$Y_i(z_1)$	$Y_i(z_1) - Y_i(z_0)$
1	1	1	1	1	0
2	1	1	1	1	0
3	0	0	0	0	0
4	0	1	1	1	0
5	1	1	1	1	0
6	0	0	0	0	0
7	1	0	0	0	0
8	0	0	0	0	0

- Test statistic: $T: \mathcal{Z} \times \mathcal{X} \times \mathcal{W} \rightarrow \mathbb{R}$
- Randomization p-value:

$$P(Z, X, W) = \mathbb{P}(T(Z', X, W) \leq T(Z, X, W) | Z, X, W)$$

where $Z' \stackrel{d}{=} Z | X, W$, $Z' \perp\!\!\!\perp Z | X, W$.

Thm $\mathbb{P}(P(Z, X, W) \leq \alpha | X, W) \leq \alpha, \forall 0 < \alpha < 1$.

Pf Probability integral transform: Let F be the C.D.F. of r.v. T . Then $F(T)$ is (almost) uniform and.

$$P(F(T) \leq \alpha) = \alpha, \quad \forall 0 < \alpha < 1. \quad \square$$

Rem Computing $P(Z, X, W)$ requires

1. Randomization: $Z \perp\!\!\!\perp W \mid X$, $\pi(Z \mid X)$ is known.
2. W is known under H_0 .

What if H_0 is not exhaustive? Example: No spillover effect for vaccinated students.

Tool: Further condition on $g(Z') = g(Z)$ for some function g .

Ex N-R model, $A = Y = \{0, 1\}$. No covariates.

Sampling w.o. replacement.

Data can be summarized by

		Y		
		0	1	
A	0	N_{00}	N_{01}	$N_{0\cdot}$
	1	N_{10}	N_{11}	$N_{1\cdot} = n_1$
		$N_{\cdot 0}$	$N_{\cdot 1}$	$N_{\cdot\cdot} = n$

Under Fisher's sharp null, the probability of observing $(N_{00}, N_{01}, N_{10}, N_{11})$ is given by the hypergeometric prob.

$$\frac{\binom{N_{0\cdot}}{N_{00}} \binom{N_{1\cdot}}{N_{10}}}{\binom{N}{N_{0\cdot}}} \Rightarrow \text{Fisher's exact test.}$$

Slide 18 (Randomization inference: Estimation)

• Randomization dist. of $\hat{\beta}_{\text{DID}} = \bar{Y}_1 - \bar{Y}_0$,

$$\left[\bar{Y}_1 = \frac{\sum A_i Y_i}{\sum A_i}, \quad \bar{Y}_0 = \frac{\sum (1-A_i) Y_i}{\sum (1-A_i)} \right]$$

consistency of p.o. $\Rightarrow = \frac{1}{n_1} \sum_{i=1}^n A_i Y_i(1) - \frac{1}{n_0} \sum_{i=1}^n (1-A_i) Y_i(0)$
sampling w.o. replacement

$$\begin{aligned} \text{So } \mathbb{E}[\hat{\beta} | W] &= \frac{1}{n_1} \sum_{i=1}^n \mathbb{E}(A_i) \cdot Y_i(1) - \frac{1}{n_0} \sum_{i=1}^n \mathbb{E}(1-A_i) Y_i(0) \\ \text{exogeneity} &= \frac{1}{n_1} \sum_{i=1}^n \frac{n_1}{n} \cdot Y_i(1) - \frac{1}{n_0} \sum_{i=1}^n \frac{n_0}{n} \cdot Y_i(0) \\ &= \frac{1}{n} \sum_{i=1}^n Y_i(1) - Y_i(0) \\ &= \beta_n. \end{aligned}$$

Exercise: Show that $\text{Var}(\hat{\beta}_{\text{DID}} | W) = \frac{1}{n_0} S^2(0) + \frac{1}{n_1} S^2(1) - \frac{1}{n} S^2(0,1)$.

Hint: $\hat{\beta}_{\text{DID}}$ is linear in W , so $\text{Var}(\hat{\beta}_{\text{DID}} | W)$ is quadratic.

Use symmetry to argue

$$\text{Var}(\hat{\beta}_{\text{DID}} | W) = c_0 S^2(0) + c_1 S^2(1) + c_{01} S^2(0,1)$$

Now consider the model

$$\begin{pmatrix} Y_i(0) \\ Y_i(1) \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_0^2 & \rho\sigma_0\sigma_1 \\ \rho\sigma_0\sigma_1 & \sigma_1^2 \end{pmatrix} \right).$$

- Variance estimators:

$$\hat{S}_1^2 = \frac{1}{n_1-1} \sum_{i=1}^n A_i (Y_i - \bar{Y}_1)^2$$

$$\hat{S}_0^2 = \frac{1}{n_0-1} \sum_{i=1}^n (1-A_i) (Y_i - \bar{Y}_0)^2.$$

Exercise $\mathbb{E}(\hat{S}_a^2 | W) = S^2(a), \quad a=0,1.$

However, $S^2(0,1)$ cannot be directly estimated.

- Finite-sample CLT: Under additional assumptions on p.o., one can show

$$\frac{\vec{\beta} - \beta_n}{\sqrt{\text{Var}(\vec{\beta} | W)}} \xrightarrow{d} N(0, I) \quad \text{as } n \rightarrow \infty.$$

Slide 19 (Randomization inference: F-test)

- Randomization dist. of S_A and S_E .

Assume: $\mathcal{A} = \{0,1\} \Rightarrow S_A = \frac{n_1 n_0}{n} (\bar{Y}_1 - \bar{Y}_0)^2.$

sampling w.o. replacement

Fisher's sharp null: $H_0: Y_i(0) = Y_i(1), \quad \forall i \Rightarrow S^2(0) = S^2(1)$
 $S^2(0,1) = 0.$

$S_0 \quad \mathbb{E}(S_A | W) = \frac{n_1 n_0}{n} \text{Var}(\hat{\beta}_{01M} | W) = \frac{n_1 n_0}{n} S^2(0) \left(\frac{1}{n_0} + \frac{1}{n_1} \right) = S^2(0).$

On the other hand,

$$S_A + S_E = S_T = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$\mathbb{E}(S_T | W) = \sum_{i=1}^n (Y_i(1) - Y_i(0))^2 = (n-1) S^2(0).$$

$$\Rightarrow \frac{\mathbb{E}(S_A | W)}{\mathbb{E}(S_E | W)} = \frac{1}{n-2}.$$

It can be further shown that the F-test is approximately valid in the randomization model, without assuming the normal linear model.

Slide 20 (Repeated sampling)

- Positivity / overlap assumption: $\mathbb{P}(A=a | X=x) > 0$, $\forall a \in A, x \in X$.
- Causal identification:

Thm Assume exogeneity (\Rightarrow stratified Bernoulli trials) and positivity.

$$\text{Then } \mathbb{P}(Y(a) \leq y | X=x) = \mathbb{P}(Y \leq y | A=a, X=x), \forall a, x, y.$$

$$\begin{aligned} \text{Pf } \mathbb{P}(Y(a) \leq y | X=x) &= \mathbb{P}(Y(a) \leq y | X=x, A=a) \\ &= \mathbb{P}(Y \leq y | X=x, A=a). \quad \square \\ &\quad \uparrow \\ &\quad \text{consistency} \end{aligned}$$

- Average treatment effect: $\beta_{ATE} = \mathbb{E}(Y(1) - Y(0)).$

- Estimators: Outcome regression $\hat{\beta}_{OR} = \frac{1}{n} \sum_{i=1}^n \hat{\mu}_1(x_i) - \hat{\mu}_0(x_i).$

Inverse-probability weighting $\hat{\beta}_{IPW} = \frac{1}{n} \sum_{i=1}^n \frac{A_i Y_i}{\pi(x_i)} - \frac{(1-A_i) Y_i}{1-\pi(x_i)}.$

Exercise $\sqrt{n} (\hat{\beta}_{IPW} - \beta_{ATE}) \xrightarrow{d} N(0, \sigma^2_{IPW}).$

Slide 21 (M-estimation)

- Taylor expansion

$$0 = \frac{1}{n} \sum_{i=1}^n \psi(\hat{\theta}; D_i)$$

$$= \frac{1}{n} \sum_{i=1}^n \psi(\theta_0; D_i) + (\hat{\theta} - \theta_0)^T \left\{ \frac{\partial}{\partial \theta} \psi(\theta_0; D_i) \right\} + \underbrace{O_P(\|\hat{\theta} - \theta_0\|^2)}_{\text{Negligible if:}}$$

1. $\frac{\partial^2}{\partial \theta^2} \psi$ bounded

2. $\hat{\theta} - \theta_0 \xrightarrow{P} 0$

$$\Rightarrow \sqrt{n}(\hat{\theta} - \theta_0) \approx \underbrace{\frac{1}{\sqrt{n}} \sum_{i=1}^n \left\{ \mathbb{E} \left[\frac{\partial}{\partial \theta} \psi(\theta_0; D_i) \right] \right\}^{-1} \psi(\theta_0; D_i)}_{\text{Influence function}}$$

Influence function

$$\xrightarrow{d} N \left(0, \left\{ \mathbb{E} \left[\frac{\partial}{\partial \theta} \psi(\theta_0) \right] \right\}^{-1} \mathbb{E}[\psi(\theta_0) \psi(\theta_0)^T] \left\{ \mathbb{E} \left[\frac{\partial}{\partial \theta} \psi(\theta_0) \right] \right\}^{-1} \right)$$

"sandwich" variance.

Apply this to $l(\theta; D) = (Y - \theta^T V)^2$:

$$\text{Let } \varepsilon = Y - \theta_0^T V$$

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N \left(0, \left\{ \mathbb{E}(V V^T) \right\}^{-1} \mathbb{E}(V V^T \varepsilon^2) \left\{ \mathbb{E}(V V^T) \right\}^{-1} \right)$$

This is robust to heteroscedasticity (i.e. $\text{Var}(\varepsilon|V)$ depends on V).

Slide 22 (Regression adjustment)

Let $(\alpha_k, \beta_k, \gamma_k, \delta_k) = \underset{Y \in P_k, \delta \in \Delta_k}{\text{argmin}} \mathbb{E} \left[\left\{ Y - \alpha - \beta A - \gamma^T X - A \cdot (\delta^T X) \right\}^2 \right]$

Suppose $X \in \mathbb{R}^P$. Then $P_1, \Delta_1 = \{0\}$. $P_2 = \mathbb{R}^P, \Delta_2 = \{0\}$. $P_3 = \Delta_3 = \mathbb{R}^P$.

Lem $\alpha_1 = \alpha_2 = \alpha_3$, $\beta_1 = \beta_2 = \beta_3 = \beta_{ATE}$.

Proof Assume ∂ and \mathbb{E} can always be interchanged.

$$\frac{\partial}{\partial \alpha} = 0 \Rightarrow \mathbb{E} \left[Y - \alpha_k - \beta_k A - \gamma_k^T X - A_i (\delta_k^T X) \right] = 0.$$

$$\frac{\partial}{\partial \beta} = 0 \Rightarrow \mathbb{E} \left[A \left(- \quad - \quad - \quad - \quad - \right) \right] = 0.$$

By using $\mathbb{E}(X) = 0$ & $A \perp X$

$$\left. \begin{aligned} \mathbb{E} \left[Y - \alpha_k - \beta_k A \right] &= 0 \\ \mathbb{E} \left[A \left(Y - \alpha_k - \beta_k A \right) \right] &= 0. \end{aligned} \right\} \begin{array}{l} (1) \\ \text{for } k=1,2,3. \\ (2) \end{array}$$

$$\begin{aligned} (2) - (1) \times \pi &\Rightarrow \beta_k = \frac{\mathbb{E}(AY) - \pi \mathbb{E}(Y)}{\pi - \pi^2} \\ &= \frac{\pi \cdot \mathbb{E}(Y|A=1) - \pi \{ \pi \mathbb{E}(Y|A=1) + (1-\pi) \mathbb{E}(Y|A=0) \}}{\pi - \pi^2} \\ &= \mathbb{E}(Y|A=1) - \mathbb{E}(Y|A=0), \quad \square \end{aligned}$$

- Applying M-estimation to

$$V_1 = \begin{pmatrix} 1 \\ A \end{pmatrix}, \quad V_2 = \begin{pmatrix} 1 \\ A \\ X \end{pmatrix}, \quad V_3 = \begin{pmatrix} 1 \\ A \\ X \\ AX \end{pmatrix}.$$

$$\text{Let } \varepsilon_k = Y - \alpha_k - \beta_k A - \gamma_k^T X - A(\delta_k^T X).$$

Thm $\sqrt{n}(\hat{\beta}_k - \beta) \xrightarrow{d} N(0, V_k), \quad V_k = \underbrace{\frac{\mathbb{E}[(A-\pi)^2 \varepsilon_k^2]}{\pi^2(1-\pi)^2}}_{\text{Verify in E.S.}}, \quad k=1,2,3$

$$\text{And } V_3 \leq \min\{V_1, V_2\}.$$

Proof $\mathbb{E}(\varepsilon_k) = 0 \Rightarrow \mathbb{E}(\varepsilon_3) = \mathbb{E}(A\varepsilon_3) = 0$

$$\mathbb{E}(\varepsilon_3 X) = \mathbb{E}(\varepsilon_3 AX) = 0.$$

$$\varepsilon_1 = \varepsilon_3 + \gamma_3^T X + A(\delta_3^T X)$$

$$\varepsilon_2 = \varepsilon_3 + (\gamma_3 - \gamma_2)^T X + A(\delta_3^T X).$$

So for $k=1,2$

$$\begin{aligned} & \mathbb{E}\{(A-\pi)^2 \varepsilon_k^2\} - \mathbb{E}\{(A-\pi)^2 \varepsilon_3^2\} \\ &= \mathbb{E}\{(A-\pi)^2 [(\gamma_3 - \gamma_k)^T X + A(\delta_3^T X)]^2\} \geq 0. \end{aligned}$$

- Remark: We do not assume linear model is correct.
- When $E(x) \neq 0$, need to center X first.

$$X_i \rightarrow X_i - \bar{X}, \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

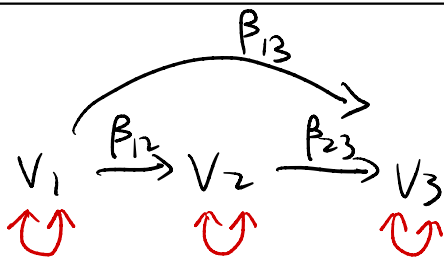
$\hat{\beta}_1, \hat{\beta}_2$ unchanged.

$$\hat{\beta}_3 \text{ becomes } \tilde{\beta}_3 = \hat{\beta}_3 + \hat{\sigma}_3 \bar{X}.$$

Can show $\tilde{\beta}_3$ is still more efficient.

Slide 28

- Example:



Bidirected self-loop often omitted.

$$V_1 = E_1$$

$$V_2 = \beta_{12} V_1 + E_2$$

$$V_3 = \beta_{13} V_1 + \beta_{23} V_2 + E_3$$

Intervention

$$\xrightarrow{V_1 = v_1}$$

$$V_2(v_1) = \beta_{12} v_1 + E_2$$

$$V_3(v_1) = \beta_{13} v_1 + \beta_{23} V_2(v_1) + E_3$$

$$= (\beta_{13} + \beta_{12} \beta_{23}) v_1 + \dots$$

- This suggests us to define the total causal effect of V_j on V_k as $\sigma(W(j \rightsquigarrow k))$.
- the set of all directed walks $j \rightarrow \dots \rightarrow k$.
-

Slide 29 (Path analysis)

- Trek rule:

$$V = \alpha + B^T V + E \Rightarrow V = (I - B)^{-T} (\alpha + E)$$

$$\Rightarrow \text{Cov}(V) = (I - B)^{-T} \Lambda (I - B)^{-1}$$

$$(I - B)^{-1} = I + B + B^2 + \dots = \sigma(I + W(V \rightsquigarrow V))$$

$$\Lambda = \sigma(W(V \leftrightarrow V))$$

$$\Rightarrow \text{Cov}(V) = \sigma(W(V \overset{t}{\rightsquigarrow} V))$$

$$\text{where } W(V \overset{t}{\rightsquigarrow} V) = \left\{ I + W(V \leftarrow V) \right\} \cdot W(V \leftrightarrow V) \cdot \left\{ I + W(V \rightsquigarrow V) \right\}$$

$$= W(V \leftrightarrow V) + W(V \leftarrow V \leftrightarrow V) + W(V \leftrightarrow V \rightsquigarrow V) + W(V \leftarrow V \leftrightarrow V \rightsquigarrow V)$$

$\overset{t}{\rightsquigarrow}$ is called a trek: a walk with 0 collider and 1 \leftrightarrow

- m -connectedness: a walk is an arc or (uncond.) m -connected if it has no collider;

$$W(V \overset{m}{\leftrightarrow} V) = W(V \overset{m}{\rightarrow} V) + W(V \overset{m}{\leftarrow} V) + W(V \overset{m}{\leftrightarrow} V)$$

Exercise How many \leftrightarrow can an arc have?

Thm Suppose AOMG \mathcal{G} has all bidirected self-loops:

$(j, j) \in \mathcal{B}$ for all $j \in V$. Then

$$W(j \overset{t}{\leftrightarrow} k) \neq \emptyset \iff P(j \overset{t}{\rightarrow} k) \neq \emptyset.$$

Pf \Leftarrow Insert a bidirected self-loop if needed.

\Rightarrow suppose $\pi \in W(j \overset{t}{\leftrightarrow} k)$

So $\pi = j \overset{t}{\rightarrow} l \leftrightarrow l' \overset{t}{\rightarrow} k$ (possibly $l=j$, $l'=k$).

If π is not already a path, suppose r is the repeated vertex closer to j .

r must appear once on $j \overset{t}{\rightarrow} l$ and once on $l' \overset{t}{\rightarrow} k$.

$$\text{So } \pi = j \leftarrow \overset{\circ}{r} \leftarrow \overset{\circ}{l} \leftrightarrow \overset{\circ}{l'} \overset{\circ}{r} \overset{\circ}{r} \rightarrow k$$

$$\underbrace{\leftarrow \overset{\circ}{r} \leftarrow \overset{\circ}{l} \leftrightarrow \overset{\circ}{l'} \overset{\circ}{r} \overset{\circ}{r} \rightarrow}_{\in W(r \overset{t}{\leftrightarrow} r)}$$

We say r is the root of π , and denote all treks from j to k with root r as $W(j \overset{t}{\leftrightarrow} k, \text{root } r)$

By assumption, $W(r \overset{t}{\leftrightarrow} r) \neq \emptyset$.

So $W(j \overset{t}{\leftrightarrow} k, \text{root } r) \neq \emptyset \Leftrightarrow P(j \overset{\circ}{\leftarrow} r \overset{\circ}{\rightarrow} k) \neq \emptyset$.

The desired conclusion follows from

$$W(j \overset{t}{\leftrightarrow} k) = P(j \overset{t}{\leftrightarrow} k) + \sum_{r \in V} W(j \overset{t}{\leftrightarrow} k, \text{root } r)$$

$$P(j \rightsquigarrow k) = P(j \overset{t}{\leftrightarrow} k) + \sum_{r \in V} P(j \overset{\circ}{\leftarrow} r \overset{\circ}{\rightarrow} k) \quad \square$$

• Wright's path analysis: As a corollary,

$$\text{Cov}(V_j, V_k) = \sigma(P(j \overset{t}{\leftrightarrow} k)) + \sum_{r \in V} \sigma(P(j \overset{\circ}{\leftarrow} r \overset{\circ}{\rightarrow} k)) \cdot \text{Var}(V_r)$$

If V is standardized so that $\text{Var}(V_j) = 1$, $\forall j \in V$, then

$$\text{Cov}(V_j, V_k) = \sigma(P(j \rightsquigarrow k))$$

Compare this to the total causal effect $\sigma(P(j \rightsquigarrow k))$.

Slide 30 (Two examples)

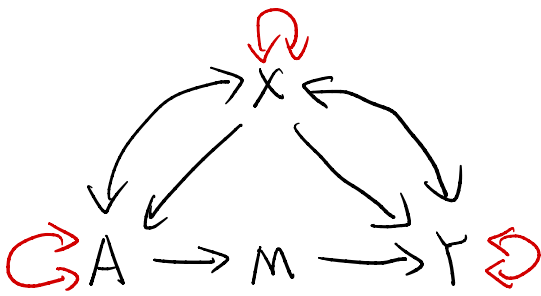


$$\text{Var}(A) = \lambda_{AA}$$

$$\text{Var}(Y) = \lambda_{YY} + \lambda_{AA} \beta_{AY} + \lambda_{AY} \beta_{AY}$$

Q: Why?

$$\text{Cov}(A, Y) = \lambda_{AY} + \lambda_{AA} \beta_{AY}$$



$$\text{Cov}(A, Y) = \beta_{AM} \beta_{AY} \cdot \text{Var}(A) + \beta_{XA} \beta_{XY} \cdot \text{Var}(X) + \beta_{XA} \lambda_{XY} + \beta_{XY} \lambda_{AX}$$

Slide 31 (Correlation vs. causation)

Examples

		Causal effect	Correlation	Partial cor.
1. Confounder:	$A \leftarrow X \rightarrow Y$	$= 0$	$\neq 0$	$\neq 0$
2. Mediator:	$A \rightarrow M \rightarrow Y$	$\neq 0$	$\neq 0$	$= 0$
3. Collider:	$A \rightarrow C \leftarrow Y$	$= 0$	$= 0$	$\neq 0$

Slide 34 (Marginalization and latent projection)

$$(I - B)^{-1} = \begin{pmatrix} ((I - B)^{-1})_{\tilde{v}, \tilde{v}} & * \\ ((I - B)^{-1})_{u, \tilde{v}} & * \end{pmatrix}$$

$$((I-B)^{-1})_{\tilde{v}, \tilde{v}} = (I - \tilde{B})^{-1}, \text{ where}$$

$$\begin{aligned} \tilde{B} &= B_{\tilde{v}, \tilde{v}} + B_{\tilde{v}, u} (I - B_{u, u})^{-1} B_{u, \tilde{v}} \\ &= \sigma(W(\tilde{v} \overset{u}{\rightsquigarrow} \tilde{v})). \end{aligned}$$

$$((I-B)^{-1})_{u, \tilde{v}} = (I - B_{u, u})^{-1} B_{u, \tilde{v}} (I - \tilde{B})^{-1}$$

$$\begin{aligned} \tilde{\Sigma} &= ((I-B)^{-1} \wedge (I-B)^{-1})_{\tilde{v}, \tilde{v}} \\ &= (I - \tilde{B})^{-1} \Lambda_{\tilde{v}, \tilde{v}} (I - \tilde{B})^{-1} \\ &\quad + (I - \tilde{B})^{-1} \Lambda_{\tilde{v}, u} ((I-B)^{-1})_{u, \tilde{v}} \\ &\quad + ((I-B)^{-1})_{\tilde{v}, u} \Lambda_{u, \tilde{v}} (I - \tilde{B})^{-1} \\ &\quad + ((I-B)^{-1})_{\tilde{v}, u} \Lambda_{u, u} ((I-B)^{-1})_{u, \tilde{v}} \\ &= (I - \tilde{B})^{-1} \tilde{\Lambda} (I - \tilde{B})^{-1}. \end{aligned}$$

$$\text{where } \tilde{\Lambda} = \Lambda_{\tilde{v}, \tilde{v}} + \Lambda_{\tilde{v}, u} (I - B_{u, u})^{-1} B_{u, \tilde{v}}$$

$$+ B_{u, \tilde{v}}^T (I - B_{u, u})^{-1} \Lambda_{u, \tilde{v}}$$

$$+ \{(I - B_{u, u})^{-1} B_{u, \tilde{v}}\}^T \Lambda_{u, u} \{(I - B_{u, u})^{-1} B_{u, \tilde{v}}\}$$

$$= \sigma(W(\tilde{v} \leftrightarrow \tilde{v})) + \sigma(W(\tilde{v} \leftrightarrow u \overset{u}{\rightsquigarrow} u \rightarrow \tilde{v}))$$

$$+ \sigma(W(\tilde{v} \leftarrow u \overset{u}{\rightsquigarrow} u \leftrightarrow \tilde{v}))$$

$$+ \sigma(W(\tilde{v} \leftarrow u \overset{u}{\rightsquigarrow} u \leftrightarrow u \overset{u}{\rightsquigarrow} u \rightarrow \tilde{v}))$$

$$= \sigma(W(\tilde{v} \overset{t, \text{via } u}{\rightsquigarrow} \tilde{v})).$$

Exercise: Show that $W(\tilde{V} \overset{\text{trivial}}{\rightleftarrows} \tilde{V}) = \emptyset$
iff $W(\tilde{V} \overset{\text{via } U}{\rightleftarrows} \tilde{V}) = \emptyset$.

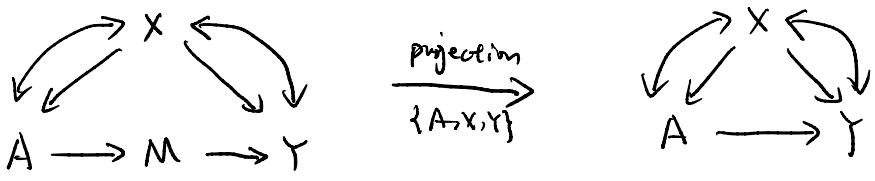
• Definition: $j \rightsquigarrow k | L \iff W(j \overset{\text{via } U}{\rightsquigarrow} k | L) \neq \emptyset$
 "j has a direct causal effect on k given L"

$j \rightleftarrows k | L \iff W(j \overset{\text{via } U}{\rightleftarrows} k | L) \neq \emptyset$
 "There is endogeneity between j and k given L"

• Corollary: $\forall \tilde{V} \subset V, a, b \in \tilde{V}, c \in \tilde{V}, a \neq b, a, b \notin c$.

$$a \left\{ \begin{array}{l} \rightsquigarrow \\ \leftarrow \\ \rightleftarrows \end{array} \right\} b | c [G] \iff a \left\{ \begin{array}{l} \rightsquigarrow \\ \leftarrow \\ \rightleftarrows \end{array} \right\} b | c [\tilde{G}]$$

• Example:



Slide 35 (A graphical criterion for conditional independence)

Definitions

- A walk is m^* -connected given L if all its colliders are in L and none of its non-colliders is in L .
- We say j and k are m -connected given L if there exists a m^* -connected walk like $j \rightsquigarrow * \leftarrow k \mid L$.
- We say j and k are confounded given L if there exists a m^* -connected walk like $j \leftarrow * \rightarrow k \mid L$.
- Unconfoundedness:

$$j \leftarrow * \rightarrow k \mid L [G] \Leftrightarrow j \leftarrow * \rightarrow k \mid L [\tilde{G}]$$

$$\Leftrightarrow j \leftrightarrow * \leftrightarrow k [G] \Rightarrow \tilde{\Lambda} = \begin{matrix} \rightarrow^j & \\ & \rightarrow^k \end{matrix} \begin{pmatrix} * & 0 \\ 0 & * \end{pmatrix}$$

Q: Why?

$$\Rightarrow (\tilde{\Lambda}^{-1})_{jk} = 0.$$

- m -separation: Recall $\tilde{\Omega} = \tilde{\Sigma}^{-1} = (I - \tilde{B}) \tilde{\Lambda}^{-1} (I - \tilde{B})^T$.

$$j \rightsquigarrow * \leftarrow k \mid L [G] \Leftrightarrow j \rightarrow * \leftarrow k [G^*]$$

$$\Rightarrow (\tilde{\Omega})_{jk} = 0$$

- Completeness: For almost all linear SEM w.r.t. G ,

$$j \rightsquigarrow * \leftarrow k \mid L [G] \Leftrightarrow (\tilde{\Omega})_{jk} = 0$$

• m -connected path: we say a path is m -connected given L if

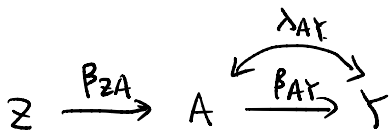
1. none of its non-colliders is in L

2. for any collider on the path, it is either in L or has a descendant in L .

It can be shown that there is no m^* -connected walk from j to k given L iff there is no m -connected path from j to k given L .

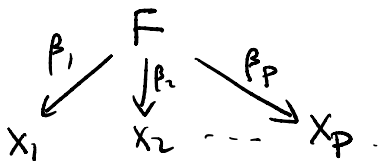
Slide 37 (Identifiability problems in linear SEM)

• IV graph



Generic identifiability: $\beta_{AY} = \frac{\text{Cov}(Z, Y)}{\text{Cov}(Z, A)}$ if $\beta_{ZA} \neq 0$.

• Factor analysis:



F is unobserved.

Thm Assume $\text{Var}(F) = 1$. Then $(\beta_1, \dots, \beta_p)$ is generically identifiable up to a sign change iff $p \geq 3$.

Pf Denote $\text{Var}(\xi_i) = \sigma_i^2$.

$$\Sigma = \text{Cov}(X) = \beta\beta^T + \text{diag}(\sigma_1^2, \dots, \sigma_p^2).$$

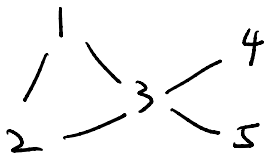
$p=3$:

$$\Sigma = \begin{pmatrix} \Sigma_{11} & & \\ \Sigma_{12} & \Sigma_{22} & \\ \Sigma_{13} & \Sigma_{23} & \Sigma_{33} \end{pmatrix} = \begin{pmatrix} \beta_1^2 + \sigma_1^2 & & \\ \beta_1\beta_2 & \beta_2^2 + \sigma_2^2 & \\ \beta_1\beta_3 & \beta_2\beta_3 & \beta_3^2 + \sigma_3^2 \end{pmatrix}$$

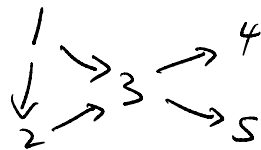
$$\Rightarrow \beta_1^2 = \frac{\Sigma_{12}\Sigma_{13}}{\Sigma_{23}}. \quad \text{Similar for } \beta_2^2, \beta_3^2. \quad \square$$

• Measurement model: In psychometrics, it is often of interest to infer causal relationships between abstract concepts that are measured by questionnaires. See example paper.

Slide 40 (Factorization)



$$f(v) = \psi(v_1, v_2, v_3) \psi(v_3, v_4) \psi(v_3, v_5)$$



$$f(v) = f(v_1) f(v_2|v_1) f(v_3|v_1, v_2) \cdot f(v_4|v_3) f(v_5|v_3)$$

Slide 44 (Global Markov \Leftrightarrow Factorization in DAG models)

WLOG suppose $(1, \dots, p)$ is a topological ordering.

It's obvious that factorization is equivalent to

$$V_j \perp\!\!\!\perp V_{[j-1] \setminus \text{pa}(j)} \mid V_{\text{pa}(j)}, \forall j \quad (\text{Ordered Markov})$$

It suffices to show Order Markov \Rightarrow Global Markov.

Below: Proof by induction.

Induction hypothesis If $J \cup K \cup L \subseteq [m]$,

$$J \perp\!\!\!\perp K \mid L \Rightarrow V_J \perp\!\!\!\perp V_K \mid V_L.$$

$$m=2 \quad \checkmark$$

Now suppose it is true up to $m-1$.

Also given: $V_m \perp\!\!\!\perp V_{[m-1] \setminus \text{pa}(m)} \mid V_{\text{pa}(m)}$.

$$J \perp\!\!\!\perp K \mid L$$

Want to prove: $V_J \perp\!\!\!\perp V_K \mid V_L$.

Idea: Apply the chain rule.

Let $N = \text{pa}(m) \setminus L$. $L_1 = L \cap \text{pa}(m)$, $L_2 = L \setminus \text{pa}(m) \setminus \{m\}$

Three cases: ① $m \in J$ ② $m \in K$ ③ $m \in L$.

① $m \in J$ $m \rightsquigarrow \not\Leftarrow K \mid L$.



So $K \rightarrow m$, $K \cap N = \emptyset$, $K \rightsquigarrow \not\Leftarrow N \mid L$.

$\Rightarrow J \cup N \rightsquigarrow \not\Leftarrow K \mid L$.

$\Rightarrow^{(1)}$ $V_{J \cup N \setminus \{m\}} \perp V_K \mid V_L$.

By
(Induction
hypothesis)

Further,

$m \rightsquigarrow \not\Leftarrow K \mid J \cup N \cup L \setminus \{m\}$

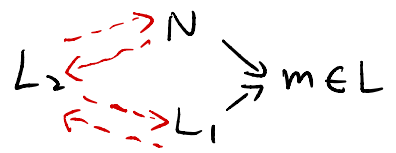
$\Rightarrow^{(2)}$ $V_m \perp V_K \mid V_L, V_{J \cup N \setminus \{m\}}$ (By
Ordered
Markov)

$\Rightarrow^{(1)}$ $V_{J \cup N} \perp V_K \mid V_L$ (chain rule)

$\Rightarrow V_J \perp V_K \mid V_L$.

② $m \in K$. symmetric.

③ $m \in L$



Claim $J \cap N = \emptyset$ and $N \rightsquigarrow^* \leftarrow K \mid L$

or

$K \cap N = \emptyset$ and $N \rightsquigarrow^* \leftarrow K \mid L$

Otherwise $J \rightsquigarrow^* \leftarrow N \rightarrow m \leftarrow N \rightsquigarrow^* \leftarrow K \mid L$.

WLOG, suppose $N \rightsquigarrow^* \leftarrow K \mid L$.

$\Rightarrow K \rightarrow m, J \cup N \rightsquigarrow^* \leftarrow K \mid L$.

Because m is the last vertex and only has edges like $\rightarrow m$, this shows

$J \cup N \rightsquigarrow^* \leftarrow K \mid L \setminus \{m\}$

$\Rightarrow^{(3)} V_{J \cup N} \perp V_K \mid V_{L \setminus \{m\}}$.

Ordered local Markov ($\text{pa}(m) = N \cup L_1$)

$\Rightarrow V_m \perp V_{[m-1] \setminus N \setminus L_1} \mid V_{N \cup L_1}$.

$\Rightarrow^{(4)} V_m \perp V_K \mid V_{J \cup N \cup L_1 \setminus \{m\}}$ (Chain rule)

(3), (4) $\Rightarrow V_K \perp V_{J \cup N \cup \{m\}} \mid V_{L \setminus \{m\}}$

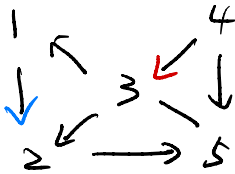
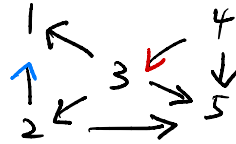
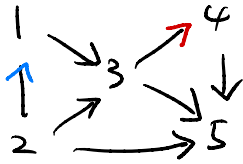
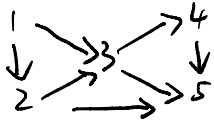
$\Rightarrow V_K \perp V_J \mid V_L$ (Chain rule)

□

Slide 47 (Markov equivalence of DAG models)

Markov equivalence class:

Ex 1



Slide 50 (NPSEM)

- Recursive substitution

- Basic potential outcomes: $V_j(V_{\setminus j}) = f_j(V_{\text{pa}(j)}, E_j), j \in V.$

- Derived potential outcomes: $V_j(V_I) = V_j(V_{\text{pa}(j)} \cap I, V_{\text{pa}(j)} \setminus I(V_I)).$

This defines a "natural counterfactual" $V_i(V_I)$ for $i \in I.$

- Example



Basic : $V_1(v_1, v_2, v_3) = V_1$, $V_2(v_1, v_2, v_3) = V_2(v_1)$

$V_3(v_1, v_2, v_3) = V_3(v_1, v_2)$.

I = {v} : $V_1(v_2) = V_1$, $V_2(v_2) = V_2$, $V_3(v_2) = V_3(V_1, v_2)$.

- Simplification of potential outcomes.

Prop Disjoint $I, I' \subseteq V$. For any $j \in V$,

$I' \not\rightarrow j \mid I \Rightarrow V_j(v_I, v_{I'}) = V_j(v_I)$.

Corollary $V_j(v_I) = V_j(v_I \cap \text{ant}(j))$.

- Consistency of potential outcomes.

It follows from the definition that $V_j = V_j(v_\emptyset) = V_j(V_{\text{pa}(j)})$.

Prop Disjoint $I, I' \subseteq V$.

$V_{I'}(v_I) = v_{I'} \Rightarrow V(v_I, v_{I'}) = V(v_I)$.

Slide 51 (Markov properties of basic p.o.)

Ex



- Single-world: $V_1 \perp\!\!\!\perp V_2(v_1) \perp\!\!\!\perp V_3(v_1, v_2)$, $\forall v_1, v_2$.

- Multiple-world: $V_1 \perp\!\!\!\perp (V_2(v_1) : v_1) \perp\!\!\!\perp (V_3(v_1, v_2) : v_1, v_2)$

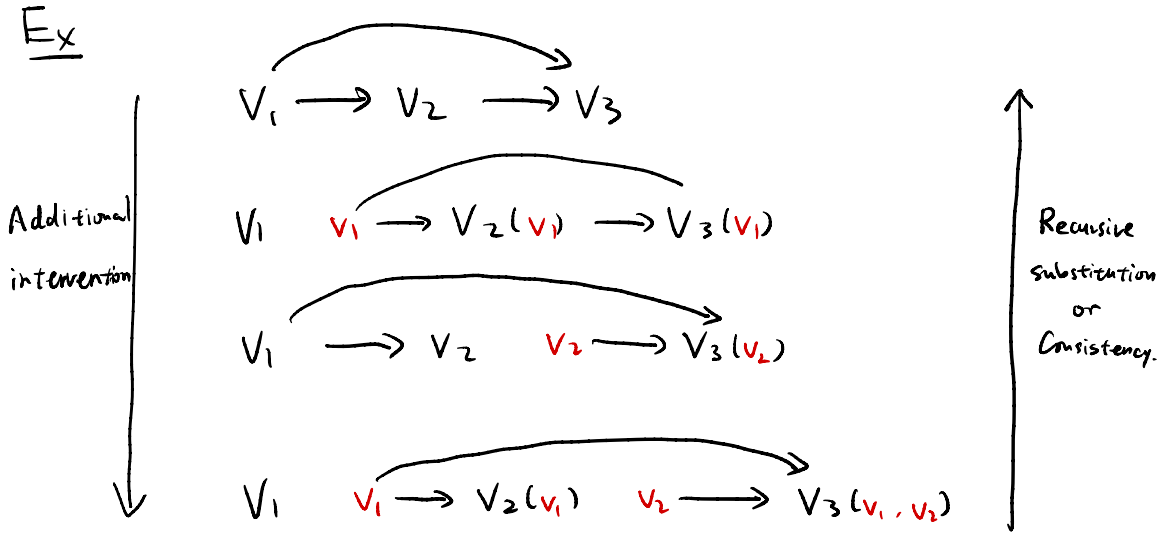
which includes all single-world independence but also cross-world

independence like

$$V_2(v_1) \perp\!\!\!\perp V_3(v_1, v_2)$$

- NPSEM via basic p.o.: $E_j = (V_j(V_{pa(j)}): V_{pa(j)})$
 f_j selects the corresponding p.o.

Slide 52 (Representing recursive substitution as a graph)



Slide 53 (Recursive substitution preserves global Markov)

Proof by induction: Suppose $del(j) \subseteq I$ and $I' = I \cup \{j\}$. If $V(v_Z)$ is Markov w.r.t. $\mathcal{E}_j[I']$, then $V(v_Z)$ is Markov w.r.t. $\mathcal{E}_j[I]$.

Claim 1 Disjoint $K, L, M \subseteq V$, $j \in L$. Then

$$V_K(v_Z) \not\perp\!\!\!\perp V_L(v_Z) \mid V_M(v_Z) \Rightarrow V_K(v_Z) \perp\!\!\!\perp V_L(v_Z) \mid V_M(v_Z)$$

Pf By assumption, $j \rightarrow K$. Then

$$V_K(v_Z) \not\perp\!\!\!\perp V_L(v_Z) \mid V_M(v_Z)$$

$$\Rightarrow V_K(v_Z) \rightsquigarrow^* \leftarrow \left(V_L(v_Z), V_{M \cap \text{ch}(j)}(v_Z) \right) \mid V_{M \cap \text{ch}(j)}(v_Z). \quad (1)$$

[Because $V_{\text{ch}(j)}(v_Z)$ has no outgoing edges like $V_{\text{ch}(j)}(v_Z) \rightarrow^*$
See ES2 Q9.]

$$\Rightarrow V_K(v_{Z'}) \rightsquigarrow^* \leftarrow \left(V_L(v_{Z'}), V_{M \cap \text{ch}(j)}(v_{Z'}) \right) \mid V_{M \cap \text{ch}(j)}(v_{Z'}).$$

$$\Rightarrow V_K(v_Z, v_j) \perp \left(V_L(v_Z, v_j), V_{M \cap \text{ch}(j)}(v_Z, v_j) \right) \mid V_{M \cap \text{ch}(j)}(v_Z, v_j).$$

$$\Rightarrow V_K(v_Z) \perp \left(V_L(v_Z), V_{M \cap \text{ch}(j)}(v_Z) \right) \mid V_{M \cap \text{ch}(j)}(v_Z).$$

[$j \in L$. consistency. Lemma 2]

Claim 2 Disjoint $K, L, M \subseteq V$, $j \notin K \cup L$. Then

$$V_K(v_Z) \rightsquigarrow^* \leftarrow V_L(v_Z) \mid V_M(v_Z) \Rightarrow V_K(v_Z) \perp V_L(v_Z) \mid V_M(v_Z).$$

Pf If $(K \cup L \cup M) \cap \text{ch}(j) = \emptyset$. Trivial.

Now suppose $j \rightarrow K \cup L \cup M$.

Observation 1 $V_K(v_Z) \rightsquigarrow^* \leftarrow V_L(v_Z) \mid V_M(v_Z), V_j(v_Z)$.

Otherwise $V_K(v_Z) \rightsquigarrow^* \leftarrow V_j(v_Z) \leftarrow^* \leftarrow V_L(v_Z) \mid V_M(v_Z), V_j(v_Z)$

This leads to contradictions when $j \rightarrow K \cup L \cup M$.

Observation 2 Either $V_K(v_Z) \rightsquigarrow^* \leftarrow V_j(v_Z) \mid V_M(v_Z)$

or $V_L(v_Z) \rightsquigarrow^* \leftarrow V_j(v_Z) \mid V_M(v_Z)$

Otherwise $V_K(v_Z) \rightsquigarrow^* \leftarrow V_j(v_Z) \rightsquigarrow^* \leftarrow V_L(v_Z) \mid V_M(v_Z)$

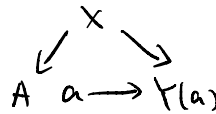
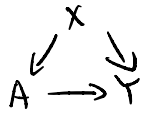
This contradicts the assumption or Observation 1.

WLOGT Suppose $V_k(V_Z) \rightsquigarrow \not\leftarrow \leftarrow V_j(V_Z) \mid V_M(V_Z)$

\Rightarrow $V_k(V_Z) \rightsquigarrow \not\leftarrow \leftarrow V_{LU\{j\}}(V_Z) \mid V_M(V_Z)$.
 claim1 \Rightarrow $V_k(V_Z) \perp\!\!\!\perp V_{LU\{j\}}(V_Z) \mid V_M(V_Z)$. \square

Slide 54 (DAG causal models)

• Example:



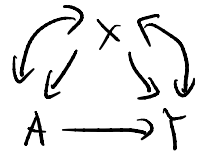
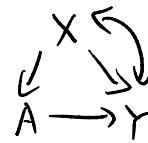
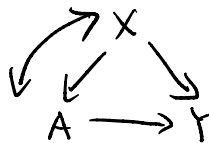
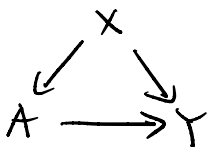
So $A \perp\!\!\!\perp Y(a) \mid X$.

$$P(A=\tilde{a}, X=\tilde{x}, Y(a)=\tilde{y}) = P(A=\tilde{a} \mid X=\tilde{x}) P(X=\tilde{x}) P(Y=\tilde{y} \mid A=\tilde{a}, X=\tilde{x})$$

$$\begin{aligned} \Rightarrow P(Y(a)=\tilde{y}) &= \sum_{\tilde{a}, \tilde{x}} P(A=\tilde{a} \mid X=\tilde{x}) P(X=\tilde{x}) P(Y=\tilde{y} \mid A=\tilde{a}, X=\tilde{x}) \\ &= \sum_{\tilde{x}} P(X=\tilde{x}) P(Y=\tilde{y} \mid A=\tilde{a}, X=\tilde{x}) \\ &= \mathbb{E} \{ P(Y=\tilde{y} \mid A=\tilde{a}, X) \}. \end{aligned}$$

Slide 55 (Back-door criterion)

Ex



Slide 56 (Front-door criterion)

$$\begin{aligned}
 \mathbb{P}(Y(a) = y) &= \mathbb{P}(Y(a, M(a)) = y) \\
 &= \mathbb{P}(Y(M(a)) = y) \quad [Y(a, m) = Y(m)] \\
 &= \sum_m \mathbb{P}(Y(m) = y | M(a) = m) \mathbb{P}(M(a) = m) \\
 &= \sum_m \mathbb{P}(Y(m) = y) \mathbb{P}(M(a) = m) \\
 &= \sum_m \left\{ \sum_{a'} \mathbb{P}(Y=y | M=m, A=a') \mathbb{P}(A=a') \right\} \mathbb{P}(M=m | A=a)
 \end{aligned}$$

Nonparametric path analysis.

Slide 57 (The fixing operator)

Pf $V = \{i\} \cup \text{de}(i) \cup \text{nd}(i)$. $\text{nd} = \text{non-descendant}$.

Claim $V_i(V_i) \rightsquigarrow \leftarrow \leftarrow V_{\text{de}(i)}(V_i) \mid V_{\text{nd}(i)}(V_i)$

$\Leftrightarrow V_i(V_i) \leftarrow \leftarrow \leftarrow V_{\text{de}(i)}(V_i) \mid V_{\text{nd}(i)}(V_i)$

$\Leftrightarrow V_i(V_i) \leftarrow \leftarrow \leftarrow V_{\text{de}(i)}(V_i) \mid V_{\text{nd}(i)}(V_i)$

because V_i is split acyclicity (so $V_{\text{de}(i)} \not\perp V_{\text{nd}(i)}$)

$\Leftrightarrow i$ fixable.

Thus $\mathbb{P}(V(V_i) = \tilde{v})$

$= \mathbb{P}(V_{\text{nd}(i)}(V_i) = \tilde{v}_{\text{nd}(i)}) \mathbb{P}(V_i(V_i) = \tilde{v}_i \mid V_{\text{nd}(i)}(V_i) = \tilde{v}_{\text{nd}(i)})$

$\cdot \mathbb{P}(V_{\text{de}(i)}(V_i) = \tilde{v}_{\text{de}(i)} \mid V_{\text{nd}(i)}(V_i) = \tilde{v}_{\text{nd}(i)}, V_i(V_i) = \tilde{v}_i)$

$$\begin{aligned}
&= \mathbb{P}(V_{nd(i)} = \tilde{V}_{nd(i)}) \cdot \mathbb{P}(V_i = \tilde{V}_i \mid V_{nd(i)} = \tilde{V}_{nd(i)}) \\
&\quad \cdot \mathbb{P}(V_{del(i)}(v_i) = \tilde{V}_{del(i)} \mid V_{nd(i)}(v_i) = \tilde{V}_{nd(i)}, V_i(v_i) = v_i) \\
&= \mathbb{P}(V_{nd(i)} = \tilde{V}_{nd(i)}) \cdot \mathbb{P}(V_i = \tilde{V}_i \mid V_{nd(i)} = \tilde{V}_{nd(i)}) \\
&\quad \cdot \mathbb{P}(V_{del(i)} = \tilde{V}_{del(i)} \mid V_{nd(i)} = \tilde{V}_{nd(i)}, V_i = v_i)
\end{aligned}$$

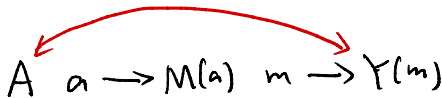
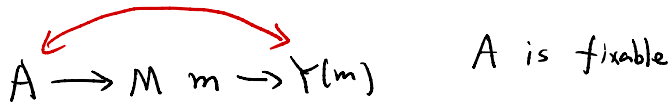
Claimed \perp
↓

$$\begin{aligned}
\text{So } \frac{\mathbb{P}(V(v_i) = \tilde{V})}{\mathbb{P}(V_i = v_i, V_{del(i)} = \tilde{V}_{del(i)})} &= \frac{\mathbb{P}(V_i = \tilde{V}_i \mid V_{nd(i)} = \tilde{V}_{nd(i)})}{\mathbb{P}(V_i = v_i \mid V_{nd(i)} = \tilde{V}_{nd(i)})} \\
&= \frac{\mathbb{P}(V_i = \tilde{V}_i \mid V_{mb(i)} = \tilde{V}_{mb(i)})}{\mathbb{P}(V_i = v_i \mid V_{mb(i)} = \tilde{V}_{mb(i)})}
\end{aligned}$$

Exercise Show that if i is fixable, then $i \rightsquigarrow \ast \leftarrow \text{nd}(i) \setminus \text{mb}(i) \mid \text{mb}(i)$. □

- Examples:

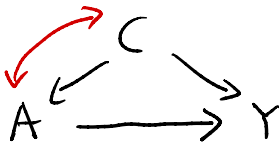
1.



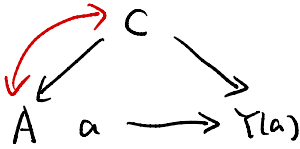
\Rightarrow Prob. dist. of $(A, M(a), Y(m))$ identifiable.

$$\begin{aligned}
\Rightarrow \mathbb{P}(Y(a) = y) &= \mathbb{E} \{ \mathbb{P}(Y(a) = y \mid M(a)) \} \\
&= \mathbb{E} \{ \mathbb{P}(Y(M(a)) = y \mid M(a)) \} \quad \text{identifiable.}
\end{aligned}$$

2.

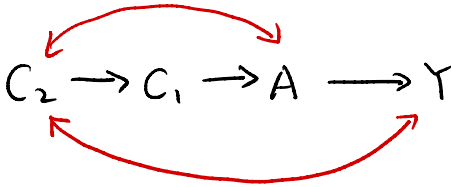


A is fixable



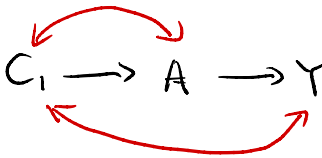
\Rightarrow Prob. dist. of $(A, C, Y(a))$ identifiable.

3.

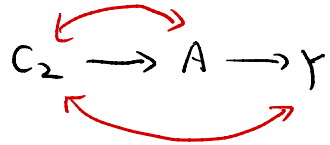


This is a challenging problem because the following are not identifiable

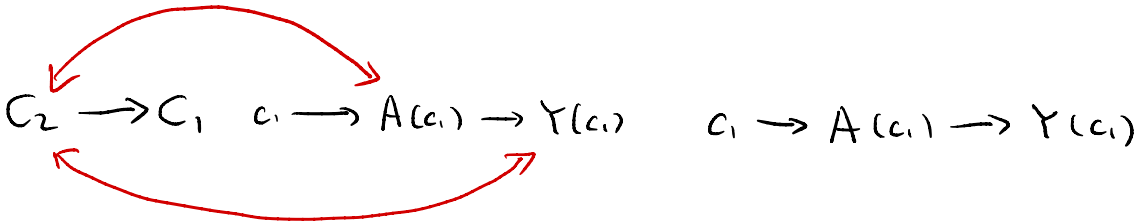
Projection on $\{C_1, A, Y\}$



$\{C_2, A, Y\}$



But we can first fix C_1 and then marginalize C_1, C_2



Now $A(c_1)$ is fixable: $c_1 \rightarrow A(c_1) \quad a \rightarrow Y(a)$

Identification formula:

$$P(A(c_1) = a, Y(c_1) = y)$$

$$= \sum_{\tilde{c}_1, \tilde{c}_2} P(C_2 = \tilde{c}_2, C_1 = c_1, A = a, Y = y) \cdot \frac{P(C_1 = \tilde{c}_1 | C_2 = \tilde{c}_2)}{P(C_1 = c_1 | C_2 = \tilde{c}_2)}$$

$$= \sum_{c_2} P(C_2 = c_2) \cdot P(A = a, Y = y | C_1 = c_1, C_2 = c_2)$$

$$\text{So } P(Y(a) = y)$$

$$= P(Y(c_1) = y | A(c_1) = a)$$

$$= \frac{E \{ P(A = a, Y = y | C_1 = c_1, C_2) \}}{E \{ P(A = a | C_1 = c_1, C_2) \}}$$

Slide 67 (van Mises expansion)

- Empirical distribution: $IP_n = \frac{1}{n} \sum_{i=1}^n \delta_{V_i} \rightarrow P$ (Glivenko-Cantelli)
at rate $\frac{1}{\sqrt{n}}$ (Donsker)

δ_v means the point mass at v .

Thus we may expect $\hat{\beta}_{\text{plug-in}} = \beta(IP_n)$ to converge to $\beta(P)$ at the same rate if β is "smooth".

- van-Mises expansion.

The definition of Gateaux differentiability implies

$$\beta(P + \varepsilon(Q - P)) - \beta(P) = \varepsilon \beta'_{IP}(Q - P) + o(\varepsilon).$$

Substituting $\varepsilon = \frac{1}{\sqrt{n}}$, $Q - P = \sqrt{n}(IP_n - P)$, we obtain

$$\beta(IP_n) - \beta(P) = \frac{1}{\sqrt{n}} \beta'_{IP}(\sqrt{n}(IP_n - P)) + \underbrace{R(IP_n, P)}_{\text{should be } o_p(\frac{1}{\sqrt{n}})}$$

$$= \frac{1}{n} \sum_{i=1}^n \underbrace{\beta'_{IP}(\delta_{V_i} - P)}_{\text{"Influence" of } V_i} + R(IP_n, P).$$

Slide 68 (Influence function)

Define $\phi_{IP}(v) = \beta'_{IP}(\delta_v - P)$.

• Properties : ① $\mathbb{E}_{\mathbb{P}}(\phi_{\mathbb{P}}(V)) = \mathbb{E}_{\mathbb{P}}\{\beta'_{\mathbb{P}}(\delta_V - \mathbb{P})\}$
 $= \beta'_{\mathbb{P}}\{\mathbb{E}_{\mathbb{P}}(\delta_V - \mathbb{P})\}$
 $= \beta'_{\mathbb{P}}(0) = 0.$

② $\mathbb{E}_{\mathbb{Q}}(\phi_{\mathbb{P}}(V)) = \mathbb{E}_{\mathbb{Q}}\{\beta'_{\mathbb{P}}(\delta_V - \mathbb{P})\}$
 $= \beta'_{\mathbb{P}}\{\mathbb{E}_{\mathbb{Q}}(\delta_V - \mathbb{P})\}$
 $= \beta'_{\mathbb{P}}(\mathbb{Q} - \mathbb{P})$
 $= \frac{\partial}{\partial \varepsilon} \beta(\mathbb{P}_{\varepsilon}) \Big|_{\varepsilon=0}.$

In other words, $\phi_{\mathbb{P}}$ is the Riesz representation of $\beta'_{\mathbb{P}}$.

• Examples :

① Mean : $\beta(\mathbb{P}) = \mathbb{E}_{\mathbb{P}}(V).$

$$\hat{\beta}_{\text{plug-in}} = \beta(\mathbb{P}_n) = \mathbb{E}_{\mathbb{P}_n}(V) = \frac{1}{n} \sum_{i=1}^n V_i.$$

Gâteaux : $\beta'_{\mathbb{P}}(\mathbb{Q} - \mathbb{P}) = \mathbb{E}_{\mathbb{Q}}(V) - \mathbb{E}_{\mathbb{P}}(V).$

IF : $\phi_{\mathbb{P}}(V) = V - \mathbb{E}_{\mathbb{P}}(V).$

$R(\mathbb{P}_n, \mathbb{P}) = 0.$ Not surprising because β is linear

② Z-estimation. $\beta = \beta(\mathbb{P})$ defined by

$$\mathbb{E}_{\mathbb{P}}(m(\beta; V)) = 0.$$

Plug-in estimator solves

$$\frac{1}{n} \sum_{i=1}^n m(\beta; V_i) = 0$$

IF: " is
$$\phi_{\mathbb{P}}(V) = \mathbb{E} \left(\left\{ \frac{\partial}{\partial \beta} m(\beta; V) \right\}^{-1} \right) m(\beta; V).$$

Slide 69 (Bias correction using IF)

• Heuristics for $\hat{\beta}_{1\text{-step}}$.

Let's restate the von Mises expansion:

$$\beta(\mathbb{Q}) - \beta(\mathbb{P}) = \mathbb{E}_{\mathbb{Q}}(\phi_{\mathbb{P}}(V)) + R(\mathbb{Q}, \mathbb{P}) \quad (*)$$

If \mathbb{Q} is "close" to \mathbb{P} , consider

$$\tilde{\mathbb{P}}_{\varepsilon} = \mathbb{P} + \varepsilon \frac{\mathbb{Q} - \mathbb{P}}{\|\mathbb{Q} - \mathbb{P}\|}$$

Then we may expect

$$|R(\mathbb{Q}, \mathbb{P})| \leq O(\varepsilon^2) = O(\|\mathbb{Q} - \mathbb{P}\|^2)$$

Now setting $\mathbb{Q} = \mathbb{P}$ and $\mathbb{P} = \hat{\mathbb{P}}$ in (4),

$$\beta(\mathbb{P}) - \beta(\hat{\mathbb{P}}) = \underbrace{\mathbb{E}_{\mathbb{P}}(\phi_{\hat{\mathbb{P}}}(V))}_{\text{- bias}} + \underbrace{R(\mathbb{P}, \hat{\mathbb{P}})}_{\text{remainder.}}$$

We may thus improve $\beta(\hat{\mathbb{P}})$ by

$$\hat{\beta}_{1\text{-step}} = \beta(\hat{\mathbb{P}}) + \mathbb{E}_{\mathbb{P}_n}(\phi_{\hat{\mathbb{P}}}(V)).$$

This is similar to take one Newton-Raphson step at $\hat{\mathbb{P}}$.

- Expansion

$$\begin{aligned} \hat{\beta}_{1\text{-step}} - \beta(\mathbb{P}) &= \mathbb{E}_{\mathbb{P}_n}(\phi_{\hat{\mathbb{P}}}(V)) - \mathbb{E}_{\mathbb{P}}(\phi_{\hat{\mathbb{P}}}(V)) - R(\mathbb{P}, \hat{\mathbb{P}}) \\ &= \underbrace{\mathbb{E}_{\mathbb{P}_n - \mathbb{P}}(\phi_{\hat{\mathbb{P}}}(V))}_{\text{CLT}} + \underbrace{\mathbb{E}_{\mathbb{P}_n - \mathbb{P}}(\phi_{\hat{\mathbb{P}}}(V) - \phi_{\mathbb{P}}(V))}_{\text{Empirical process}} - \underbrace{R(\mathbb{P}, \hat{\mathbb{P}})}_{\text{second-order}} \end{aligned}$$

where $\mathbb{E}_{\mathbb{Q} - \mathbb{P}}(\cdot) = \mathbb{E}_{\mathbb{Q}}(\cdot) - \mathbb{E}_{\mathbb{P}}(\cdot)$.

- Cross-fitting: the empirical process term shows up because

$\hat{\mathbb{P}}$ depends on \mathbb{P}_n . This prevents us from claiming

$$\mathbb{E}_{\mathbb{P}_n - \mathbb{P}}(\phi_{\hat{\mathbb{P}}}(V) - \phi_{\mathbb{P}}(V)) = O\left(\frac{1}{\sqrt{n}} \|\phi_{\hat{\mathbb{P}}} - \phi_{\mathbb{P}}\|\right).$$

$$\text{Let } P_n^{(1)} = \frac{1}{n/2} \sum_{i=1}^{n/2} \delta_{V_i}, \quad P_n^{(2)} = \frac{1}{n/2} \sum_{i=n/2+1}^n \delta_{V_i}.$$

$$\text{So } P_n = \frac{1}{2} (P_n^{(1)} + P_n^{(2)}).$$

Let $\hat{P}^{(k)}$ be a smooth estimator of P based on $P_n^{(k)}$, $k=1, 2$. Denote $P_n^{(-1)} = P_n^{(2)}$ and $P_n^{(-2)} = P_n^{(1)}$.

$$\text{Define } \hat{\beta}^{(k)} = \beta(\hat{P}^{(-k)}) + \mathbb{E}_{P_n^{(k)}}(\phi_{\hat{P}^{(-k)}}(V)), \quad k=1, 2$$

$$\text{and } \hat{\beta} = \frac{1}{2} (\hat{\beta}^{(1)} + \hat{\beta}^{(2)}).$$

$$\begin{aligned} \text{Then } \hat{\beta} - \beta &= \frac{1}{2} \sum_{k=1}^2 \left\{ \mathbb{E}_{P_n^{(k)} - P}(\phi_P(V)) \right. \\ &\quad + \mathbb{E}_{P_n^{(k)} - P}(\phi_{\hat{P}^{(-k)}}(V) - \phi_P(V)) \\ &\quad \left. + R(P, \hat{P}^{(-k)}) \right\} \\ &= \mathbb{E}_{P_n - P}(\phi_P(V)) + \frac{1}{2} \sum_{k=1}^2 R(P, \hat{P}^{(-k)}) \\ &\quad + \underline{o_p\left(\frac{1}{\sqrt{n}}\right)} \end{aligned}$$

$$\text{If } \phi_{\hat{P}^{(-k)}} \xrightarrow{P} \phi_P, \quad k=1, 2.$$

Slide 70 (Calculus of IF)

Let $V = (X, Y)$. Consider the submodel indexed by ε :

$$P_{\varepsilon}(V=v) = (1-\varepsilon) P_{\varepsilon}(V=v) + \varepsilon \delta_v$$

Then

$$\begin{aligned} \frac{\partial}{\partial \varepsilon} P_{\varepsilon}(X=x) \Big|_{\varepsilon=0} &= \left(\frac{\partial}{\partial \varepsilon} (1-\varepsilon) \cdot P(X=x) + \varepsilon I\{x=\tilde{x}\} \right) \Big|_{\varepsilon=0} \\ &= I\{x=\tilde{x}\} - P(X=x) \end{aligned}$$

So the IF is $\delta_x - P(X=x)$.

$$\begin{aligned} & \frac{\partial}{\partial \varepsilon} E_{\varepsilon}(Y|X=x) \Big|_{\varepsilon=0} \\ &= \frac{\partial}{\partial \varepsilon} \sum_y y \cdot \frac{(1-\varepsilon) \cdot P(V=v) + \varepsilon \cdot I\{v=\tilde{v}\}}{(1-\varepsilon) \cdot P(X=x) + \varepsilon \cdot I\{x=\tilde{x}\}} \Big|_{\varepsilon=0} \\ &= \sum_y y \cdot \frac{(I\{v=\tilde{v}\} - P(V=v)) P(X=x) - P(V=v) \cdot (I\{x=\tilde{x}\} - P(X=x))}{P(X=x)^2} \\ &= \sum_y y \cdot \frac{I\{v=\tilde{v}\} - I\{x=\tilde{x}\} \cdot P(Y=y|X=x)}{P(X=x)} \\ &= \frac{(\tilde{y} - E(Y|X=x)) \cdot I\{x=\tilde{x}\}}{P(X=x)} \end{aligned}$$

So the IF is $\{Y - E(Y|X=x)\} \cdot \frac{\delta_x}{P(X=x)}$

Now consider $\beta = \mathbb{E}(\mathbb{E}(Y|A=1, X))$.

$$= \sum_x \mathbb{E}(Y|A=1, X=x) \cdot P(X=x)$$

$$IC(\beta) = \sum_x IC(\mathbb{E}(Y|A=1, X=x)) \cdot P(X=x)$$

$$+ IC(P(X=x)) \cdot \mathbb{E}(Y|A=1, X=x)$$

$$= \sum_x \frac{\mathbb{1}_{\{X=x, A=1\}}}{P(X=x, A=1)} \cdot (Y - \mathbb{E}(Y|X=x, A=1)) \cdot P(X=x)$$

$$+ (\mathbb{1}_{\{X=x\}} - P(X=x)) \cdot \mathbb{E}(Y|A=1, X=x)$$

$$= \frac{A}{P(A=1|X)} \cdot (Y - \mathbb{E}(Y|X, A=1)) + \mathbb{E}(Y|A=1, X) - \beta$$

$$= \frac{A}{\pi(X)} \cdot (Y - \mu_1(X)) + \mu_1(X) - \beta$$

• AIPW/DR estimator

Initial estimator: $\hat{\beta}_{OR} = \frac{1}{n} \sum_{i=1}^n \hat{\mu}_1(X_i)$

One-step estimator:
(AIPW or DR) $\hat{\beta}_{DR} = \hat{\beta}_{OR} + \frac{1}{n} \sum_{i=1}^n IC(\beta; \pi_i, \hat{\mu}_1)$
 $= \frac{1}{n} \sum_{i=1}^n \frac{A_i}{\pi(X_i)} \cdot (Y_i - \hat{\mu}_1(X_i)) + \hat{\mu}_1(X_i)$

Remainder term:

$$\begin{aligned}R(P, \hat{P}) &= \beta(P) - \beta(\hat{P}) - \mathbb{E}_P(\phi_{\hat{P}}(V)) \\&= \mathbb{E}_P(\mu_1(X)) - \mathbb{E}_P\left(\frac{A}{\pi(X)} \cdot (Y - \hat{\mu}_1(X)) + \hat{\mu}_1(X)\right) \\&= \mathbb{E}_P(\mu_1(X) - \hat{\mu}_1(X)) - \mathbb{E}_P\left(\frac{\pi(X)}{\hat{\pi}(X)} \cdot (\mu_1(X) - \hat{\mu}_1(X))\right) \\&= -\mathbb{E}_P\left(\left(\frac{1}{\hat{\pi}(X)} - \frac{1}{\pi(X)}\right)(\hat{\mu}_1(X) - \mu_1(X)) \cdot \pi(X)\right).\end{aligned}$$

Define: $MSE(\hat{\pi}) = \mathbb{E}_P((\hat{\pi}(X) - \pi(X))^2)$

So the remainder term $\xrightarrow{P} 0$ if $\hat{\pi} \rightarrow \pi$ or $\hat{\mu}_1 \rightarrow \mu_1$,

"double robustness":

_____ is negligible if $MSE(\hat{\pi}) \cdot MSE(\hat{\mu}_1) = o(\frac{1}{n})$.
and $\pi(x) \geq c > 0$ for some c .

This theory can be easily extended to estimating the ATE, ATT...

Slide 72 (Entropy balancing)

• Lagrangian dual

$$L(w, \lambda, \nu) = \sum_{i=n+1}^n w_i \log w_i + \sum_{i=n+1}^n \lambda_i w_i + \nu^T \left\{ \sum_{i=1}^{n_1} X_i - \sum_{i=n+1}^n w_i X_i \right\}$$

$\lambda_i \leq 0.$

$$\frac{\partial}{\partial w_i} L = \log w_i + 1 + \lambda_i - (\nu^T X_i)$$

$$\Rightarrow w_i^* = e^{(\nu^T X_i) - \lambda_i - 1}$$

Lagrangian dual function $g(\lambda, \nu) = \inf_w L(w, \lambda, \nu)$

$$= L(w^*, \lambda, \nu)$$

$$\begin{aligned} \frac{\partial}{\partial \lambda_i} g &= (\log w_i^* + 1) \cdot \frac{\partial}{\partial \lambda_i} w_i^* + w_i^* + (\lambda_i - \nu^T X_i) \cdot \frac{\partial}{\partial \lambda_i} w_i^* \\ &= -w_i^* (\log w_i^* + \lambda_i - \nu^T X_i) \\ &= w_i^* > 0. \end{aligned}$$

So optimal dual $\lambda_i^* = 0$.

Lagrangian dual optimization: Denote $\pi_i = \frac{e^{\nu^T X_i - 1}}{1 + e^{\nu^T X_i - 1}}$

$$\begin{aligned} &\sup_{\nu} L(w^*, \lambda^*, \nu) \\ &= \sup_{\nu} \left\{ \sum_{i=1}^n \frac{\pi_i}{1 - \pi_i} (\nu_i^T X_i - 1) + \nu^T \sum_{i=1}^n X_i - \sum_{i=1}^n \frac{\pi_i}{1 - \pi_i} X_i \right\} \\ &= \sup_{\nu} \left\{ - \sum_{i=1}^n \frac{\pi_i}{1 - \pi_i} + \nu^T \sum_{i=1}^n X_i \right\} \\ &= \inf_{\nu} \sum_{i=1}^n (1 - A_i) \frac{\pi_i}{1 - \pi_i} - A_i \log \frac{\pi_i}{1 - \pi_i}. \end{aligned}$$

Compare to the maximum likelihood problem:

$$\max_{\nu} \sum_{i=1}^n [(1 - A_i) \log(1 - \pi_i) + A_i \log \pi_i]$$

Thus, the dual problem is fitting a logistic regression with a different loss function.

- Double robustness

A benefit of this is that when we do the one-step correction to $\hat{\beta} = \sum_{i=1}^n (1-A_i) w_i Y_i$, we get

$$\begin{aligned} \hat{\beta}_{DR} &= \frac{1}{n} \sum_{i=1}^n \cancel{(1-A_i)} w_i Y_i \\ &\quad + \frac{1}{n} \sum_{i=1}^n (1-A_i) w_i (Y_i - \hat{\mu}_0(X_i)) + A_i \hat{\mu}_0(X_i) - \frac{1}{n} \sum_{i=1}^n \cancel{(1-A_i)} w_i Y_i \\ &= \frac{1}{n} \sum_{i=1}^n (1-A_i) w_i Y_i + \underbrace{\frac{1}{n} \sum_{i=1}^n A_i \hat{\mu}_0(X_i) - \frac{1}{n} \sum_{i=1}^n \cancel{(1-A_i)} w_i Y_i}_{=0 \text{ if } \hat{\mu}_0(X_i) \text{ is linear in } X_i} \\ &= \hat{\beta} \end{aligned}$$

Slide 75 (Linear SEMs: Multiple IVs)

- Linear SEM: Assume all variables have mean 0.

$$A = Z^T \beta_{ZA} + X \beta_{XA} + \varepsilon_A$$

$$Y = A \beta_{AY} + X \beta_{XY} + \varepsilon_Y$$

$\varepsilon_A, \varepsilon_Y$ may be correlated. $(X, Z) \perp (\varepsilon_A, \varepsilon_Y)$.

$$\text{Then } \mathbb{E}(Y|Z, X) = \beta_{AY} \mathbb{E}(A|Z, X) + X\beta_{XY}$$

• This motivates the two-stage least squares estimator:

1. Regress A on Z, X to obtain $\hat{\mathbb{E}}(A|Z, X)$.

2. Regress Y on $\hat{\mathbb{E}}(A|Z, X)$.

Slide 77 (Semiparametric estimation with IV)

• Plug-in estimator: Let $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$ and $\bar{A} = \frac{1}{n} \sum_{i=1}^n A_i$.

So $\hat{\alpha}(\beta) = \bar{Y} - \beta \bar{A}$. By solving the empirical estimating equation, we obtain

$$\hat{\beta}_g = \frac{\frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y}) g(Z_i)}{\frac{1}{n} \sum_{i=1}^n (A_i - \bar{A}) g(Z_i)} \rightarrow \frac{\text{Cov}(Y, g(Z))}{\text{Cov}(A, g(Z))} = \beta.$$

For fixed $g(\cdot)$ and under regularity conditions, $\hat{\beta}_g$ is asymptotically linear with the following IF:

$$\boxed{\text{ES}} \quad \psi_{g|Z, A, Y} = \frac{(\{Y - \mathbb{E}(Y)\} - \beta\{A - \mathbb{E}(A)\}) [g(Z) - \mathbb{E}\{g(Z)\}]}{\text{Cov}(A, g(Z))}.$$

So $\sqrt{n}(\hat{\beta}_g - \beta) \xrightarrow{d} N(0, \sigma_g^2)$ as $n \rightarrow \infty$.

$$\text{where } \sigma_g^2 = \frac{\text{Var}(Y - \beta A) \cdot \text{Var}(g(Z))}{\text{Cov}(A, g(Z))^2}$$

By Cauchy-Schwarz inequality, this is minimised at

$$g(Z) = g^*(Z) = \mathbb{E}(A|Z). \quad \text{"optimal instrument"}$$

To improve efficiency, it is common to estimate $g^*(\cdot)$ first and then plug it in $\hat{\beta}_g$. Let the resulting estimator be

$\hat{\beta}_{\hat{g}}$. This is essentially the two-stage least squares.

Rem 1. As long as \hat{g} is not too complex, empirical process theory suggests $\hat{\beta}_{\hat{g}} - \hat{\beta}_g$ is negligible (g is the limit of \hat{g}), so $\hat{\beta}_{\hat{g}}$ is still asymptotically normal.

2. A robustness property: $\hat{\beta}_g$ is \sqrt{n} -consistent as long as $\text{Cov}(A, g(Z)) > 0$. Do not need a consistent estimator of $g^*(\cdot)$. This is similar to regression adjustment methods for randomised experiments.

Slide 79 (IV identification: Complier average causal effect)

$$\text{PF} \quad \beta = \frac{\text{Cov}(Z, Y)}{\text{Cov}(Z, A)} = \frac{\mathbb{E}(Y|Z=1) - \mathbb{E}(Y|Z=0)}{\mathbb{E}(A|Z=1) - \mathbb{E}(A|Z=0)}.$$

By assumption, $Z \perp\!\!\!\perp C$ and $P(C=de) = 0$.

$$\begin{aligned}
\mathbb{E}(Y|Z=1) &= \sum_{C \in \{at, nt, \omega, de\}} \mathbb{E}\{Y|Z=1, C=c\} P(C=c|Z=1) \\
&= \mathbb{E}\{Y(1)|C=at\} P(C=at) \\
&\quad + \mathbb{E}\{Y(0)|C=nt\} P(C=nt) \\
&\quad + \mathbb{E}\{Y(1)|C=\omega\} P(C=\omega) \\
&\quad + \mathbb{E}\{Y(0)|C=de\} P(C=de)
\end{aligned}$$

Similarly $\mathbb{E}(Y|Z=0) = \mathbb{E}\{Y(1)|C=at\} P(C=at)$
 $+ \mathbb{E}\{Y(0)|C=nt\} P(C=nt)$
 $+ \mathbb{E}\{Y(0)|C=\omega\} P(C=\omega)$
 $+ \mathbb{E}\{Y(1)|C=de\} P(C=de)$

So $\mathbb{E}(Y|Z=1) - \mathbb{E}(Y|Z=0) = \mathbb{E}\{Y(1) - Y(0)|C=\omega\} P(C=\omega)$.
Similarly, $\mathbb{E}(A|Z=1) - \mathbb{E}(A|Z=0) = P(C=\omega)$.

The desired conclusion then follows. □