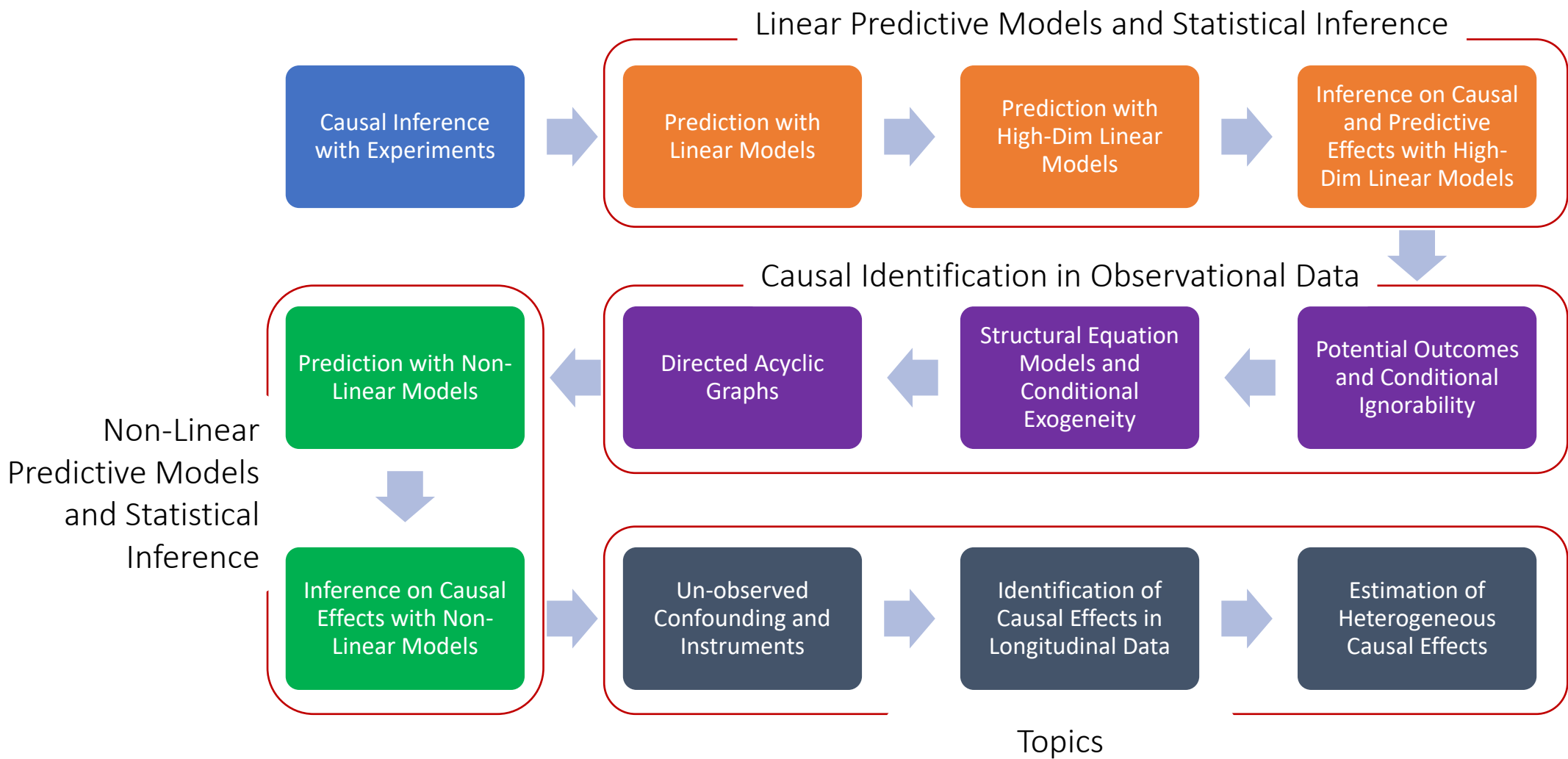
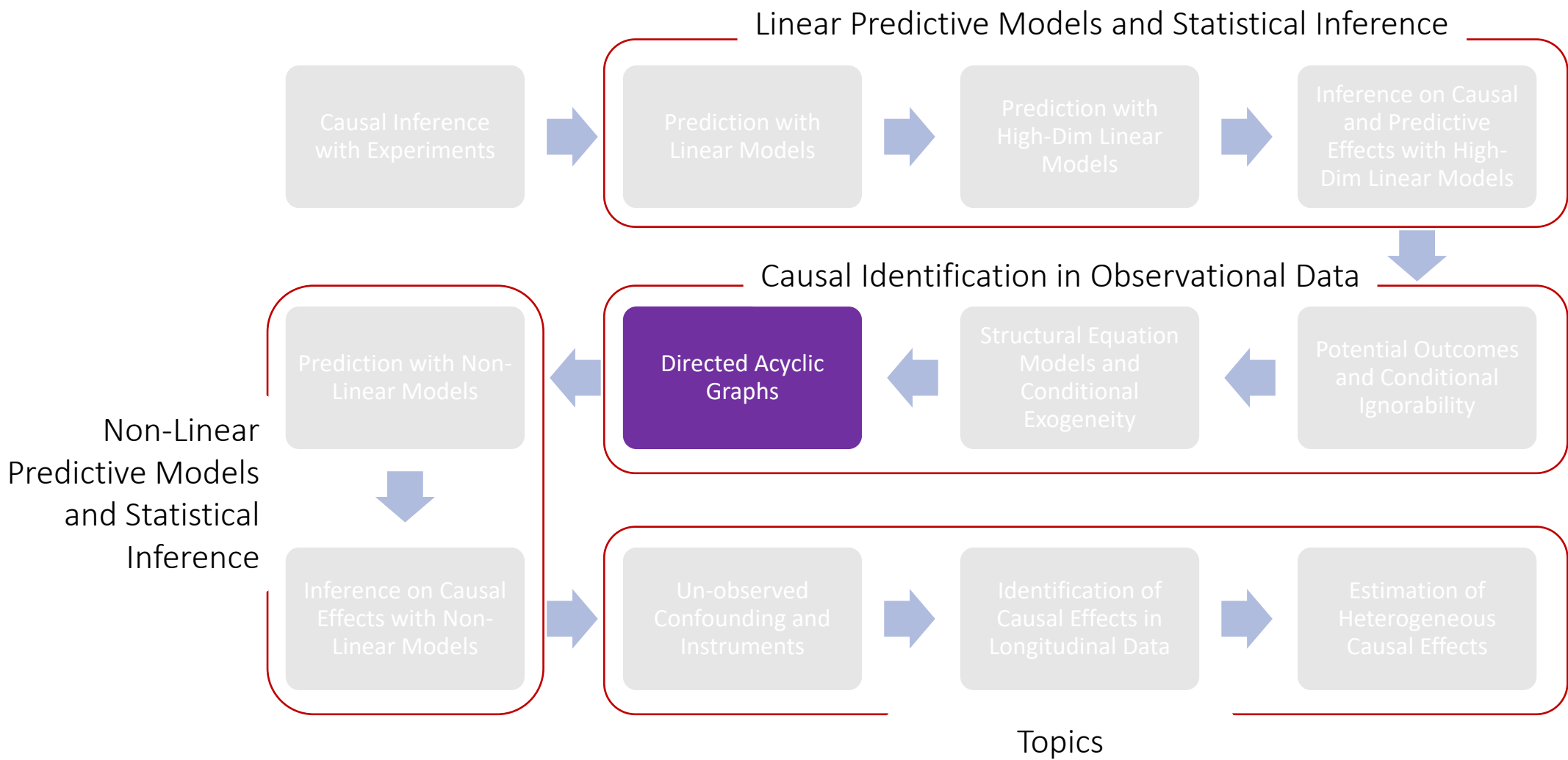


MS&E 228: Directed Acyclic Graphs and Non-Linear SEMs

Vasilis Syrgkanis

MS&E, Stanford





Goals for Today

- Learn the “language” of Directed Acyclic Graphs (DAGs) and their associated non-linear structural equation models (SEMs)
- Introduce “intervention” concepts “do” and “fix”
- Introduce d-separation and conditional independence in DAGs
- Proof sketch of fundamental theorem d-separation \Rightarrow conditional ind.

Next lecture

- Graphical criteria for selection of adjustment set
- Crash course on good and bad “controls”

DAGs

Judea Pearl. 'Causal diagrams for empirical research'. In: *Biometrika* 82.4 (1995), pp. 669–688 (cited on page 30).

Trygve Haavelmo. 'The probability approach in econometrics'. In: *Econometrica: Journal of the Econometric Society* 12 (1944), pp. iii–vi+1–115 (cited on pages 30, 32).

James Heckman and Rodrigo Pinto. 'Causal analysis after Haavelmo'. In: *Econometric Theory* 31.1 (2015 (NBER 2013)), pp. 115–151 (cited on pages 30, 35).

James Robins. 'A new approach to causal inference in mortality studies with a sustained exposure period—application to control of the healthy worker survivor effect'. In: *Mathematical modelling* 7.9-12 (1986), pp. 1393–1512 (cited on page 54).

Thomas S. Richardson and James M. Robins. *Single world intervention graphs (SWIGs): A unification of the counterfactual and graphical approaches to causality*. Working Paper

Non-Linear SEMs

Non-Linear SEMs

- Last time we looked at linear SEMs
- The language of SEMs does not really rely on the linearity assumption
- For example, the Triangular Structural Equation (TSEM)

$$Y := \delta P + \beta' X + \epsilon_Y$$

$$P := \nu' X + \epsilon_P$$

$$X := \epsilon_X$$

Non-Linear SEMs

- Last time we looked at linear SEMs
- The language of SEMs does not really rely on the linearity assumption
- For example, the Triangular Structural Equation (TSEM)
- Can be made non-linear

$$Y := f_Y(P, X, \epsilon_Y)$$

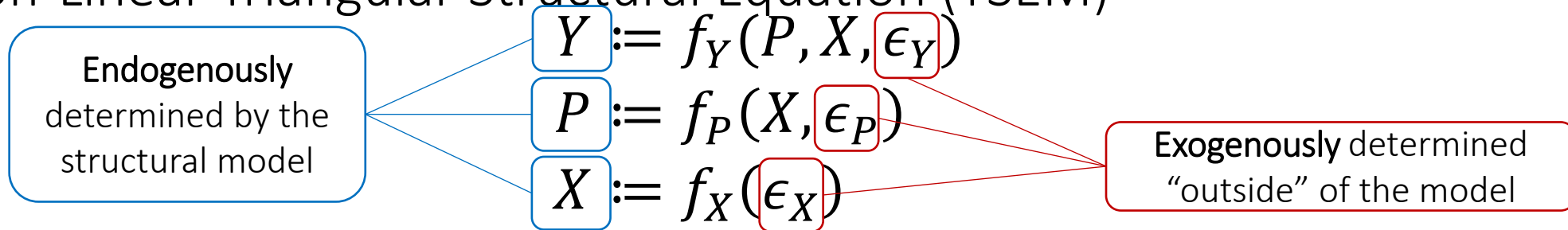
$$P := f_P(X, \epsilon_P)$$

$$X := f_X(\epsilon_X)$$

- While we still maintain that $\epsilon_Y, \epsilon_P, \epsilon_X$ are independent “exogenous” shocks

Non-Linear SEMs

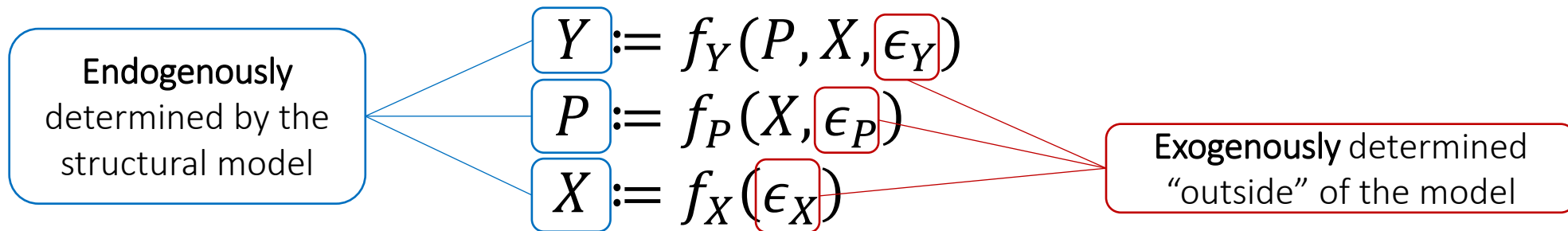
- Non-Linear Triangular Structural Equation (TSEM)



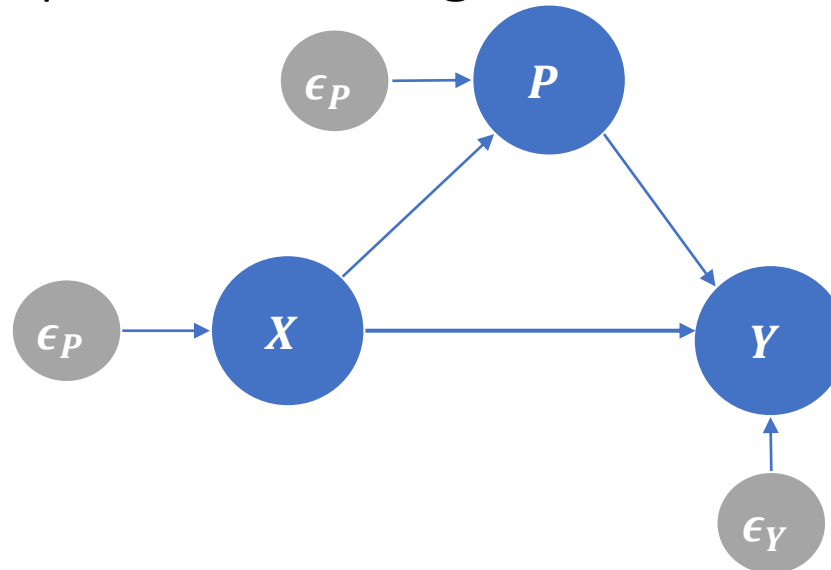
- $\epsilon_Y, \epsilon_P, \epsilon_X$ are independent “exogenous” shocks
- The functions f_Y, f_P, f_X are deterministic “structural functions”
- Instead of “structural parameters” we now have “structural functions”
- Moreover, the dimension of exogenous shocks is un-restricted
- Note that the TSEM implies: $\epsilon_Y \perp\!\!\!\perp P, X$ and $\epsilon_P \perp\!\!\!\perp X$

Non-Linear SEMs

- Non-Linear Triangular Structural Equation (TSEM)

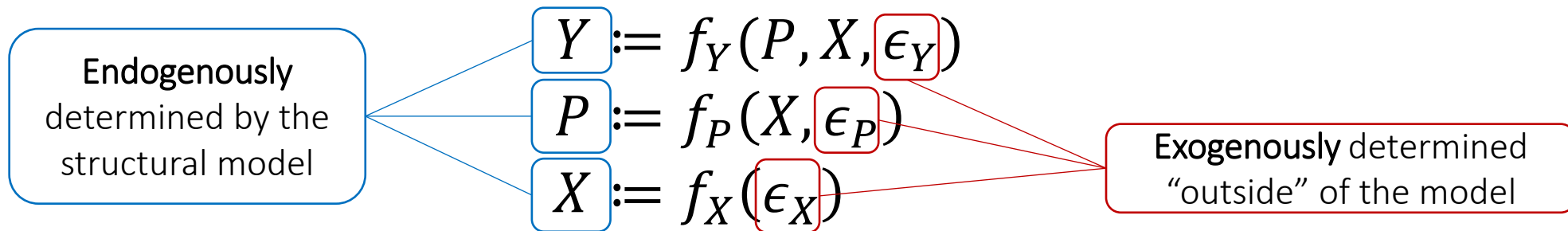


- $\epsilon_Y, \epsilon_P, \epsilon_X$ are independent "exogenous" shocks

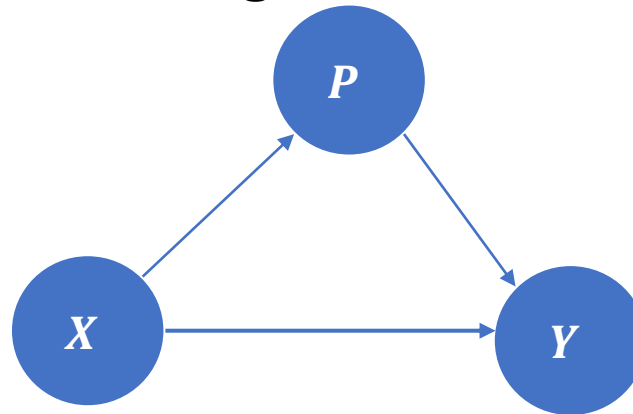


Non-Linear SEMs

- Non-Linear Triangular Structural Equation (TSEM)

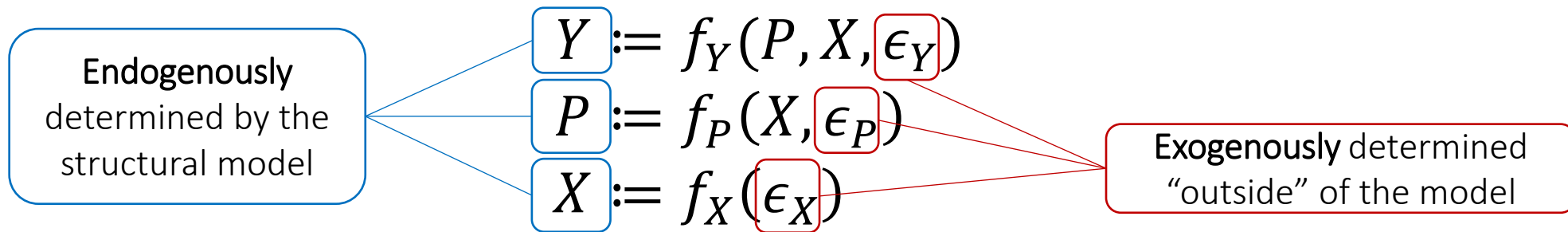


- $\epsilon_Y, \epsilon_P, \epsilon_X$ are independent “exogenous” shocks; typically omitted from DAG visualization



Non-Linear SEMs

- Non-Linear Triangular Structural Equation (TSEM)



- $\epsilon_Y, \epsilon_P, \epsilon_X$ are independent “exogenous” shocks; typically omitted from DAG visualization
- A TSEM is simply a statistical “generative” model that determines a distribution over observed random variables (*c.f. Neural-Causal Models in further reading)

Structural Form

- TSEM is “structural” in that it is endowed with the following properties
- Made up of a collection of stochastic potential outcome processes indexed by (p, x)

$$Y(p, x) := f_Y(p, x, \epsilon_Y)$$

$$P(x) := f_P(x, \epsilon_P)$$

$$X := f_X(\epsilon_X)$$

- **Exogeneity:** $\epsilon_P, \epsilon_Y, \epsilon_X$ are independent “shock” variables generated outside of the model
- **Consistency:** endogenous variables (Y, P, X) generated by recursive substitutions
$$Y := Y(P, X), \quad P := P(X), \quad X := \epsilon_X$$
- **Invariance:** structure remains invariant to changes of distributions of shocks

Link to Potential Outcomes

- Define arbitrary random potential outcomes

$$Y(p, x), P(x), X(.)$$

- Which are independent random processes

$$\{Y(p, x)\}_{p, x} \perp \{P(x)\}_x \perp X(.)$$

- Each potential outcome process can be represented as a SEM

- We can define $\epsilon_Y := \epsilon_Y(p, x) = Y(p, x)$ and

$$f_Y(p, x, \epsilon_Y) = \epsilon_Y(p, x) = Y(p, x)$$

The Language of Interventions and Intervention Counterfactuals

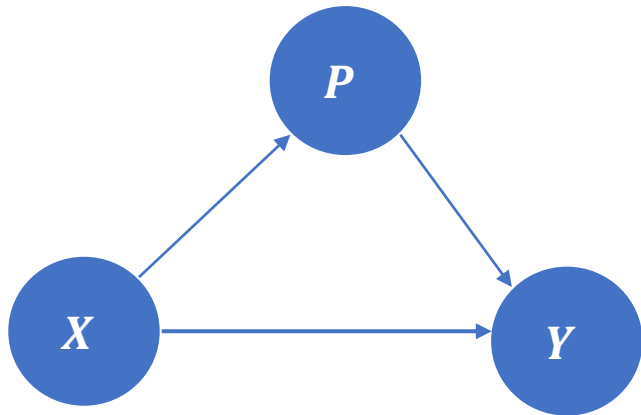
Do Interventions: $do(P = p)$

Original Data Generative Model

$$Y := f_Y(P, X, \epsilon_Y)$$

$$P := f_P(X, \epsilon_P)$$

$$X := \epsilon_X$$

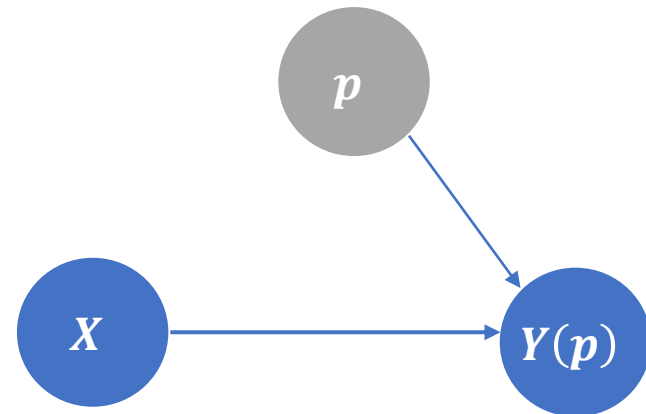


Data Generative Model under $do(P = p)$

$$Y \mid \begin{matrix} P \\ X \end{matrix} \quad do(P = p) := f_Y(p, X, \epsilon_Y)$$

$$P \mid \begin{matrix} Y \\ X \end{matrix} \quad do(P = p) := p$$

$$X \mid \begin{matrix} Y \\ P \end{matrix} \quad do(P = p) := \epsilon_X$$



Interventions

- Do-interventions is only one way of defining counterfactuals
- We can define any type of counterfactual by simply changing one of the equations to something else
- Wright in his seminal work in '28 defined an intervention where the demand equation was replaced by another one that reflects a tax hike
- We can also define “soft-interventions”: increase price by 10% of its current value
- Another useful variant of do-interventions does not replace the treatment equation are “fix” interventions

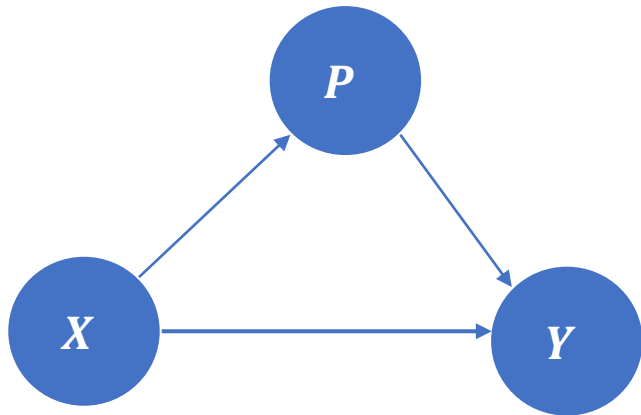
Fix Interventions: $\text{fix}(P = p)$

Original Data Generative Model

$$Y := f_Y(P, X, \epsilon_Y)$$

$$P := f_P(X, \epsilon_P)$$

$$X := \epsilon_X$$

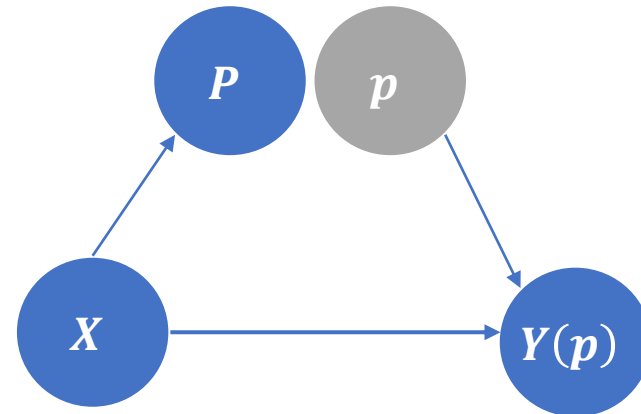


Data Generative Model under $\text{fix}(P = p)$

$$Y \mid \text{fix}(P = p) := f_Y(p, X, \epsilon_Y)$$

$$P \mid \text{fix}(P = p) := f_P(X, \epsilon_P)$$

$$X \mid \text{fix}(P = p) := \epsilon_X$$



Fix Interventions

- A fix intervention is a form of “localized” do intervention
- We are only fixing the value of P in the structural equation for Y
- The random variables generated by the fix intervention are the triplets
 $(Y(p), P, X)$
- The intervention does not affect the P, X equations nor the distribution of the exogenous shock ϵ_Y in the outcome equation

Conditional Ignorability and Mean Intervention Counterfactuals

- Identification by conditioning and conditional ignorability extends to average intervention counterfactuals

- Recall to do identification by conditioning we need for a set S

$$Y(p) \perp\!\!\!\perp P \mid S$$

- Then predictive response equals interventional response

$$E[Y(p) \mid S] = E[Y(p) \mid P = p, S] = E[Y(P) \mid P = p, S] = E[Y \mid P = p, S]$$

- Average predictive response equals average interventional response

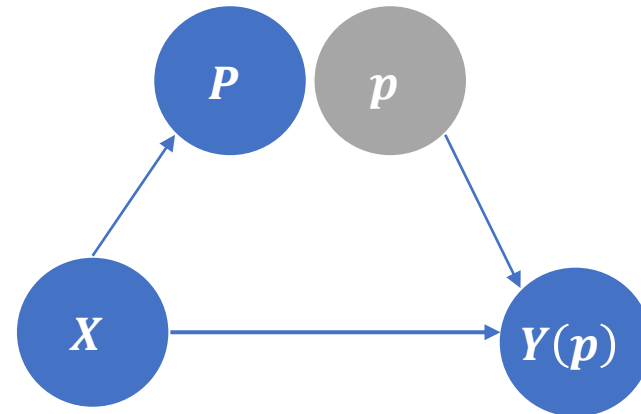
$$E[Y(p)] = E[E[Y \mid P = p, S]]$$

Single World Intervention Graphs

- The graphs that represent the generative model under a fix intervention
- Easy to verify visually that
$$Y(p) \perp\!\!\!\perp P \mid X$$
- Then we can do identification based on conditional ignorability

Data Generative Model under $\text{fix}(P = p)$

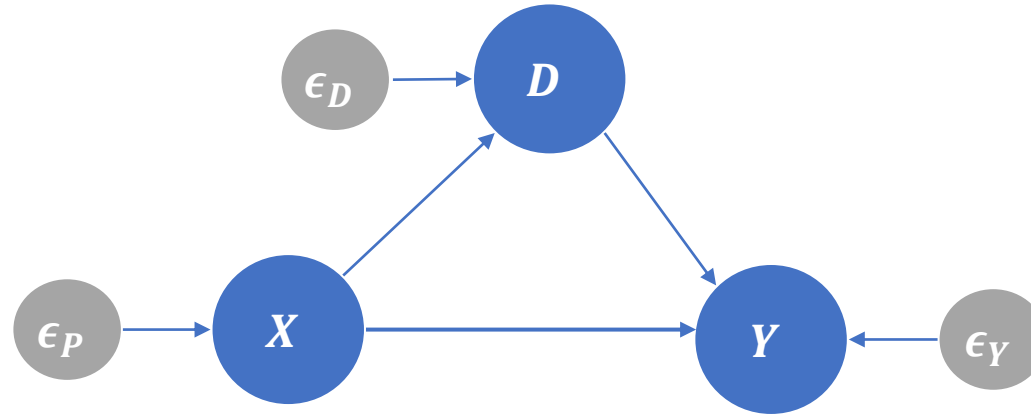
$$\begin{array}{l|l} Y & := f_Y(p, X, \epsilon_Y) \\ P & \text{fix}(P = p) \quad f_P(X, \epsilon_P) \\ X & := \epsilon_X \end{array}$$



Single World Intervention Graph



Non-Linear versions of structural equation models are equivalent to Directed Acyclic Graphs



Exogenously determined
“outside” of the model

Endogenously
determined by
the structural
model

For any DAG, we can write ASEM

$$X_j := f_j(\text{Parents}_j, \epsilon_j) = f_j(\text{Pa}_j, \epsilon_j)$$

Shocks ϵ_j are jointly independent and independent of $\{X_j\}$



Corresponding structural response functions

$$X_j(pa_j) := f_j(pa_j, \epsilon_j)$$

Shocks can be multi-dimensional
e.g. separate shock variable per
parental value

Potential/Counterfactual
Outcome Processes

Structural Response
Function

Potential values of
parents

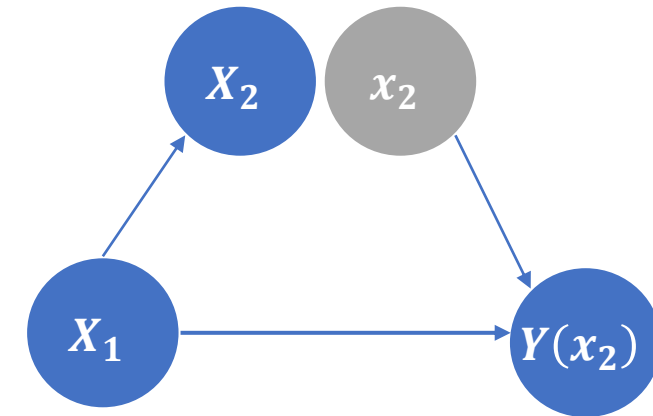
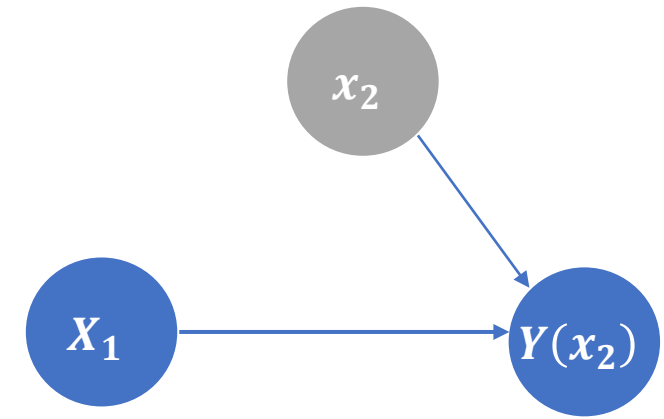
Fix Interventions $\text{fix}(X_j = x_j)$

Locally replace X_j in every RHS of a structural equation with x_j . Leave as-is structural response of X_j . Also measures potential outcome $Y(x_j)$

Fix intervention visually represented as **SWIG** $\tilde{G}(x_j)$. Depicts potential outcome $Y(x_j)$ and original variable X_j on the same graph



If we can check $Y(x_j) \perp X_j \mid S$ based on the SWIG, we can identify $E[Y(x_j)]$ via conditioning



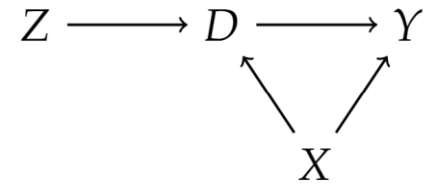
Graphical Criteria for Conditional Independence

DAGs Encode Factorization of Probability

- Graph implies factorization of the probability law
$$p(y, d, x, z) = p(y|x, d)p(d|x, z)p(x)p(z)$$

Proof.

- By Bayes rule: $p(y, d, x, z) = p(y|d, x, z)p(d, x, z)$
- From graph: $p(y|d, x, z) = p(y|d, x)$
- By Bayes rule: $p(d, x, z) = p(d|x, z)p(x, z)$
- By Bayes rule: $p(x, z) = p(x|z)p(z)$
- From graph: $p(x|z) = p(x)$



General DAGs and Factorization

- The probability law factorizes as:

$$p(\{x_\ell\}_{\ell \in V}) = \prod_{\ell \in V} p(x_\ell | pa_\ell)$$

DAGs Encode Conditional Independencies

- Any two variables X, Y are independent conditional on a set S if they are D(irected)-separated in the graph

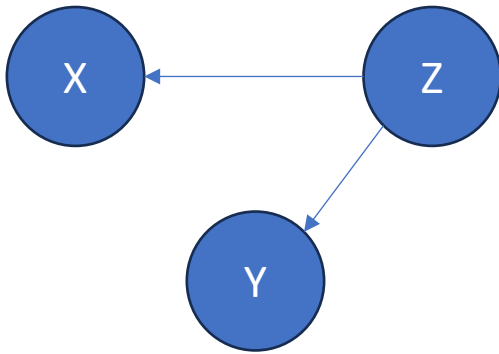
$$(X \perp\!\!\!\perp_d Y \mid S)_G \Rightarrow X \perp\!\!\!\perp Y \mid S$$

- Need to define the concept of D-separation

Graph Separation and Conditional Independence

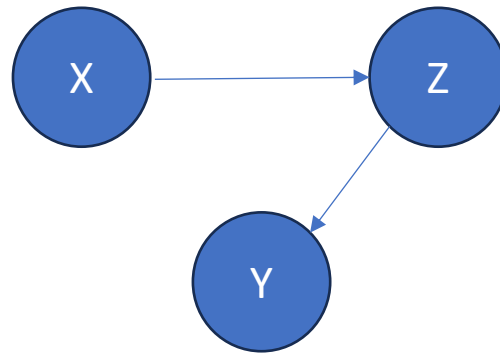
- Looking at a graph, when can we conclude that $X \perp\!\!\!\perp Y \mid Z$

Case 1:
Z is common cause



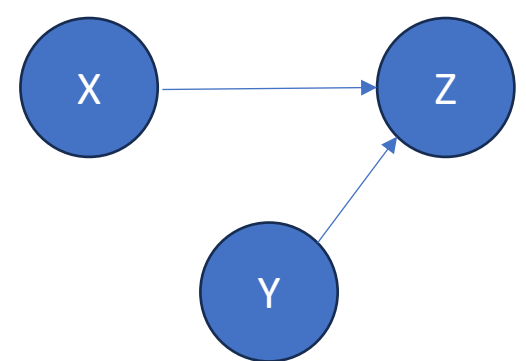
$$\begin{aligned} p(x, y, z) &= p(x \mid z) \cdot p(y \mid z) \cdot p(z) \\ &= f(x, z) \cdot g(y, z) \end{aligned}$$

Case 2:
Z is mediator



$$\begin{aligned} p(x, y, z) &= p(z \mid x) \cdot p(x) \cdot p(y \mid z) \\ &= f(x, z) \cdot g(y, z) \end{aligned}$$

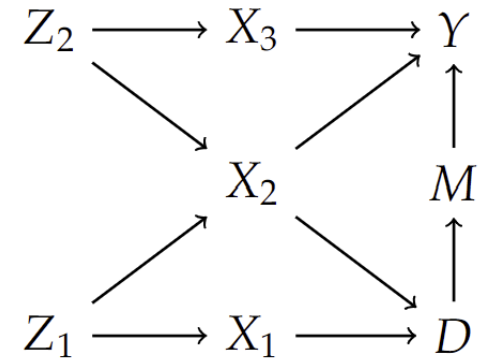
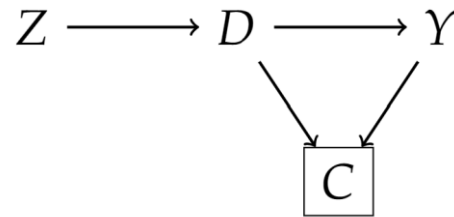
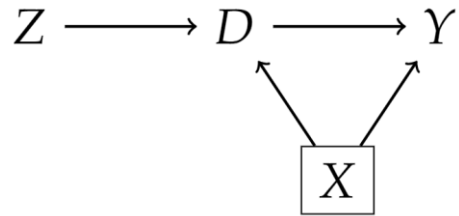
Case 3:
Z is collider



$$\begin{aligned} p(x, y, z) &= p(z \mid x, y) \cdot p(x) \cdot p(y) \\ &\neq f(x, z) \cdot g(y, z) \end{aligned}$$

Some Graph Definitions

- A path π in a graph is blocked by a set of nodes S if
 - Either π contains a chain $i \rightarrow m \rightarrow j$ or a fork $i \leftarrow m \rightarrow j$ and $m \in S$
 - Or π contains a collider $i \rightarrow m \leftarrow j$ and neither m nor its descendants are in S



D-Separation

- In a DAG G , two nodes X, Y are D-separated by a set of nodes S if S blocks all paths between X and Y
- We denote it as:

$$(X \perp\!\!\!\perp_d Y \mid S)_G$$

D-separation
implies
conditional
independency

$$(X \perp\!\!\!\perp_d Y \mid S)_G \Rightarrow X \perp\!\!\!\perp Y \mid S,$$

(Verma, Pearl, '88)



DAGs encode conditional independencies: S d-separates X from Y in DAG G implies $X \perp\!\!\!\perp Y \mid S$

$$(X \perp\!\!\!\perp_d Y \mid S)_G \Rightarrow X \perp\!\!\!\perp Y \mid S$$



Implies testable restrictions we can use to refute DAG from data;
e.g. for linear ASEM, BLP of Y using X, S should have zero on X

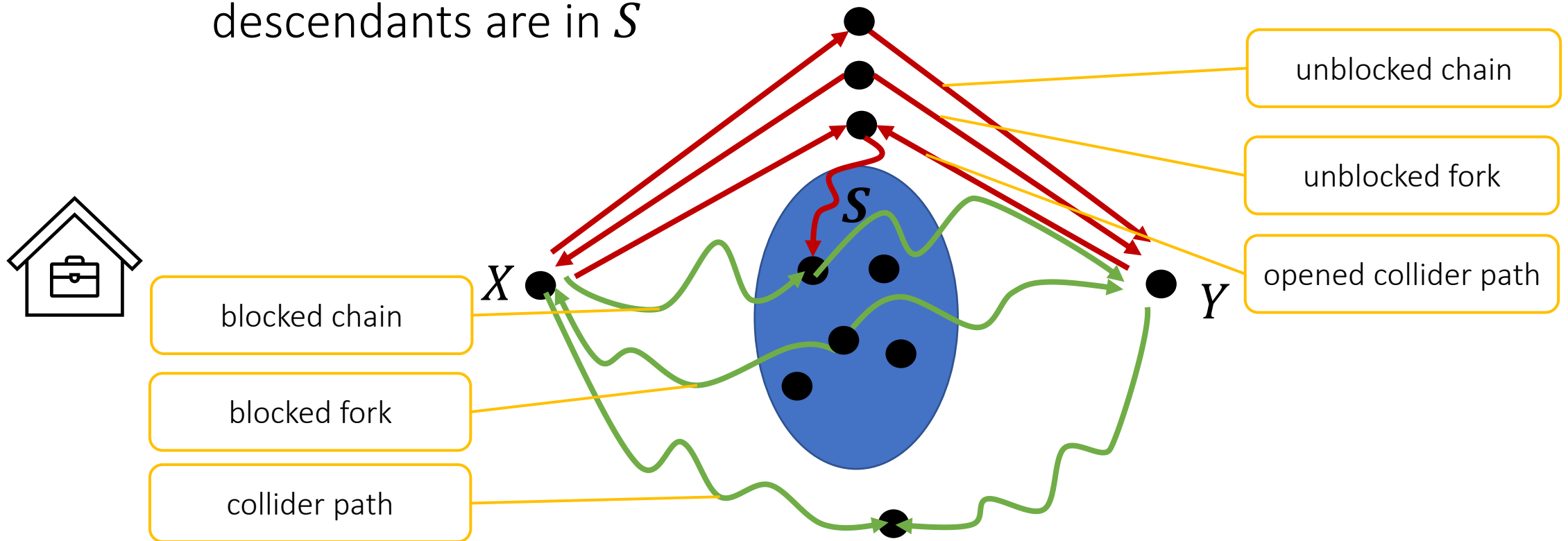
$$Y = \alpha X + \beta' S + \epsilon, \quad \epsilon \perp (X, S)$$

Test whether it
is non-zero!

X is d-separated from Y by S if **every path** from X to Y is **blocked**.

S **blocks a path** if one of the following holds:

- path contains chain $X \rightarrow M \rightarrow Y$ or fork $X \leftarrow M \rightarrow Y$ and $M \in S$
- path contains collider $X \rightarrow M \leftarrow Y$ and neither M nor its descendants are in S



Graphical Criteria for Valid Adjustment Sets

Conditional Ignorability

- Recall to do identification by conditioning we need for a set S

$$Y(d) \perp\!\!\!\perp D \mid S$$

- Then predictive response equals structural response

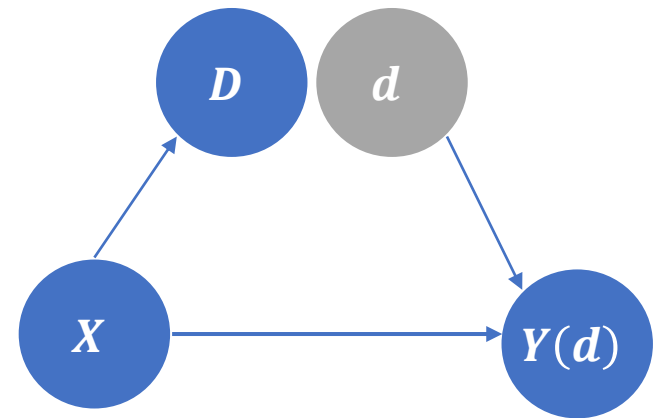
$$E[Y(d) \mid S] = E[Y \mid D = d, S]$$

- Average predictive response equals average structural response

$$E[Y(d)] = E[E[Y \mid D = d, S]]$$

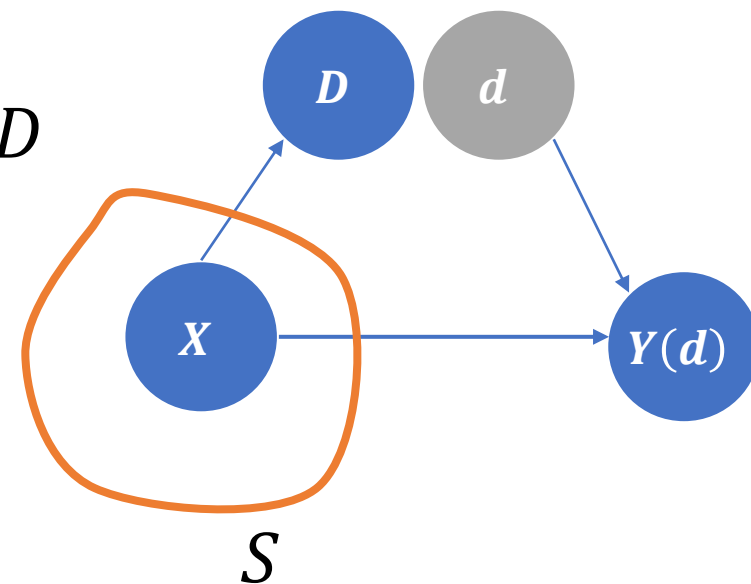
How can we check Conditional Ignorability

- Given a DAG, can we visually inspect if conditional ignorability holds
- Note that the SWIG graph contains both $Y(d)$ and D !
- We can simply check if $Y(d)$ is independent of D conditional on S on the SWIG graph!
- This is just a conditional independence statement on a DAG
- We can use d-separation!

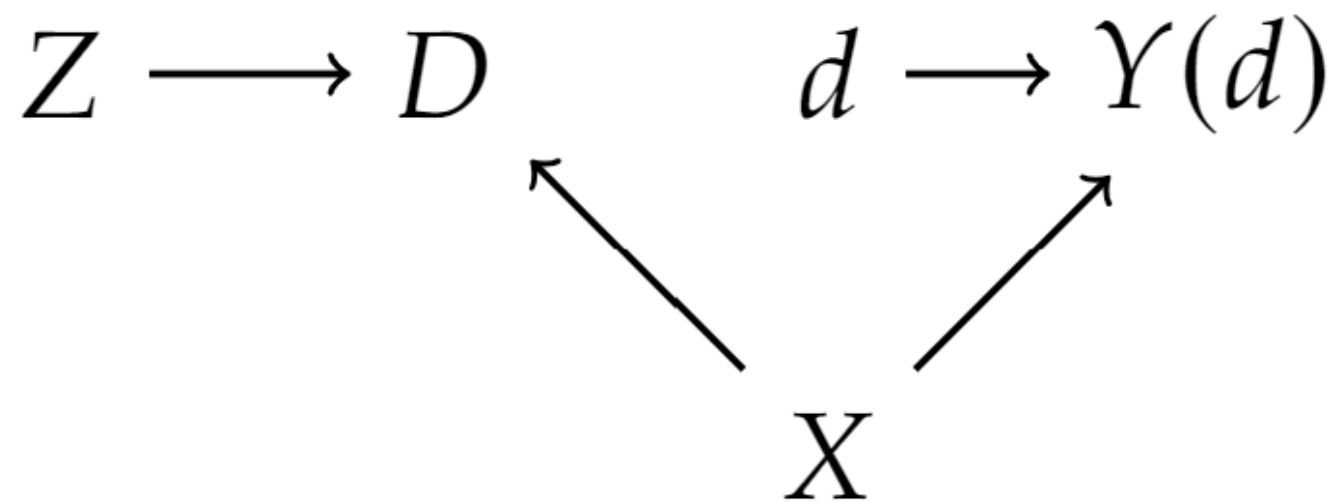




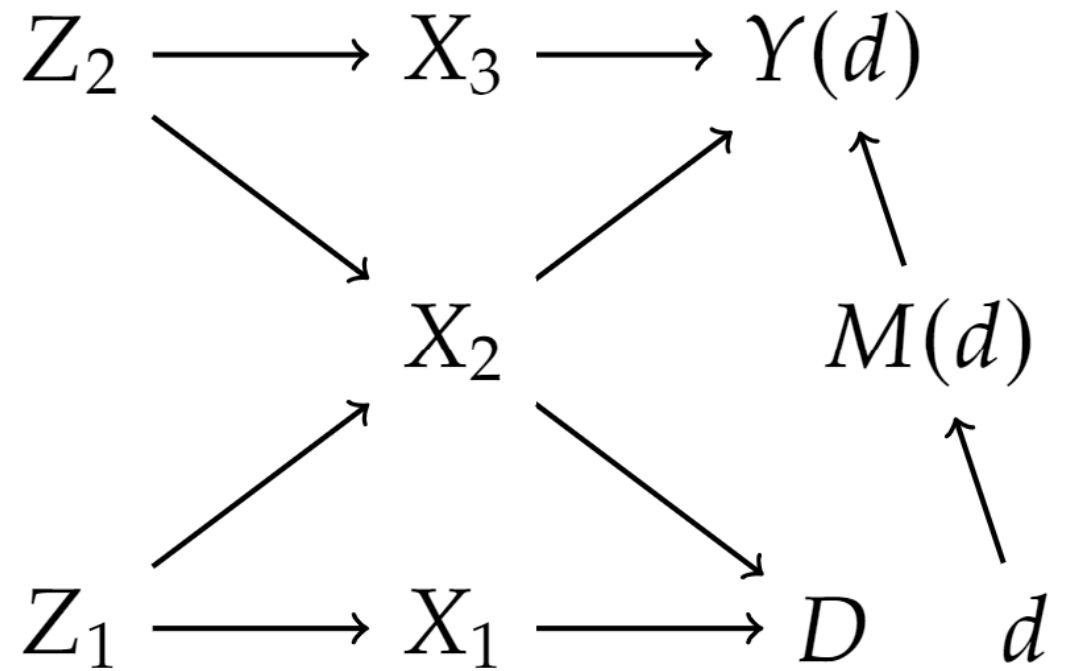
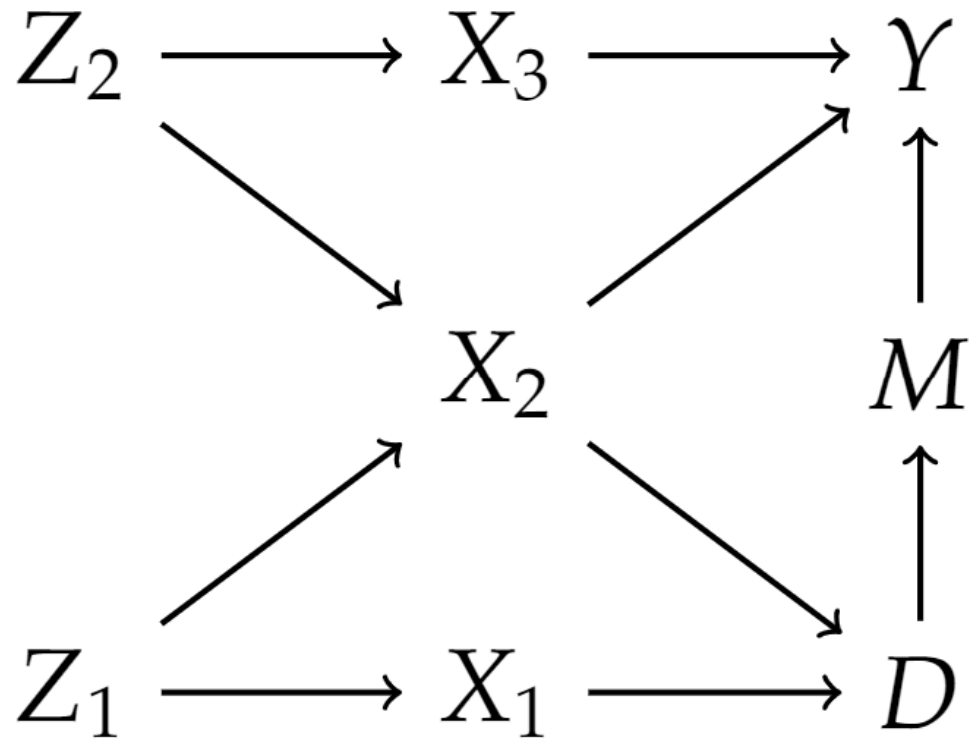
Conditional ignorability between treatment D and outcome Y conditional on set S holds if $Y(d)$ is d-separated from D on SWIG $\tilde{G}(d)$ induced by $\text{fix}(D = d)$ by the set S



Example



Example



Useful Adjustment Strategies

Adjustment Strategies

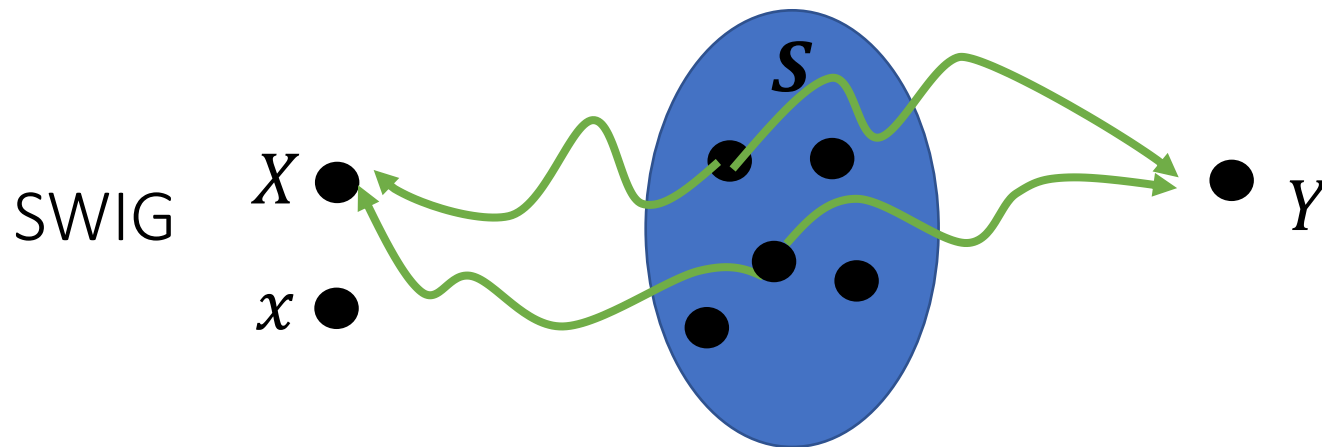
Assuming Y is a descendant of D (otherwise effect is zero)

- Condition on all parents of Y (that are not descendants of D)
- Condition on all parents of D
- Condition on the union of the above

- Condition using backdoor blocking criterion (minimal adjustment)
- Condition on all common causes of D and Y

Conditioning on Parents

- Empirically widely used strategy
- Requires only partial knowledge of the graph
- If we only know the parents of D or the parents of Y we are ok
- Adding any further set to the parent strategy maintains validity



Parents of X

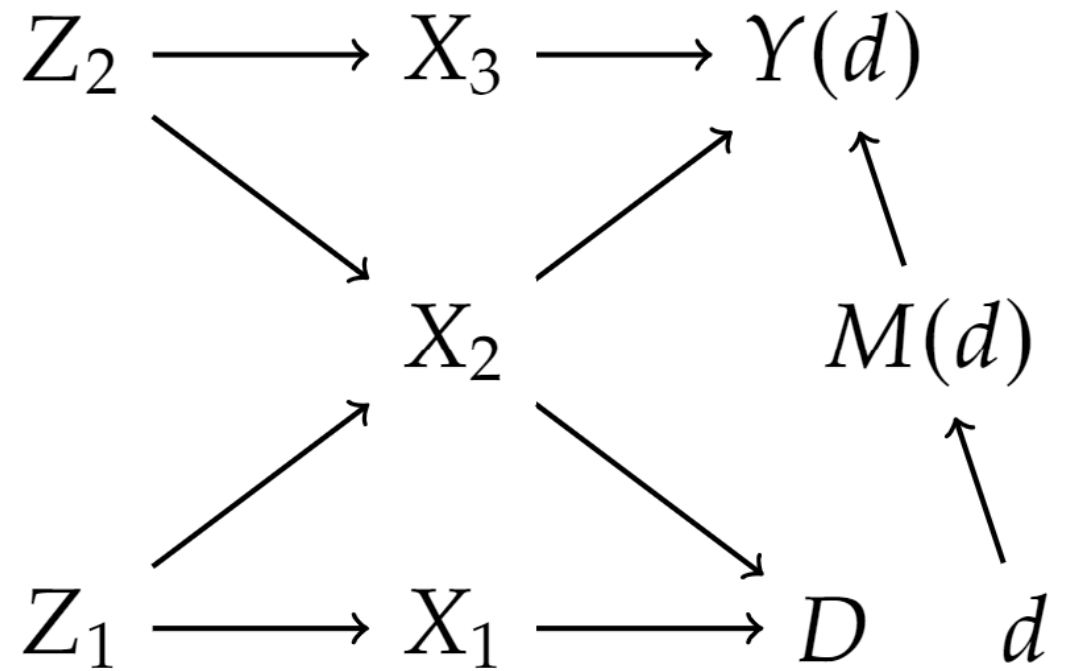
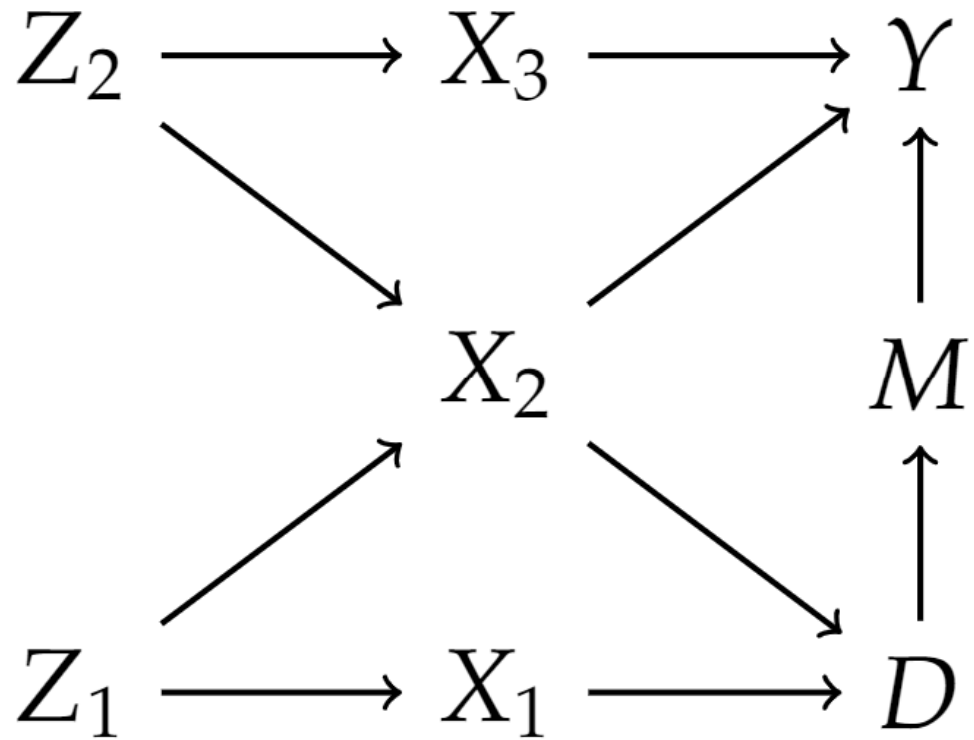
In SWIG, X has no descendants
All paths $X - Y$ must pass from parent

Parent has to be a fork or chain

Backdoor Blocking

- Non-directed path: a path from Y to D that has some reverse edge \leftarrow
- Backdoor path: a non-directed path from Y to D that ends with a \rightarrow
- S is a valid adjustment set if it contains no descendant of D and all backdoor paths from Y to D are blocked
- Allows us to find minimal adjustment sets (small size)

Example



All Common Causes

- Consider all common ancestors of D and Y
 - These are called the “common causes of D and Y ”
 - Very common empirical practice
 - More conservative version: union of ancestors of D and Y
-
- Any backdoor path from Y to D must either contain a common ancestor or contain a collider
 - Conditioning on common ancestors blocks all paths that don't contain a collider

Good and Bad Controls

High Level Categorization of Variables

- **Pre-treatment variables:** variables whose value is determined before the assignment of the treatment
- **Post-treatment variables:** variables whose value is determined after the assignment of the treatment

Good and Bad Controls: High Level

- Pre-treatment variables that are ancestors of either D or Y are ok controls (worst-case can hurt precision/variance)
- There exist pre-treatment variables not of that sort, that can lead to wrong answer (M-bias)
- Most post-treatment variables are bad controls (or don't do much)
- Typically introduce either mediation bias or collider bias (Heckman selection)

“Pre-Treatment Variables”

Examples: Good Controls

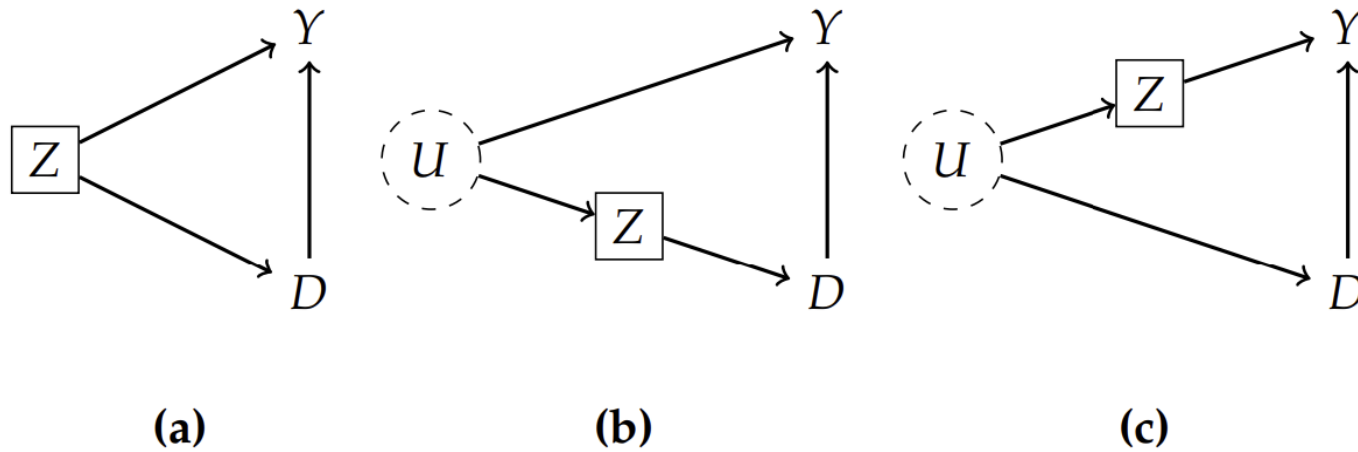
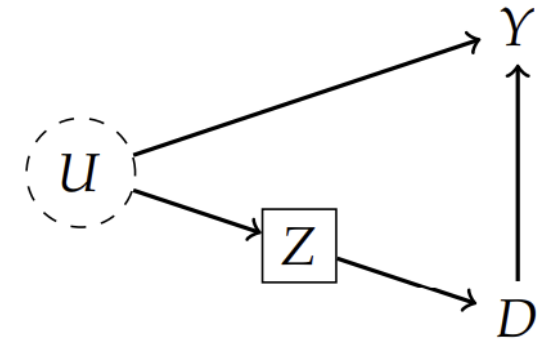


Figure 8.5: Good controls: (a) observed common cause, (b) complete treatment proxy control of unobserved common cause, (c) complete outcome proxy control of unobserved common cause.

Example of Proxy

- Effect of Prenatal Multivitamin consumption (D) on birth defects (Y)
 - Prior family history of birth defects (Z) can influence mother's decision for multivitamin consumption
 - Unmeasured genetic factors (U) cause family history (Z) of birth defects
 - Unmeasured genetic factors (U) have direct influence on birth defects (Y)
-
- Family history of birth defects is a good proxy control for unmeasured genetic factors



Examples: Good Controls

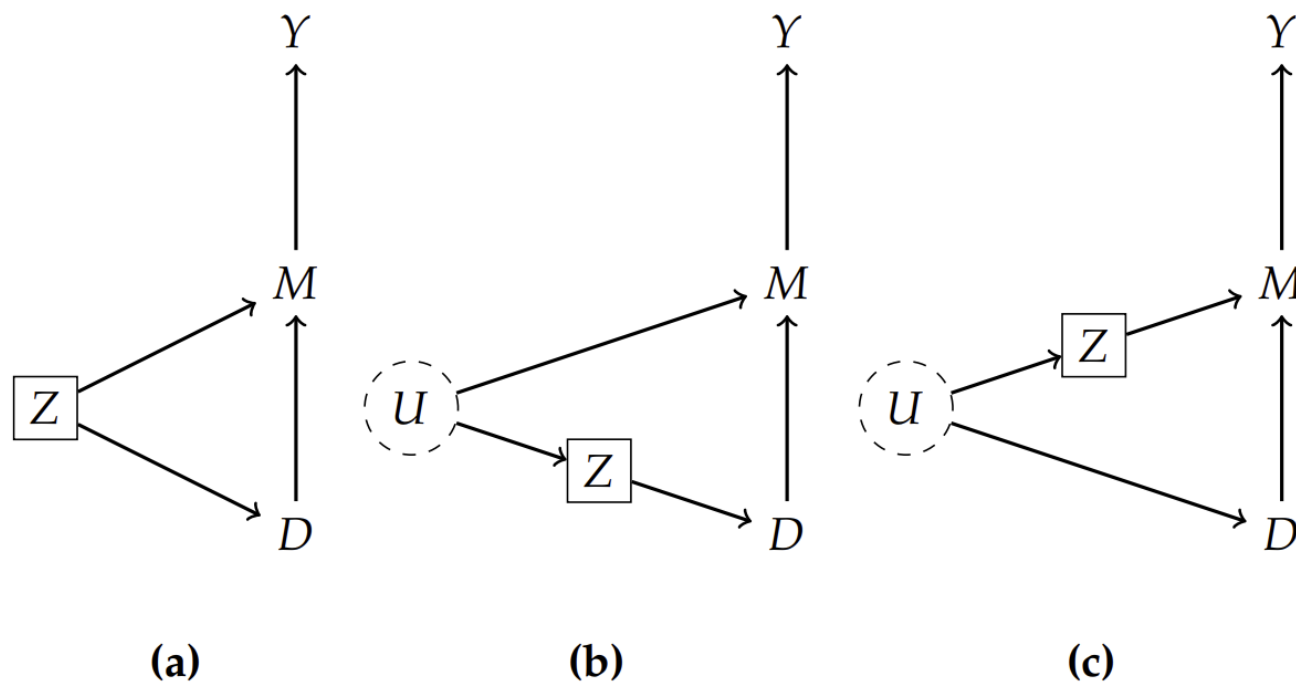


Figure 8.6: Good controls: (a) confounded mediator with observed common cause, (b) confounded mediator, with observed complete treatment proxy control of unobserved common cause, (c) confounded mediator with observed complete outcome proxy control of unobserved common cause.

Examples: Neutral Controls

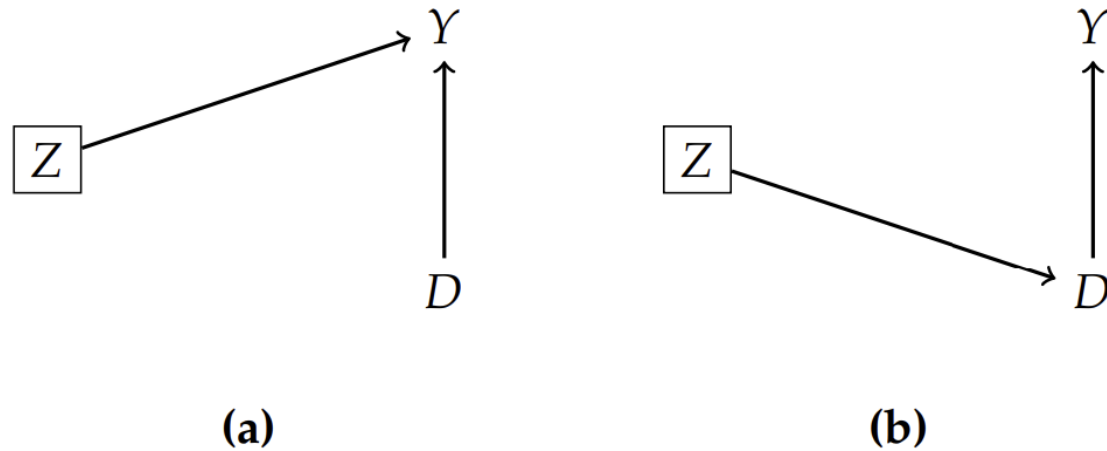


Figure 8.7: Neutral controls: **(a)** Outcome-only cause. Can improve precision; decrease variance. **(b)** Treatment-only cause. Can decrease precision; introduce variance.

Examples: Pre-Treatment Variable that is Bad Control

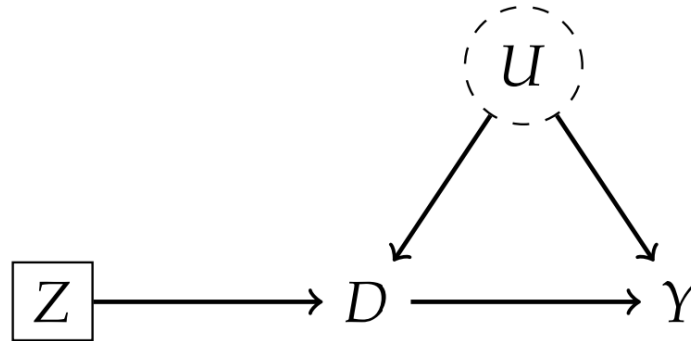


Figure 8.8: Bad control. Bias amplification by adjusting for an *instrument*. Treatment-only cause (*instrument*) that can amplify unobserved confounding bias.

Example: Pre-Treatment Variable that is Bad Control

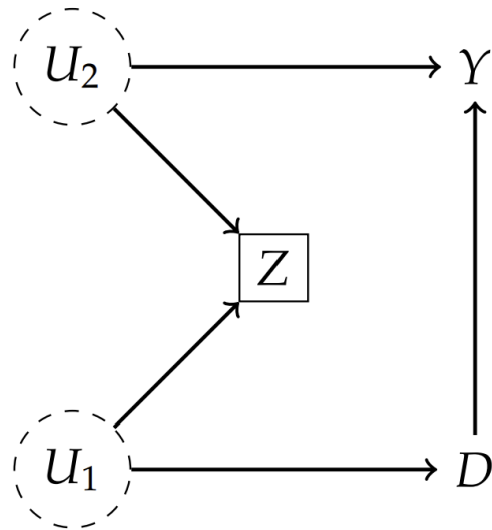
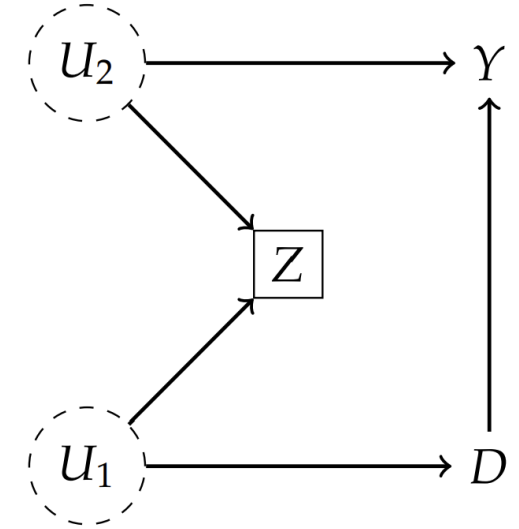


Figure 8.9: Bad control. M-Bias. Pre-treatment variable that introduces Heckman selection bias between two un-correlated unobserved causes.

Example: Homophily Bias in Peer Effects

- Peer effects on civic engagement level
- Consider samples of pairs of friends
- D is civic engagement of one friend at time t
- Y is civic engagement of other friend at time $t+1$
- Civic engagement can be driven by personal traits U_1 and U_2 that are independently drawn well before friendship and affect civic engagement (e.g. level of altruism)
- Friendship Z driven also by personal traits (homophily)
- By looking at pairs of friends we are implicitly conditioning on Z



Example: M-Bias not robust to perturbation

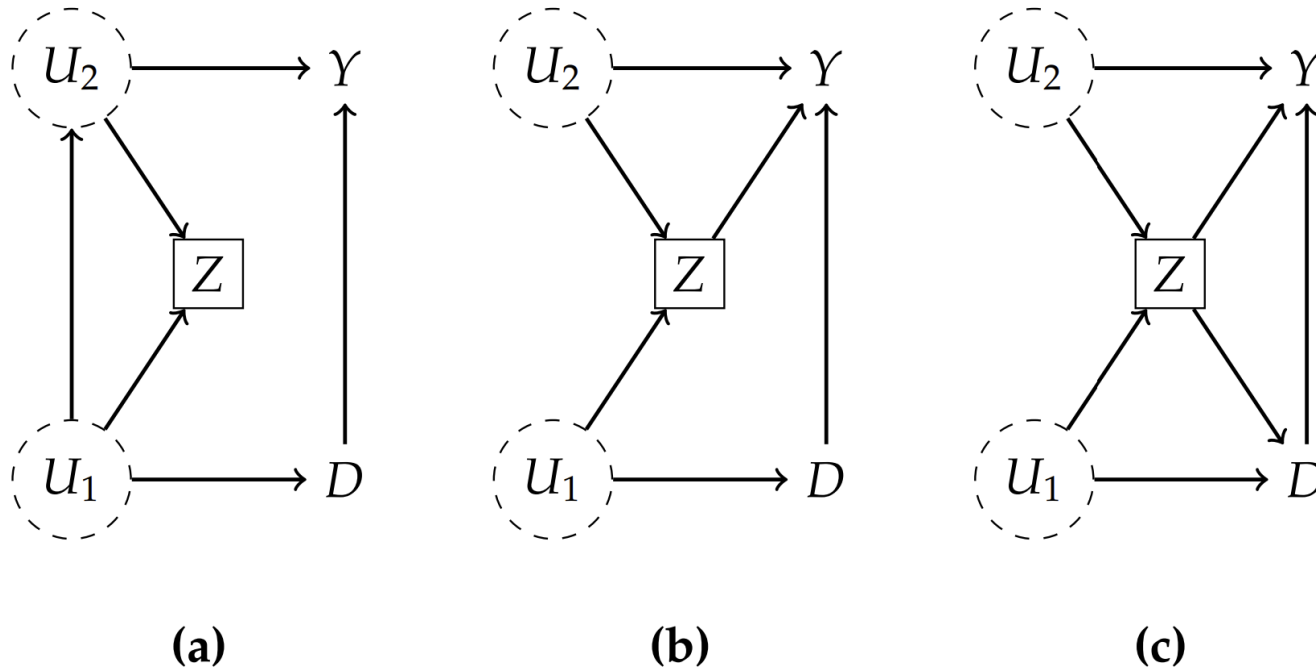


Figure 8.10: Unclear controls: (a) M-bias with correlated unobserved factors, (b) M-Bias with confounding. Pre-treatment variable that introduces Heckman selection between two un-correlated unobserved causes, but also is a confounder itself. No solution is perfect. (c) Butterfly Bias. M-bias with direct confounding.

“Post-Treatment Variables”

Example: Bad Controls, Mediation Bias

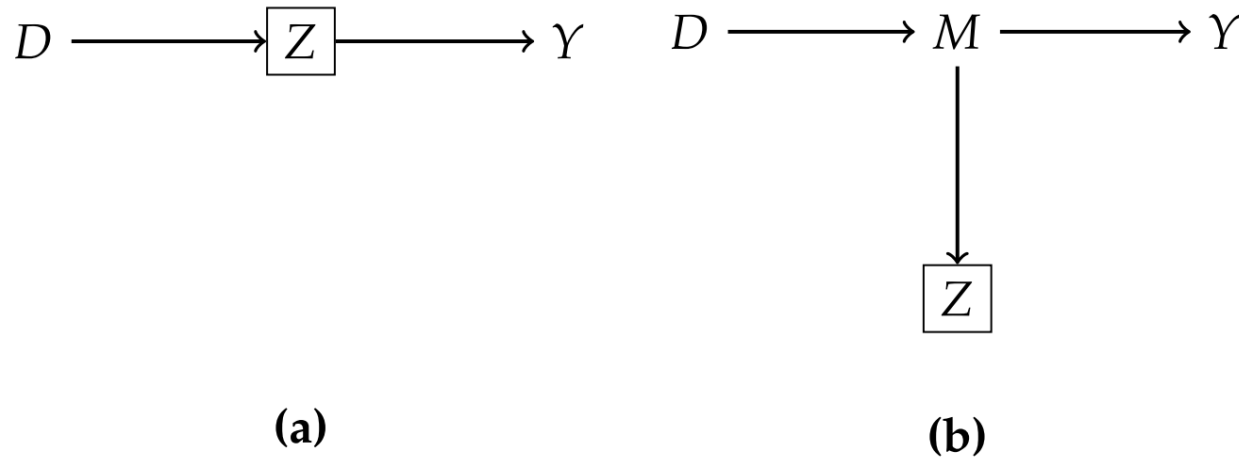


Figure 8.11: Bad controls: **(a)** over-control bias, by controlling on a mediator, **(b)** over-control bias, by controlling on some outcome caused by a mediator.

Example: Exception to Mediation Bias

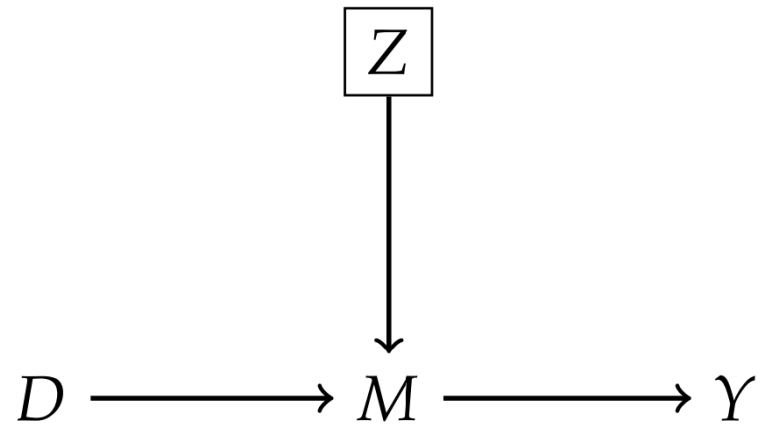


Figure 8.12: Neutral control. Cause of a mediator. Can potentially improve precision.

Example: Controlled Direct Effect Gone Wrong

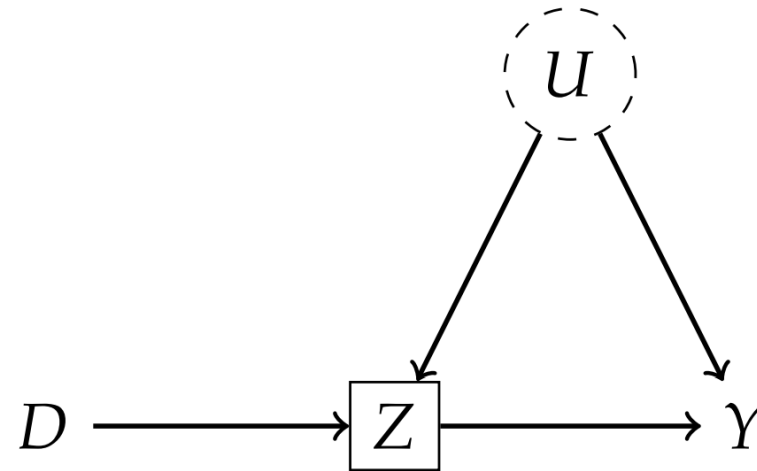


Figure 8.13: Bad control even for *controlled direct effect*. Confounded mediator bias.

Example: Bad Controls, Collider (Heckman Selection) bias

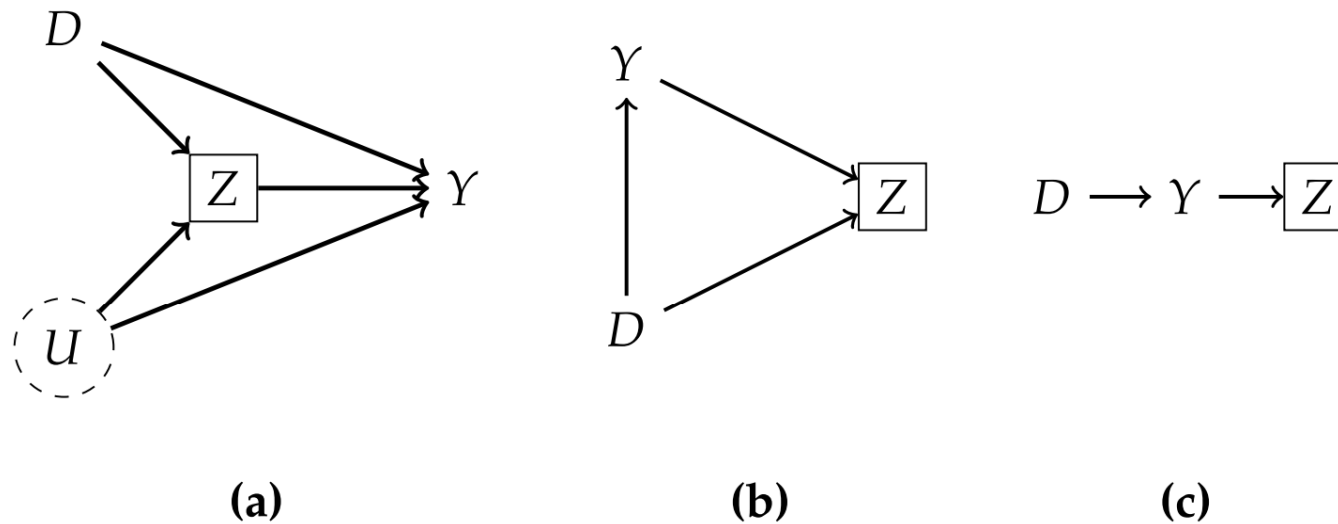


Figure 8.14: Bad controls: (a) collider stratification bias (e.g. low birth-weight “paradox” example), (b) collider stratification bias, (c) case control bias, by controlling on another outcome of the outcome of interest.

Concluding Remarks

- Any study that claims it is estimating the effect of D on Y by conditioning on S must be based on a rigorous thought process
- The DAG/ASEM framework is a rigorous form of this process
- Enables explicit incorporation of domain knowledge in a domain expert friendly manner
- Automatic identification arguments and testable restrictions
- Effective in communicating assumptions of observational study

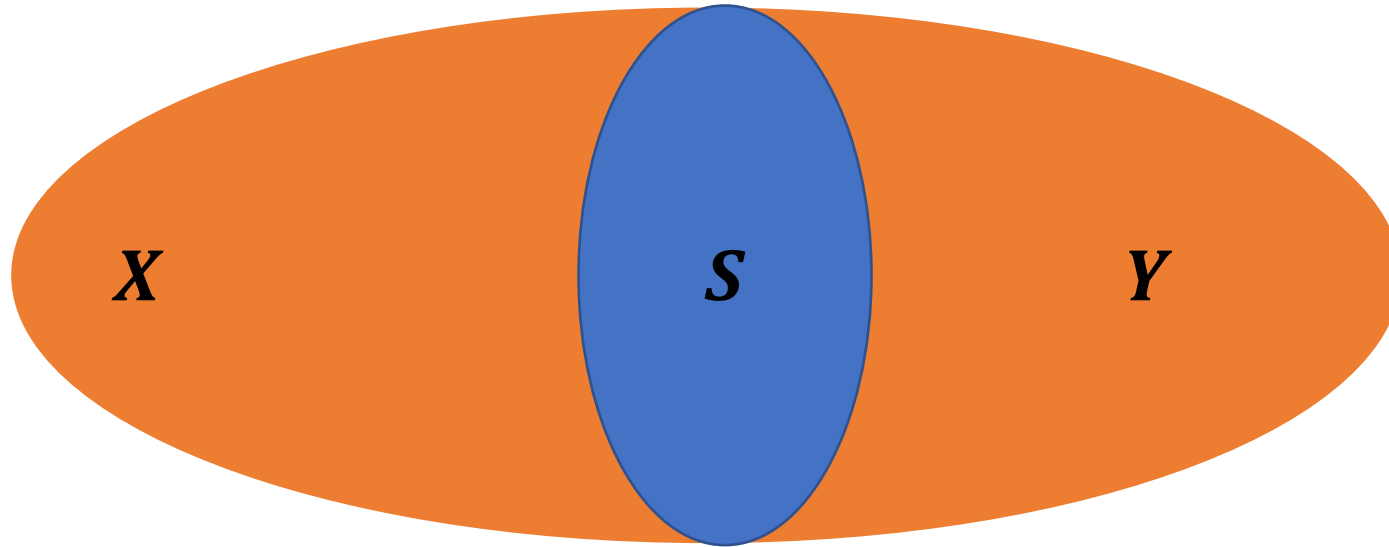
Proving the Main Theorem!

Proof Step 1

- A set of nodes \mathbf{X} is called ancestral if all ancestors of \mathbf{X} are in \mathbf{X}
- Removing all nodes outside of an ancestral set and looking at the resulting graph and ASEM, the probability law is the same as the probability law of \mathbf{X} in the original graph (exercise)

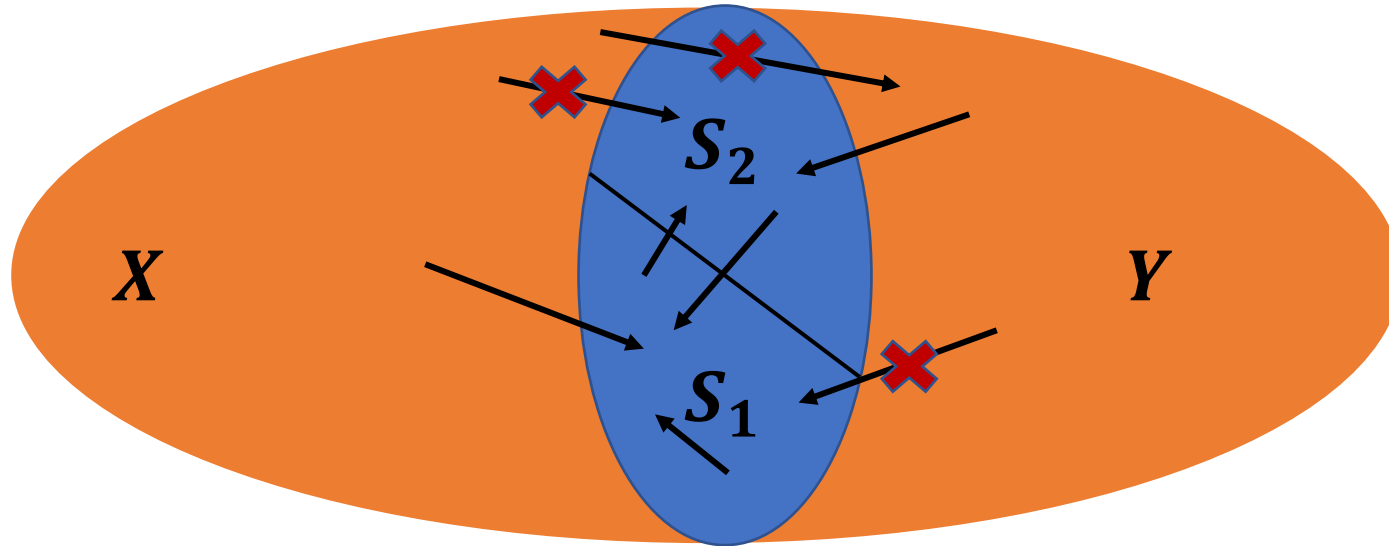
Proof Step 2

- Suppose that a set of nodes \mathbf{X} is D-separated from a set of nodes \mathbf{Y} by a set of nodes \mathbf{S}
- And that $\mathbf{X} \cup \mathbf{Y} \cup \mathbf{S}$ is the set of all nodes



Proof Step 2

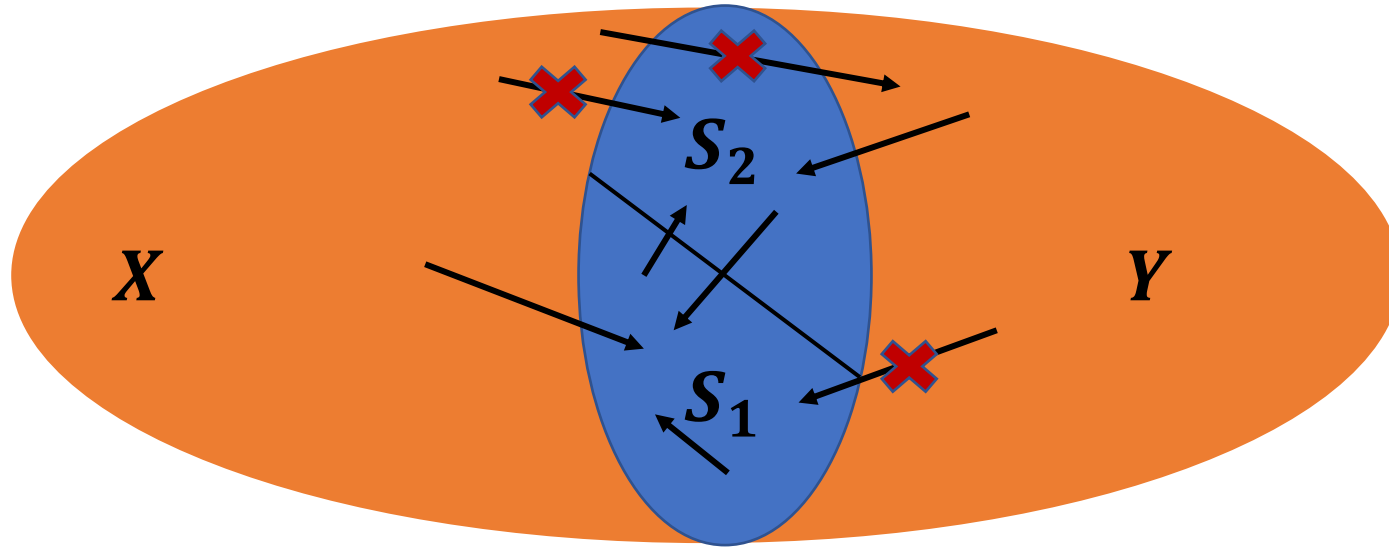
- Let S_1 the subset of S that has a parent in X . Let S_2 the remainder.
- It has to be that $Pa(X \cup S_1) \in X \cup S$
- It has to be that $Pa(Y \cup S_2) \in Y \cup S$



Proof Step 2

- We can factorize:

$$p(x, y, s) = \prod_{W \in X \cup S_1} p(w|pa_w) \prod_{W \in Y \cup S_2} p(w|pa_w) = f(x, s_1)g(y, s_2)$$



Proof Step 2

- We can factorize:

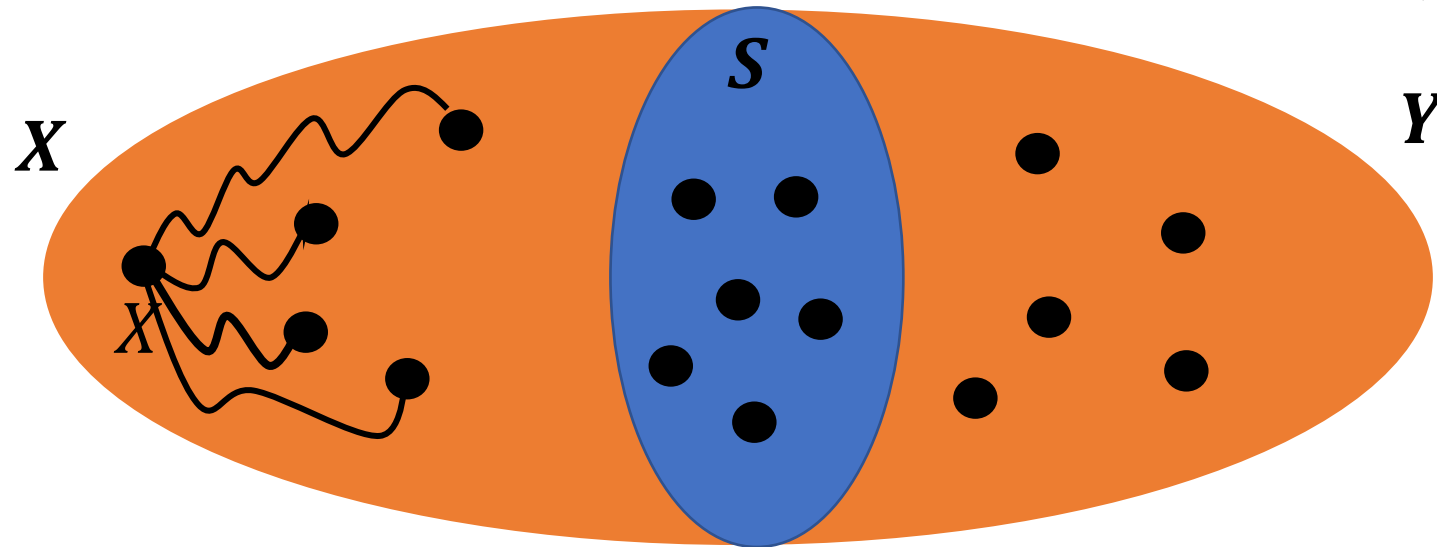
$$p(x, y, s) = \prod_{W \in X \cup S_1} p(w|pa_w) \prod_{W \in Y \cup S_2} p(w|pa_w) = f(x, s_1)g(y, s_2)$$

- Implies that:

$$X \perp\!\!\!\perp Y \mid S$$

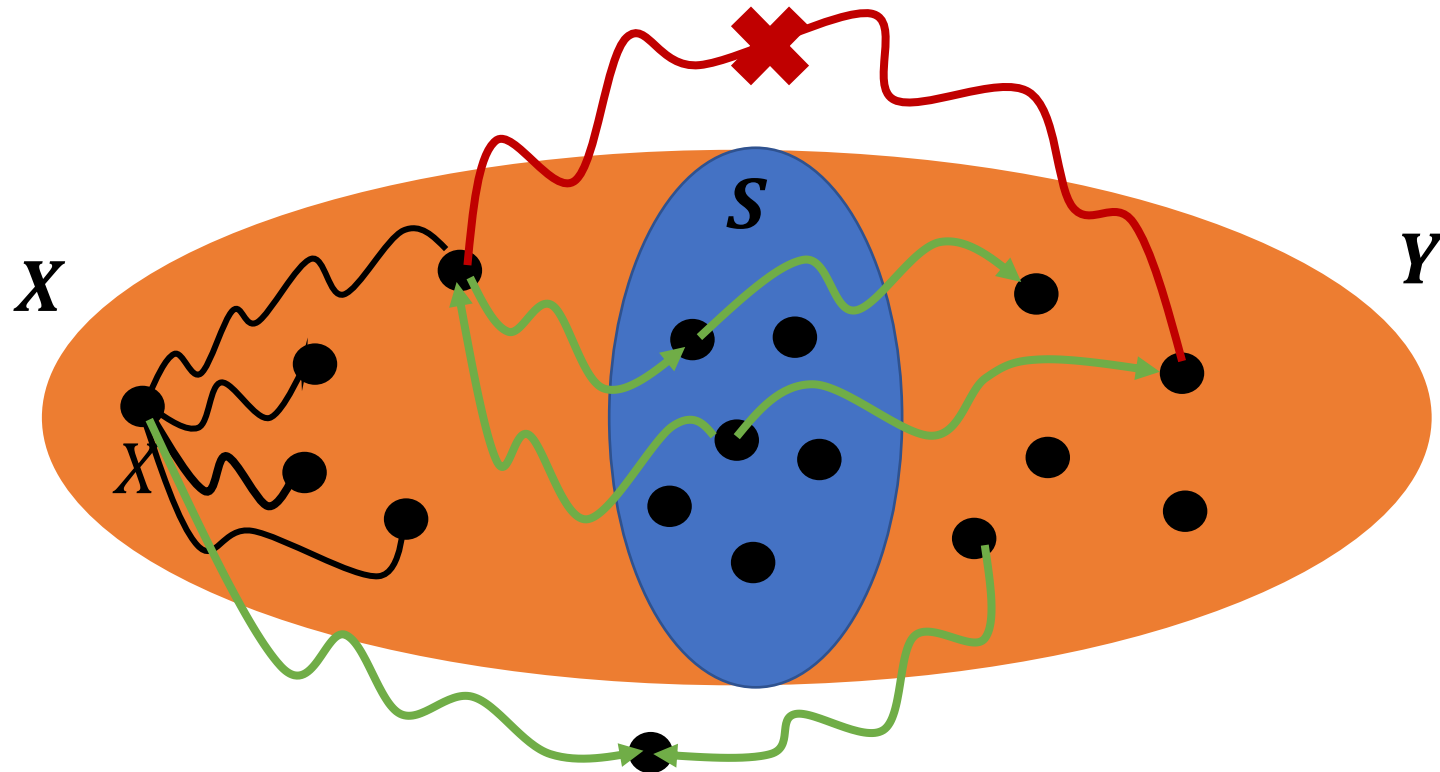
Final Step

- By first step, we can restrict to ancestral set of $X \cup Y \cup S$
- Does not change conditional independence relations (exercise)
- Does not change d-separation relations (exercise)
- Define **X** nodes in ancestral set of $X \cup Y \cup S$ not d-separated from X
- Define **Y** the remainder of nodes in ancestral set not in **X, S** .



Final Step

- By definition of d-separation, S must d-separate X from Y (exercise)
- We can invoke previous critical lemma



Final Step

- By marginalization

$$p(x, y, \mathbf{s}) = \int \int p(x, \mathbf{x}', y, \mathbf{y}', \mathbf{s}) d\mathbf{x}' d\mathbf{y}'$$

- By step 2

$$p(x, y, \mathbf{s}) = \int \int f(x, \mathbf{x}', \mathbf{s}) g(y, \mathbf{y}', \mathbf{s}) d\mathbf{x}' d\mathbf{y}'$$

- We can split integrals

$$p(x, y, \mathbf{s}) = \int f(x, \mathbf{x}', \mathbf{s}) d\mathbf{x}' \int g(y, \mathbf{y}', \mathbf{s}) d\mathbf{y}'$$

- Thus

$$p(x, y, \mathbf{s}) = \bar{f}(x, \mathbf{s}) \bar{g}(y, \mathbf{s}) \Rightarrow X \perp\!\!\!\perp Y \mid S$$