

# Causal Graphical Methods For Handling Nonignorable Missing Data

**Razieh Nabi, Ph.D.**

Department of Biostatistics and Bioinformatics  
Rollins School of Public Health, Emory University  
✉ [razieh.nabi@emory.edu](mailto:razieh.nabi@emory.edu)

**Daniel Scharfstein, Sc.D.**

Department of Population Health Sciences  
University of Utah School of Medicine  
✉ [daniel.scharfstein@hsc.utah.edu](mailto:daniel.scharfstein@hsc.utah.edu)

Upon using the materials, please cite appropriately.

# Motivation

- ▶ In statistics, missing data or missing values occur when no data value is stored for one or multiple variables in an observational unit.
- ▶ We are often interested in a functional of an underlying distribution:

- ▶ Population-level outcome  $Y$ :

$$\psi_1 = \mathbb{E}[Y]$$

- ▶ Outcome mean in a sub-population  $X = x$ :

$$\psi_2 = \mathbb{E}[Y \mid X = x]$$

- ▶ Average causal effect of binary treatment  $T$  on outcome  $Y$ :

$$\psi_3 = \mathbb{E}\left[\mathbb{E}[Y \mid X, T = 1] - \mathbb{E}[Y \mid X, T = 0]\right].$$

(Note that  $\psi_3$  can only be interpreted as a causal effect under certain assumptions, such as *consistency*, *positivity*, and *conditional ignorability*.)

- ▶ The question is how to compute these estimands?
  - ▶ We need some replicates of the underlying distribution, e.g.,  $\{X, T, Y\} \sim p(X, T, Y)$

# Missing data indicators

- ▶ We look at the observed sample and it looks like the following:

$X^*$	$T^*$	$Y^*$
$x_1$	$t_1$	?
$x_2$	?	$y_2$
$x_3$	$t_3$	$y_3$
$\vdots$	$\vdots$	$\vdots$
?	$t_n$	?

- ▶ **Missingness indicator:** each variable with missing values can have an underlying missingness indicator:
  - ▶  $R_V = 1$  if variable  $V$  is observed and  $R_V = 0$  if  $V$  is "?".

$X^*$	$T^*$	$Y^*$	$R_X$	$R_T$	$R_Y$
$x_1$	$t_1$	?	1	1	0
$x_2$	?	$y_2$	1	0	1
$x_3$	$t_3$	$y_3$	1	1	1
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
?	$t_n$	?	0	1	0

# Missing data challenges

- ▶ How can we use a sample with missing values to estimate the parameter of interest?
  - ▶ Should we ignore rows with missing values?
  - ▶ Should we impute the missing values? If so, how?
  - ▶ Should we do something else?
- ▶ Our choices may affect our data analysis by
  - ▶ introducing **bias** resulting from differences between missing and complete data, and/or
  - ▶ losing **efficiency** when we ignore part of the observed sample.
- ▶ Before we choose what method to use, we need to know **why we have missing data** in the first place!

# Sources of missingness

We encounter missing data for a variety of reasons:

- ▶ A survey was conducted and values were just randomly missed when being entered in the computer.
- ▶ A respondent chooses not to respond to a question like "Have you ever recreationally used opioids?"
- ▶ You decide to start collecting a new variable (due to new actions: like a pandemic) partway through the data collection of a study.
- ▶ You want to measure the speed of meteors, and some observations are just "too quick" to be measured properly.

The source of missing values in data leads to three distinct missingness mechanisms: MCAR, MAR, MNAR (Rubin, 1976).

# Rubin's hierarchy of missingness

1. **Missing Completely at Random (MCAR)** - the probability of missingness in a variable is the same for all units. Like randomly poking holes in a data set.
2. **Missing at Random (MAR)** - the probability of missingness in a variable depends only on available information (in other predictors).
3. **Missing Not at Random (MNAR)** - the probability of missingness depends on information that has not been recorded and this information also predicts the missing values.

# Rubin's hierarchy of missingness

Let  $Z$  : variables with no missingness

$X$  : variables that are sometimes missing

$X_{\text{obs}}$ : observed entries,  $X_{\text{miss}}$ : missing entries

$R$  : missingness indicators

---

1. **MCAR**:  $R \perp\!\!\!\perp X_{\text{miss}}, X_{\text{obs}}, Z, \quad p(R \mid X, Z) = p(R).$

- ▶ probability that any observation is missing is independent of all data values, regardless of whether they are observed or unobserved.

2. **MAR**:  $R \perp\!\!\!\perp X_{\text{miss}} \mid Z, X_{\text{obs}}, \quad p(R \mid X, Z) = p(R \mid Z, X_{\text{obs}})$

- ▶ probability that any observation is missing depends only on elements that are observed.

3. **MNAR**:  $R \not\perp\!\!\!\perp X_{\text{miss}}$  neither marginally nor conditionally

- ▶ probability that any observation is missing depends on elements that are missing themselves – a missingness mechanism that is neither MCAR nor MAR.

# Beyond the traditional missingness hierarchy

- ▶ Rubin's categorization does not determine the best approach for handling missing data in multiple variables.
- ▶ How best to handle missing data?
  - ▶ It depends on the assumed causal relationships between variables and their missingness, and
  - ▶ What these relationships imply in terms of the identification/recoverability of the target estimand.
- ▶ Main takeaway: encourage the use of *missing data DAGs* in data analysis to
  - ▶ Make underlying **assumptions** more explicit, and
  - ▶ Use **identification** procedures as a guide for **estimation** methods.



# Course outline and objectives

## Part I. Missing data DAGs

- ▶ Representing missingness mechanisms graphically; transferring expert knowledge into a concise graphical representation; interpret a missing data DAG model as a class of distributions with a set of independence restrictions.

## Part II. Non-parametric identification

- ▶ Identifying a given query, write it down as a function of observed data, or prove the given query of interest is not identified without making further assumptions.

## Part III. Non/Semi-parametric estimation

- ▶ Given an identified query, derive the non-parametric influence functions; Given a missing data DAG model, derive the tangent space of the underlying full data and observed data distributions.

## Part IV. Sensitivity analysis

- ▶ Assessing deviations from assumptions.

## **Part 1. Missing Data DAG Models**

# Missing data and causal inference

- ▶ Causal inference and missing data are analogous in terminology, theory of identification, and statistical inference.
- ▶ Causal inference has been viewed as a missing data problem:
  - ▶ Responses to some (hypothetical) treatment interventions are not observed.
  - ▶ Given the treatment vs placebo option, we only observe the potential outcome under treatment received or the potential outcome under placebo received, but not both.
- ▶ Missing data can be viewed as a causal inference problem:
  - ▶ Missingness indicators can be treated as intervenable treatments.
  - ▶ We can view variable  $X$  as a potential outcome had the missingness indicator  $R_X$  been set to 1 (had there been no missingness).
- ▶ In this part: get inspired by developments in causal graphical models to reason about missing data models.

# A causal workflow

1. Define a causal estimand in terms of counterfactuals.
2. Define a causal model that links counterfactuals to factual variables.
  - ▶ Impose assumptions about the distribution over counterfactual and factual variables.
3. Identify the causal estimand as a function of observed data in the assumed causal model.
4. Define a statistical model to estimate the identifiable causal estimand.
  - ▶ Perform statistical inference which includes testing and estimating the magnitude of a causal estimand given the observed data.
5. Assess assumptions with sensitivity analysis.

## Example: a causal workflow

1. **Average causal effect:**  $ACE := \mathbb{E}[Y^{(1)} - Y^{(0)}]$

- ▶  $Y^{(t)}$ : potential outcome  $Y$  when binary treatment  $T$  is assigned to  $t = \{0, 1\}$ .

2.  $\mathcal{M}$ : a causal model relating counterfactuals to factials

- ▶ **Consistency**: observed outcome  $Y$  is equal to the potential outcome  $Y^{(t)}$  when the treatment received is  $T = t$ ,
- ▶ **Positivity**:  $p(T = t \mid X = x) > 0$  for all  $x$  in the state space of  $X$ ,
- ▶ **Conditional ignorability**:  $Y^{(t)} \perp T \mid X$

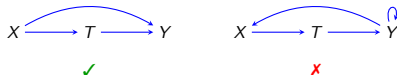
3. Under the above causal model, we can identify ACE via the following functional, known as **adjustment formula** or **g-formula**:

$$\mathbb{E}[Y^{(1)} - Y^{(0)}] = \mathbb{E}\left[\mathbb{E}[Y \mid T = 1, X] - \mathbb{E}[Y \mid T = 0, X]\right].$$

4, 5. Luckily, we are not short of any estimation or sensitivity analysis techniques!

# Directed acyclic graph (DAG)

- ▶ The second step in the causal workflow is what distinguishes causal analyses from traditional statistical analyses.
  - ▶ Graphical models like directed acyclic graphs (DAGs) are often used to encode assumptions about the causal model.
- ▶ A graph  $\mathcal{G}(V, E)$  is a set of vertices/nodes  $V$  that correspond to random variables and a set  $E$  that contains the set of edges between variables.
- ▶ The graph  $\mathcal{G}(V, E)$  is said to be **directed** and **acyclic** if:
  - ▶ There are only directed edges ( $V_i \rightarrow V_j$ )
  - ▶ There are no directed cycles – for any  $V_i \in V$  there is no sequence of directed edges in  $G$  such that  $V_i \rightarrow \dots \rightarrow V_i$



- ▶ For notational convenience, we often refer to  $\mathcal{G}(V, E)$  as  $\mathcal{G}(V)$  or simply  $\mathcal{G}$ .

# Statistical model of a DAG

- ▶ DAG  $\mathcal{G}(V)$  encodes a set of independence restrictions on the joint distribution  $p(V)$ .
- ▶ The joint distribution  $p(V)$  corresponding to DAG  $\mathcal{G}(V)$  has three **equivalent** characterizations:
  - ▶ **Factorization** (writes the distribution as a set of small factors.)
  - ▶ **Local Markov property** (lists a complete set of independence constraints.)
  - ▶ **Global Markov property** (lists all independence constraints in the model.)

## Statistical model of a DAG ctd.

The joint distribution  $p(V)$  satisfies the **factorization property** wrt DAG  $\mathcal{G}(V)$  if:

$$p(V) = \prod_{V_i \in V} p(V_i \mid \text{pa}_{\mathcal{G}}(V_i)),$$

- ▶  $\text{pa}_{\mathcal{G}}(V_i) = \{V_j \in V \mid V_j \rightarrow V_i\}$  denotes parents of  $V_i$  in  $\mathcal{G}(V)$ .
- 

The joint distribution  $p(V)$  satisfies the **local Markov property** wrt DAG  $\mathcal{G}(V)$  if:

$$V_i \perp\!\!\!\perp \text{nd}_{\mathcal{G}}(V_i) \setminus \text{pa}_{\mathcal{G}}(V_i) \mid \text{pa}_{\mathcal{G}}(V_i), \quad \forall V_i \in V$$

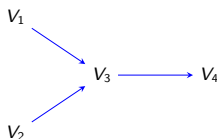
(variable  $V_i$  is independent of its non-descendant non-parents given its parents)

- ▶  $\text{deg}_{\mathcal{G}}(V_i) = \{V_j \in V \mid V_i \rightarrow \dots \rightarrow V_j\}$  denotes descendants of  $V_i$  in  $\mathcal{G}(V)$ ,
- ▶  $\text{nd}_{\mathcal{G}}(V_i) = \{V_j \in V \mid V_j \notin \text{deg}_{\mathcal{G}}(V_i)\}$  denotes non-descendants of  $V_i$  in  $\mathcal{G}(V)$ .



## Example: statistical model of a DAG

Consider the following DAG:



- ▶ According to the **DAG factorization**, the statistical model of this DAG is a set of distributions  $p(V)$ , where  $V = \{V_1, V_2, V_3, V_4\}$  s.t.,

$$\left\{ p(V) = p(V_1) \times p(V_2) \times p(V_3 \mid V_1, V_2) \times p(V_4 \mid V_3) \right\}.$$

- ▶ According to the **local Markov property**, the statistical model of this DAG is a set of distributions  $p(V)$  s.t.,

$$\left\{ p(V) \text{ s.t. } V_1 \perp\!\!\!\perp V_2 \text{ and } V_4 \perp\!\!\!\perp V_1, V_2 \mid V_3 \right\}.$$

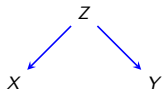
The above list implies a larger set of independence restrictions, e.g.,  $V_4 \perp\!\!\!\perp V_1 \mid V_3$  or  $V_4 \perp\!\!\!\perp V_1 \mid V_2, V_3$ . (graphoid axioms)

# Global Markov property: d-separation

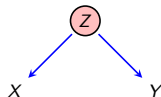
- ▶ Given a DAG  $\mathcal{G}(V)$ , we need to be able to answer arbitrary independence queries of the form  $X \perp\!\!\!\perp Y \mid Z$  in  $p(V)$ , where  $X, Y, Z$  are disjoint subsets of  $V$ .
- ▶ **d-separation** (directed-separation) is a graphical criterion that allows one to answer such queries in an automated fashion.
- ▶ Here are the three types of triplets that define d-separation:
  - ▶ Forks  $X \leftarrow Z \rightarrow Y$
  - ▶ Chains  $X \rightarrow Z \rightarrow Y$  or  $X \leftarrow Z \leftarrow Y$
  - ▶ Colliders  $X \rightarrow Z \leftarrow Y$

# Fork triplets

In a fork  $X \leftarrow Z \rightarrow Y$ , the variables  $X$  and  $Y$  are marginally dependent, but conditionally independent given  $Z$ .



$$X \not\perp Y$$



$$X \perp Y \mid Z$$

- ▶ Intuition:  $X$  and  $Y$  **share a common cause** and thus dependent.
  - ▶ Upon observing the common cause  $Z$ , the two effects  $X$  and  $Y$  are no longer related.
- ▶ Example: warmer weather draw more people to the beach and also drive up ice cream sales.
  - ▶  $X$ : shark attacks,  $Z$ : warm whether, and  $Y$ : ice cream sales.

# Chain triplets

In a chain  $X \rightarrow Z \rightarrow Y$ , the variables  $X$  and  $Y$  are marginally dependent, but conditionally independent given  $Z$ .



$$X \not\perp\!\!\!\perp Y$$

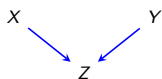


$$X \perp\!\!\!\perp Y \mid Z$$

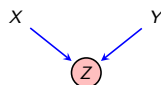
- ▶ Intuition: if  $Z$  is a noisy version of  $X$  and  $Y$  is a noisy version of  $Z$ , then  $Y$  is a noisy version of  $X$ .
  - ▶  $Z$  **screens off the effect** of  $X$  on  $Y$ .
  - ▶ Upon observing  $Z$ ,  $X$  holds no extra information about  $Y$ .
- ▶ Example: blood sugar causes hunger, but only indirectly through increasing the stomach acidity.
  - ▶  $X$ : blood sugar,  $Z$ : stomach acidity, and  $Y$ : hunger.

# Collider triplets

In a collider  $X \rightarrow Z \leftarrow Y$ , the variables  $X$  and  $Y$  are marginally independent, but conditionally dependent given  $Z$ .



$$X \perp\!\!\!\perp Y$$

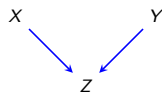


$$X \not\perp\!\!\!\perp Y \mid Z$$

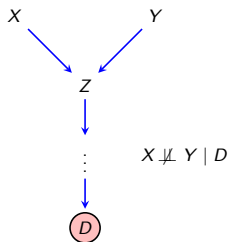
- ▶ Intuition: if  $X$  and  $Y$  only **share a common effect**, they are independent. That is the common effect has two independent sources of causes.
  - ▶ Upon observing the common effect, the two causes become dependent.
  - ▶ This is often referred to as a **Berkson's paradox**.
- ▶ Example: if we observe that the car fails to start, then knowing something about the fuel status tells us something about the battery status & vice versa
  - ▶  $X$ : battery,  $Z$ : car starts, and  $Y$ : fuel.
- ▶ Conditioning can induce dependence, not just remove the dependence.

# Collider extensions

In a collider  $X \rightarrow Z \leftarrow Y$ , the variables  $X$  and  $Y$  are marginally independent, but conditionally dependent given a *descendant* of  $Z$ .



$X \perp\!\!\!\perp Y$



$X \not\perp\!\!\!\perp Y \mid D$

- ▶ Example: extend the previous example where the car was taken to a mechanic.
  - ▶  $X$ : battery,  $Z$ : car starts,  $Y$ : fuel, and  $D$ : taken to mechanic

# Summary of the (in)dependence rules

$$X \leftarrow Z \rightarrow Y$$

$$X \not\perp Y$$

$$X \leftarrow (Z) \rightarrow Y$$

$$X \perp Y \mid Z$$

$$X \rightarrow Z \rightarrow Y$$

$$X \not\perp Y$$

$$X \rightarrow (Z) \rightarrow Y$$

$$X \perp Y \mid Z$$

$$X \rightarrow Z \leftarrow Y$$

$$X \perp Y$$

$$X \rightarrow (Z) \leftarrow Y$$

$$X \not\perp Y \mid Z$$

$$X \rightarrow Z \leftarrow Y$$



...



$$X \not\perp Y \mid D$$

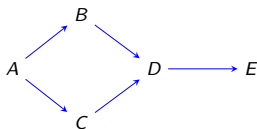
- ▶ Forks/chains are *open*, but become *blocked* upon conditioning.
- ▶ Colliders are *blocked*, but become *open* upon conditioning.

# From blocked triplets to d-separation

- ▶ A **path** from  $X$  to  $Y$  is a sequence of consecutive edges connecting  $X$  and  $Y$  such that no vertex (and consequently no edge) appears more than once in the sequence.
- ▶ A path from  $X$  to  $Y$  is blocked by  $Z$  if there is a **blocking triplet** on the path.
  - ▶ There exists a blocked chain or fork on the path, or
  - ▶ A collider that is not open.
- ▶ Dependence is like water flow and paths are pipes. A single block is enough to block the whole path.
- ▶  $X$  and  $Y$  are **d-separated** given  $Z$  if all paths from  $X$  to  $Y$  are blocked by  $Z$ , and is denoted by  $X \perp\!\!\!\perp_d Y \mid Z$ .



## Examples: d-separation



- Is  $A \perp\!\!\!\perp_d E \mid C$ ?

No! because  $A \rightarrow B \rightarrow D \rightarrow E$  is still open.

- Is  $B \perp\!\!\!\perp_d C \mid A$ ?

Yes!  $B \leftarrow A \rightarrow C$  and  $B \rightarrow D \leftarrow C$  are both blocked.

- Is  $B \perp\!\!\!\perp_d C \mid E, A$ ?

No!  $B \rightarrow D \leftarrow C$  is now open (condition on  $E$  opens up the collider at  $D$ ).

# Global Markov property

For any distribution  $p(V)$  that satisfies the DAG factorization wrt  $\mathcal{G}(V)$ , the following **Global Markov property** holds: for all disjoint subsets  $X, Y, Z$  of  $V$  we have,

$$(X \perp\!\!\!\perp_d Y \mid Z) \Big|_{\mathcal{G}(V)} \implies (X \perp\!\!\!\perp Y \mid Z) \Big|_{p(V)}$$

where  $(\perp\!\!\!\perp_d) \Big|_{\mathcal{G}}$  denotes d-separation in  $\mathcal{G}$  and  $(\perp\!\!\!\perp) \Big|_p$  denotes independence in  $p$ .

- ▶ We can apply a purely graphical criterion to a DAG  $\mathcal{G}(V)$  to tell us about conditional independence facts in the joint distribution  $p(V)$ .
- ▶ The above is a one way implication!
  - ▶ We could indeed have extra independence restrictions in  $p(V)$  that we cannot read them by d-separation (and that happens in *unfaithful* distributions).

# Equivalence of DAG properties

The statistical model of a DAG is characterized with three definitions:

- ▶ **Factorization** (writes the distribution as a set of small factors).
- ▶ **Local Markov property** (lists a small but complete set of independence constraints).
- ▶ **Global Markov property** (lists all independence constraints in the model).

A distribution  $p(V)$  factorizes according to a DAG  $\mathcal{G}(V)$  if and only if it obeys the local Markov property according to  $\mathcal{G}(V)$  if and only if it obeys the global Markov property according to  $\mathcal{G}(V)$  (Verma and Pearl, 1990).

DAG factorization  $\iff$  Local Markov property  $\iff$  Global Markov property

# Causal model of a DAG

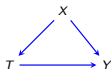
- ▶ The causal model of a DAG can be formally defined in terms of a **Nonparametric Structural Equation Model (NPSEM)**.
  - ▶ It describes how “nature” assigns values to each variable in the model.
- ▶ For every  $V_i \in V$ :
$$V_i \leftarrow f_{V_i}(\text{pa}_{\mathcal{G}}(V_i), \epsilon_{V_i})$$
  - ▶  $\epsilon_{V_i}$  denotes the error term (all external unmeasured causes of  $V_i$ ).
  - ▶  $f_{V_i}$  is nonparametric. It does not constrain the dependence of  $V_i$  on its parents and  $\epsilon_{V_i}$  in any way.
  - ▶ This is an imperative assignment, not an equality! which means the model is not “reversible.”
- ▶ Unmeasured factors are typically assumed to be independent – a reflection of the fact that all unmeasured common causes have been measured.
  - ▶ This is known as the **NPSEM with independent errors (NPSEM-IE)**
  - ▶ The explicit assumption is that  $\perp\!\!\!\perp \{\epsilon_{V_i}, \forall V_i \in V\}$ , and thus
$$p(\epsilon) = \prod_{V_i \in V} p(\epsilon_{V_i})$$

# Intervention in causal models

- ▶ Let  $T$  be the variable that we would like to (hypothetically) intervene on and set it to  $t$ .
- ▶ An intervention that sets  $T = t$ , entails the following three changes:
  - I. **Structural** changes to the SEM,
    - ▶ In the corresponding NPSEM, replace  $T \leftarrow f_T(\text{pa}_{\mathcal{G}}(T), \epsilon_T)$  with  $T \leftarrow t$ .
    - ▶ The change in the structural equation for other variables depend on their genealogical relations to  $T$ .
  - II. **Graphical** changes the DAG  $\mathcal{G}$ ,
    - ▶ In the corresponding DAG  $\mathcal{G}$ , delete all incoming edges into  $T$  and fix the  $T$  node to take value  $t$ .
  - III. **Probabilistical** changes the joint law  $p(V)$ .
    - ▶ In the corresponding joint law  $p(V)$ , drop the term  $p(T \mid \text{pa}_{\mathcal{G}}(T))$  from the factorization of  $p(V)$ , and evaluate all terms at  $T = t$ .

## Example: intervention in causal models

- Let  $T$  be the treatment of interest in the following DAG:



- Structural operation of intervening on  $T$  and setting it to  $t$ :

$$X \leftarrow f_X(\epsilon_X)$$

$$T \leftarrow f_T(X, \epsilon_T)$$

$$Y \leftarrow f_Y(X, T, \epsilon_Y)$$

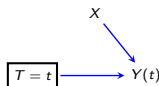
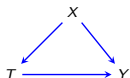
$$X \leftarrow f_X(\epsilon_X)$$

$$T \leftarrow t$$

$$Y^{(t)} \leftarrow f_Y(X, t, \epsilon_Y)$$

NPSEM-IE implies the conditional ignorability assumption:  $Y^{(t)} \perp\!\!\!\perp T \mid X$ .

- Graphical operation of this intervention is illustrated as follows:



- Probabilistic operation of this intervention is a truncated factorization:

$$p(X, Y^{(t)}) = \frac{p(X, Y, T)}{p(T \mid X)} \Big|_{T=t} = \frac{p(X, Y, T=t)}{p(T=t \mid X)}.$$

Do-calculus notation of Pearl:  $p(X, Y^{(t)}) \equiv p(X, Y \mid \text{do}(T = t))$ .

# Missing data notation

- ▶  $X = (X_1, \dots, X_K)^T$ : a vector of random variables
- ▶ Given a finite sample:
  - ▶  $R = (R_1, \dots, R_K)^T$ : binary missingness indicators  
 $R_k = 1$  if  $X_k$  is observed, and  $R_k = 0$  otherwise
  - ▶  $X^* = (X_1^*, \dots, X_K^*)^T$ : coarsened version of  $X$   
 $X_k^* = X_k$  if  $R_k = 1$ , and  $X_k^* = ?$  otherwise
- ▶ Causal interpretation of the tuple  $(X_k, R_k, X_k^*)$ :
  - ▶  $R_k$ : a treatment variable that can be intervened on
  - ▶  $X_k$ : a counterfactual – had we intervened and set  $R_k = 1$
  - ▶  $X_k^*$ : a factual variable
- ▶ Switching notation to emphasize counterfactual connotation:

$$X, R, X^* \mapsto L^{(1)}, R, L$$

$$L^{(1)} = (L_1^{(1)}, \dots, L_K^{(1)})^T \text{ and } L = (L_1, \dots, L_K)^T.$$

- ▶  $Z$ : completely observed variables

# Missing data models

- ▶ A missing data model  $\mathcal{M}$  is a set of distributions defined over variables in  $\{Z, L^{(1)}, R, L\}$ .
- ▶ By chain rule of probability, we can factorize  $p(Z, L^{(1)}, R, L)$  as follows:

$$\underbrace{\underbrace{p(Z, L^{(1)})}_{\text{target law}} \times \underbrace{p(R \mid L^{(1)}, Z)}_{\text{missingness mechanism}}}_{\text{full law } p(L^{(1)}, Z, R)} \times \underbrace{p(L \mid L^{(1)}, R)}_{\text{deterministic terms}}.$$

- ▶ Consistency assumption:  $L_k = \begin{cases} L_k^{(1)} & \text{if } R_k = 1 \\ ? & \text{if } R_k = 0 \end{cases}$
- ▶ Observed data margin is  $p(Z, R, L)$ , where counterfactuals are marginalized out.

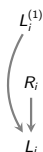


# A missing data workflow

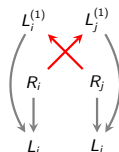
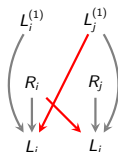
1. Define the **estimand** (often done in the absence of missing data).
  - ▶ Some function of target law  $p(Z, L^{(1)})$  or full law  $p(Z, L^{(1)}, R)$ .
2. Assume a **model** that links the counterfactual, factual, and missingness indicator variables.
  - ▶ Use **Directed Acyclic Graphs (DAGs)** to encode the modeling assumptions.
3. Determine whether the estimand is **identifiable** in the assumed model.
  - ▶ Focus on **identification of the target and full laws**.
4. If estimand is identifiable, find the best **estimation** strategy, and if it is not, perhaps stronger assumptions are needed (or alternatively obtaining bounds).
5. Conduct **sensitivity analysis** to reflect on the assumptions.

# Introducing missing data DAGs

- ▶ Define missing data models via restrictions on the full data distribution, that can be represented by a DAG (similar to causal inference).
- ▶ In missing data DAGs: (Mohan et al., 2013)
  1. Observed and counterfactual variables appear on the same graph
  2. There are certain edge restrictions: (marked in red)



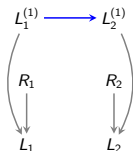
deterministic edges



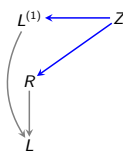
"no interference"

- ▶ The "no interference" assumption can be relaxed (Srinivasan et al., 2023).

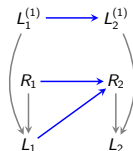
# Examples: missing data DAGs



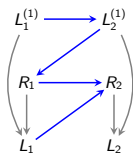
MCAR



MAR

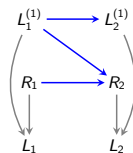


MAR



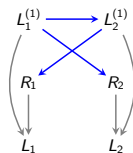
MNAR: permutation

(Robins, 1997)



MNAR: block-conditional

(Zhou et al., 2010)



MNAR: block-parallel

(Mohan et al., 2013)

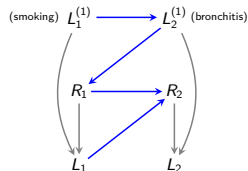
## Examples: missing data DAGs ctd.

How models differ in telling a story about the missingness mechanisms?

- ▶  $L_1^{(1)}$ : true smoking status of an individual
- ▶  $L_2^{(1)}$ : diagnosis of bronchitis.
- ▶  $R_1, R_2$ : encode whether these variables have been measured or not.

- ▶  $L_2^{(1)} \rightarrow R_1$   
A doctor inquires about the patient's smoking status on a suspected diagnosis of bronchitis before administering the test.

- ▶  $R_1 \rightarrow R_2 \leftarrow L_1$   
Whether the true bronchitis status is measured via a diagnostic test depends on the doctor's awareness of the individual's smoking status ( $R_1$ ) and their observed value of smoking ( $L_1$ ).



MNAR: permutation

(Robins, 1997)

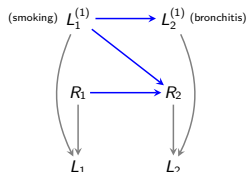
## Examples: missing data DAGs ctd.

How models differ in telling a story about the missingness mechanisms?

- ▶  $L_1^{(1)}$ : true smoking status of an individual
- ▶  $L_2^{(1)}$ : diagnosis of bronchitis.
- ▶  $R_1, R_2$ : encode whether these variables have been measured or not.

- ▶  $R_1$  has no parent  
Inquiry into smoking status is random (e.g., as in random screening programs or surveys).

- ▶  $R_1 \rightarrow R_2 \leftarrow L_1^{(1)}$   
Administration of a diagnostic test depends on the inquiry into smoking, as well as the potentially unobserved past history of smoking.



MNAR: block-conditional

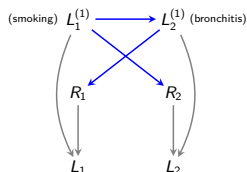
(Zhou et al., 2010)

## Examples: missing data DAGs ctd.

How models differ in telling a story about the missingness mechanisms?

- ▶  $L_1^{(1)}$ : true smoking status of an individual
- ▶  $L_2^{(1)}$ : diagnosis of bronchitis.
- ▶  $R_1, R_2$ : encode whether these variables have been measured or not.

- ▶  $R_1 \leftarrow L_2^{(1)}$   
Inquiry into smoking status depends on a suspected diagnosis of bronchitis.
- ▶  $R_2 \leftarrow L_1^{(1)}$   
Administration of the diagnostic test depends on the suspected smoking status of an individual.



MNAR: block-parallel

(Mohan et al., 2013)

# Missing data DAG models

- ▶ Denote the missing data DAG (m-DAG) defined over  $V = (Z, L^{(1)}, R, L)$  via  $\mathcal{G}(V)$ .
- ▶ The statistical model of a m-DAG is a set of distributions that factorize as:

$$\begin{aligned} p(Z, L^{(1)}, R, L) &= \prod_{V_i \in V} p(V_i \mid \text{pa}_{\mathcal{G}}(V_i)) \\ &= \prod_{V_i \in V \setminus L} p(V_i \mid \text{pa}_{\mathcal{G}}(V_i)) \times \prod_{L_i \in L} p(L_i \mid L_i^{(1)}, R_i). \end{aligned}$$

- ▶ Familiar concepts like d-separation and Markov properties carry over.
  - ▶ **Factorization:** probability distribution as a set of small factors.
  - ▶ **Local Markov property:** a small but complete set of indep. constraints.

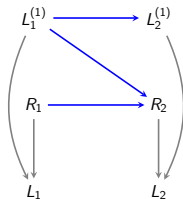
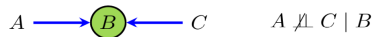
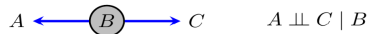
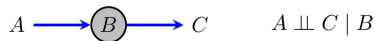
$$V_i \perp\!\!\!\perp \text{nd}_{\mathcal{G}}(V_i) \setminus \text{pa}_{\mathcal{G}}(V_i) \mid \text{pa}_{\mathcal{G}}(V_i), \quad \forall V_i \in V.$$

- ▶ **Global Markov property:** all independence constraints in the model.

$$\text{Given } X, Y, Z \in V : \quad (X \perp\!\!\!\perp_{\text{d-sep}} Y \mid Z)_{\mathcal{G}(V)} \implies (X \perp\!\!\!\perp Y \mid Z)_{p(V)}.$$

- ▶ All three properties are equivalent.

# d-separation refresher

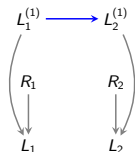


$$R_1 \perp\!\!\!\perp L_1^{(1)}$$

$$R_1 \not\perp\!\!\!\perp L_1^{(1)} \mid R_2$$



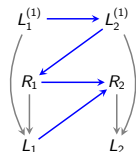
# Examples: m-DAG models



MCAR

$$R_1 \perp\!\!\!\perp R_2, L_1^{(1)}, L_2^{(1)}$$

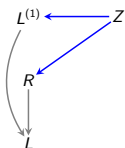
$$R_2 \perp\!\!\!\perp R_1, L_1^{(1)}, L_2^{(1)}$$



Permutation (Robins, 1997)

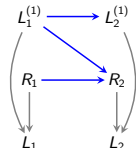
$$R_1 \perp\!\!\!\perp L_1^{(1)} \mid L_2^{(1)}$$

$$R_2 \perp\!\!\!\perp L_1^{(1)}, L_2^{(1)} \mid R_1, L_1$$



MAR

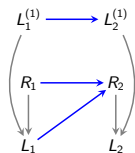
$$R \perp\!\!\!\perp L^{(1)} \mid Z$$



Block-conditional (Zhou et al., 2010)

$$R_1 \perp\!\!\!\perp L_1^{(1)}, L_2^{(1)}$$

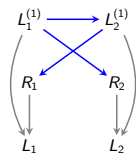
$$R_2 \perp\!\!\!\perp L_1^{(1)} \mid R_1, L_1^{(1)}$$



MAR

$$R_1 \perp\!\!\!\perp L_1^{(1)}, L_2^{(1)}$$

$$R_2 \perp\!\!\!\perp L_1^{(1)}, L_2^{(1)} \mid R_1, L_1$$



Block-conditional (Mohan et al., 2013)

$$R_1 \perp\!\!\!\perp R_2, L_1^{(1)} \mid L_2^{(1)}$$

$$R_2 \perp\!\!\!\perp R_1, L_2^{(1)} \mid L_1^{(1)}$$

# Graphical representations of MCAR, MAR, MNAR mechanisms

- ▶ **MCAR**: if  $\prod_{R_k \in R} p(R_k \mid \text{pa}_{\mathcal{G}}(R_k))$  is not a function of variables in  $L^{(1)}$  and  $L$ .
  - ▶ Graphically speaking, there are no edges that point to variables in  $R$ .
- ▶ **MAR**: if  $\prod_{R_k \in R} p(R_k \mid \text{pa}_{\mathcal{G}}(R_k))$  is not a function of variables in  $L^{(1)}$ .
  - ▶ Graphically speaking, there are no edges from variables in  $L^{(1)}$  to variables in  $R$ .
- ▶ **MNAR** otherwise.

## **Part 2. Nonparametric Identification**

# Identification in missing data models

- ▶ Let  $\psi := \mathbb{E}[h(p(Z, L^{(1)}))]$  denote the parameter (estimand) of interest.
- ▶ Let the full law  $p(Z, L^{(1)}, R)$  be Markov relative to the m-DAG  $\mathcal{G}(V)$ .
- ▶ To do *inference* on  $\psi$ , we first need to argue whether  $\psi$  is *identified* as a function of the observed data law in the assumed m-DAG or not?
- ▶ The estimand  $\psi$  is identified in the assumed m-DAG  $\mathcal{G}$ , if it can be expressed as a unique function of the observed data law  $p(Z, L, R)$ . This means:
  - ▶ A parameter is identified under a particular collection of assumptions if these assumptions imply that the distribution of the **observed data is compatible with a single value** of the parameter.
  - ▶ If there exists no unique mapping between the counterfactual distribution and the observed data law, then the parameter is **not identified**.
- ▶ Instead of  $\psi$ , we might be interested in identification of the entire target law  $p(Z, L^{(1)})$  or the entire full law  $p(Z, L^{(1)}, R)$ .

## Example: identification

Is the target law  $p(Z, L^{(1)})$  identified as a function of the observed data law  $p(Z, R, L)$ ?

- ▶ Under MCAR missingness, target law is identified:

$$\begin{aligned} p(Z, L^{(1)}) &= p(Z, L^{(1)} \mid R = 1) && R \perp\!\!\!\perp Z, L^{(1)} \\ &= p(Z, L \mid R = 1). && \text{consistency} \end{aligned}$$

- ▶ Under MAR missingness, target law is identified:

$$\begin{aligned} p(Z, L^{(1)}) &= p(Z) \times p(L^{(1)} \mid Z) \\ &= p(Z) \times p(L^{(1)} \mid Z, R = 1) && R \perp\!\!\!\perp L^{(1)} \mid Z \\ &= p(Z) \times p(L \mid Z, R = 1). && \text{consistency} \end{aligned}$$

- ▶ MNAR model:

$$p(Z, L^{(1)}) = ???$$

MNAR models  $\equiv$  causal models with unmeasured confounding

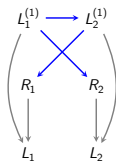
- ▶ Sometimes we succeed and sometimes we fail!

# Nonparametric identification theory in causal inference

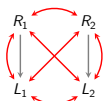
- ▶ Identification questions in causal inference: given an arbitrary DAG with hidden/unmeasured variables, is  $p(Y^{(t)})$  identified?
- ▶ **Sound** and **complete** algorithms exist for identification of causal effects
  - ▶ Soundness: functionals of identified effects are correct.
  - ▶ Completeness: no-identifiability of the causal estimand is *provable*.  
(Shpitser and Pearl, 2006; Huang and Valtorta, 2006; Richardson et al., 2017; Bhattacharya et al., 2020)
- ▶ Similarly, given that assumptions/restrictions in a missing data model are encoded via a m-DAG, we would like to know whether the underlying full/target law is identified or not.
- ▶ **Causal identification theory is *incomplete* for missing data identification!**
  - ▶ There are indeed identified MNAR models for which causal identification theory fails.

# Incompleteness of causal identification theory for m-DAGs

- Causal identification theory is **incomplete** for missing data identification.



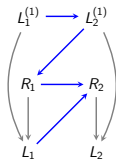
Block-parallel



Observed margin

One-line ID (Richardson et al., 2017)

- ▶  $Y^* = \{L_1, L_2\}$
- ▶  $G_{Y^*} = L_1 \leftrightarrow L_2$
- ▶ District:  $\{L_1, L_2\}$
- ▶ Need to fix  $R_1, R_2$  and fail.



Permutation



Observed margin

One-line ID (Richardson et al., 2017)

- ▶  $Y^* = \{L_1, L_2\}$
- ▶  $G_{Y^*} = L_1 \leftrightarrow L_2$
- ▶ District:  $\{L_1, L_2\}$
- ▶ Need to fix  $R_1, R_2$  and fail.

# Nonparametric identification in m-DAGs

- ▶ The target law is identified **if and only if** the missingness mechanism  $p(R = 1 \mid L^{(1)}, Z)$  is identified. Using Bayes rule:

$$p(R = 1 \mid L^{(1)}, Z) = \frac{p(Z, L^{(1)}, R = 1)}{p(L^{(1)}, Z)} \rightarrow p(Z, L^{(1)}) = \frac{p(Z, L^{(1)}, R = 1)}{p(R = 1 \mid Z, L^{(1)})}.$$

- ▶ The full law is identified **if and only if** the missingness mechanism  $p(R = r \mid L^{(1)}, Z)$  is identified, for all possible missingness pattern  $R = r$ . Using chain rule:

$$\begin{aligned} p(Z, L^{(1)}, R = r) &= p(Z, L^{(1)}) \times p(R = r \mid L^{(1)}, Z) \\ &= \frac{p(Z, L^{(1)}, R = 1)}{p(R = 1 \mid L^{(1)}, Z)} \times p(R = r \mid L^{(1)}, Z). \end{aligned}$$

- ▶ Game plan: focus on identification of the missingness mechanism  $p(R \mid L^{(1)}, Z)$  in a given m-DAG.



# Np-identification of missingness mechanisms in m-DAGs

Given an m-DAG, is the missingness mechanism identified or not? We look at two different parameterizations of  $p(R \mid L^{(1)}, Z)$ :

(i) m-DAG factorization: (Bhattacharya et al., 2019)

$$p(R \mid \text{pa}_{\mathcal{G}}(R)) = \prod_{R_k \in R} p(R_k \mid \text{pa}_{\mathcal{G}}(R_k))$$

Identify each **propensity score**  $p(R_k \mid \text{pa}_{\mathcal{G}}(R_k))$ , for all  $R_k \in R$ .

(ii) Odds ratio parameterization: (Chen, 2007; Malinsky et al., 2021)

$$\prod_{k=1}^K p(R_k \mid R_{-k} = 1, \text{pa}_{\mathcal{G}}(R)) \times \prod_{k=2}^K \text{OR}(R_k, R_{\prec k} \mid R_{\succ k} = 1, \text{pa}_{\mathcal{G}}(R)),$$

where  $R_{-k} = R \setminus R_k$ ,  $R_{\prec k} = \{R_1, \dots, R_{k-1}\}$ ,  $R_{\succ k} = \{R_{k+1}, \dots, R_K\}$ .

Identify each **univariate conditionals** and **pairwise odds ratio** terms.

# Identification arguments

1. m-DAG factorization of the missingness mechanism
2. Odds ratio parameterization of the missingness mechanism

# Identification via m-DAG factorization

- ▶ We would like to know whether and how full/target law is identified in a given m-DAG  $\mathcal{G}$ .
  - ▶ Full law: argue for identification of  $p(R = r \mid \text{pa}_{\mathcal{G}}(R))$ ,  $\forall r$ .
  - ▶ Target law: argue for identification of  $p(R = 1 \mid \text{pa}_{\mathcal{G}}(R))$ .
- ▶ Target law is identified if each propensity score  $p(R_k \mid \text{pa}_{\mathcal{G}}(R_k))$  is identified when evaluated at  $R = 1$

$$p(R = 1 \mid \text{pa}_{\mathcal{G}}(R)) = \prod_{R_k \in R} p(R_k \mid \text{pa}_{\mathcal{G}}(R_k)) \Big|_{R=1}.$$

- ▶ Full law is identified if each propensity score  $p(R_k \mid \text{pa}_{\mathcal{G}}(R_k))$  is identified at all levels of  $R_i \in \text{pa}_{\mathcal{G}}(R_k)$ .

$$p(R = r \mid \text{pa}_{\mathcal{G}}(R)) = \prod_{R_k \in R} p(R_k \mid \text{pa}_{\mathcal{G}}(R_k)) \Big|_{R=r}$$

- ▶ There are two major ideas for propensity scores identification:
  1. *Associational* irrelevancy: d-separation
  2. *Causal* irrelevancy: invariance property

# 1. Associational irrelevancy

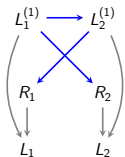
- ▶ In order to identify the propensity score of  $R_k$ ,  $p(R_k \mid \text{pa}_{\mathcal{G}}(R_k))$ , we need to **select on** the following missingness indicators:

$$R_k^s = \left\{ R_i \in R \mid L_i^{(1)} \in \text{pa}_{\mathcal{G}}(R_k) \right\} \quad (\text{selection set for } R_k)$$

- ▶ For any  $R_i \in R_k^s$ ,  $R_i$  is either a **descendant** of  $R_k$  or a **non-descendant** of  $R_k$ .
- ▶ If  $R_i$  is a **non-descendant** of  $R_k$ , then we can apply the local Markov property which states  $R_k \perp\!\!\!\perp \text{nd}_{\mathcal{G}}(R_k) \setminus \text{pa}_{\mathcal{G}}(R_k) \mid \text{pa}_{\mathcal{G}}(R_k)$ .
  - ▶ So we can include  $R_i = 1$  in the conditioning set and replace  $L_i^{(1)}$  with  $\{L_i, R_i = 1\}$ .
- ▶ Formally, for any  $R_i \in R_k^s \cap \text{nd}_{\mathcal{G}}(R_k)$ , we can write:

$$p(R_k \mid \text{pa}_{\mathcal{G}}(R_k)) \Big|_{R=1} = p(R_k \mid \underbrace{\text{pa}_{\mathcal{G}}(R_k)}_{\text{includes } L_i^{(1)}}, R_i = 1) \Big|_{R=1}.$$

## Example 1/2: associational irrelevancy (block-parallel)



- Is the full/target law identified in the **block-parallel** model?

$$p(R \mid \text{pa}_{\mathcal{G}}(R)) = p(R_1 \mid \text{pa}_{\mathcal{G}}(R_1)) \times p(R_2 \mid \text{pa}_{\mathcal{G}}(R_2))$$

Identification of  $p(R_1 \mid \text{pa}_{\mathcal{G}}(R_1))$

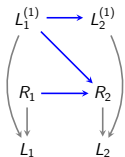
$$\begin{aligned} p(R_1 \mid \text{pa}_{\mathcal{G}}(R_1)) &= p(R_1 \mid L_2^{(1)}) \\ &= p(R_1 \mid R_2 = 1, L_2^{(1)}) && R_1 \perp\!\!\!\perp R_2 \mid L_2^{(1)} \\ &= p(R_1 \mid R_2 = 1, L_2) && \text{consistency} \end{aligned}$$

Identification of  $p(R_2 \mid \text{pa}_{\mathcal{G}}(R_2))$

$$\begin{aligned} p(R_2 \mid \text{pa}_{\mathcal{G}}(R_2)) &= p(R_2 \mid L_1^{(1)}) \\ &= p(R_2 \mid R_1 = 1, L_1^{(1)}) && R_2 \perp\!\!\!\perp R_1 \mid L_1^{(1)} \\ &= p(R_2 \mid R_1 = 1, L_1) && \text{consistency} \end{aligned}$$

- So  $p(R \mid \text{pa}_{\mathcal{G}}(R))$  is ID, which means the full and target laws are both ID.

## Example 2/2: associational irrelevancy (block-conditional MAR)



- Is the full/target law identified in the **block-conditional MAR** model?

$$p(R \mid \text{pa}_{\mathcal{G}}(R)) = p(R_1 \mid \text{pa}_{\mathcal{G}}(R_1)) \times p(R_2 \mid \text{pa}_{\mathcal{G}}(R_2))$$

Identification of  $p(R_2 \mid \text{pa}_{\mathcal{G}}(R_2))$

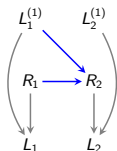
$$\begin{aligned} p(R_2 \mid \text{pa}_{\mathcal{G}}(R_2)) &= p(R_2 \mid R_1, L_1^{(1)}) \\ &= ??? \end{aligned}$$

$$\begin{aligned} p(R_2 \mid \text{pa}_{\mathcal{G}}(R_2)) \Big|_{R=1} &= p(R_2 = 1 \mid R_1 = 1, L_1^{(1)}) \\ &= p(R_2 = 1 \mid R_1 = 1, L_1) \quad \text{consistency} \end{aligned}$$

- So only  $p(R_2 \mid R_1 = 1, L_1^{(1)})$  is identified.
- Since  $p(R_1 \mid \text{pa}_{\mathcal{G}}(R_1)) = p(R_1)$  is also identified, we conclude that the target law is ID.
- However, the full law might **NOT** be identified. We have to prove that the full law is not identified.

# A non-identified structure

- ▶ Claim:  $p(R_2 \mid R_1 = 0, L_1^{(1)})$  is not ID.



- 
- ▶ Assume binary data. So the full law  $p(R, L^{(1)})$  would have 7 parameters.
  - ▶ 5 identified parameters:
    - ▶  $p(L_1^{(1)} = h_1) = p(L_1^{(1)} = h_1 \mid R_1 = 1) = p(L_1 = h_1 \mid R_1 = 1)$
    - ▶  $p(L_2^{(1)} = h_2) = p(L_2^{(1)} = h_2 \mid R_2 = 1) = p(L_2 = h_2 \mid R_2 = 1)$
    - ▶  $p(R_1 = r_1)$
    - ▶  $p(R_2 = r_2 \mid R_1 = 1, L_1^{(1)} = h_1) = p(R_2 = r_2 \mid R_1 = 1, L_1 = h_1)$
  - ▶ 2 unidentified parameters:
    - ▶  $p(R_2 = r_2 \mid R_1 = 0, L_1^{(1)} = h_1), h_1 \in \{0, 1\}$

# Proofing non-identifiability claims

$R_1$	$p(R_1)$	$L_1^{(1)}$	$p(L_1^{(1)})$	$L_2^{(1)}$	$p(L_2^{(1)})$	$R_2$	$R_1$	$L_1^{(1)}$	$p(R_2   R_1, L_1^{(1)})$	
0	$a$	0	$b$	0	$c$	0	0	0	$d$	
1	$1 - a$	1	$1 - b$	1	$c$	1	0	0	$1 - d$	
						0	1	0	$e$	
						1	1	0	$1 - e$	
						0	0	1	$f$	
						1	0	1	$1 - f$	
						0	1	1	$g$	
						1	1	1	$1 - g$	

$R_1$	$R_2$	$L_1^{(1)}$	$L_2^{(1)}$	$p(\text{FULL LAW})$	$L_1$	$L_2$	$p(\text{OBSERVED LAW})$
0	0	0	0	$abcd$	?	?	$a[db + f(1 - b)]$
		1	0	$af(1 - b)c$			
		0	1	$adb(1 - c)$			
		1	1	$af(1 - b)(1 - c)$			
1	0	0	0	$(1 - a)ebc$	0	?	$(1 - a)eb$
		1	0	$(1 - a)g(1 - b)c$			
		0	1	$(1 - a)eb(1 - c)$			
		1	1	$(1 - a)g(1 - b)(1 - c)$			
0	1	0	0	$a(1 - d)bc$	?	0	$ac[1 - (db + f(1 - b))]$
		1	0	$a(1 - f)(1 - b)c$			
		0	1	$a(1 - d)b(1 - c)$			
		1	1	$a(1 - f)(1 - b)(1 - c)$			
1	1	0	0	$(1 - a)(1 - e)bc$	0	0	$(1 - a)(1 - e)bc$
		1	0	$(1 - a)(1 - g)(1 - b)c$			
		0	1	$(1 - a)(1 - e)b(1 - c)$			
		1	1	$(1 - a)(1 - g)(1 - b)(1 - c)$			

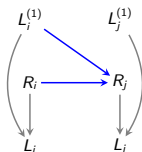
We can pick any  $\{d, f\}$  as long as  $bd + (1 - b)f$  stays the same.



# Collider: non-identified structure

## Definition (Collider)

If  $\exists R_i, R_j \in R$  such that  $R_i \rightarrow R_j \leftarrow L_i^{(1)}$ , then a special collider structure forms at  $R_j$ , referred to as collider.

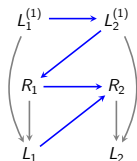


## Lemma (Collider non-identification)

If  $\exists R_i, R_j \in R$  such that  $R_i \rightarrow R_j \leftarrow L_i^{(1)}$  then  $p(R_j \mid \text{pa}_{\mathcal{G}}(R_j) \setminus R_i, R_i = 0)$  is not identified (Bhattacharya et al., 2019).

The above result means that whenever we spot a collider on the m-DAG, we can immediately conclude that the full law is not identified.

# Associational irrelevancy: limitations



- Is the full/target law identified in the **permutation** model?

$$p(R \mid L^{(1)}) = p(R_1 \mid \text{pa}_{\mathcal{G}}(R_1)) \times p(R_2 \mid \text{pa}_{\mathcal{G}}(R_2))$$

## Identification of $p(R_i \mid \text{pa}_{\mathcal{G}}(R_i))$

$$p(R_2 \mid \text{pa}_{\mathcal{G}}(R_2)) = p(R_2 \mid R_1, L_1) \quad \checkmark$$

$$p(R_1 \mid \text{pa}_{\mathcal{G}}(R_1)) = p(R_1 \mid L_2^{(1)}) \quad R_1 \not\perp R_2 \mid L_2^{(1)}$$

- What does this mean? Is the propensity score of  $R_1$  not identified? which would then imply the full law is not ID.
- To answer this question, we either need to prove the propensity score is not identified or find a way to identify it.

## 2. Causal irrelevancy

### Invariance property:

- ▶ Given the propensity score for  $R_k \in R$ , the conditioning set  $\text{pa}_{\mathcal{G}}(R_k)$  captures all the direct causes of  $R_k$ . Hence, it remains invariant to any set of interventions that disrupts other parts of the full law.
- ▶ Formally, given  $R^* \subseteq R \setminus R_k$ , we have

$$p(R_k \mid \text{pa}_{\mathcal{G}}(R_k)) = p(R_k \mid \text{pa}_{\mathcal{G}}(R_k), \text{do}(R^* = 1)).$$

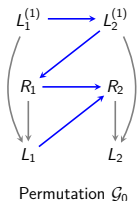
- ▶ Due to *invariance* property of the propensity scores, we can sometimes succeed in propensity score identification by exploring interventional distributions where a subset of variables are intervened on.

# Intervention operation on missingness indicators

An intervention that sets  $R_k = 1$ , entails the following changes:

- ▶ **Graphical** changes to the missing data DAG  $\mathcal{G}(V)$ 
  - ▶ In the corresponding m-DAG  $\mathcal{G}$ , delete all the incoming edges into  $R_k$  and fix  $R_k$  to take value 1, and
  - ▶ Treat the counterfactual variable  $L_k^{(1)}$  as  $L_k$  (by consistency).
- ▶ **Probabilistical** changes to the joint distribution  $p(V)$ 
  - ▶ In the corresponding joint law  $p(V)$ , drop the term  $p(R_k \mid \text{pa}_{\mathcal{G}}(R_k))$  from the factorization of  $p(V)$ , and evaluate all terms at  $R_k = 1$ .

## Example: causal irrelevancy

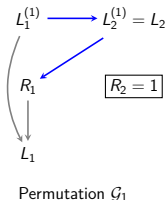


$$\mathcal{G}_0 : R_1 \not\perp R_2 \mid L_2^{(1)}$$

Invariance property:

$$p(R_1 \mid L_2^{(1)}) = p(R_1 \mid L_2^{(1)}, \text{do}(R_2 = 1))$$

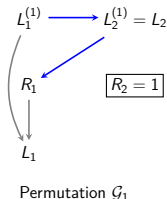
Graphical and probabilistical changes after intervening on  $R_2$ :



$$p(R_1, L_1^{(1)}, L_2^{(1)}, L_1 \mid \text{do}(R_2 = 1)) = \frac{p(R_1, L_1^{(1)}, L_2^{(1)}, L_1, R_2)}{p(R_2 \mid R_1, L_1)} \Big|_{R_2=1}$$

- The propensity score of  $R_1$  is identified from the above intervention dist.

## Example ctd. causal irrelevancy



$$p(R_1, L_1^{(1)}, L_2^{(1)}, L_1 \mid \text{do}(R_2 = 1)) = \frac{p(R_1, L_1^{(1)}, L_2, L_1, R_2 = 1)}{p(R_2 = 1 \mid R_1, L_1)}$$

$$p(R_1, L_2^{(1)} \mid \text{do}(R_2 = 1)) = \sum_{l_1} \frac{p(R_1, l_1, L_2, R_2 = 1)}{p(R_2 = 1 \mid R_1, l_1)}.$$

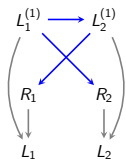
- The propensity score of  $R_1$  is identified from the above intervention dist.

$$\begin{aligned} p(R_1 \mid L_2^{(1)}) &= p(R_1 \mid L_2^{(1)}, \text{do}(R_2 = 1)) \\ &= \frac{p(R_1, L_2^{(1)} \mid \text{do}(R_2 = 1))}{p(L_2^{(1)} \mid \text{do}(R_2 = 1))}. \end{aligned}$$

First equality holds by the invariance property and second holds by Bayes rule.

# Order of interventions

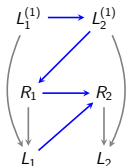
- Target law is ID via **parallel interventions** on  $R_1$  and  $R_2$ .



Block-parallel

$$\begin{aligned}
 p(L^{(1)}) &= \frac{p(L^{(1)}, R)}{p(R_1 \mid \text{pa}_G(R_1)) \times p(R_2 \mid \text{pa}_G(R_2))} \Big|_{R=1} \\
 &= \frac{p(L, R=1)}{p(R_1=1 \mid R_2=1, L_2) \times p(R_2=1 \mid R_1=1, L_1)}.
 \end{aligned}$$

- Target law ID is obtained via **sequential interventions** on first  $R_2$  and then  $R_1$ .

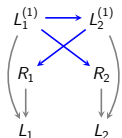


Permutation

$$\begin{aligned}
 p(L^{(1)}) &= \frac{p(L^{(1)}, R)}{p(R_1 \mid \text{pa}_G(R_1)) \times p(R_2 \mid \text{pa}_G(R_2))} \Big|_{R=1} \\
 &= \frac{p(L, R=1)}{p(R_1=1 \mid L_2^{(1)}, \text{do}(R_2=1)) \times p(R_2=1 \mid R_1=1, L_1)}.
 \end{aligned}$$

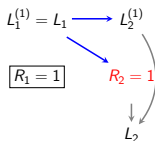
# Main identification challenge: selection bias

Can we apply the causal-irrelevancy idea to the block-parallel model?



$$\begin{aligned}
 p(R_2 \mid L_1^{(1)}) &= p(R_2 \mid L_1^{(1)}, \text{do}(R_1 = 1)) \\
 &= \frac{p(R_2, L_1^{(1)} \mid \text{do}(R_1 = 1))}{p(L_1^{(1)} \mid \text{do}(R_1 = 1))}
 \end{aligned}$$

An intervention on  $R_1$  implies:



$$p(R_2, L_1^{(1)}, L_2^{(1)}, L_2 \mid \text{do}(R_1 = 1)) = \frac{p(L_1, R_1 = 1, L_2^{(1)}, L_2, R_2)}{p(R_1 = 1 \mid L_2^{(1)})}$$

We can only evaluate the above expression when  $R_2 = 1$ :

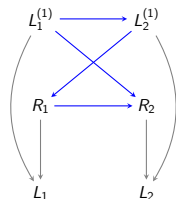
$$p(L_1^{(1)}, \textcolor{red}{R_2} = 1 \mid \text{do}(R_1 = 1)) = \sum_{L_2^{(1)}} \frac{p(L_2^{(1)}, R_1 = 1, \textcolor{red}{R_2} = 1)}{p(R_1 = 1 \mid L_2^{(1)}, \textcolor{red}{R_2} = 1)}$$

Intervening on  $R_1$  **induces a selection** on  $R_2$ .



# Inevitable selection bias

- ▶ The lesson is that **sequential interventions** do not help with identification arguments in the *block-parallel* model. Indeed, we need **parallel interventions** to dodge the selection bias issue.
- ▶ Sometimes we cannot avoid the selection bias and end up with unidentified distributions. An example of this is the so-called *criss-cross* model,



Criss-cross (Nabi and Bhattacharya, 2022)

- ▶ Full law  $p(L^{(1)}, R)$  is not identified because of the collider at  $R_2$ .
- ▶ Target law  $p(L^{(1)})$  is also **provably** not identified (Nabi and Bhattacharya, 2022).

# Partial orders of interventions

- ▶ Sufficient rules for identification: finding **valid partial orders** of interventions that avoid the issue of selection bias.
  - ▶ That is a combination of sequential and parallel interventions (as opposed to a *total order* in causal inference).
- ▶ **Dodging selection bias** requires:
  - ▶ Set interventions
  - ▶ Intervening on variables other than  $R$
  - ▶ Interventions on margins of  $\mathcal{G}$  (pseudo-propensity scores)
- ▶ See Bhattacharya et al. (2019) and Nabi et al. (2022) for more discussions.

# Identification arguments

1. m-DAG factorization of the missingness mechanism
2. Odds ratio parameterization of the missingness mechanism

# Odds ratio parameterization (Chen, 2007)

Given disjoint sets of variables  $A, B, C$ , and reference values  $a_0, b_0$  :

$$p(A, B \mid C) = \frac{1}{\mathcal{Z}(C)} \times p(A \mid B = b_0, C) \times p(B \mid A = a_0, C) \times \text{OR}(A, B \mid C),$$

where

$$\begin{aligned} \text{OR}(A = a, B = b \mid C) &= \frac{p(A = a \mid B = b, C)}{p(A = a_0 \mid B = b, C)} \times \frac{p(A = a_0 \mid B = b_0, C)}{p(A = a \mid B = b_0, C)} \\ &= \frac{p(B = b \mid A = a, C)}{p(B = b_0 \mid A = a, C)} \times \frac{p(B = b_0 \mid A = a_0, C)}{p(B = b \mid A = a_0, C)} \end{aligned}$$

$$\mathcal{Z}(C) = \sum_{a,b} p(A = a \mid B = b_0, C) \times p(B = b \mid A = a_0, C) \times \text{OR}(A = a, B = b \mid C).$$

- 
- ▶ It is symmetric:  $\text{OR}(A, B \mid C) = \text{OR}(B, A \mid C)$
  - ▶  $\text{OR}(A = a_0, B \mid C) = \text{OR}(A, B = b_0 \mid C) = \text{OR}(A = a_0, B = b_0 \mid C) = 1$

# Identification via odds ratio parameterization

- ▶ m-DAG factorization view for identification:

$$p(R_i, R_j \mid \text{pa}_{\mathcal{G}}(R_i, R_j)) = p(R_i \mid \text{pa}_{\mathcal{G}}(R_i)) \times p(R_j \mid \text{pa}_{\mathcal{G}}(R_j))$$

Identify:

- ▶  $p(R_i \mid \text{pa}_{\mathcal{G}}(R_i))$ ,
- ▶  $p(R_j \mid \text{pa}_{\mathcal{G}}(R_j))$ .

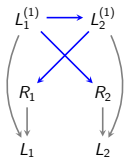
- ▶ Odds ratio parameterization view for identification:

$$p(R_i, R_j \mid L^{(1)}) = \frac{1}{\mathcal{Z}(L^{(1)})} \times p(R_i \mid R_j = 1, L^{(1)}) \times p(R_j \mid R_i = 1, L^{(1)}) \times \text{OR}(R_i, R_j \mid L^{(1)})$$

Identify:

- ▶  $p(R_i \mid R_j = 1, L^{(1)})$ ,
- ▶  $p(R_j \mid R_i = 1, L^{(1)})$ ,
- ▶  $\text{OR}(R_i = 0, R_j = 0 \mid L^{(1)})$ .

## Example: block-parallel model



► Is the full/target law identified in the **block-parallel** model?

$$p(R \mid L_1^{(1)}) = \frac{1}{\mathcal{Z}(L^{(1)})} \times \underbrace{p(R_1 \mid R_2 = 1, L^{(1)})}_{(1)} \times \underbrace{p(R_2 \mid R_1 = 1, L^{(1)})}_{(2)} \times \underbrace{\text{OR}(R_1, R_2 \mid L^{(1)})}_{(3)}$$

Note that  $R_1 \perp\!\!\!\perp L_1^{(1)} \mid R_2, L_2^{(1)}$  and  $R_2 \perp\!\!\!\perp L_2^{(1)} \mid R_1, L_1^{(1)}$ :

$$(1) : p(R_1 \mid R_2 = 1, L^{(1)}) = p(R_1 \mid R_2 = 1, L_2^{(1)}) = p(R_1 \mid R_2 = 1, L_2)$$

$$(2) : p(R_2 \mid R_1 = 1, L^{(1)}) = p(R_2 \mid R_1 = 1, L_1^{(1)}) = p(R_2 \mid R_1 = 1, L_1)$$

$$(3) : \text{OR}(R_1, R_2 \mid L^{(1)}) = 1.$$

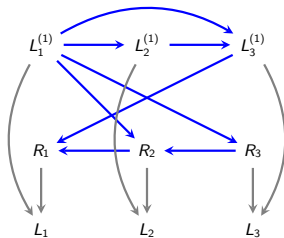
Yes, both full and target laws are identified.

# Why the two parameterizations?

There exist identified MANR models where:

- ▶ *m-DAG factorization* approach fails to identify the model, but the *odds ratio parameterization* approach succeeds.
- ▶ *Odds ratio parameterization* approach fails to identify the model, but the *m-DAG factorization* approach succeeds.

## Example: m-DAG factorization fails!

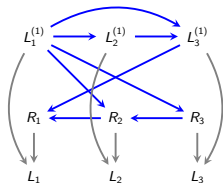


Is the target/full law identified?

$$\begin{aligned} p(R \mid L^{(1)}) &= \prod_{R_i \in R} p(R_i \mid \text{pa}_{\mathcal{G}}(R_i)) \\ &= p(R_1 \mid R_2, L_3^{(1)}) \times p(R_2 \mid R_3, L_1^{(1)}) \times p(R_3 \mid L_1^{(1)}). \end{aligned}$$



## Example ctd. identification of the propensity score of $R_1$

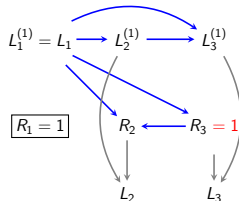
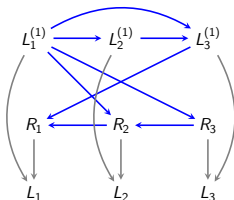


$$\begin{aligned}
 p(R \mid L^{(1)}) &= \prod_{R_i \in R} p(R_i \mid \text{pa}_{\mathcal{G}}(R_i)) \\
 &= p(R_1 \mid R_2, L_3^{(1)}) \times p(R_2 \mid R_3, L_1^{(1)}) \times p(R_3 \mid L_1^{(1)}).
 \end{aligned}$$

### Nonparametric identification of $p(R_1 \mid \text{pa}_{\mathcal{G}}(R_i))$

$$\begin{aligned}
 p(R_1 \mid \text{pa}_{\mathcal{G}}(R_1)) &= p(R_1 \mid R_2, L_3^{(1)}) \\
 &= p(R_1 \mid R_3 = 1, R_2, L_3^{(1)}) && R_1 \perp\!\!\!\perp R_3 \mid R_2, L_3^{(1)} \\
 &= p(R_1 \mid R_3 = 1, R_2, L_3) && \text{consistency.}
 \end{aligned}$$

## Example ctd. identification of the propensity score of $R_2$

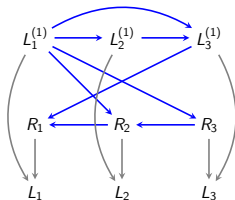


$$p^* = p(L^{(1)}, R_1, R_3 \mid \text{do}(R_1 = 1)) = \frac{p(L^{(1)}, R)}{p(R_1 = 1 \mid R_3 = 1, R_2, L_3^{(1)})}.$$

Nonparametric identification of  $p(R_2 \mid \text{pa}_{\mathcal{G}}(R_i))|_{R=1}$

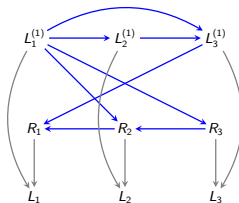
$$\begin{aligned} p(R_2 = 1 \mid \text{pa}_{\mathcal{G}}(R_2))|_{R=1} &= p(R_2 = 1 \mid R_3 = 1, L_1^{(1)}) \\ &= p(R_2 = 1 \mid R_3 = 1, L_1^{(1)}, \text{do}(R_1 = 1)) && \text{causal irrelevance} \\ &= p^*(R_2 = 1 \mid R_3 = 1, L_1^{(1)}) \\ &= p^*(R_2 = 1 \mid R_1 = 1, R_3 = 1, L_1^{(1)}) && R_2 \perp\!\!\!\perp_{\mathcal{G}^*} R_1 \mid R_3, L_1^{(1)} \\ &= p^*(R_2 = 1 \mid R_1 = 1, R_3 = 1, L_1) && \text{consistency} \end{aligned}$$

## Example ctd. identification of the propensity score of $R_3$



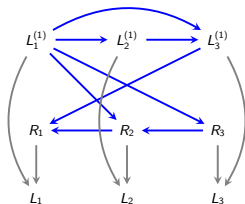
- ▶ Unfortunately, similar tricks do not help with identification of  $p(R_3 \mid \text{pa}_G(R_3))$  due to selection bias on  $R_3$  from intervening on either  $R_1$  or  $R_2$ .
- ▶ It seems that the missingness mechanism is not identified. Thus, it seems neither the full law nor the target law are identified. Can we prove the non-identification?!
  - ▶ The answer is no, because the model is indeed identified. We can prove identification using a different parameterization of the missingness mechanism (via odds ratio parameterization).

## Example ctd. odds ratio parameterization



- ▶  $p(R_1 \mid R_2 = 1, R_3 = 1, L^{(1)}) = p(R_1 \mid R_2 = 1, R_3 = 1, L_3)$
- ▶  $p(R_2 \mid R_1 = 1, R_3 = 1, L^{(1)}) = p(R_2 \mid R_1 = 1, R_3 = 1, L_1, L_3)$
- ▶  $p(R_3 \mid R_1 = 1, R_2 = 1, L^{(1)}) = p(R_3 \mid R_2 = 1, R_1 = 1, L_1)$
- ▶  $\text{OR}(R_1, R_2 \mid R_3 = 1, L^{(1)})$
- ▶  $\text{OR}(R_2, R_3 \mid R_1 = 1, L^{(1)})$
- ▶  $\text{OR}(R_1, R_3 \mid R_2 = 1, L^{(1)})$

## Example ctd. identification of the odds ratio terms



$$\begin{aligned}
 & \text{OR}(R_2 = 0, R_3 = 0 \mid R_1 = 1, L^{(1)}) \\
 &= \frac{p(R_3 = 0 \mid R_2 = 0, R_1 = 1, L^{(1)})}{p(R_3 = 1 \mid R_2 = 0, R_1 = 1, L^{(1)})} \times \frac{p(R_3 = 1 \mid R_2 = 1, R_1 = 1, L^{(1)})}{p(R_3 = 0 \mid R_2 = 1, R_1 = 1, L^{(1)})} \\
 &= \frac{p(R_3 = 0 \mid R_2 = 0, R_1 = 1, L_1^{(1)})}{p(R_3 = 1 \mid R_2 = 0, R_1 = 1, L_1^{(1)})} \times \frac{p(R_3 = 1 \mid R_2 = 1, R_1 = 1, L_1^{(1)})}{p(R_3 = 0 \mid R_2 = 1, R_1 = 1, L_1^{(1)})} \\
 &= \frac{p(R_3 = 0 \mid R_2 = 0, R_1 = 1, L_1)}{p(R_3 = 1 \mid R_2 = 0, R_1 = 1, L_1)} \times \frac{p(R_3 = 1 \mid R_2 = 1, R_1 = 1, L_1)}{p(R_3 = 0 \mid R_2 = 1, R_1 = 1, L_1)}.
 \end{aligned}$$

$\text{OR}(R_2, R_3 \mid R_1 = 1, L^{(1)})$  and  $\text{OR}(R_1, R_3 \mid R_2 = 1, L^{(1)})$  can be similarly identified. Thus, the missingness mechanism, full law, and target law are all identified.

# Odds ratio parameterization of $p(R \mid L^{(1)})$

$$p(R \mid L^{(1)}) = \frac{1}{\mathcal{Z}(L^{(1)})} \times \prod_{k=1}^K p(R_k \mid R_{-k} = 1, L^{(1)}) \times \prod_{k=2}^K \text{OR}(R_k, R_{\prec k} \mid R_{\succ k} = 1, L^{(1)}),$$

where  $R_{-k} = R \setminus R_k$ ,  $R_{\prec k} = \{R_1, \dots, R_{k-1}\}$ ,  $R_{\succ k} = \{R_{k+1}, \dots, R_K\}$ , and

$$\text{OR}(R_k, R_{\prec k} \mid R_{\succ k} = 1, L^{(1)}) = \frac{p(R_k \mid R_{\succ k} = 1, R_{\prec k}, L^{(1)})}{p(R_k = 1 \mid R_{\succ k} = 1, R_{\prec k}, L^{(1)})} \times \frac{p(R_k = 1 \mid R_{-k} = 1, L^{(1)})}{p(R_k \mid R_{-k} = 1, L^{(1)})}.$$

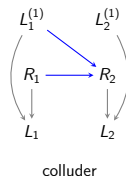
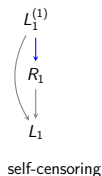
Need to identify:

- ▶ Conditional distributions:  $p(R_k \mid R_{-k} = 1, L^{(1)})$
- ▶ Odds ratio terms:  $\text{OR}(R_k, R_{\prec k} \mid R_{\succ k} = 1, L^{(1)})$
- ▶ When can we succeed?
- ▶ When do we fail?

# Full law identification theory in m-DAGs

Full law  $p(R, L^{(1)}, Z)$  that is Markov relative to a missing data DAG  $\mathcal{G}$  is identified **if and only if**  $\mathcal{G}$  does not contain the following two structures: (Nabi et al., 2020)

- ▶ self-censoring edge:  $L_i^{(1)} \rightarrow R_i$ ,
- ▶ collider:  $L_j^{(1)} \rightarrow R_i \leftarrow R_j$ .



- ▶ These graphical conditions are **sound** and **complete** for full law ID.
- ▶ Identification functional is given by the odds ratio parameterization of  $p(R \mid Z, L^{(1)})$ .

# Proof sketch

- ▶ Absence of self-censoring edges and colluders imply: (without loss of generality assume  $Z = \emptyset$ , and let  $V_{-k} = V \setminus V_k$ )

$$R_k \perp\!\!\!\perp L_k^{(1)} \mid R_{-k}, L^{(1)} \setminus L_k^{(1)}$$

- ▶ Odds ratio parameterization: ( $R_{\prec k} = \{R_1, \dots, R_{-k}\}$ )

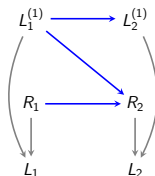
$$p(R \mid L^{(1)}) = \frac{1}{\mathcal{Z}(L^{(1)})} \times \prod_{k=1}^K p(R_k \mid R_{-k} = 1, L^{(1)}) \times \prod_{k=2}^K \text{OR}(R_k, R_{\prec k} \mid R_{\succ k} = 1, L^{(1)}),$$

- ▶  $p(R_k \mid R_{-k} = 1, L^{(1)}) = p(R_k \mid R_{-k} = 1, L_{-k})$ .
- ▶  $\text{OR}(R_i, R_j \mid R_{-\{i,j\}} = 1, L^{(1)})$  is identified via “symmetric argument,”
  - ▶ It is not a function of  $L_i^{(1)}$  and it is not a function of  $L_j^{(1)}$
- ▶ DAGs with no self-censoring edges and no colluders are submodels of *Itemwise Conditionally Independence Nonresponse* model (Sadinle and Reiter, 2017; Shpitser, 2016; Malinsky et al., 2021)



## Example: odds ratio parameterization fails!

- ▶ Even though the odds ratio parameterization led to completeness results for identification for the full law, the m-DAG factorization is still the only tool for arguing about the target law identification.
- ▶ If there exists a collider in the m-DAG, then the full law is not identifiable. However, the target law might still be identified.
- ▶ An example of this is the block-conditional MAR model.



- ▶ As we saw earlier,  $p(R = 1 \mid \text{pa}_{\mathcal{G}}(R))$  is easily identifiable as follows:

$$p(R = 1 \mid \text{pa}_{\mathcal{G}}(R)) = p(R_1 = 1) \times p(R_2 = 1 \mid R_1 = 1, L_1).$$

Therefore, the target law is identified.

## **Part 3. Non/Semi-parametric Estimation**

# General Theory of Estimation

- ▶ Let  $L^{(1)}$  be the *full* data
- ▶ Let  $R$  be missingness indicators ( $R = 1$  is no missingness)
- ▶  $C = (R, L^{(1)})$  be the *complete* data
- ▶ Let  $\varphi_r(L^{(1)})$  be the components of  $L^{(1)}$  when  $R = r$ 
  - ▶  $\varphi_1(L^{(1)}) = L^{(1)}$
- ▶ Let  $O = (R, \varphi_R(L^{(1)}))$  be the *observed* data
- ▶ We assume  $P[R = 1|L^{(1)}] > 0$
- ▶ Let  $F_{L^{(1)}}$  represent the model for  $L^{(1)}$ , indexed by parameters  $\mu$  (scalar target parameter) and  $\theta$  (nuisance)
- ▶ Let  $F_{R|L^{(1)}}$  represent the model for  $R|L^{(1)}$ , indexed by parameters  $\eta$  (nuisance)

# General Theory of Estimation

- ▶ Let  $\Lambda_1 = \Lambda(F_{L^{(1)}})$  be the collection of nuisance scores for  $\theta$  based on observation of  $C$
- ▶ Let  $\Lambda_2 = \Lambda(F_{R|L^{(1)}})$  be the collection of nuisance scores for  $\eta$  based on observation of  $C$
- ▶ Let  $\Lambda_j^O = \overline{R(g \cdot \Pi_j)}$  where
  - ▶  $\Pi_j$  is a projection operator that maps from  $C$  to  $\Lambda_j$  and
  - ▶  $g(\cdot) = E[\cdot|O]$ .
- ▶  $\Lambda^O = \overline{\Lambda_1^O + \Lambda_2^O}$
- ▶  $\Lambda^{O,\perp} = \Lambda_1^{O,\perp} \cap \Lambda_2^{O,\perp}$
- ▶ Influence functions (up to a normalizing constant) for regular and asymptotically linear (RAL) estimator of  $\mu$  live in  $\Lambda^{O,\perp}$
- ▶ A RAL estimator  $\hat{\mu}$  of  $\mu^*$  has the property that

$$\sqrt{n}(\hat{\mu} - \mu^*) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \underbrace{\psi(O_i; \mu^*, \theta^*, \eta^*)}_{\text{Influence Function}} + o_P(1)$$

# General Theory of Estimation

- ▶ Consider  $W(O; \mu^*, \theta^*, \eta^*) \in \Lambda^{O, \perp}$ .
- ▶ Let  $\hat{\theta}$  and  $\hat{\eta}$  be estimators for  $\theta$  and  $\eta$ , respectively.
- ▶ Estimate  $\hat{\mu}$  as the solution to

$$\sum_{i=1}^n W(O_i; \mu, \hat{\theta}, \hat{\eta}) = 0$$

- ▶ Under regularity conditions on the rate of convergence of  $\hat{\theta}$  and  $\hat{\eta}$ ,  $\hat{\mu}$  will be RAL with influence function

$$\left\{ \frac{\partial E[W(O; \mu^*, \theta^*, \eta^*)]}{\partial \mu} \right\}^{-1} W(O; \mu^*, \theta^*, \eta^*)$$

- ▶ Could also use profile estimators  $\hat{\theta}(\mu)$  and  $\hat{\eta}(\mu)$  and estimate  $\hat{\mu}$  as the solution to

$$\sum_{i=1}^n W(O_i; \mu, \hat{\theta}(\mu), \hat{\eta}(\mu)) = 0$$

$$\Lambda_1^{O,\perp} = \left\{ \frac{I(R=1)}{P[R=1|L^{(1)}]} a(L^{(1)}) + b(O) : a(L^{(1)}) \in \Lambda_1^\perp, E[b(O)|L^{(1)}] = 0 \right\}$$

$$\Lambda_2^{O,\perp} = \{b(O) : b(O) \in \Lambda_2^\perp\}$$

Plan:

- ▶ Write out an expression for the elements of  $\Lambda_1^{O,\perp}$
- ▶ Find restrictions on these elements to ensure orthogonality to  $\Lambda_2$

# Application of General Theory of Estimation

- ▶ Let  $L^{(1)} = (L_1^{(1)}, L_2^{(1)})$
- ▶ Let  $R = (R_1, R_2)$
- ▶ Let  $\pi_{ij}(L^{(1)}) = P[R_1 = i, R_2 = j | L^{(1)}]$

Any observed data random variable can be written as

$$b(O) = R_1 R_2 c_{11}(L^{(1)}) + R_1 (1 - R_2) c_{10}(L_1^{(1)}) + (1 - R_1) R_2 c_{01}(L_2^{(1)}) + (1 - R_1)(1 - R_2) c_{00}$$

What restrictions are  $c_{11}(L_1^{(1)}, L_2^{(1)})$ ,  $c_{10}(L_1^{(1)})$ ,  $c_{01}(L_2^{(1)})$  and  $c_{00}$  ensure that the  $E[b(O) | L^{(1)}] = 0$

$$c_{11}(L^{(1)}) = \frac{-\pi_{10}(L^{(1)})c_{10}(L_1^{(1)}) - \pi_{01}(L^{(1)})c_{01}(L_2^{(1)}) - \pi_{00}(L^{(1)})c_{00}}{\pi_{11}(L^{(1)})}$$

# Application of General Theory of Estimation

Any observed data random variable that has mean zero given  $L^{(1)}$  can be expressed as

$$\begin{aligned} & \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} c_{10}(L_1^{(1)}) + \\ & \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1)R_2 \right\} c_{01}(L_2^{(1)}) + \\ & \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00} \end{aligned}$$



# Application of General Theory of Estimation

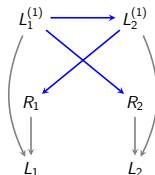
- ▶ Let  $\mu = E[h(L^{(1)})]$  for a specified function  $h(\cdot)$
- ▶ Suppose no restrictions are placed on the distribution of  $L^{(1)}$
- ▶ Let  $\Lambda_{1,\dagger}$  be the nuisance tangent space under no-restrictions.
- ▶ The elements of  $\Lambda_{1,\dagger}^\perp$  will be proportional to  $h(L^{(1)}) - \mu$
- ▶ If there are restrictions placed on the distribution of  $L^{(1)}$ , then  $\Lambda_{1,\dagger}^\perp \subset \Lambda_1^\perp$ .
- ▶ Thus,  $h(L^{(1)}) - \mu \in \Lambda_1^\perp$ .

# Application of General Theory of Estimation

We will work with

$$\Lambda_{1,\dagger}^{O,\perp} = \left\{ \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} c_{10}(L_1^{(1)}) + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1)R_2 \right\} c_{01}(L_2^{(1)}) + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00} : \right. \\ \left. c_{10}(L_1^{(1)}), c_{01}(L_2^{(1)}), c_{00} \right\}$$

# Block-Parallel Model



$$\pi_{jk}(L^{(1)}) = \pi_1(L_2^{(1)})^j \{1 - \pi_1(L_2^{(1)})\}^{1-j} \pi_2(L_1^{(1)})^k \{1 - \pi_2(L_1^{(1)})\}^{1-k}$$

$$\Lambda_2 = \left\{ \{R_1 - \pi_1(L_2^{(1)})\} g_1(L_2^{(1)}) + \{R_2 - \pi_2(L_1^{(1)})\} g_2(L_1^{(1)}) : g_1(L_2^{(1)}), g_2(L_1^{(1)}) \right\}$$

# Block-Parallel Model

What choices of  $c_{10}(L_1^{(1)})$ ,  $c_{01}(L_2^{(1)})$ ,  $c_{00}$  ensure orthogonality with all elements of  $\Lambda_2$ ?

$$\begin{aligned} E \left[ \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} \{R_1 - \pi_1(L_2^{(1)})\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} c_{10}(L_1^{(1)}) \{R_1 - \pi_1(L_2^{(1)})\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1)R_2 \right\} c_{01}(L_2^{(1)}) \{R_1 - \pi_1(L_2^{(1)})\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00} \{R_1 - \pi_1(L_2^{(1)})\} \middle| L_2^{(1)} \right] = 0 \quad (1) \end{aligned}$$

$$\begin{aligned}
 E \left[ \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} \{R_2 - \pi_2(L_1^{(1)})\} + \right. \\
 \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} c_{10}(L_1^{(1)}) \{R_2 - \pi_2(L_1^{(1)})\} + \right. \\
 \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1)R_2 \right\} c_{01}(L_2^{(1)}) \{R_2 - \pi_2(L_1^{(1)})\} + \right. \\
 \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00} \{R_2 - \pi_2(L_1^{(1)})\} \right| L_1^{(1)} \Big] = 0 \quad (2)
 \end{aligned}$$

# Block-Parallel Model

(1) implies that

$$c_{01}(L_2^{(1)}; c_{00}) = \underbrace{\frac{E \left[ h(L^{(1)}) - \mu \mid L_2^{(1)} \right]}{E \left[ \pi_2(L_1^{(1)}) \mid L_2^{(1)} \right]}}_{c_{01}(L_2^{(1)})} - \underbrace{\frac{E \left[ 1 - \pi_2(L_1^{(1)}) \mid L_2^{(1)} \right]}{E \left[ \pi_2(L_1^{(1)}) \mid L_2^{(1)} \right]}}_{c'_{01}(L_2^{(1)})} c_{00}$$

(2) implies that

$$c_{10}(L_1^{(1)}; c_{00}) = \underbrace{\frac{E \left[ h(L^{(1)}) - \mu \mid L_1^{(1)} \right]}{E \left[ \pi_1(L_2^{(1)}) \mid L_1^{(1)} \right]}}_{c_{10}(L_1^{(1)})} - \underbrace{\frac{E \left[ 1 - \pi_1(L_2^{(1)}) \mid L_1^{(1)} \right]}{E \left[ \pi_1(L_2^{(1)}) \mid L_1^{(1)} \right]}}_{c'_{10}(L_1^{(1)})} c_{00}$$

So,  $\Lambda^{O,\perp}$  contains a collection of elements indexed by  $c_{00}$ .

We will work with

$$\begin{aligned} & \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} + \\ & \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} \left\{ c_{10}(L_1^{(1)}) - c'_{10}(L_1^{(1)})c_{00} \right\} + \\ & \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1)R_2 \right\} \left\{ c_{01}(L_2^{(1)}) - c'_{01}(L_2^{(1)})c_{00} \right\} + \\ & \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00} \end{aligned}$$

Need estimators for

- ▶  $\pi_1(L_2^{(1)})$
- ▶  $\pi_2(L_1^{(1)})$
- ▶  $E[\pi_2(L_1^{(1)})|L_2^{(1)}]$
- ▶  $E[\pi_1(L_2^{(1)})|L_1^{(1)}]$
- ▶  $E[h(L^{(1)})|L_1^{(1)}]$
- ▶  $E[h(L^{(1)})|L_2^{(1)}]$



To find the optimal choice of  $c_{00}$ , minimize

$$E \left[ \left\{ \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} + \right. \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} \left\{ c_{10}(L_1^{(1)}) - c'_{10}(L_1^{(1)})c_{00} \right\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1)R_2 \right\} \left\{ c_{01}(L_2^{(1)}) - c'_{01}(L_2^{(1)})c_{00} \right\} + \right. \\ \left. \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00} \right\}^2 \right]$$

# Block-Parallel Model

Set derivative with respect to  $c_{00}$  equal to zero.

$$\begin{aligned} E \left[ \left\{ \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \left\{ h(L^{(1)}) - \mu \right\} + \right. \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} \left\{ c_{10}(L_1^{(1)}) \right\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} \left\{ -c'_{10}(L_1^{(1)}) c_{00} \right\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1) R_2 \right\} \left\{ c_{01}(L_2^{(1)}) \right\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1) R_2 \right\} \left\{ -c'_{01}(L_2^{(1)}) c_{00} \right\} + \right. \\ \left. \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00} \right\} \times \right. \\ \left. \left\{ \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} \left\{ -c'_{10}(L_1^{(1)}) \right\} + \right. \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1) R_2 \right\} \left\{ -c'_{01}(L_2^{(1)}) \right\} + \right. \\ \left. \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} \right\} \right] = 0 \end{aligned}$$

# Block-Parallel Model

$$c_{00} = \frac{a}{b}$$

where

$$\begin{aligned} b = & -E \left[ \left\{ \frac{\pi_{10}(L^{(1)})}{\pi_{11}(L^{(1)})} + 1 \right\} \pi_{10}(L^{(1)}) \{c'_{10}(L^{(1)})\}^2 \right] - \\ & E \left[ \left\{ \frac{\pi_{10}(L^{(1)})\pi_{01}(L^{(1)})}{\pi_{11}(L^{(1)})} \right\} \{2c'_{10}(L^{(1)})c'_{01}(L^{(1)})\} \right] + \\ & E \left[ \left\{ \frac{\pi_{10}(L^{(1)})\pi_{00}(L^{(1)})}{\pi_{11}(L^{(1)})} \right\} \{2c'_{10}(L^{(1)})\} \right] - \\ & E \left[ \left\{ \frac{\pi_{01}(L^{(1)})}{\pi_{11}(L^{(1)})} + 1 \right\} \pi_{01}(L^{(1)}) \{c'_{01}(L^{(1)})\}^2 \right] + \\ & E \left[ \left\{ \frac{\pi_{01}(L^{(1)})\pi_{00}(L^{(1)})}{\pi_{11}(L^{(1)})} \right\} \{2c'_{01}(L^{(1)})\} \right] + \\ & E \left[ \left\{ \frac{\pi_{00}(L^{(1)})}{\pi_{11}(L^{(1)})} + 1 \right\} \pi_{00}(L^{(1)}) \right] \end{aligned}$$

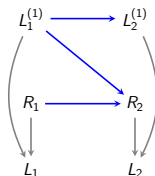
# Block-Parallel Model

$$c_{00} = \frac{a}{b}$$

where

$$\begin{aligned} a = & E \left[ \left\{ \frac{\pi_{10}(L^{(1)})}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} \right\} \{c'_{10}(L^{(1)}_1)\} \right] + \\ & E \left[ \left\{ \frac{\pi_{01}(L^{(1)})}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} \right\} \{c'_{01}(L^{(1)}_2)\} \right] - \\ & E \left[ \left\{ \frac{\pi_{00}(L^{(1)})}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} \right\} \right] - \\ & E \left[ \left\{ \frac{\pi_{10}(L^{(1)})}{\pi_{11}(L^{(1)})} + 1 \right\} \pi_{10}(L^{(1)}) \{c_{10}(L^{(1)}_1)\} \{c'_{10}(L^{(1)}_1)\} \right] - \\ & E \left[ \left\{ \frac{\pi_{10}(L^{(1)})\pi_{01}(L^{(1)})}{\pi_{11}(L^{(1)})} \right\} \{c_{10}(L^{(1)}_1)c'_{01}(L^{(1)}_2) + c_{01}(L^{(1)}_2)c'_{10}(L^{(1)}_1)\} \right] + \\ & E \left[ \left\{ \frac{\pi_{10}(L^{(1)})\pi_{00}(L^{(1)})}{\pi_{11}(L^{(1)})} \right\} \{c_{10}(L^{(1)}_1)\} \right] - \\ & E \left[ \left\{ \frac{\pi_{01}(L^{(1)})}{\pi_{11}(L^{(1)})} + 1 \right\} \pi_{01}(L^{(1)}) \{c_{01}(L^{(1)}_2)\} \{c'_{01}(L^{(1)}_2)\} \right] + \\ & E \left[ \left\{ \frac{\pi_{01}(L^{(1)})\pi_{00}(L^{(1)})}{\pi_{11}(L^{(1)})} \right\} \{c_{01}(L^{(1)}_2)\} \right] \end{aligned}$$

# Block-Conditional Model



$$\pi_{jk}(L^{(1)}) = \pi_1^j \{1 - \pi_1\}^{1-j} \times \\ \pi_2(1, L_1^{(1)})^{jk} \{1 - \pi_2(1, L_1^{(1)})\}^{j(1-k)} \pi_2(0, L_1^{(1)})^{(1-j)k} \{1 - \pi_2(0, L_1^{(1)})\}^{(1-j)(1-k)}$$

$$\Lambda_2 = \left\{ \{R_1 - \pi_1\}g_1 + R_1\{R_2 - \pi_2(1, L_1^{(1)})\}g_2(1, L_1^{(1)}) + \right. \\ \left. (1 - R_1)\{R_2 - \pi_2(0, L_1^{(1)})\}g_2(0, L_1^{(1)}) : g_1, g_2(1, L_1^{(1)}), g_2(0, L_1^{(1)}) \right\}$$

# Block-Conditional Model

What choices of  $c_{10}(L_1^{(1)})$ ,  $c_{01}(L_2^{(1)})$ ,  $c_{00}$  ensure orthogonality with all elements of  $\Lambda_2$ ?

$$\begin{aligned} E \left[ \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} R_1 \{R_2 - \pi_2(1, L_1^{(1)}) + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} c_{10}(L_1^{(1)}) R_1 \{R_2 - \pi_2(1, L_1^{(1)}) + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1)R_2 \right\} c_{01}(L_2^{(1)}) R_1 \{R_2 - \pi_2(1, L_1^{(1)}) + \right. \\ \left. \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00} R_1 \{R_2 - \pi_2(1, L_1^{(1)}) \right\} \middle| L_1^{(1)} \right] = 0 \end{aligned} \quad (3)$$

## Block-Conditional Model

$$\begin{aligned} E \left[ \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} (1 - R_1) \{R_2 - \pi_2(0, L_1^{(1)})\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} c_{10}(L_1^{(1)}) (1 - R_1) \{R_2 - \pi_2(0, L_1^{(1)})\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1) R_2 \right\} c_{01}(L_2^{(1)}) (1 - R_1) \{R_2 - \pi_2(0, L_1^{(1)})\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00}(1 - R_1) \{R_2 - \pi_2(0, L_1^{(1)})\} \right] \Bigg| L_1^{(1)} \Bigg] = 0 \end{aligned} \quad (4)$$

$$\begin{aligned}
 E \left[ \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} \{R_1 - \pi_1\} + \right. \\
 \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} c_{10}(L_1^{(1)}) \{R_1 - \pi_1\} + \right. \\
 \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1)R_2 \right\} c_{01}(L_2^{(1)}) \{R_1 - \pi_1\} + \right. \\
 \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00} \{R_1 - \pi_1\} \right] = 0 \quad (5)
 \end{aligned}$$



(4) implies

$$E[c_{01}(L_2^{(1)})|L_1^{(1)}] = c_{00}$$

(5) implies

$$E[\pi_2(0, L_1^{(1)})c_{01}(L_2^{(1)})] + c_{00}E[1 - \pi_2(0, L_1^{(1)})] = 0$$

(4) and (5) imply that

$$c_{00}E[\pi_2(0, L_1^{(1)})] + c_{00}E[1 - \pi_2(0, L_1^{(1)})] = c_{00} = 0$$

# Block-Conditional Model

What choices of  $c_{01}(L_2^{(1)})$  make  $E[c_{01}(L_2^{(1)})|L_1^{(1)}] = 0$ ?

- ▶ Fredholm integral equation of the first kind.
- ▶ Obviously,  $c_{01}(L_2^{(1)}) = 0$
- ▶ Non-trivial choices may or may not exist depending on the conditional distribution of  $L_2^{(1)}$  given  $L_1^{(1)}$ .
- ▶ If the conditional distribution of  $L_2^{(1)}$  given  $L_1^{(1)}$  is from a canonical exponential family, then  $c_{01}(L_2^{(1)}) = 0$  a.s.

With  $c_{01}(L_2^{(1)}) = c_{00} = 0$ , (3) implies

$$c_{10}(L_1^{(1)}) = \frac{E[h(L^{(1)}) - \mu | L_1^{(1)}]}{\pi_1}$$

# Block-Conditional Model

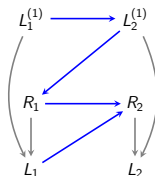
We will work with

$$\begin{aligned} & \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} + \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} \frac{E[h(L^{(1)}) - \mu | L_1^{(1)}]}{\pi_1} \\ &= \frac{R_1}{\pi_1} \left\{ \frac{R_2}{\pi_2(1, L_1^{(1)})} h(L^{(1)}) + \left( 1 - \frac{R_2}{\pi_2(1, L_1^{(1)})} \right) E[h(L^{(1)}) | L_1^{(1)}] - \mu \right\} \end{aligned}$$

Need estimators for

- ▶  $\pi_1$
- ▶  $\pi_2(1, L_1^{(1)})$
- ▶  $E[h(L^{(1)}) | L_1^{(1)}]$

# Permutation Model



$$\pi_{jk}(L^{(1)}) = \pi_1(L_2^{(1)})^j \{1 - \pi_1(L_2^{(1)})\}^{1-j} \times \\ \pi_2(1, L_1^{(1)})^{jk} \{1 - \pi_2(1, L_1^{(1)})\}^{j(1-k)} \pi_2(0)^{(1-j)k} \{1 - \pi_2(0)\}^{(1-j)(1-k)}$$

$$\Lambda_2 = \left\{ \{R_1 - \pi_1(L_2^{(1)})\} g_1(L_2^{(1)}) + R_1 \{R_2 - \pi_2(1, L_1^{(1)})\} g_2(1, L_1^{(1)}) + \right. \\ \left. (1 - R_1) \{R_2 - \pi_2(0)\} g_2(0) : g_1(L_2^{(1)}), g_2(1, L_1^{(1)}), g_2(0) \right\}$$

# Permutation Model

What choices of  $c_{10}(L_1^{(1)})$ ,  $c_{01}(L_2^{(1)})$ ,  $c_{00}$  ensure orthogonality with all elements of  $\Lambda_2$ ?

$$\begin{aligned} E \left[ \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} \{R_1 - \pi_1(L_2^{(1)})\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} c_{10}(L_1^{(1)}) \{R_1 - \pi_1(L_2^{(1)})\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1)R_2 \right\} c_{01}(L_2^{(1)}) \{R_1 - \pi_1(L_2^{(1)})\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00} \{R_1 - \pi_1(L_2^{(1)})\} \middle| L_2^{(1)} \right] = 0 \quad (6) \end{aligned}$$

# Permutation Model

$$\begin{aligned} E \left[ \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} R_1 \{R_2 - \pi_2(1, L_1^{(1)})\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} c_{10}(L_1^{(1)}) R_1 \{R_2 - \pi_2(1, L_1^{(1)})\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1) R_2 \right\} c_{01}(L_2^{(1)}) R_1 \{R_2 - \pi_2(1, L_1^{(1)})\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00} R_1 \{R_2 - \pi_2(1, L_1^{(1)})\} \middle| L_1^{(1)} \right] = 0 \end{aligned} \quad (7)$$

# Permutation Model

$$\begin{aligned} E \left[ \frac{R_1 R_2}{\pi_{11}(L^{(1)})} \{h(L^{(1)}) - \mu\} (1 - R_1) \{R_2 - \pi_2(0)\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} c_{10}(L_1^{(1)}) (1 - R_1) \{R_2 - \pi_2(0)\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1) R_2 \right\} c_{01}(L_2^{(1)}) (1 - R_1) \{R_2 - \pi_2(0)\} + \right. \\ \left. \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} c_{00}(1 - R_1) \{R_2 - \pi_2(0)\} \right] = 0 \end{aligned} \quad (8)$$



# Permutation Model

(8) implies

$$c_{00} = \frac{E[\{1 - \pi_1(L_2^{(1)})\}c_{01}(L_2^{(1)})]}{E[\{1 - \pi_1(L_2^{(1)})\}]}$$

(6) implies

$$c_{01}(L_2^{(1)}) = \frac{E[\{h(L^{(1)}) - \mu\} | L_2^{(1)}]}{\pi_2(0)} - \frac{1 - \pi_2(0)}{\pi_2(0)} c_{00}$$

Together, this implies that

$$c_{00} = \frac{E[h(L^{(1)})\{1 - \pi_1(L_2^{(1)})\}]}{E[1 - \pi_1(L_2^{(1)})]} - \mu$$

and

$$c_{01}(L_2^{(1)}) = \frac{E[h(L^{(1)}) | L_2^{(1)}]}{\pi_2(0)} - \frac{1 - \pi_2(0)}{\pi_2(0)} \left\{ \frac{E[h(L^{(1)})\{1 - \pi_1(L_2^{(1)})\}]}{E[1 - \pi_1(L_2^{(1)})]} \right\} - \mu$$

Adding (7) implies

$$c_{10}(L_1^{(1)}) = \frac{E[h(L^{(1)}) - E[h(L^{(1)})|L_2^{(1)}]|L_1^{(1)}]}{E[\pi_1(L_2^{(1)})|L_1^{(1)}]} + \frac{E[\pi_1(L_2^{(1)})E[h(L^{(1)})|L_2^{(1)}]|L_1^{(1)}]}{E[\pi_1(L_2^{(1)})|L_1^{(1)}]} - \mu$$

# Permutation Model

We will work with

$$\begin{aligned} & \frac{R_1 R_2}{\pi_{11}(L^{(1)})} h(L^{(1)}) + \\ & \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{10}(L^{(1)}) + R_1(1 - R_2) \right\} \times \\ & \left\{ \frac{E[h(L^{(1)}) - E[h(L^{(1)})|L_2^{(1)}]|L_1^{(1)}]}{E[\pi_1(L_2^{(1)})|L_1^{(1)}]} + \frac{E[\pi_1(L_2^{(1)})E[h(L^{(1)})|L_2^{(1)}]|L_1^{(1)}]}{E[\pi_1(L_2^{(1)})|L_1^{(1)}]} \right\} + \\ & \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{01}(L^{(1)}) + (1 - R_1)R_2 \right\} \times \\ & \left\{ \frac{E[h(L^{(1)})|L_2^{(1)}]}{\pi_2(0)} - \frac{1 - \pi_2(0)}{\pi_2(0)} \left\{ \frac{E[h(L^{(1)})\{1 - \pi_1(L_2^{(1)})\}]}{E[1 - \pi_1(L_2^{(1)})]} \right\} \right\} + \\ & \left\{ -\frac{R_1 R_2}{\pi_{11}(L^{(1)})} \pi_{00}(L^{(1)}) + (1 - R_1)(1 - R_2) \right\} \times \\ & \left\{ \frac{E[h(L^{(1)})\{1 - \pi_1(L_2^{(1)})\}]}{E[1 - \pi_1(L_2^{(1)})]} \right\} - \mu \end{aligned}$$

# Permutation Model

Need estimators for

- ▶  $\pi_2(0)$
- ▶  $\pi_2(1, L_1^{(1)})$
- ▶  $\pi_1(L_2^{(1)})$
- ▶ Conditional means of functions of  $L_2^{(1)}$  given  $L_1^{(1)}$  and conditional means of functions of  $L_1^{(1)}$  given  $L_2^{(1)}$ 
  - ▶  $E[h(L^{(1)})|L_1^{(1)}]$
  - ▶  $E[h(L^{(1)})|L_2^{(1)}]$
  - ▶  $E[\pi_1(L_2^{(1)})|L_1^{(1)}]$
  - ▶  $E[\pi_1(L_2^{(1)})E[h(L^{(1)})|L_2^{(1)}]|L_1^{(1)}]$
  - ▶  $E[E[h(L^{(1)})|L_2^{(1)}]|L_1^{(1)}]$

**Wrap up**

# Revisiting course outline and objectives

## Part I. Missing data DAGs

- ▶ Represented missingness mechanisms graphically; interpreted a missing data DAG model as a class of distributions with a set of independence restrictions.

## Part II. Non-parametric identification

- ▶ Discussed identification tricks for full and target laws, and showed how non-identification proofs go.

## Part III. Non/Semi-parametric estimation

- ▶ Given and identified query, derived the non-parametric influence functions; Given three types of m-DAGs with MNAR missingness, derived the tangent space of the underlying full data and observed data distributions.

## Part IV. Sensitivity analysis

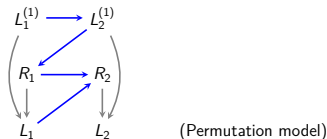
- ▶ Assessing deviations from assumptions.

# General questions and thoughts

1. How can identification results guide estimation strategies?
2. What if there exist variables that are not just missing but completely unobserved?
3. Similar to DAGs, absence of an edge in an m-DAG implies a restriction of the form  $A \perp B \mid C$ . Is this restriction testable from observed finite samples?
4. What if the model is not nonparametrically identified?

# I. Estimation strategies (weighted estimating equations)

- How identification results can guide estimation strategies?



- Let our parameter of interest be  $\beta_h = \mathbb{E}[h(L_1^{(1)}, L_2^{(1)})]$ , which can be rewritten as:

$$\beta_h = \mathbb{E} \left[ \frac{R_1 R_2}{p(R_1 = 1 | L_2^{(1)}) p(R_2 = 1 | R_1 = 1, L_1)} \times h(L_1^{(1)}, L_2^{(1)}) \right].$$

$p(R_1 = 1 | L_2^{(1)})$  is identified from an intervention distribution where  $R_2$  is intervened on:

$$p(R_1 = 1 | L_2^{(1)}) = p(R_1 = 1 | L_2, \text{do}(R_2 = 1)).$$

- How does identification of  $p(R_1 = 1 | L_2^{(1)})$  help with estimation?

Assume  $p(R_1 = 1 | L_2^{(1)}) = p(R_1 = 1 | L_2^{(1)}; \alpha)$ ,  $\alpha \in \mathbb{R}^p$  and  $\mathbb{E}[U(R_1, L_2^{(1)}; \alpha)] = 0$ , then:

$$\mathbb{E} \left[ \frac{R_2}{p(R_2 | R_1, L_1)} \times U(R_1, L_2^{(1)}; \alpha) \right] = 0.$$



# I. Estimation strategies (EIF derivation)

$$\begin{aligned}\beta_h &= \mathbb{E} \left[ h(L_1^{(1)}, L_2^{(1)}) \right] \\ &= \mathbb{E} \left[ \frac{R_1 R_2}{p(R_1 = 1 \mid L_2^{(1)}) p(R_2 = 1 \mid R_1 = 1, L_1)} \times h(L_1^{(1)}, L_2^{(1)}) \right]\end{aligned}$$

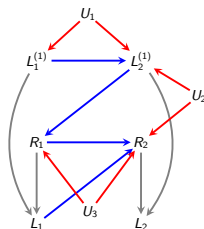
- ▶ A more reasonable estimator for  $\beta_h$  is the estimator derived based on the nonparametric efficient influence function (EIF).
- ▶ The core idea of deriving the EIF for  $\beta_h$  is to use an intermediate variable that first takes care of the missingness of  $L_1^{(1)}$ , and then  $L_2^{(1)}$  in a sequential manner

$$\begin{aligned}\tilde{\beta}_h(L_2^{(1)}) &= \frac{R_1}{p(R_1 = 1 \mid L_2^{(1)})} \times h(L_1^{(1)}, L_2^{(1)}), \\ \beta_h &= \mathbb{E} \left[ \frac{R_2}{p(R_2 = 1 \mid R_1 = 1, L_1)} \times \tilde{\beta}_h(L_2^{(1)}) \right].\end{aligned}$$

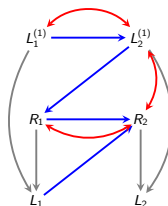
(Robins, 1997)

## II. Missing data DAGs with hidden variables

- ▶ What if there exist variables that are not just missing but completely unobserved?
- ▶ Summarize the observed data distribution with a missing data acyclic directed mixed graph (ADMG).



(a)  $G(V, U)$



(b)  $G(V)$

$L_1^{(1)}$ : smoking,  $L_2^{(1)}$ : lung cancer

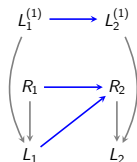
$U_1$ : genotypic traits,  $U_2$ : occupation,  $U_3$ : ethnicity

- ▶ Identification results are extended to missing data ADMGs, and results remain complete for full law identification (Nabi et al., 2020).

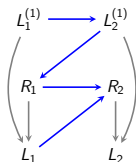
### III. Testable implications in m-DAGs

- ▶ Similar to DAGs, absence of an edge in an m-DAG implies a restriction of the form  $A \perp B \mid C$ . Is this restriction testable from observed finite samples?
- ▶ If all the restrictions encoded in a missing data DAG are provably untestable (i.e., no restriction on the observed data law), the full law Markov relative to the DAG is said to be **non-parametric saturated** (Robins; 1997)
  - ▶ A self-censoring mechanism imposes no restriction on the observed data law
  - ▶ Another example of a non-parametric saturated model is the permutation model.
- ▶ Submodels of a non-parametric saturated model can still be tested using partially observed data (Nabi and Bhattacharya, 2022).

### III. Testable implications



MAR model



Permutation supermodel

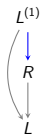
- ▶ Is  $R_1 \perp\!\!\!\perp L_2^{(1)}$ ?
- ▶ Fit  $p(R_1)$  and  $p(R_1 \mid L_2^{(1)})$  and compare the goodness of fits
- ▶ Use a weighted estimating equation to fit  $p(R_1 \mid L_2^{(1)}; \alpha)$

$$\mathbb{P}_n \left[ \frac{R_2}{p(R_2 \mid R_1, L_1)} \times U(., \alpha) \right] = 0,$$

where  $\mathbb{P}_n[U(., \alpha)] = 0$  with respect to the full law.

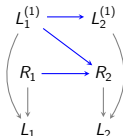
## IV. Missing data non-identification

- ▶ What if a parameter of interest is not identified from the observed data law?



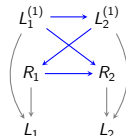
Self-censoring

$p(L^{(1)})$  not ID.



Collider

$p(L^{(1)}, R)$  not ID.

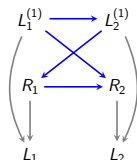


Criss-cross

$p(L^{(1)})$  not ID.

- ▶ Generally speaking, we have two general options:
  - ▶ Restrict the missing data model by posing **extra** assumptions on the full law.
  - ▶ Obtain bounds, conduct sensitivity analysis, etc.  
(Rotnitzky et al., 1998; Robins et al., 2000; Scharfstein and Irizarry, 2003)

## IV. Partial/parametric identification



$p(L_1^{(1)}, L_2^{(1)})$  is not identified.

Criss-cross structure

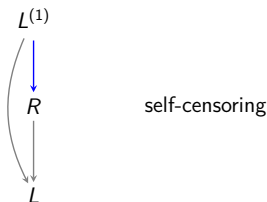
- ▶ Partial identification:  $p(L_1^{(1)} \mid L_2^{(1)})$  and  $\text{OR}(L_1^{(1)}, L_2^{(1)})$  are identified.
  - ▶ We can test  $L_1^{(1)} \perp\!\!\!\perp L_2^{(1)}$  without further assumptions.
- ▶ Under what conditions  $p(L_1^{(1)}, L_2^{(1)})$  is identified? assume  $p(L_1^{(1)})$  and  $p(L_2^{(1)} \mid L_1^{(1)})$  follow exponential family distributions:

$$L_1^{(1)} \sim \exp \left\{ \frac{l_1^{(1)} \eta_h - b_h(\eta_h)}{\Phi_h} + c_h(l_1; \Phi_h) \right\}$$

$$L_2^{(1)} \mid L_1^{(1)} \sim \exp \left\{ \frac{l_2^{(1)} \eta - b(\eta)}{\Phi} + c(l_2^{(1)}; \Phi) \right\}, \quad g(\mu(\eta)) = \alpha + \beta l_1^{(1)}.$$

- ▶ What are sufficient conditions for target law ID in the above class of distributions? (Guo et al., 2023)

## IV. Alternative handling of non-identified models



$$p(L^{(1)}) = p(L^{(1)} \mid R = 0) \times p(R = 0) + p(L^{(1)} \mid R = 1) \times p(R = 1).$$

$$p(L^{(1)} \mid R = 0) \propto p(L^{(1)} \mid R = 1) \times \exp(\gamma S(L^{(1)}))$$

- ▶ The relation is controlled by the sensitivity parameter  $\gamma$ .
- ▶  $S(L^{(1)})$  is a specified function of  $L^{(1)}$ .

# Many interesting open problems

- ▶ Missing data DAGs with or without unmeasured confounding:
  - ▶ A concise and precise representation of MNAR mechanisms.
- ▶ **Identification:**
  - ▶ Complete characterization of target law ID remains an open problem while such characterizations for full law ID exist.
  - ▶ Partial identification.
- ▶ **Estimation:**
  - ▶ Intuitive estimation strategies: IPW-style estimators
  - ▶ An understudied research area: influence-function based estimators in m-DAGs.
- ▶ **Testable implications:**
  - ▶ Data-driven structure learning approaches.



# Ananke: Software for causal inference

## ***ananke-causal:***

A python package for causal inference using the language of graphical models.

## **Highlights:**

- Identification algorithms
- Semiparametric inference
- Surrogate experiments
- Missing data



## **Links:**

- ▶ Documentation: <https://ananke.readthedocs.io/en/latest/index.html>
- ▶ Gitlab repository: <https://gitlab.com/causal/ananke>
- ▶ Python package index: <https://pypi.org/project/ananke-causal/>

**Contributors:** Rohit Bhattacharya, Jaron Lee, and RN

# Acknowledgments

In alphabetical order:

- ▶ Rohit Bhattacharya, Ph.D., Williams College
- ▶ Anna Guo, Emory University
- ▶ James Robins, Ph.D., Harvard School of Public Health
- ▶ Dan Scharfstein, Sc.D., University of Utah
- ▶ Ilya Shpitser, Ph.D., Johns Hopkins University
- ▶ Jiwei Zhao, Ph.D., University of Wisconsin-Madison

# References I

- R. Bhattacharya, R. Nabi, I. Shpitser, and J. Robins. Identification in missing data models represented by directed acyclic graphs. In *Proceedings of the Thirty Fifth Conference on Uncertainty in Artificial Intelligence (UAI-35th)*. AUAI Press, 2019.
- R. Bhattacharya, R. Nabi, and I. Shpitser. Semiparametric inference for causal effects in graphical models with hidden variables. *arXiv preprint arXiv:2003.12659*, 2020.
- H. Y. Chen. A semiparametric odds ratio model for measuring association. *biometrics*, 63:413–421, 2007.
- A. Guo, J. Zhao, and R. Nabi. Semiparametric identifiability and estimation under missing not at random mechanisms, 2023.
- Y. Huang and M. Valtorta. Pearl’s calculus of interventions is complete. In *Twenty Second Conference On Uncertainty in Artificial Intelligence*, 2006.
- D. Malinsky, I. Shpitser, and E. J. Tchetgen Tchetgen. Semiparametric inference for nonmonotone missing-not-at-random data: the no self-censoring model. *Journal of the American Statistical Association*, pages 1–9, 2021.
- K. Mohan, J. Pearl, and J. Tian. Graphical models for inference with missing data. In *Advances in Neural Information Processing Systems 26*, pages 1277–1285. Curran Associates, Inc., 2013.

## References II

- R. Nabi and R. Bhattacharya. On testability and goodness of fit tests in missing data models. *arXiv preprint arXiv:2203.00132*, 2022.
- R. Nabi, R. Bhattacharya, and I. Shpitser. Full law identification in graphical models of missing data: Completeness results. In *Proceedings of the Twenty Seventh International Conference on Machine Learning (ICML-20)*, 2020.
- R. Nabi, R. Bhattacharya, I. Shpitser, and J. Robins. Causal and counterfactual views of missing data models. *arXiv preprint arXiv:2210.05558*, 2022.
- T. S. Richardson, R. J. Evans, J. M. Robins, and I. Shpitser. Nested markov properties for acyclic directed mixed graphs. *arXiv preprint arXiv:1701.06686*, 2017.
- J. M. Robins. Non-response models for the analysis of non-monotone non-ignorable missing data. *Statistics in Medicine*, 16:21–37, 1997.
- J. M. Robins, A. Rotnitzky, and D. O. Scharfstein. Sensitivity analysis for selection bias and unmeasured confounding in missing data and causal inference models. In *Statistical models in epidemiology, the environment, and clinical trials*, pages 1–94. Springer, 2000.
- A. Rotnitzky, J. M. Robins, and D. O. Scharfstein. Semiparametric regression for repeated outcomes with nonignorable nonresponse. *Journal of the american statistical association*, 93(444):1321–1339, 1998.

## References III

- D. B. Rubin. Causal inference and missing data (with discussion). *Biometrika*, 63: 581–592, 1976.
- M. Sadinle and J. P. Reiter. Itemwise conditionally independent nonresponse modelling for incomplete multivariate data. *Biometrika*, 104(1):207–220, 2017.
- D. O. Scharfstein and R. A. Irizarry. Generalized additive selection models for the analysis of studies with potentially nonignorable missing outcome data. *Biometrics*, 59(3):601–613, 2003.
- I. Shpitser. Consistent estimation of functions of data missing non-monotonically and not at random. In *Proceedings of the Thirtieth Annual Conference on Neural Information Processing Systems (NIPS-16)*. Curran Associates, Inc., 2016.
- I. Shpitser and J. Pearl. Identification of joint interventional distributions in recursive semi-Markovian causal models. In *Proceedings of the Twenty-First National Conference on Artificial Intelligence (AAAI-06)*. AAAI Press, Palo Alto, 2006.
- R. Srinivasan, R. Bhattacharya, R. Nabi, E. L. Ogburn, and I. Shpitser. Graphical models of entangled missingness. *arXiv preprint arXiv:2304.01953*, 2023.
- T. S. Verma and J. Pearl. Equivalence and synthesis of causal models. Technical Report R-150, Department of Computer Science, University of California, Los Angeles, 1990.

Y. Zhou, R. J. A. Little, and K. J. D. Block-conditional missing at random models for missing data. *Statistical Science*, 25(4):517–532, 2010.

Upon using the materials, please cite appropriately.

**Razieh Nabi, PhD**

Rollins Assistant Professor

Department of Biostatistics and Bioinformatics

Rollins School of Public Health

Emory University

🐦 @razielnabi

✉ razieh.nabi@emory.edu

🏠 <https://razielnabi.com>

**Daniel Scharfstein, ScD**

Chief of the Division of Biostatistics

Department of Population Health Sciences

School of Medicine

University of Utah

🐦 @dscharf3

✉ Daniel.Scharfstein@hsc.utah.edu

🏠 <https://medicine.utah.edu/dscharf>