

# Exercises MLSS 2019: *Causality*

Joris Mooij

August 27, 2019

## 1 In-class exercises

### 1.1 Simpson's Paradox

You are investigating the effectiveness of a drug against a deadly disease. You are given access to data collected by health insurance companies about their customers. You divide the diseased customers into two groups: those that took the drug, and those that didn't take the drug. Some of the customers recovered, others unfortunately didn't recover. The reasons why some patients were treated and others were not, are unknown to you. You find the following numbers:

	Recovery	No recovery	Total	Recovery rate
Drug	20	20	40	...%
No drug	16	24	40	...%
Total	36	44	80	

- Calculate the recovery rates (in %) for both groups ("drug" and "no drug").
- If you were diseased, would you take the drug, or not?

Upon closer inspection of the data, you notice something peculiar when you group patients according to gender:

<b>Males</b>	Recovery	No recovery	Total	Recovery rate
Drug	18	12	30	...%
No drug	7	3	10	...%
Total	25	15	40	

<b>Females</b>	Recovery	No recovery	Total	Recovery rate
Drug	2	8	10	...%
No drug	9	21	30	...%
Total	11	29	40	

- Calculate the recovery rates (in %) for both groups ("drug" and "no drug"), for each subpopulation (males and females) separately.
- In light of these numbers, would you take the drug if you were diseased, or not?
- What would be your advice to a diseased patient with unknown gender?

This phenomenon is known as "Simpson's paradox". A lot has been written about this paradox, but it dissolves once you recognize that you should not make the mistake of interpreting correlations as causations, as we'll see later today.

## 1.2 Paths, colliders, d-blocked paths and d-separation

**Definition 1 (Paths, Ancestors)** Let  $\mathcal{G}$  be a directed mixed graph.

- A **path**  $q$  in  $\mathcal{G}$  is a sequence of adjacent edges in  $\mathcal{G}$  in which no node occurs more than once.
- A path consisting of directed edges  $X_{i_1} \rightarrow X_{i_2} \rightarrow X_{i_3} \rightarrow \dots \rightarrow X_{i_k}$  that all point in the same direction is called a **directed path**.
- If there is a directed path from  $X$  to  $Y$  (or if  $X = Y$ ),  $X$  is called a **ancestor** of  $Y$ .
- The ancestors of  $Y$  are denoted  $\text{ang}(Y)$ , and include  $Y$ .

**Definition 2 (Colliders, Blocked Paths, d-separation)** Let  $\mathcal{G}$  be a directed mixed graph, and  $q$  a path on  $\mathcal{G}$ .

- A **collider** on  $q$  is a (non-endpoint) node  $X$  on  $q$  with precisely two arrowheads pointing towards  $X$  on the adjacent edges:

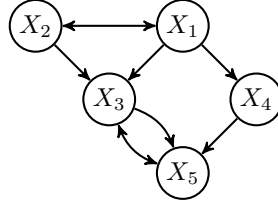
$$\rightarrow X \leftarrow, \quad \rightarrow X \leftrightarrow, \quad \leftrightarrow X \leftarrow, \quad \leftrightarrow X \leftrightarrow$$

- A **non-collider** on  $q$  is any node on the path which is not a collider.

A set of nodes  $\mathbf{S}$  in  $\mathcal{G}$  is said to **d-block**  $q$  if  $q$  contains a non-collider which is in  $\mathbf{S}$ , or a collider which is not an ancestor of  $\mathbf{S}$ .

For three sets  $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$  of nodes in  $\mathcal{G}$ , we say that  $\mathbf{X}$  and  $\mathbf{Y}$  are **d-separated by  $\mathbf{Z}$**  iff all paths between a node in  $\mathbf{X}$  and a node in  $\mathbf{Y}$  are d-blocked by  $\mathbf{Z}$ , and write  $\mathbf{X} \perp_{\mathcal{G}} \mathbf{Y} \mid \mathbf{Z}$ .

Consider the following directed mixed graph  $\mathcal{G}$ :



- Is  $X_3 \rightarrow X_5 \leftrightarrow X_3$  a path? Is it a directed path?
- Is  $X_3 \leftrightarrow X_5$  a path? Is it a directed path?
- Is  $X_5 \leftarrow X_3 \leftarrow X_1$  a path? Is it a directed path?
- What are the ancestors of  $X_4$ ?

Consider the path  $X_2 \leftrightarrow X_1 \rightarrow X_3 \leftrightarrow X_5 \leftarrow X_4$  on  $\mathcal{G}$ .

- Which nodes on the path are colliders?
- Which nodes on the path are non-colliders?
- Does  $\{X_3\}$  d-block this path? Does  $\{X_5\}$  d-block this path? Does  $\{X_3, X_5\}$  d-block this path?
- Does  $X_1$  d-separate  $X_2$  from  $X_4$ ?
- Is  $X_1 \perp_{\mathcal{G}} X_5 \mid \{X_3, X_4\}$ ?