

DAGs as a tool for understanding psychosocial processes



Briana Mezuk, PhD
MCUAAAR Analysis Core
Mini-Session
14 February 2024

Inspiration for today's session



Journal of Clinical Epidemiology 142 (2022) 264–267

**Journal of
Clinical
Epidemiology**

ORIGINAL ARTICLE

Tutorial on directed acyclic graphs

Jean C. Digitale, Jeffrey N. Martin, Medellena Maria Glymour*

Department of Epidemiology and Biostatistics, University of California, 550 16th St, 2nd Floor, San Francisco, CA 94158

Received 8 May 2021; Received in revised form 28 July 2021; Accepted 2 August 2021; Available online 8 August 2021

Abstract

Directed acyclic graphs (DAGs) are an intuitive yet rigorous tool to communicate about causal questions in clinical and epidemiologic research and inform study design and statistical analysis. DAGs are constructed to depict prior knowledge about biological and behavioral systems related to specific causal research questions. DAG components portray who receives treatment or experiences exposures; mechanisms by which treatments and exposures operate; and other factors that influence the outcome of interest or which persons are included in an analysis. Once assembled, DAGs — via a few simple rules — guide the researcher in identifying whether the causal effect of interest can be identified without bias and, if so, what must be done either in study design or data analysis to achieve this. Specifically, DAGs can identify variables that, if controlled for in the design or analysis phase, are sufficient to eliminate confounding and some forms of selection bias. DAGs also help recognize variables that, if controlled for, bias the analysis (e.g., mediators or factors influenced by both exposure and outcome). Finally, DAGs help researchers recognize insidious sources of bias introduced by selection of individuals into studies or failure to completely observe all individuals until study outcomes are reached. DAGs, however, are not infallible, largely owing to limitations in prior knowledge about the system in question. In such instances, several alternative DAGs are plausible, and researchers should assess whether results differ meaningfully across analyses guided by different DAGs and be forthright about uncertainty. DAGs are powerful tools to guide the conduct of clinical research. © 2021 Elsevier Inc. All rights reserved.

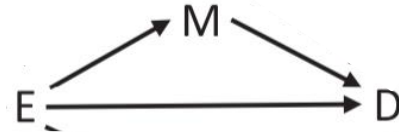
What is a Directed Acyclical Graph (DAG)?

DAGs are an “**intuitive** yet **rigorous** tool to **communicate about causal questions** in clinical and epidemiologic research and **inform study design and statistical analysis.**”

- What makes them **intuitive**?
- What makes them **rigorous**?
- How do they **communicate causal information**?
- How do they **inform study design and analysis**?

What makes DAGs intuitive?

- They are a way to **visually represent your hypotheses or assumptions** about the biopsychosocial processes that are relevant to your research question.
- In this diagram there are **2 causal paths**:
 - Exposure \rightarrow Outcome
 - Exposure \rightarrow Mediator \rightarrow Outcome



To identify the causal effect of E on D, we must block all non-causal paths and none of the causal paths between the two variables.

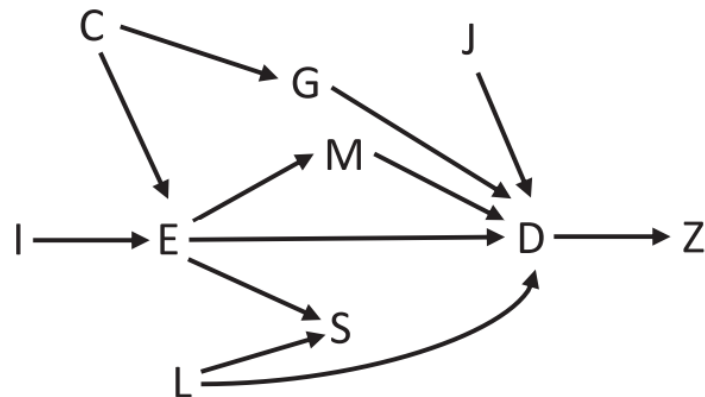
Causal paths linking E and D:

$E \rightarrow M \rightarrow D$

$E \rightarrow D$

What makes DAGs intuitive?

- **But the world is more complex!**
- DAGs they also allow you to depict:
 - What influences exposure (**I=instrument**) and **C (confounder)**?
 - All proposed mechanisms – including **non-causal paths** - linking exposure to the outcome
 - **The consequences of controlling (or failing to control) for different types of variables that are relevant to your research question.**
 - **The consequences of selecting a sample with a specific status/from a particular clinic, etc.**



To identify the causal effect of E on D, we must block all non-causal paths and none of the causal paths between the two variables.

Causal paths linking E and D:

$E \rightarrow M \rightarrow D$

$E \rightarrow D$

Non-causal paths linking E and D and how to block them:

$E \leftarrow C \rightarrow G \rightarrow D$ (block by controlling for C or G)

$E \rightarrow S \leftarrow L \rightarrow D$ (blocked provided we do not control for S)

Key terms:

- C confounds the association of E and D.
- G can be controlled to block the confounding path between E and D.
- M partially mediates the effect of E on D.
- S is a collider on a non-causal path between E and L, and therefore a collider on a non-causal path between E and D. Controlling or restricting on S will create a biased association between E and D.
- Z is a descendant of D.
- I is an instrumental variable, such as randomization, for the effect of E on D.
- J causes D and will therefore be an effect modifier of any other cause of D on at least one scale (additive or multiplicative).

Ingredients for creating a DAG

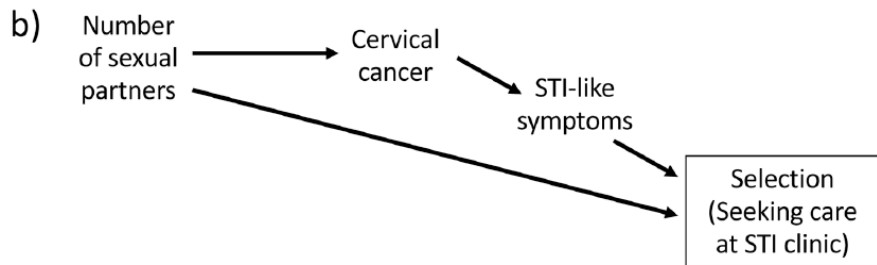
1. Specify a *causal* question, including the exposure (E) and outcome (D)
2. Specify variables (Confounders (C), Moderators (J) etc) that may influence the $E \rightarrow D$ relationship, either through their association with E, mediators of $E \rightarrow D$ (M) or D
3. Specify whether there are discrepancies between the *constructs* you are testing and their *measurement* (e.g., are the measures you have ideally what you want, or are they just proxies?)
4. Specify selection factors that influence entry into your sample
5. Specify the relationships between these variables

NOTE: Just because you don't have a specific measure in your dataset, that doesn't mean you shouldn't be represented in your DAG if it is in the list above.

This is how DAGs can help you understand the potential influence of unmeasured confounders.

One bit of epid-specific vocabulary

COLLIDER (aka a way to visualize selection bias)



- **E=Number of sexual partners**
- **D=Cervical cancer**
- **J=DECEDENT of D** (symptoms that develop after the cancer has occurred that prompt seeking care at STI clinic)
- **S=Collider** (not part of the causal relationship of interest, but a “status” that is caused by both E and D)

- **Example: You want to estimate the relationship between number of sexual partners on risk of cervical cancer.**
- **You decide to recruit your sample from a local STI-clinic.**
- **Q: How might this study design decision (e.g., where to recruit your sample) bias your analytic inferences?**
- **A: It will likely bias the E→D association towards the null (under-estimate a true effect).**
- **Why? Who is not represented and/or under-represented in the sampling frame? How does that effect the range of values for E?**

What makes a DAG a *causal* diagram?

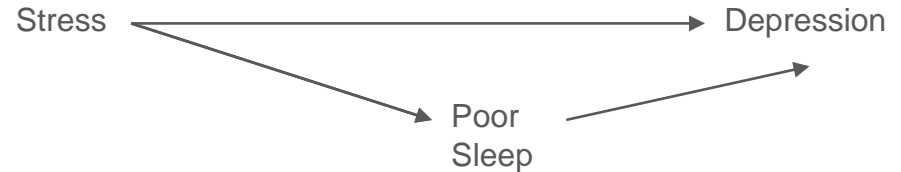
- In order to specify that $E \rightarrow D$ [or $E \rightarrow M \rightarrow D$] is a ***causal relationship***, you must demonstrate that there are no other explanations (e.g., no non-causal pathways) for why E is correlated with D
 - That is, you must “block” all other potential paths linking E and D.
 - How do you BLOCK a path?
 - By *controlling* for it (e.g., controlling for a common cause or an intermediate mechanism) in your regression models
 - By *matching* on it (e.g., case control, case-crossover, family-based designs)
 - By *not controlling* for a COLLIDER (or not influencing the $E \rightarrow D$ relationships of interest through selection bias (e.g., our STI clinic example))

Causal Hypothesis: Poor sleep the primary mechanism linking stress and depression in middle-aged adults.

Exposure (E): Stress

Outcome (O): Depression

What is the hypothesized mediator (M)?



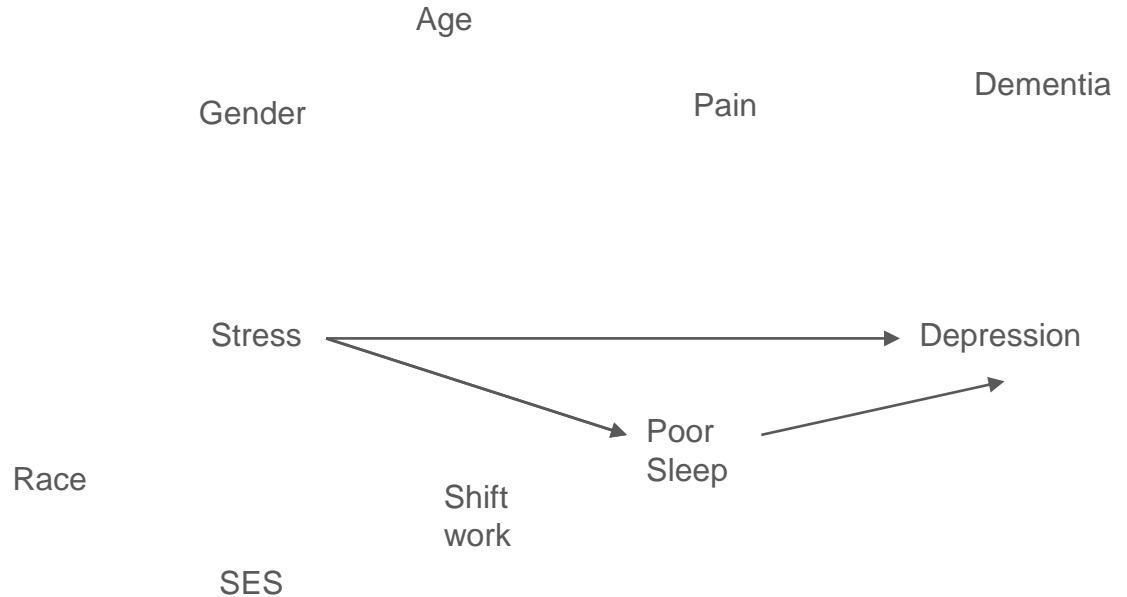
Causal Hypothesis: Poor sleep the primary mechanism linking stress and depression in middle-aged adults.

Exposure (E): Stress

Outcome (O): Depression

What is the hypothesized mediator (M)?

On your own: connect the dots!



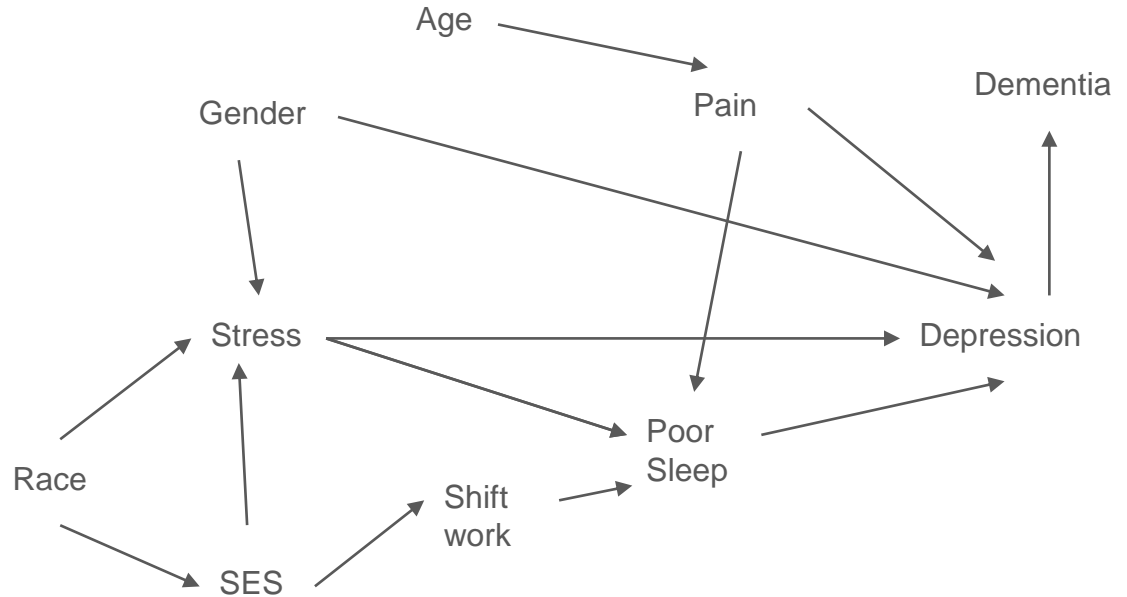
Causal Hypothesis: Poor sleep the primary mechanism linking stress and depression in middle-aged adults.

Exposure (E): Stress

Outcome (O): Depression

What is the hypothesized mediator (M)?

On your own: connect the dots!

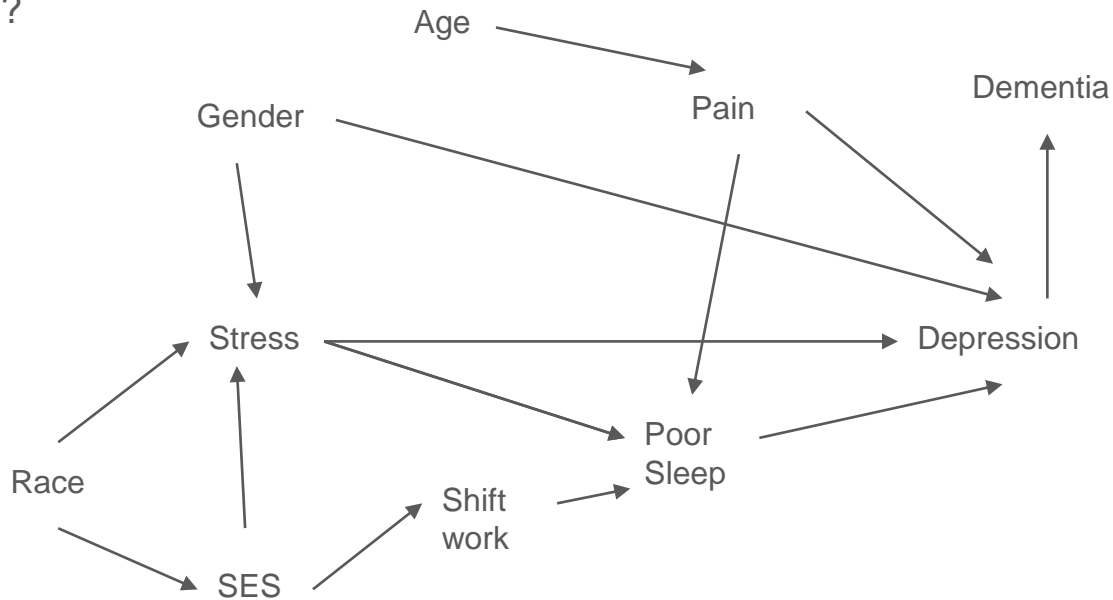


Causal Hypothesis: Poor sleep the primary mechanism linking stress and depression in middle-aged adults.

Should I control for **PAIN** to estimate the causal link between stress→depression? Why or why not?

What will happen if I recruit my sample from a factory that runs 3 shifts of work? What if I **only** recruit from the 3rd shift?

What will happen if I recruit my sample from a **memory clinic**?



Limitations of DAGs

- DAGs forces us to admit that, often, because of limitations in our prior knowledge, we may not know which of several possible DAGs is correct.
 - They can still guide our analyses and help us consider alternative “thought experiments” that we can potentially test in different ways to enhance the rigor of our analysis
 - Example: Identified a relationship in a clinic sample? See if it replicates in a general population one.
- DAGs do not convey information about magnitude or functional form of causal relationships
 - This means they are not great for visualizing effect-measure modification or moderators.

Summary

- Drawing DAGs can help us understand a wide range of **potential threats to drawing valid inferences** about causal relationships from observational data.
- Particularly useful at the start of a project to **inform study design and analysis**:
 - Sampling source and frame
 - Variables to measure (if collecting your own data)
 - Variables to include in your models (if using existing data)
- Also useful as a **tool for building “What if...” scenarios** that you may be able to test to assess the robustness of your inferences.
 - With external data
 - With other specifications of the data you have already analyzed (e.g., testing different cutoffs for binary variables, testing different forms of interactions)
 - Simulated data
 - Other sensitivity analyses