

ECON 672

Week 3: Causal Diagrams and Directed Acyclic Graphs

Samuel Rowe, PhD 12/15/2022

Overview

- Directed Acyclic Graphs (DAG)
- Confounders
- Colliders
- Examples

Directed Acyclic Graphs

Directed Acyclic Graphs

- Directed Acyclic Graphs (DAGs) are a chain of causal effects in graphical form
- One of many contributions to causal inference from Judea Pearl (2009)
- DAGs are a model
 - Based upon unobserved structural process
 - Equilibrium values of a system of behavioral equations

Directed Acyclic Graphs

- Upfront statements about DAGs
- 1) Causality usually runs in one direction
 - No time cyclic in DAGs (hence acyclic)
- 2) Reverse causality should be handled in a different way
 - Simultaneity, such as supply and demand models, would require multiple nodes and are not best handled with a DAG
- 3) DAGs explain causality in terms of counterfactuals
 - Causal effects through two potential states

Directed Acyclic Graphs

- Arrows and Nodes
 - Nodes are circles and they represent random variables
 - Arrows between nodes represent direction of causality
- Causal effects can occur in two ways
 - 1) Direct
 - 2) Indirect

Directed Acyclic Graphs

- Direct Causal Effect
 - One variable directly affects another variable
 - $D \rightarrow Y$
- Indirect Causal Effect
 - Our treatment variable D is mediated by a third variable, such as X
 - $D \rightarrow X \rightarrow Y$

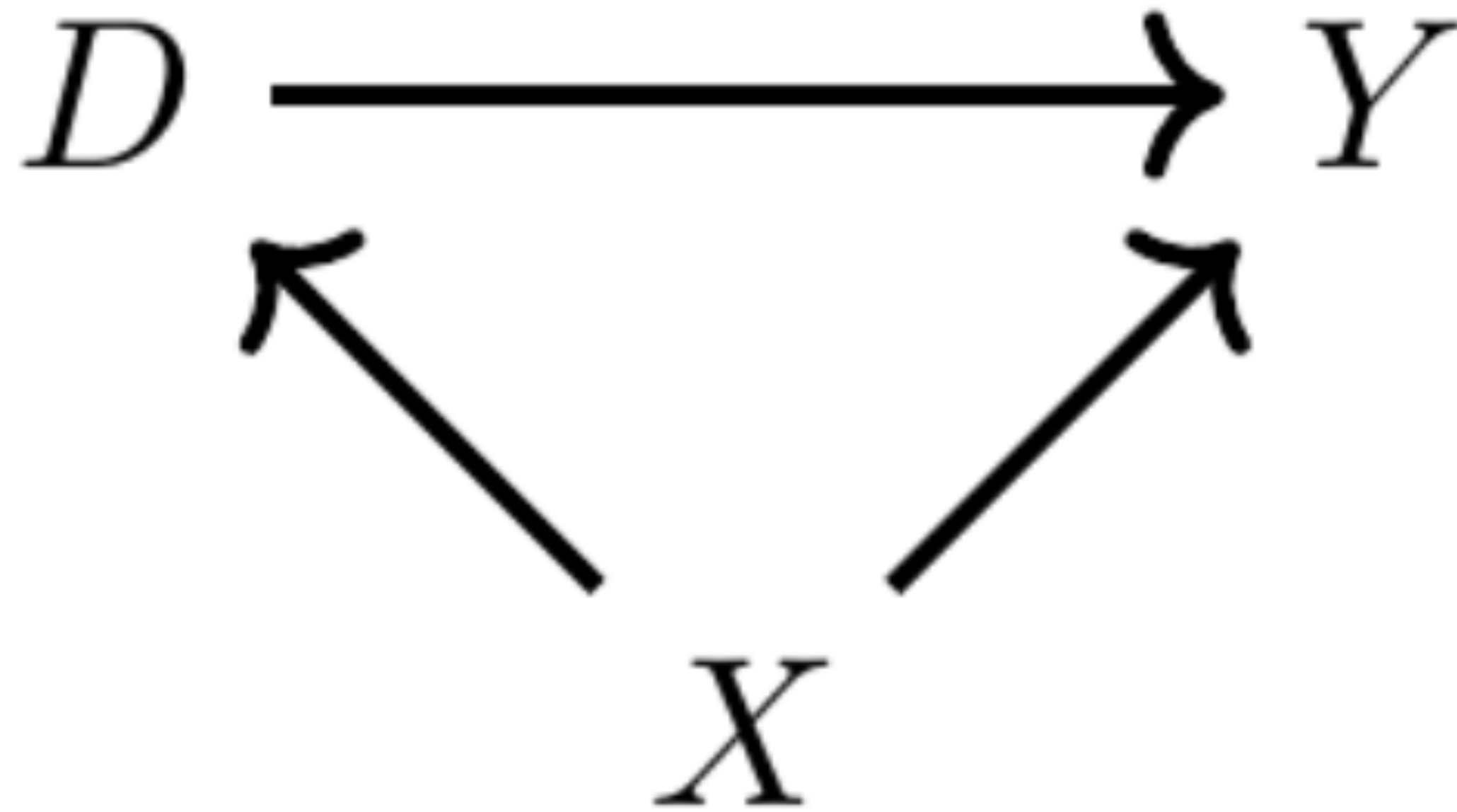
Directed Acyclic Graphs

- DAGs are meant to show all causal relationship of the model
 - These are theoretical representations of phenomena you are interested in studying
 - These can be develop through theory or prior literature
- What is included is just as important as what is not included
 - Direction of arrow between nodes imply causality
 - A lack of an arrow implies no causal relationship
- A complete DAG will have all direct effects among the variables including common causes of any pair of variables

Why Use Directed Acyclic Graphs

- A useful state-of-the-art knowledge of the phenomena you are interested in studying
 - Shows theory, literature, and institutional/prior knowledge
- They provide a picture representation of your model
 - Visually displays your research design and identification strategy
- Corroborates your research design
 - Shows causal effect of intervention by showing backdoor criterion and collider biases
- Provides your assumptions

Simple DAG



Simple DAG

- We have 3 random variables
 - D is our treatment
 - Y is our outcome of interest
 - X is a variable that affects both D and Y

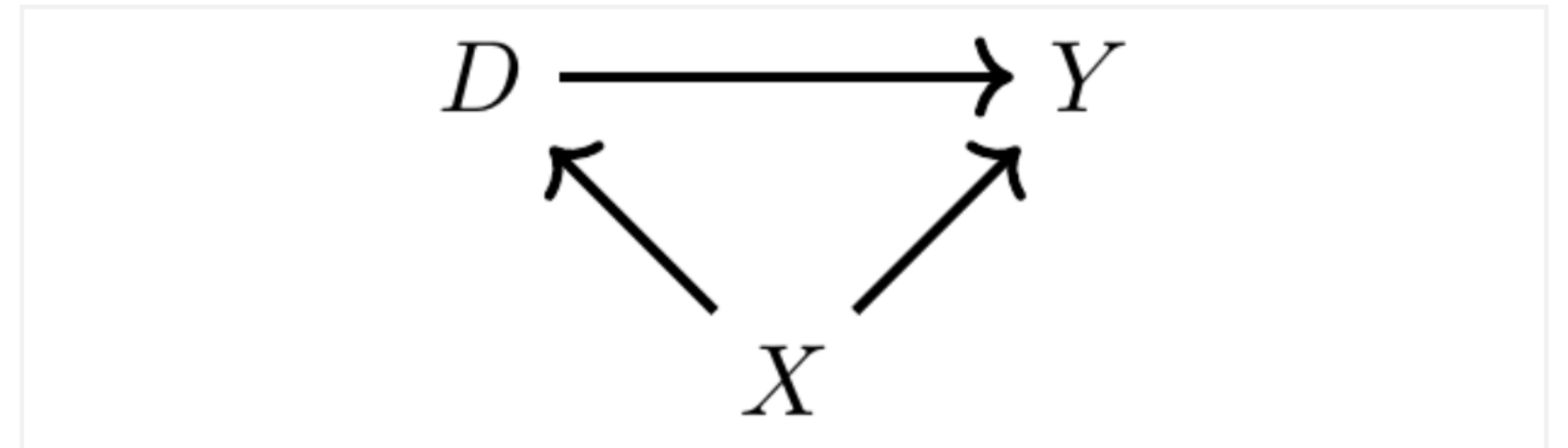
Simple DAG

- There are two pathways affecting Y
- 1) Direct Pathway: $D \rightarrow Y$
 - This is our causal effect of interest
- 2) Indirect Pathway: $D \leftarrow X \rightarrow Y$
 - A backdoor pathway shows that D and Y take on different values when X takes on different values
 - This means that part of the correlation between D and Y is spurious due to X

Confounders

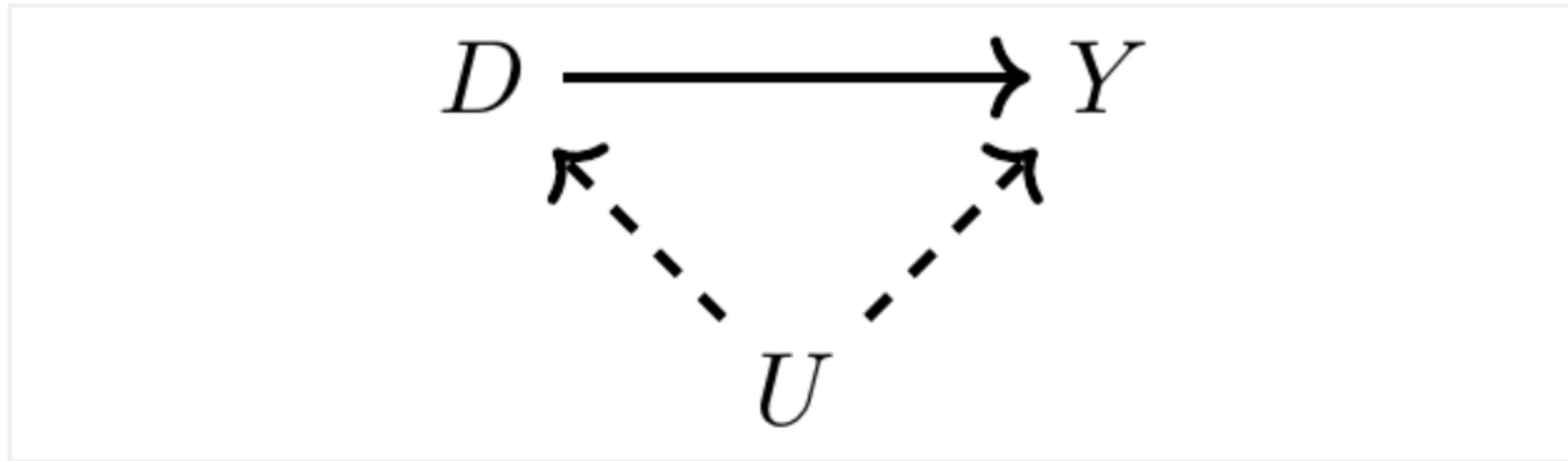
Confounders

- The backdoor pathway is one of the most important concepts with DAGs
- When X mediates the values of D and Y , then X is considered a **confounder**
- Use the Simple DAG example
 - X is a confounder
 - It jointly impacts D and Y
 - There are two pathways contained in the correlation between D and Y
- Backdoor paths are similar to Omitted Variable Bias
 - Leaving a backdoor open is similar to not controlling for a variable



Unobserved Confounders

- Here is an example of a DAG with unobserved confounder



- Unlike X, U is unobserved and jointly impact D and y
 - We will use dash lines to indicate that the variable is unobserved

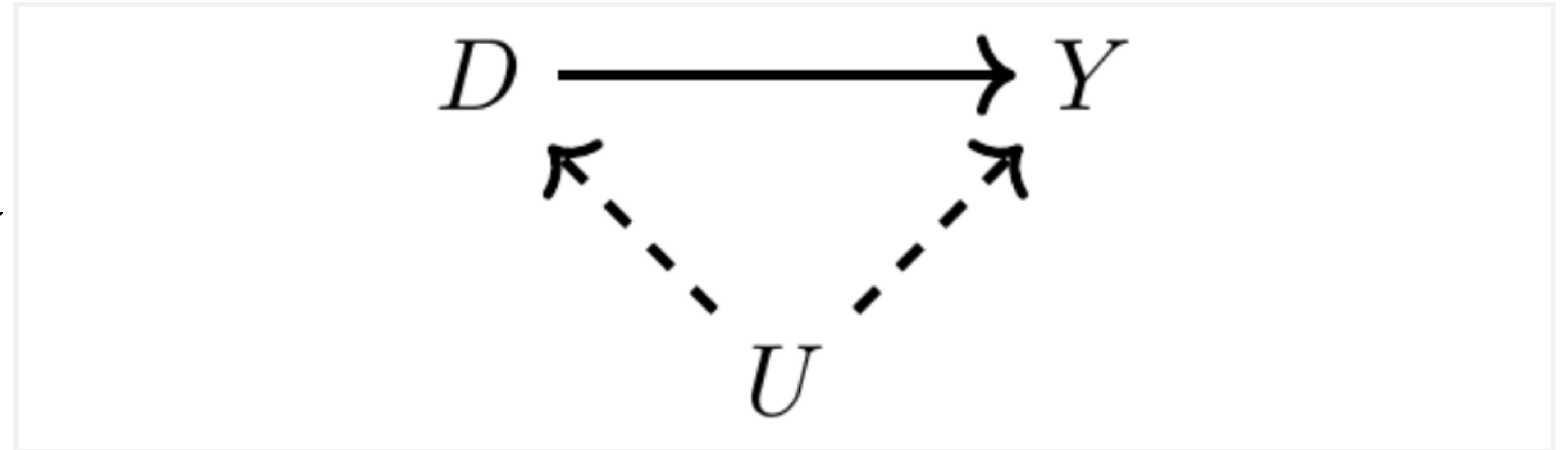
Unobserved Confounder

- There are two pathways similar to our simple DAG

- Direct: $D \rightarrow Y$

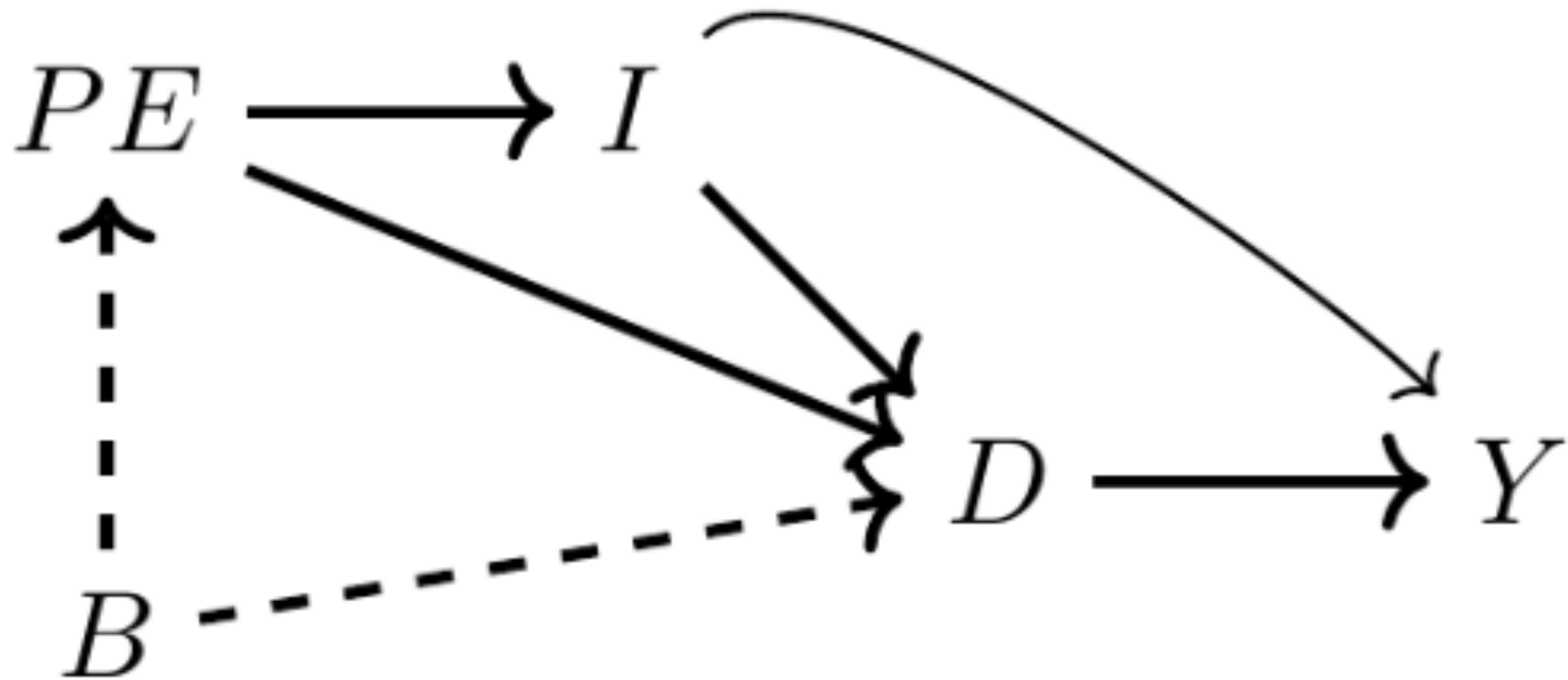
- Indirect: $D \leftarrow - U - \rightarrow Y$

- Since U is unobserved



- The pathway remains open since we cannot control for it

DAG Example: College Education and Earnings



DAG Example

- The model shows us two things
 - Pathways
 - Assumptions
- There are several variable in this DAG
 - Y is observed earnings
 - D is observed college education
 - PE is observed parental education
 - I is observed family education
 - B is unobserved background characteristics, such as genetics, family environment, mental ability, etc.

DAG Example

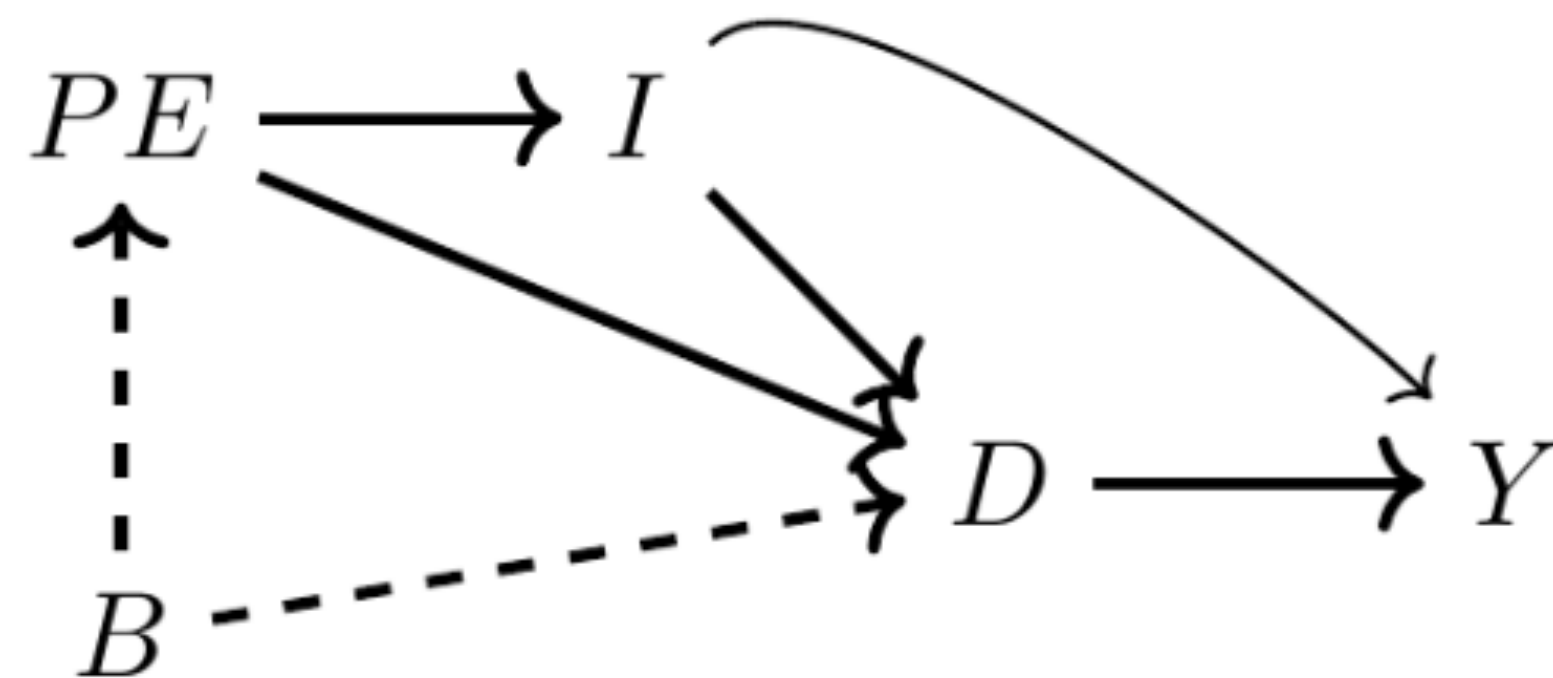
- There are 4 pathways here
- Direct: $D \rightarrow Y$
- Indirect Backdoor 1: $D \leftarrow I \rightarrow Y$
- Indirect Backdoor 2: $D \leftarrow PE \rightarrow I \rightarrow Y$
- Indirect Backdoor 3: $D \leftarrow - B - \rightarrow PE \rightarrow I \rightarrow Y$

DAG Example

- Narrative of this Example DAG
 - College education $[0,1]$ affects child's earnings ($D \rightarrow Y$)
 - Family income affects child's income due to bequests, transfers, etc
 - Parent's education affects family income and child's choice of college education
 - Family background characteristics affect parents education choice and child's education choice
- Assumption
 - Background characteristics do not directly affect child's earnings
 - Assumption is the background characteristics work indirectly through parent's and child's education choices

Example DAG

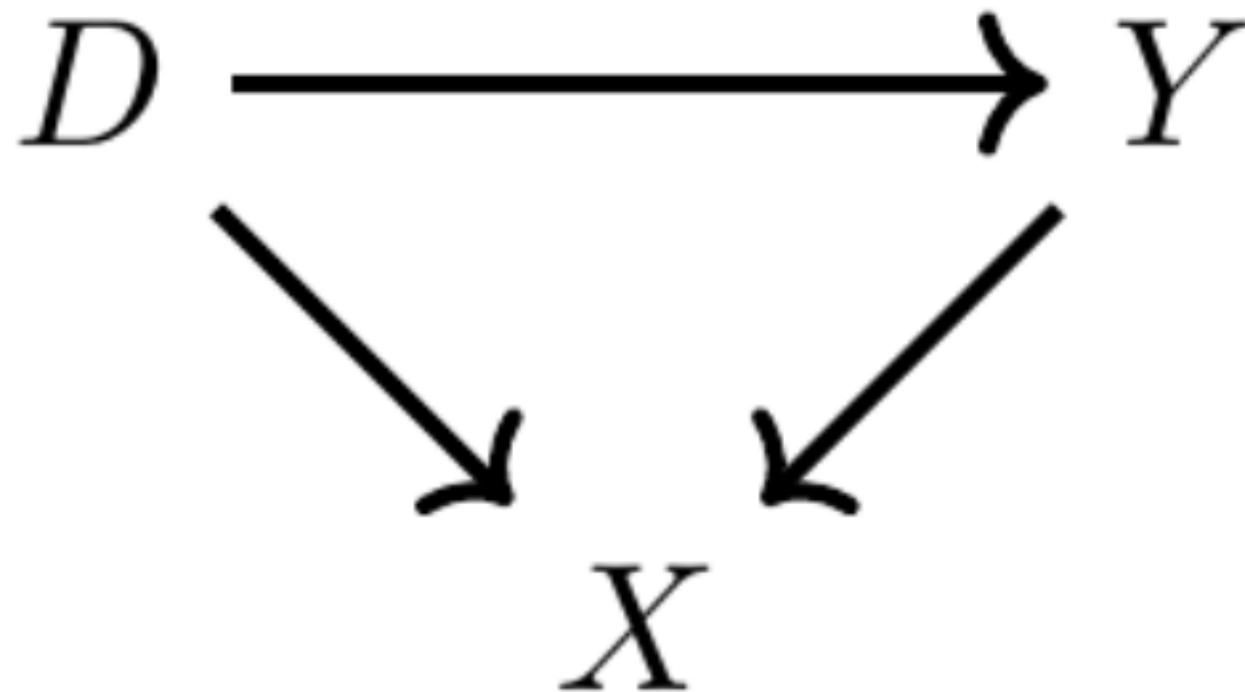
- If we naively compare outcomes of college education to no college education
 - It will be biased due to several backdoors not being closed



Colliders

Colliders

- Colliders are different from confounders
 - Colliders occur when two variables cause a third variable
 - They are a bit more complex than confounders
- Unlike a confounder, controlling for a collider will introduce bias
 - Angrist and Pischke (2009) refer to these as bad controls



Colliders

- We have three variables, D, X, and Y
 - D and Y collide at $D \rightarrow X \leftarrow Y$, such that changes in D and Y cause changes in X
- We have two pathways
 - Direct: $D \rightarrow Y$
 - Indirect (Backdoor): $D \rightarrow X \leftarrow Y$
- **Key Point**
 - Leaving a collider alone closes the backdoor pathway
 - If you control for the collider, you reopen the backdoor pathway

Backdoor Criterion

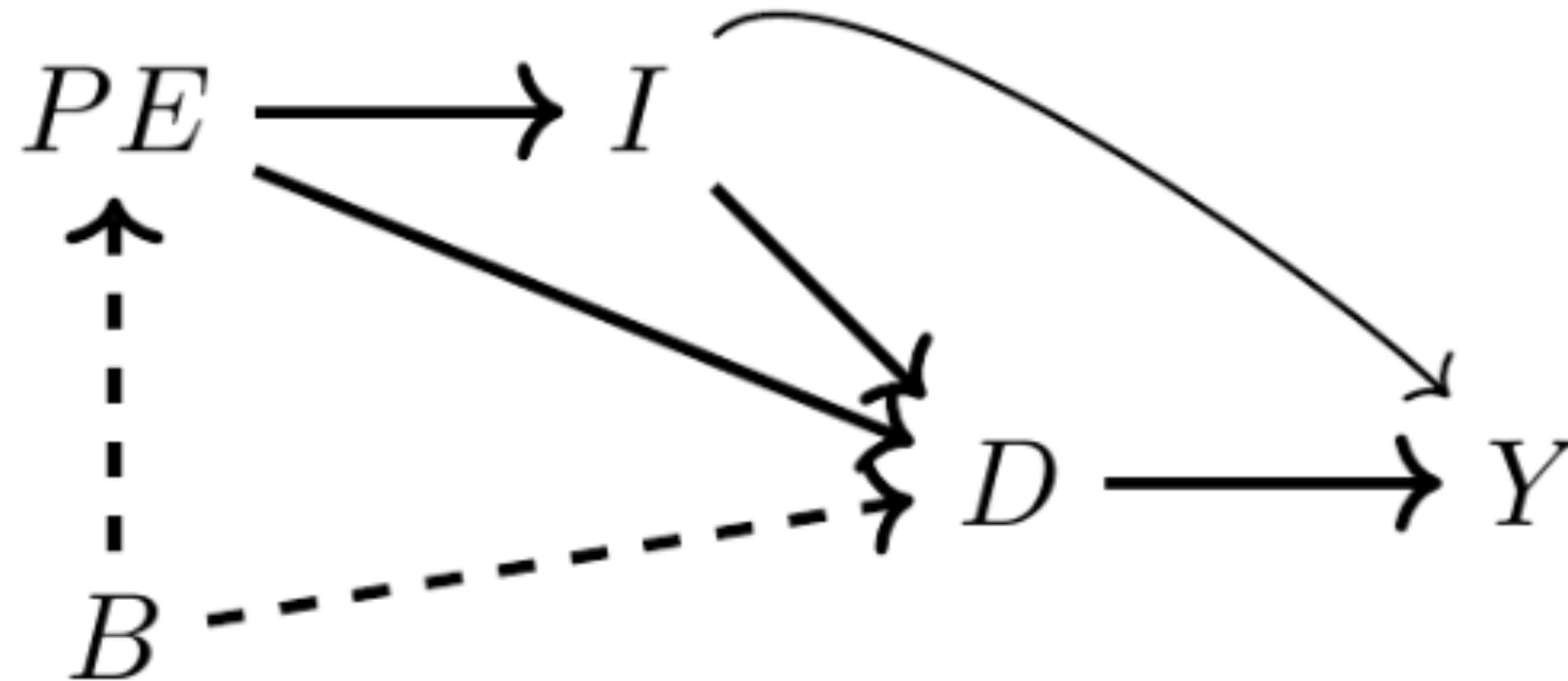
- Backdoor pathways are indirect pathways between D and Y
 - These matter since they create systematic noncausal relationships between our treatment of interest (D) and outcome of interest (Y)
- Open backdoor pathways
 - These are omitted variable bias
- Our Goal
 - Close all backdoor pathways
 - Identify the direct pathway between D and Y
- We use an identification strategy to achieve our goal in the DAG

Backdoor Criterion

- Two ways to close the backdoor criterion
 - 1) Controlling or conditioning on the confounder
 - 2) The appearance of a backdoor collider
- ***Backdoor Criterion***
 - Backdoor criterion is satisfied when all backdoor pathways are closed

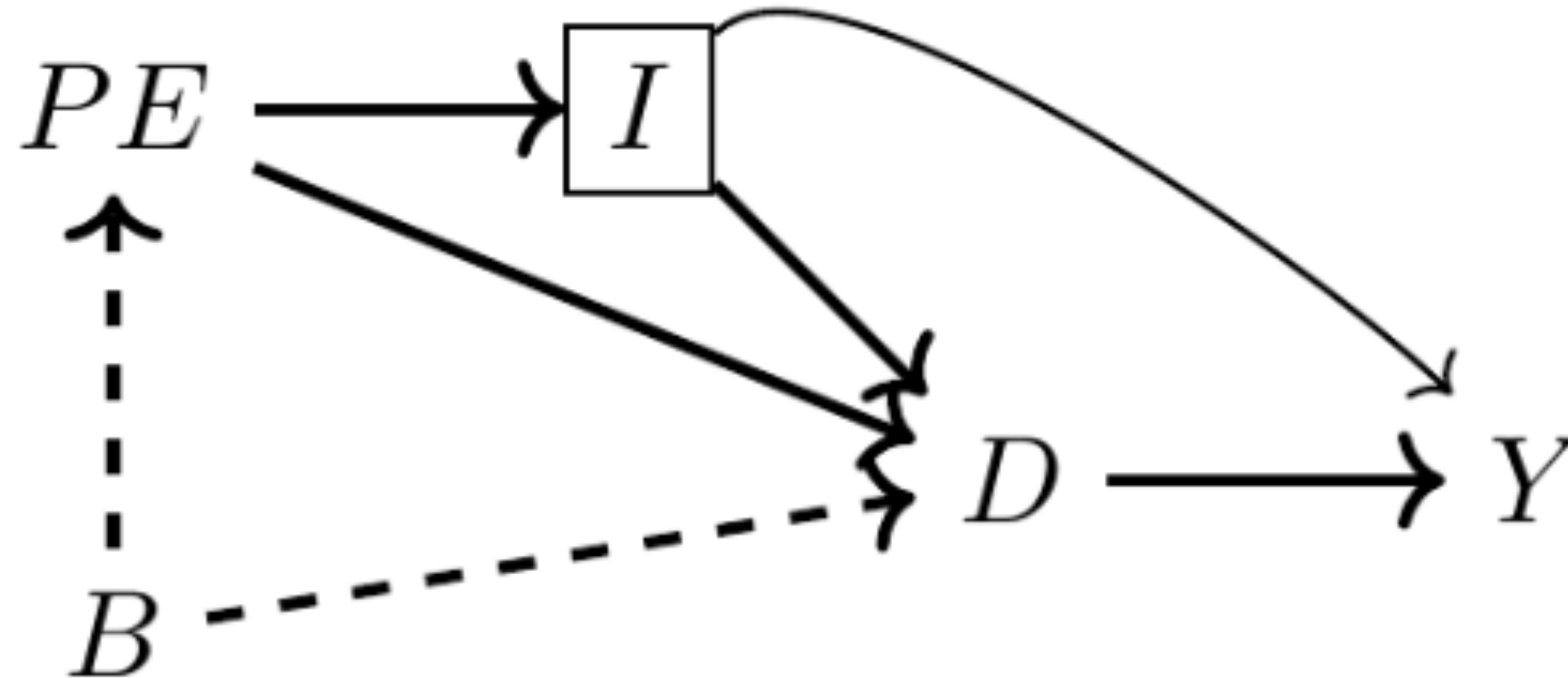
DAG Example: College Education and Earnings

- How do we satisfy the backdoor criterion in our college education and earnings example?



DAG Example: College Education and Earnings

- How do we satisfy the backdoor criterion in our college education and earnings example?

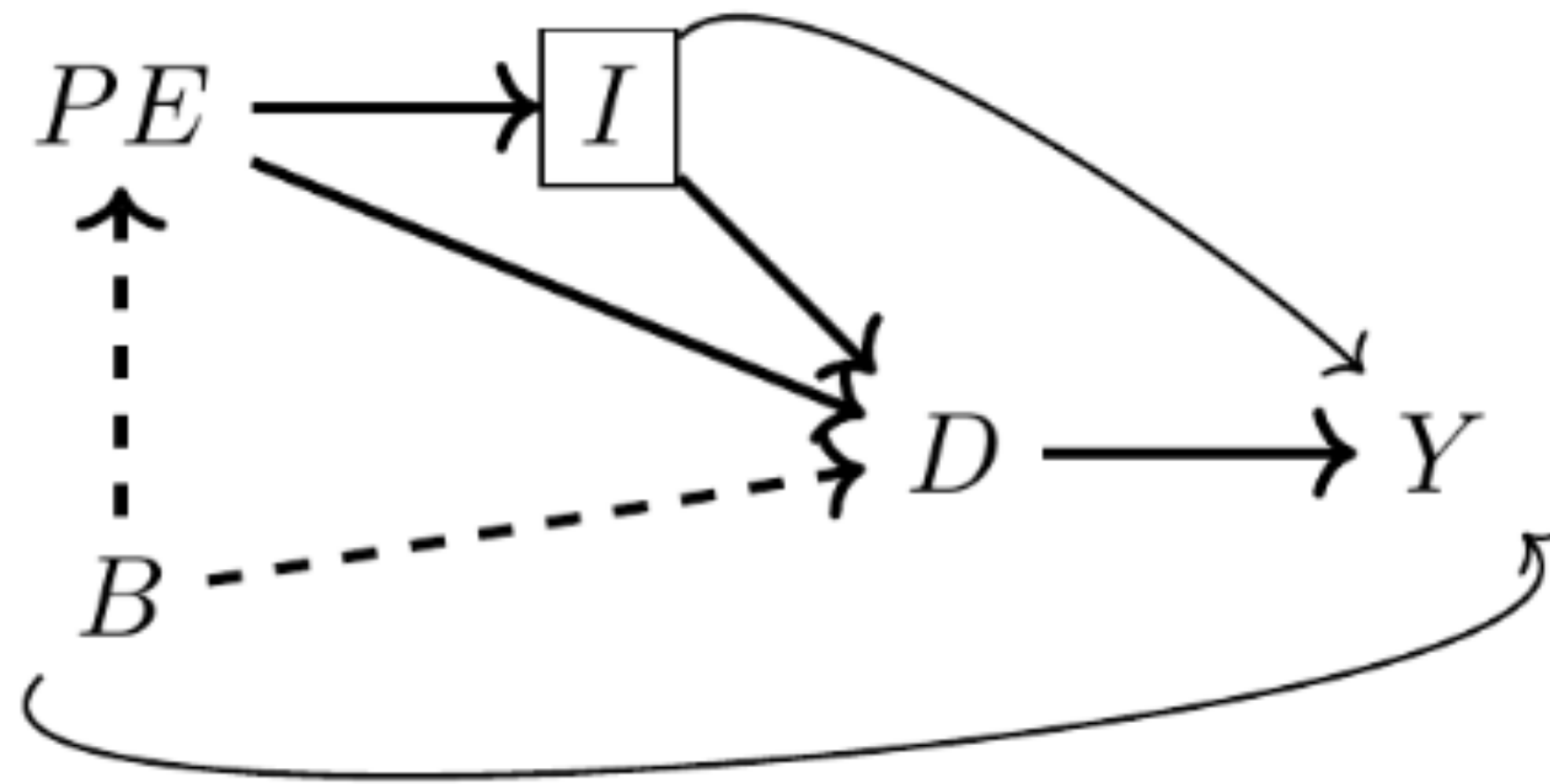


DAG Example: College Education and Earnings

- We can close all backdoor pathways when we control for Family Income
 - Family income runs along all backdoor pathways
 - This is the minimally sufficient strategy to satisfy the backdoor criterion
- $Y_i = \alpha + \hat{\delta}D_i + \beta I_i + \varepsilon_i$
- $\hat{\delta}$ takes a causal interpretation when we satisfy the backdoor criterion

Skepticism of DAG Strategy

- Is the assumption that family background doesn't affect sufficient?
 - This is the importance of theory, literature, and prior knowledge
- If family background does impact the child's college choice, then controlling for family income is an insufficient strategy, since $D \leftarrow - B - \rightarrow Y$



Bias Examples

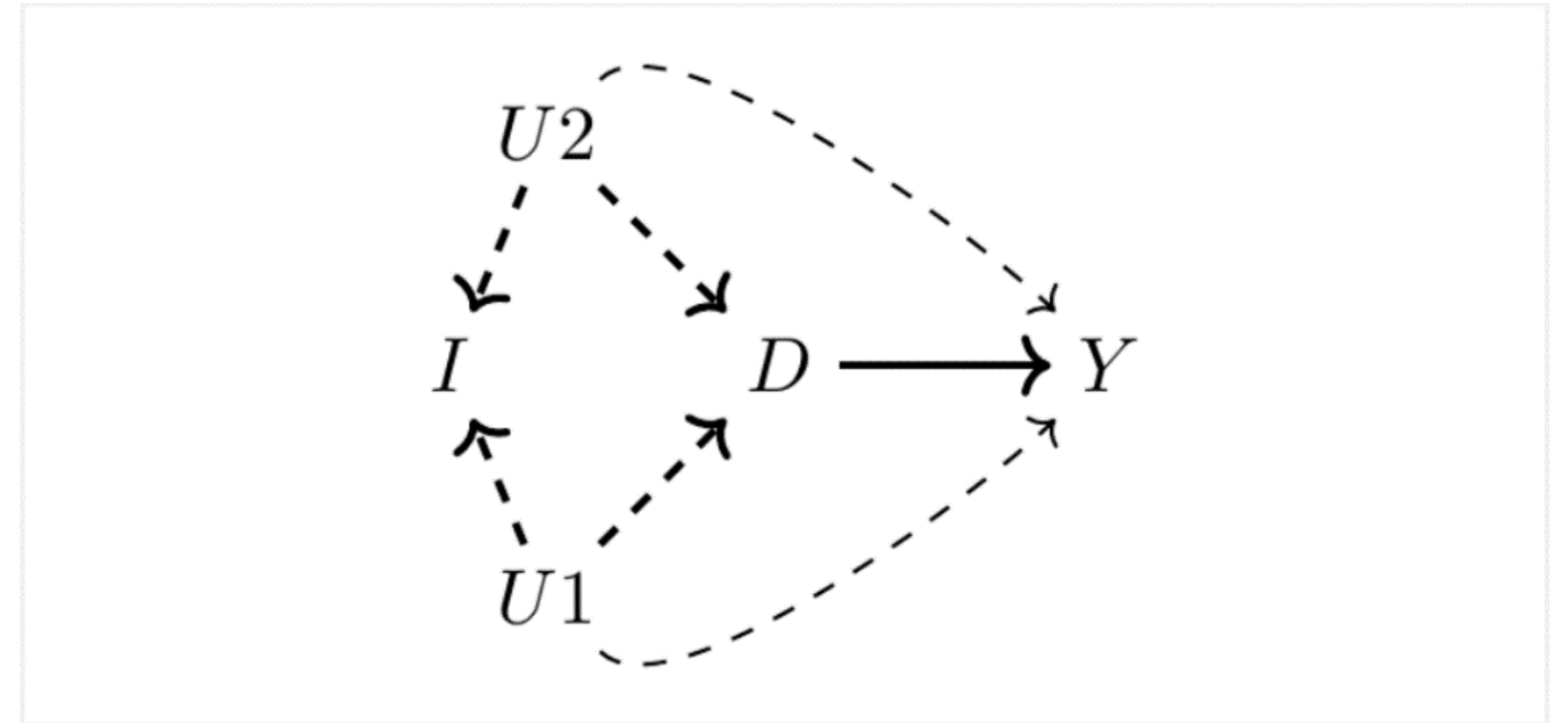
- There are **no flags** for confounders or colliders in a data set
 - We need knowledge of the *data-generation process, theory, prior literature, and logic* to assign confounders and colliders
- As mentioned colliders are a bit weird
 - When we condition or control for a collider, **we introduce bias**
 - We open the backdoor pathways when we control for a collider
- We should be familiar with confounders that we have seen in econometrics
 - When we do not condition or control for a confounder, **we introduce bias**

Bias Examples

- Setting up a DAG: You need theory, prior literature, and prior knowledge of data-generation process
 - These flag colliders and confounders
 - These establish your assumptions

Collider Bias Example 1: College and Earnings

- We revisit the child's college choice
 - D is child's college choice
 - Y is child's earnings
 - I is family income
 - $U1$ is mother's unobserved ability
 - $U2$ is father's unobserved ability



Collider Bias Example 1: College and Earnings

- We have a few pathways

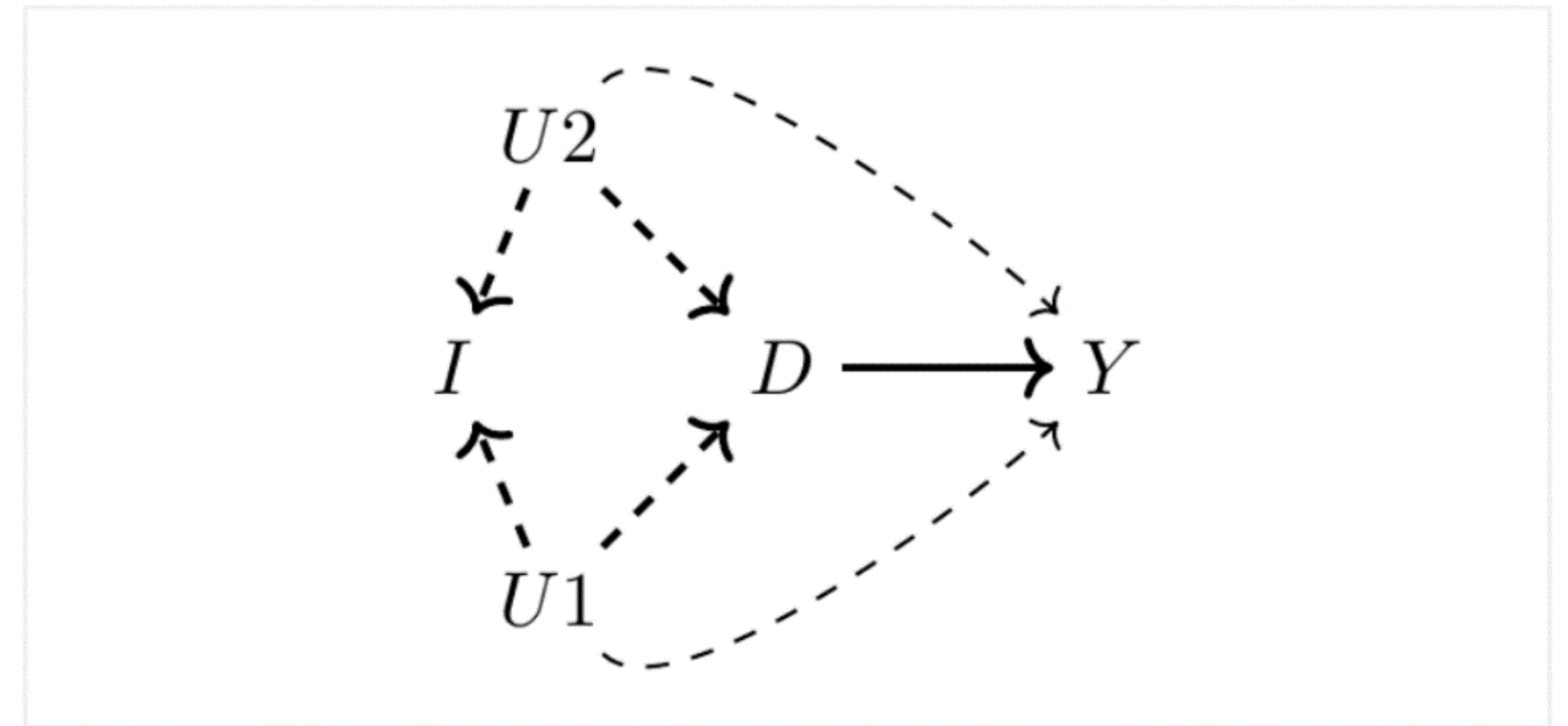
- Direct: $D \rightarrow Y$

- Backdoor 1: $D \leftarrow - U_1 - \rightarrow Y$

- Backdoor 2: $D \leftarrow - U_2 - \rightarrow Y$

- Backdoor 3: $D \leftarrow - U_1 - \rightarrow I \leftarrow - U_2 - \rightarrow Y$

- Backdoor 4: $D \leftarrow - U_2 - \rightarrow I \leftarrow - U_1 - \rightarrow Y$



- If we condition on family income, I , we will introduce collider bias

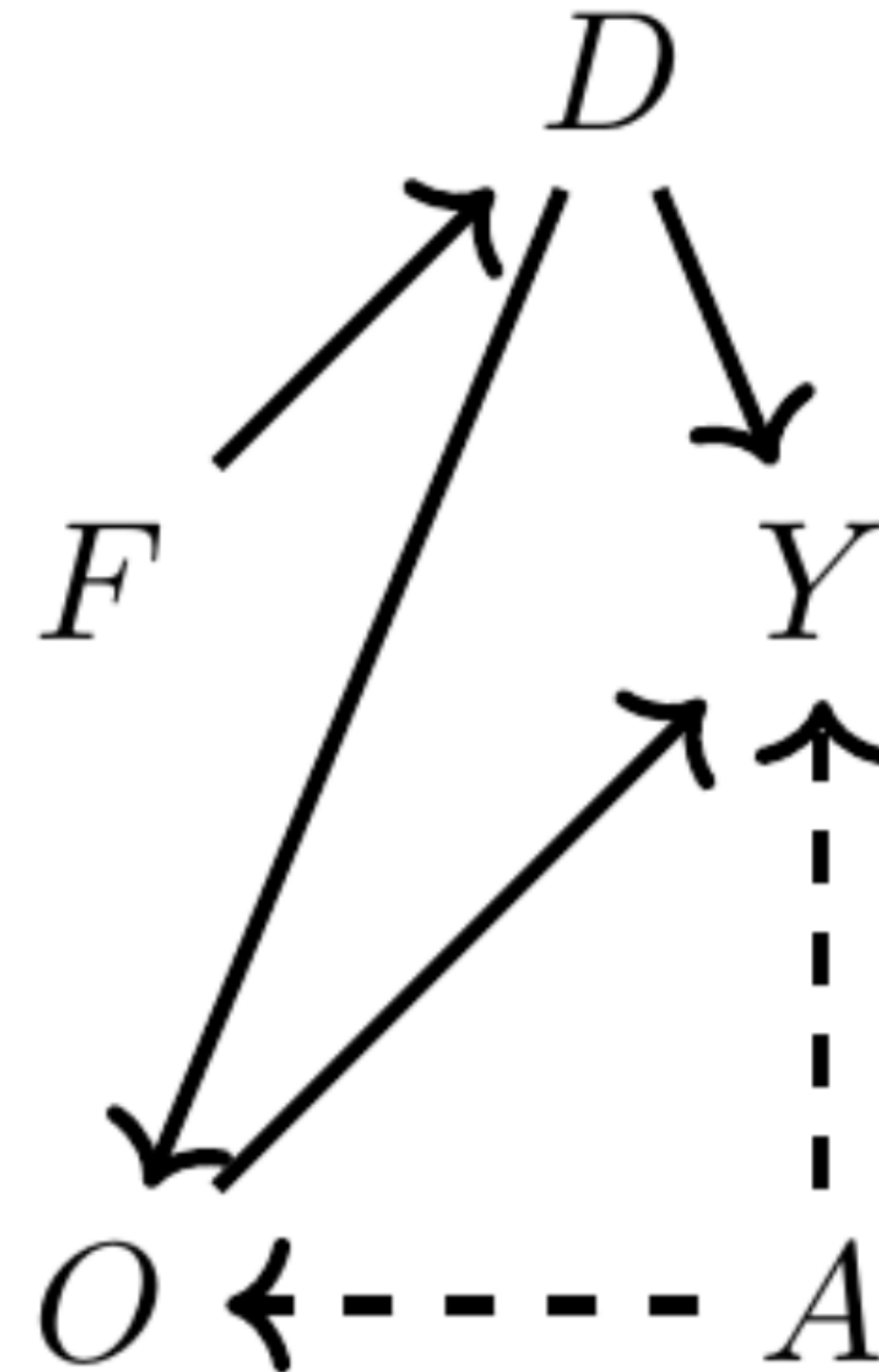
- The backdoor pathways of 3 and 4 are closed as long as we do not condition on family income, I
 - We need an identification strategy that takes care of time-invariant ability of parents

Collider Bias Example 2: Discrimination

- It is common to hear that wage disparity reduces or disappears when you control for occupation
 - An example is when an internal Google study showed that wage disparity was eliminated when they controlled for occupations within Google
- If discrimination come from job/occupational sorting, then controlling for occupation introduces collider bias
 - Worsens the bias of the estimate

Collider Bias Example 2: Discrimination

- Set up the DAG
 - Y is earnings
 - D is “treatment” of discrimination
 - O is occupation
 - F is female
 - A is unobserved ability



Collider Bias Example 2: Discrimination

- Pathways

- Direct Pathway: $D \rightarrow Y$

- Mediated Pathway: $D \rightarrow O \rightarrow Y$

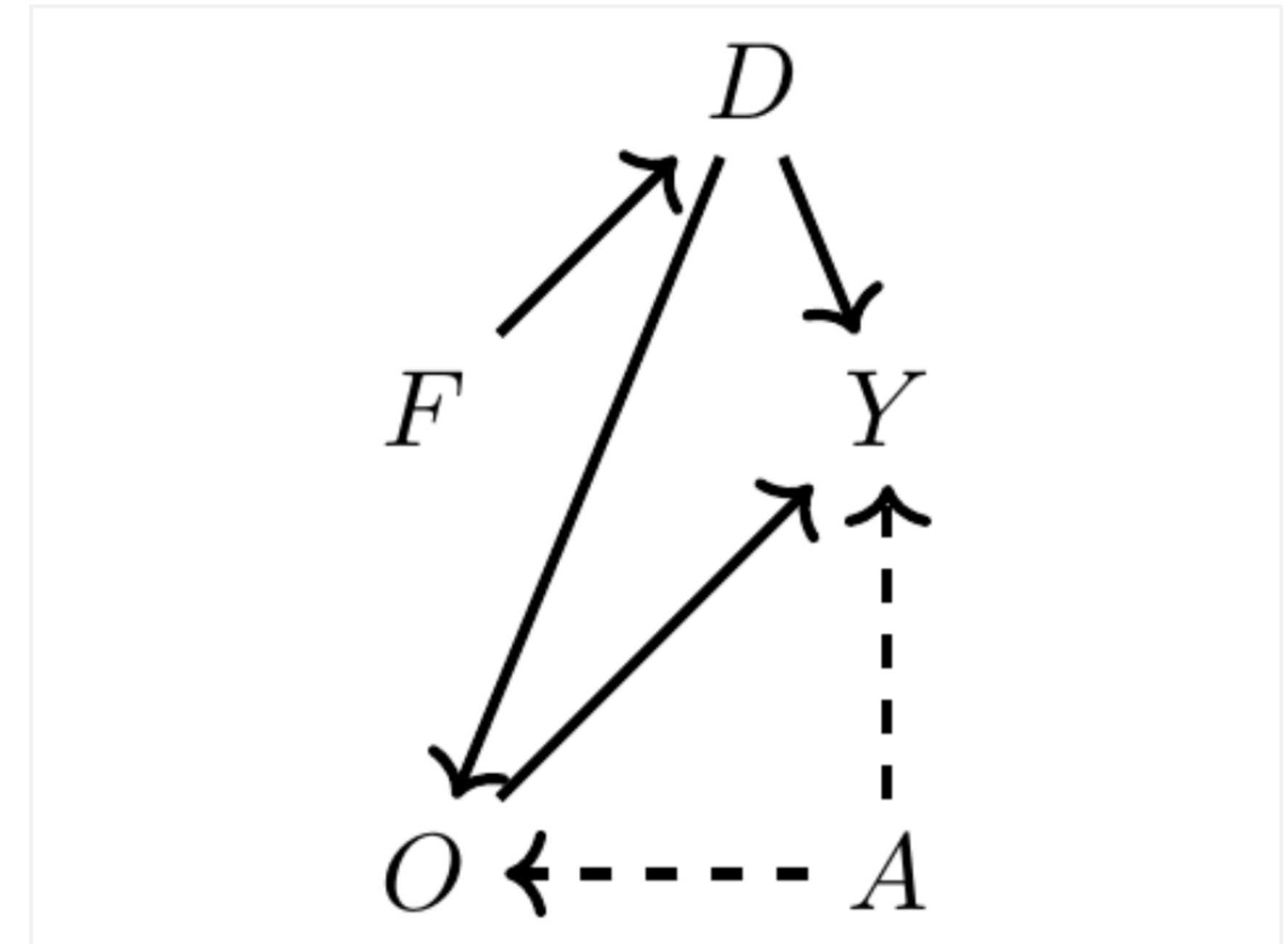
- Backdoor Pathway: $D \rightarrow O \leftarrow A \rightarrow Y$

- Mediated Pathway

- This means that discrimination is mediated by occupation

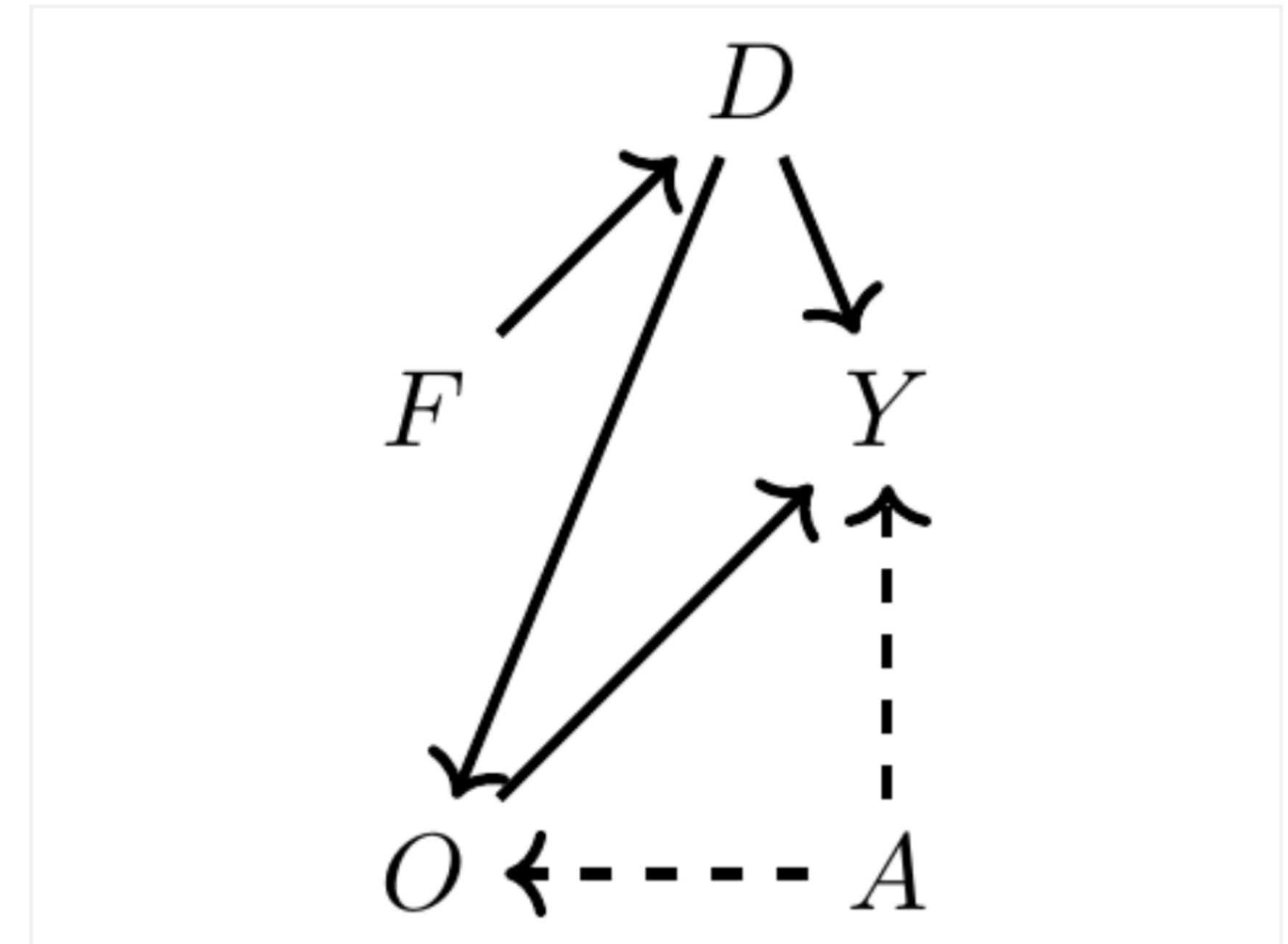
- It implies that discrimination affects the jobs or occupations that female can hold

- Discrimination means that women have fewer opportunities for higher paying jobs



Collider Bias Example 2: Discrimination

- Assumptions (what is not shown)
 - Female status has no direct impact on earnings
- Direct Pathway
 - Implies discrimination impact earnings
- Mediated Pathway
 - Implies women are discriminated against by what kind of jobs they are offered
- Backdoor Pathway
 - Implies ability affects earnings and occupations they sort into

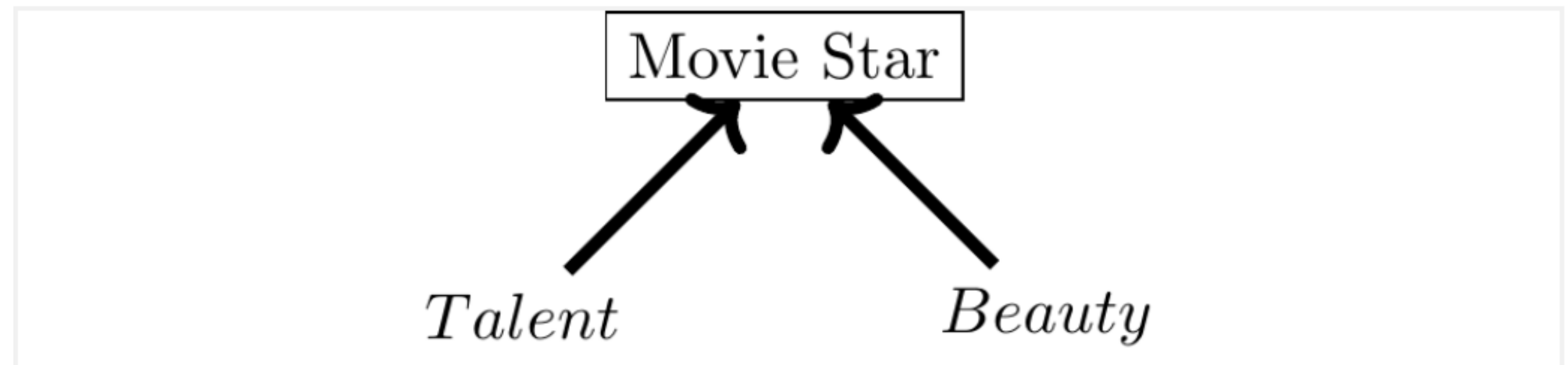


Collier Bias Example 2: Discrimination

- Stata Example
- We can get a total effect of discrimination onto earnings
 - Direct and Mediated
- When we control for occupation
 - It closes the mediated pathway, but opens up the backdoor pathway and introduces bias
 - This happens since ability does not affect D directly
 - Ironically, controlling for occupation makes the bias worse
- We need an identification strategy that controls for ability
 - We have ability in our example, but in real life we don't

Collider Bias Example 3: Sample Selection

- Collider bias can be baked into the sample
 - If the sample itself is a collider
- There is a story about beauty and talent being inversely related or negatively correlated for actors
 - *Talent* \rightarrow *MovieStar*
 - *Beauty* \rightarrow *MovieStar*



Collider Bias Example 3: Sample Selection

- Stata Example
- If there is a cutoff between all actors that separates movie stars and aspiring actors, then the frontier has a negative relationship
 - Movie Star status creates a collider bias when there is no relationship between beauty and talent
 - Movie Star status introduces a negative correlation between talent and beauty

Collider Bias Example 4: Policing and Admin Data

- DAGs can help spot subtle forms of conditioning on colliders and collider bias
- For example, admin data may be rife with collider bias
- Main problem with admin data
 - Admin data may be condition on an interaction occurring
 - Police interactions is the example used, but there could be many types of conditional interactions, such as wage violations, health violations, etc.

Collider Bias Example 4: Policing and Admin Data

- What is the data-generation process?
 - For the police admin data, data generation is conditional on a police interaction
 - The data-generation process is a ***function*** of police interactions
 - This means admin data are endogenous

Collider Bias Example 4: Policing and Admin Data

- Fryer (2019) wanted to study police force and racial bias
- He uses several public-use databases to study this problem
 - NYC Stop and Frisk database
 - This contained data on police stops and questioning of pedestrians
 - Police-Public Contact Survey
 - This was a survey of civilians describing interactions with police including the use of force from Houston
- The main problem is that these data are condition on police-civilian interactions

Collider Bias Example 4: Policing and Admin Data

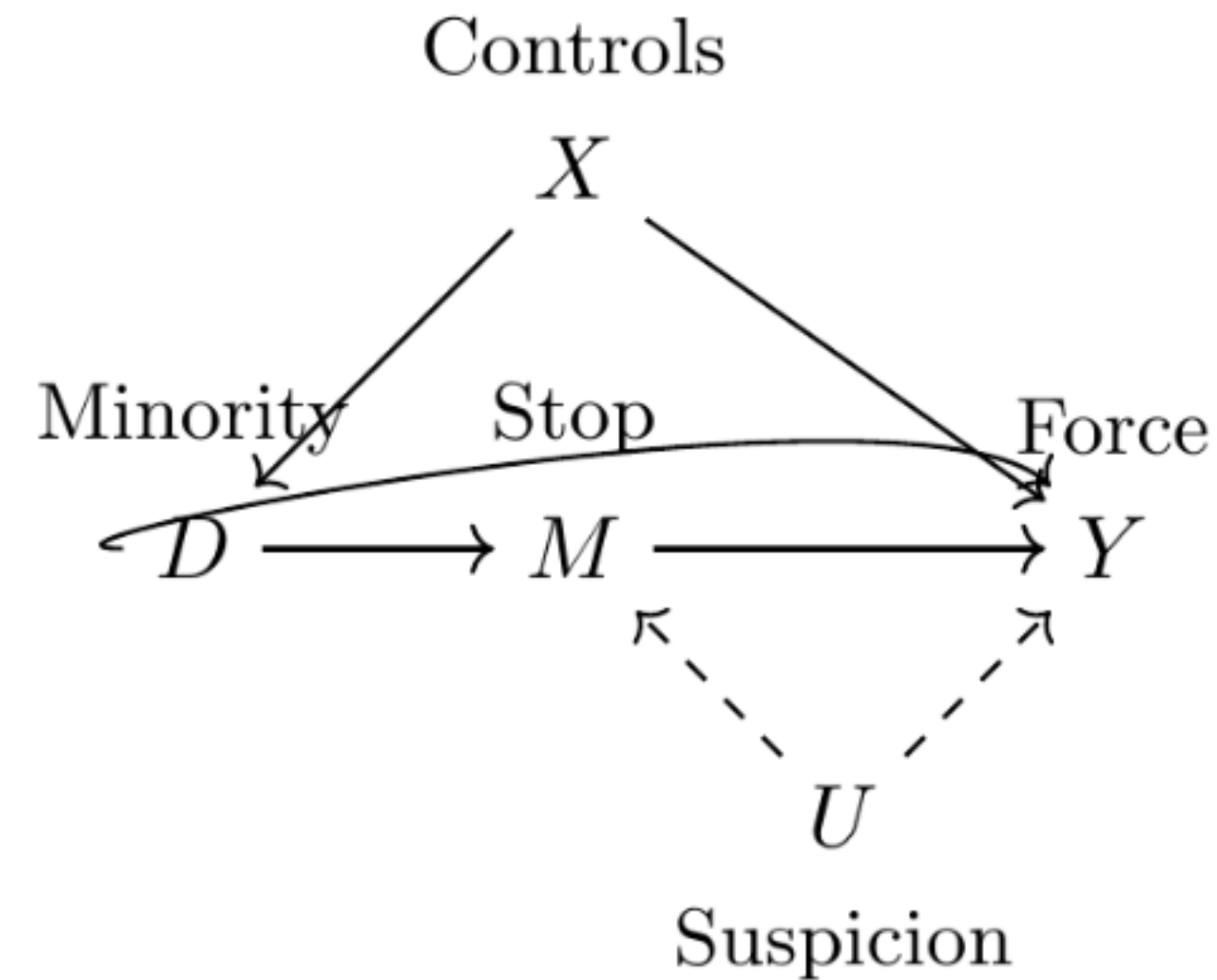
- Fryer (2019) finds in the NYC data
 - Blacks and Hispanic/Latinos were 50 percent more likely to have a interaction with police in the data
 - Blacks were 21 percent more likely than Whites to have an interaction with police in which a weapon was drawn
- Fryer (2019) finds a similar result in the Houston data
 - The racial differences are larger in the Houston data

Collider Bias Example 4: Policing and Admin Data

- Fryer (2019), however, concludes that there is no racial difference in officer-involved shootings
 - He controls for suspect demographics, officer demographics, encounter characteristics, suspect weapon, and year fixed effects
- With his model, Fryer (2019) finds
 - Blacks were 27 percent less likely to be shot by police than non-Black Non-Hispanics
- Strength of the study
 - Gathering the labor-intensive process of compiling the admin data
 - He is able to gather observed confounders that would have been unobserved without compiling the admin data

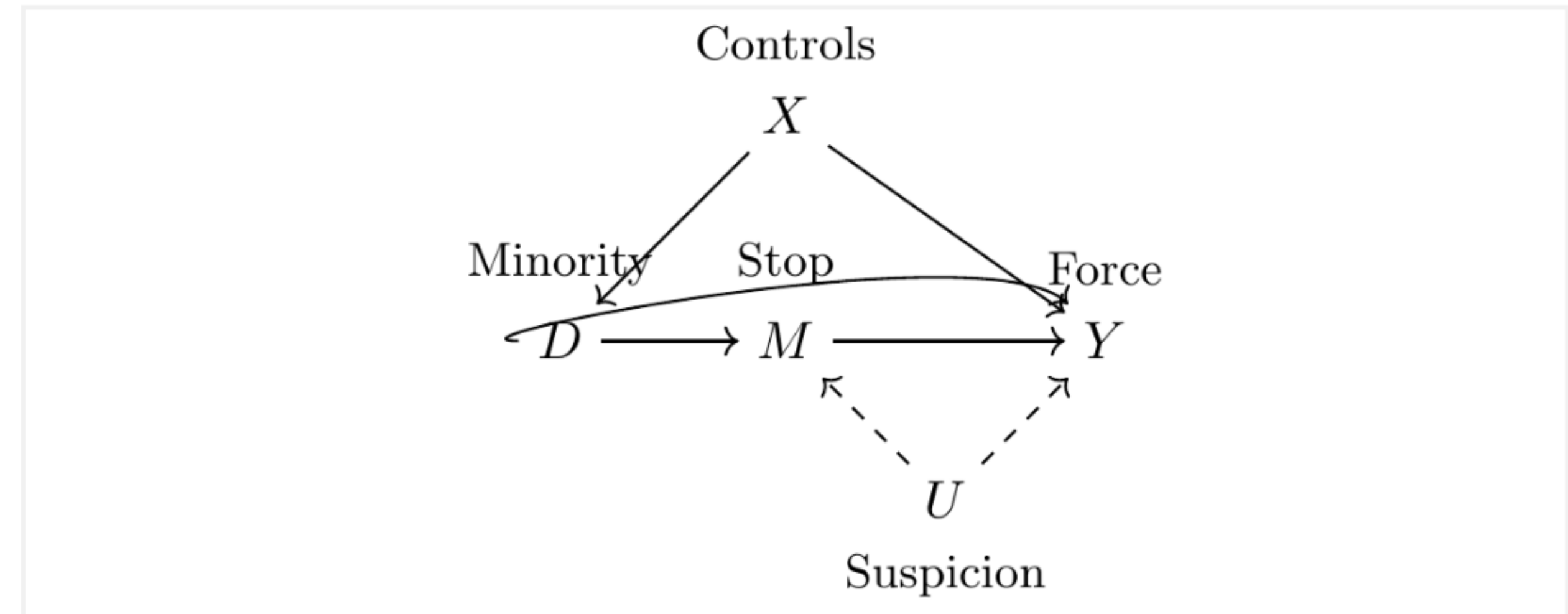
Collider Bias Example 4: Policing and Admin Data

- A critic of the study
 - Administrative data were endogenous since they were conditioned on a collider
- DAG
 - Y is forced used
 - D is Black or Hispanic/Latino
 - M is police stop
 - X are other controls
 - U is unobserved suspicion



Collider Bias Example 4: Policing and Admin Data

- Pathways
- Direct Pathway: $D \rightarrow M$
- Mediated Pathway: $D \rightarrow M \rightarrow Y$
- Backdoor Pathway 1: $D \leftarrow X \rightarrow Y$
- Backdoor Pathway 2: $D \rightarrow M \leftarrow U \rightarrow Y$



Collider Bias Example 4: Policing and Admin Data

- Fryer's (2019) data collection by compiling X controls closes that backdoor
- The direct pathway from Black/Hispanic onto Stops exists within the literature
 - M is a collider along the second backdoor pathway
- Fryer's (2019) results are conditional on police stops or interactions
 - Understanding potential selection into police data due to bias in who the police interacts with is a difficult endeavor
- Knox, Lowe, and Mummolo (2020) revisit Fryer's (2019) question and find that after applying bias correction
 - Lower bound estimates of police violence against civilians were 5 times higher than traditional approaches that ignores the sample selection problem

Concluding Thoughts on DAGs

- DAGs are a useful tool to clarify relationships among variables
 - Guides you to a credible identification strategy
- Atheoretical approaches to empiricism are subject to fail
 - Knowledge is essential for the establishing a credible identification strategy
- More data or “Big Data” do not solve the problem of potential outcomes
 - More data is an insufficient substitute for theory and literature
 - More data is an insufficient substitute for deep institutional knowledge
 - More data is an insufficient substitute for knowledge of data-generation process