

# Lecture 23: 22 April, 2026

Madhavan Mukund

<https://www.cmi.ac.in/~madhavan>

Data Mining and Machine Learning  
January–April 2026

# Correlation vs causality

- Apply flame to an old cloth

# Correlation vs causality

- Apply flame to an old cloth
- Use heavy gloves when we do this

# Correlation vs causality

- Apply flame to an old cloth
- Use heavy gloves when we do this
- Cloth burning is correlated to both applying a flame and using heavy gloves

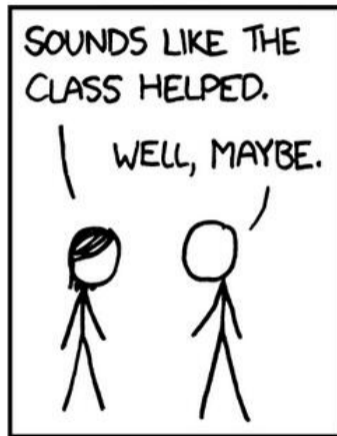
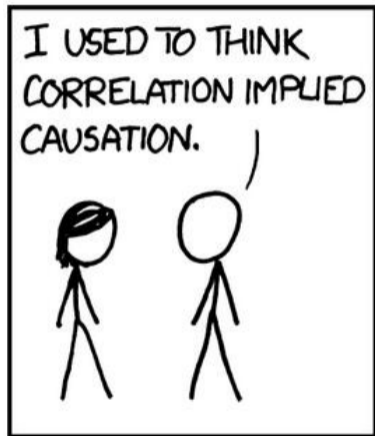
# Correlation vs causality

- Apply flame to an old cloth
- Use heavy gloves when we do this
- Cloth burning is correlated to both applying a flame and using heavy gloves
- Which of these correlations implies causation?

# Correlation vs causality

- Apply flame to an old cloth
- Use heavy gloves when we do this
- Cloth burning is correlated to both applying a flame and using heavy gloves
- Which of these correlations implies causation?
- Explain using **counterfactuals**

# Correlation vs causality



<https://xkcd.com/552>

# Simpson's paradox

Should we prefer Treatment A or Treatment B?

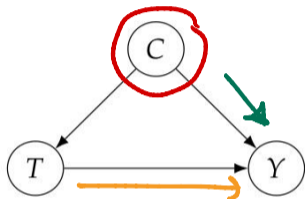
2 treatments for a disease  
Availability is 75%, 25%

Mortality rate

	Condition		
	Mild	Severe	Total
A	15% (210/1400)	30% (30/100)	16% (240/1500) ✓
B	10% ✓ (5/50)	20% ✓ (100/500)	19% (105/550) ✓

# Simpson's paradox

Condition determines the treatment



Mortality rate

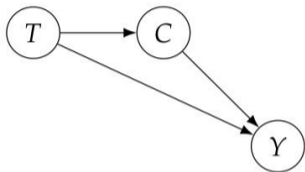
		Condition		
		Mild	Severe	Total
Treatment	A	15% (210/1400)	30% (30/100)	<b>16%</b> (240/1500)
	B	<b>10%</b> (5/50)	<b>20%</b> (100/500)	19% (105/550)

- B is preferred for severe cases
- Choose B

→ clearly "causal"

# Simpson's paradox

Treatment determines the condition



Mortality rate

		Condition		
		Mild	Severe	Total
Treatment	A	15% (210/1400)	30% (30/100)	<b>16%</b> (240/1500)
	B	<b>10%</b> (5/50)	<b>20%</b> (100/500)	19% (105/550)

- B is in short supply, delays increase severity
- Choose A

# Confounding

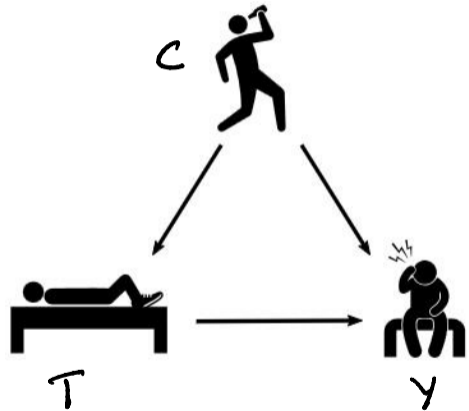
- Observed correlation:  
Go to bed with shoes on, wake up with a headache

# Confounding

- Observed correlation:  
Go to bed with shoes on, wake up with a headache
- Does wearing shoes in bed cause headaches?

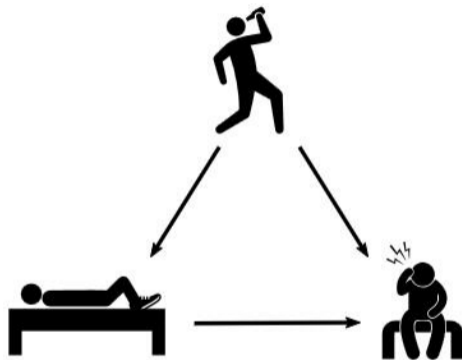
# Confounding

- Observed correlation:  
Go to bed with shoes on, wake up with a headache
- Does wearing shoes in bed cause headaches?
- Hidden cause:  
Both are a consequence of heavy partying the night before!



# Confounding

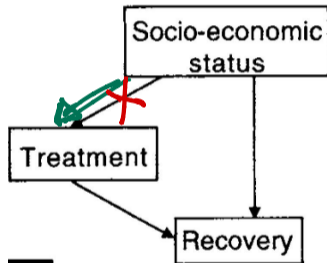
- Observed correlation:  
Go to bed with shoes on, wake up with a headache
- Does wearing shoes in bed cause headaches?
- Hidden cause:  
Both are a consequence of heavy partying the night before!
- Party **confounds** the association between sleeping with shoes and waking up with a headache



# Impact of treatment

## Confounding

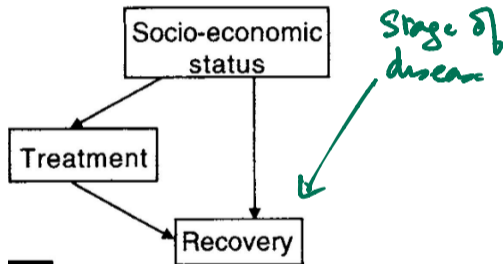
Uncontrolled conditions



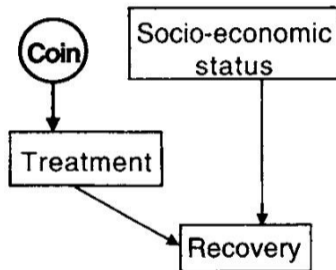
# Impact of treatment

## Confounding

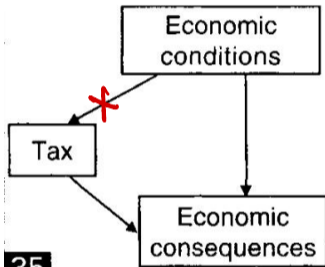
Uncontrolled conditions



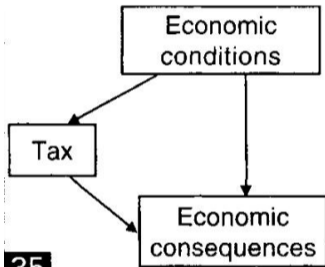
Experimental conditions



Model underlying data

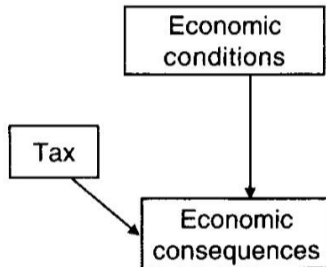


Model underlying data



35

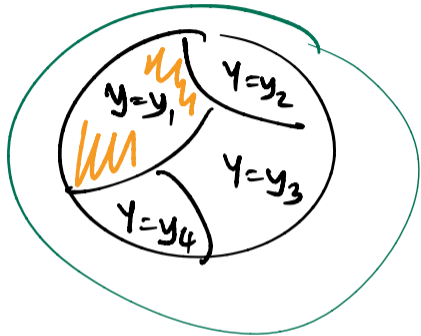
Model for policy evaluation



$$P(x \mid y = y_i)$$

Set  $y = y_1$

Intervention



# Conditioning vs intervention

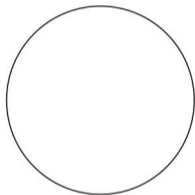
$P(Y | T = 0)$  vs  $P(Y | \text{do}(T = 0))$   
*~~~~~*

Judea Pearl

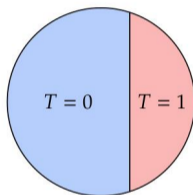
# Conditioning vs intervention

$P(Y | T = 0)$  vs  $P(Y | do(T = 0))$

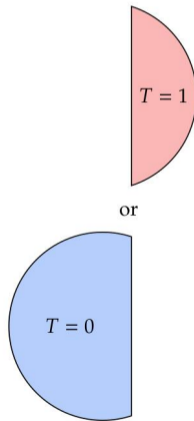
Population



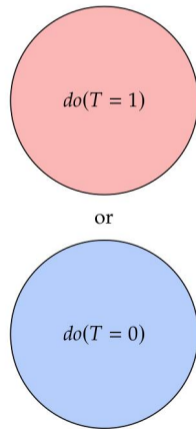
Subpopulations



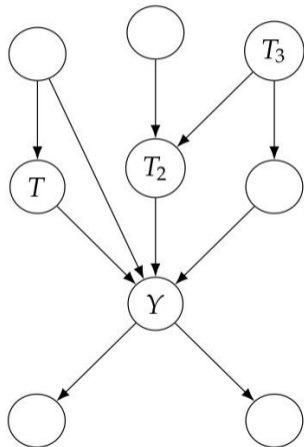
Conditioning



Intervening

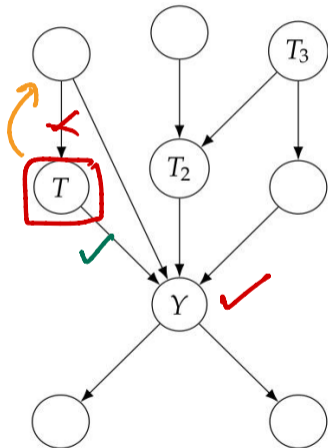


# Conditioning vs intervention



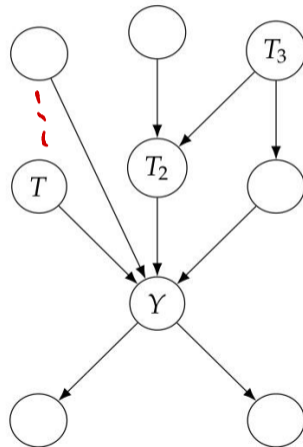
# Conditioning vs intervention

$P(Y \mid \text{do}(T = 0))$



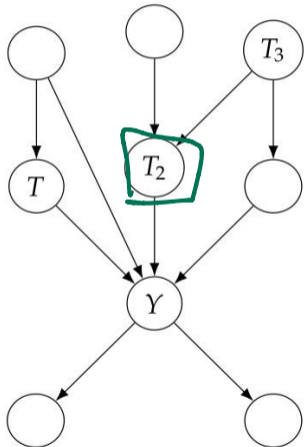
Intervene on  $T$

Delete incoming edges

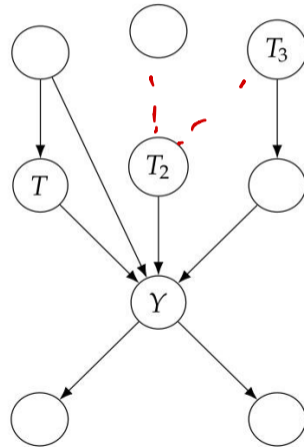


# Conditioning vs intervention

$P(Y \mid \text{do}(T_2 = 1))$

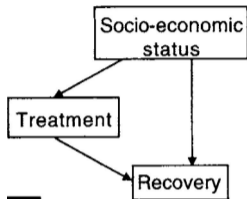


Intervene on  $T_2$

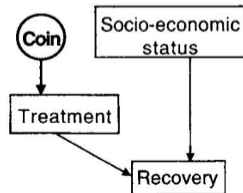


# Conditioning vs intervention

Uncontrolled conditions

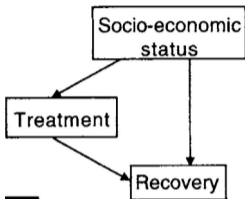


Experimental conditions

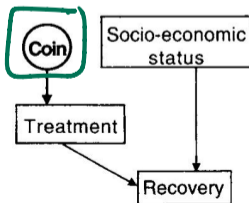


# Conditioning vs intervention

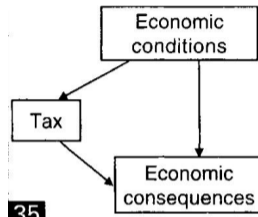
Uncontrolled conditions



Experimental conditions



Model underlying data



Model for policy evaluation

