

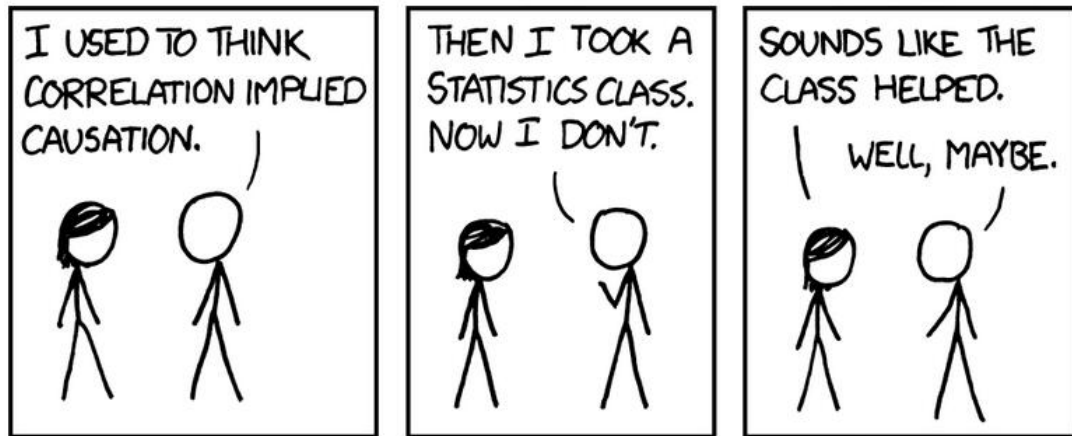
# Lecture 25: 24 April, 2025

Madhavan Mukund

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

Data Mining and Machine Learning  
January–April 2025

# Correlation vs causality



<https://xkcd.com/552>

# Simpson's paradox

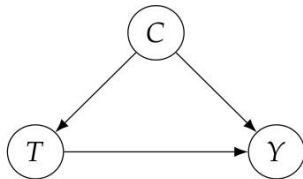
Should we prefer Treatment A or Treatment B?

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)

# Simpson's paradox

Condition determines the treatment



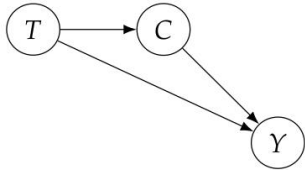
- B is preferred for severe cases
- Choose B

Mortality rate

		Condition		
Treatment		Mild	Severe	Total
	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)

# Simpson's paradox

Treatment determines the condition



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

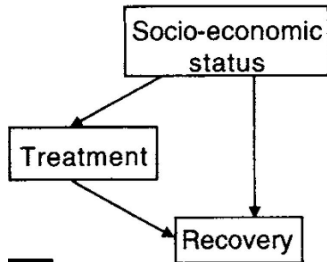
Mortality rate

		Condition		
Treatment		Mild	Severe	Total
	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)

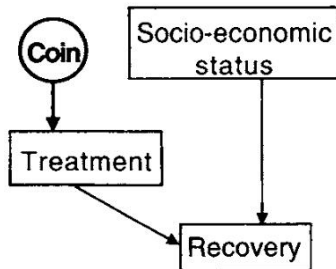
# Impact of treatment

## Confounding

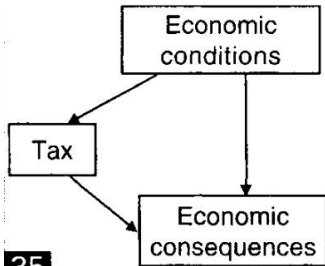
Uncontrolled conditions



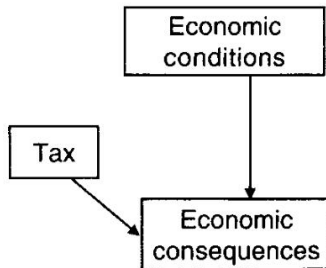
Experimental conditions



Model underlying data



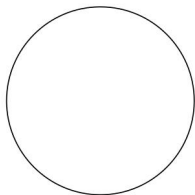
Model for policy evaluation



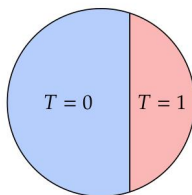
# Conditioning vs intervention

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

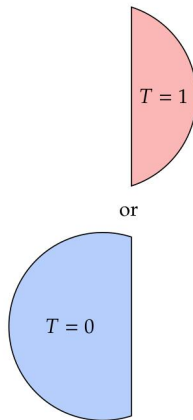
Population



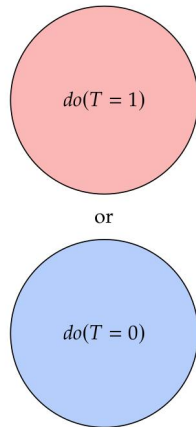
Subpopulations



Conditioning

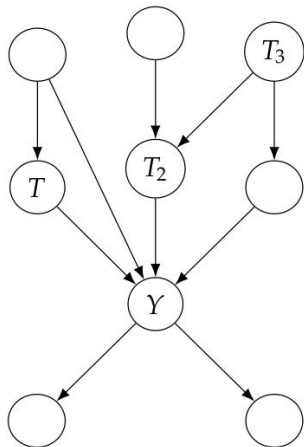


Intervening



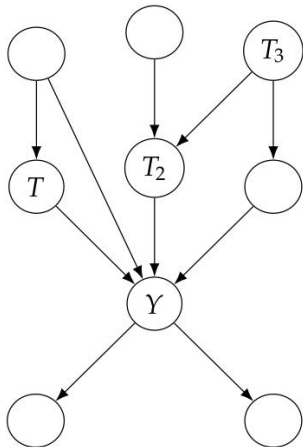


# Conditioning vs intervention

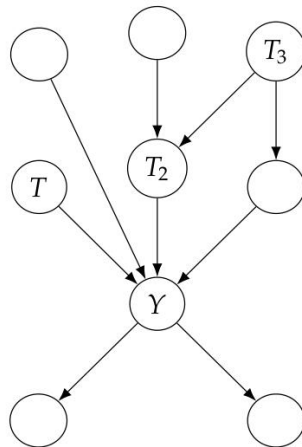


# Conditioning vs intervention

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

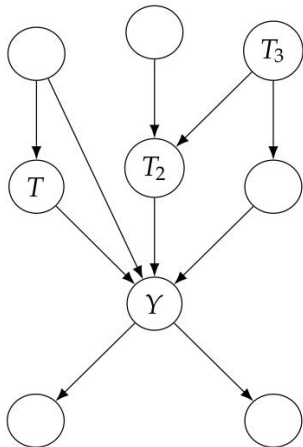


Intervene on  $T$

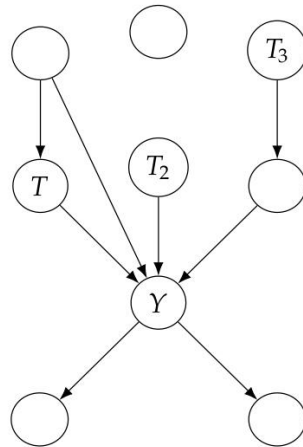


# Conditioning vs intervention

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

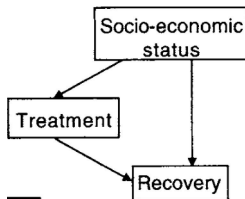


Intervene on  $T_2$

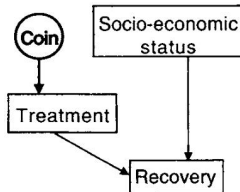


# Conditioning vs intervention

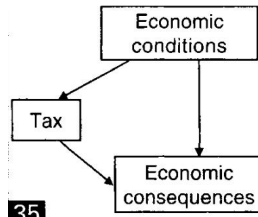
Uncontrolled conditions



Experimental conditions



Model underlying data



Model for policy evaluation

