

Counterfactuals in Artificial Intelligence

Sargur N. Srihari

University at Buffalo, The State University of New York
USA

Topics in Counterfactuals

- What are Counterfactuals?
- Three-layer Causal Hierarchy
- Counterfactual Explanations
- Mechanisms and Response Variables
- Twinned Networks

What are Counterfactuals?

- Contrary to facts
- Counterfactual thoughts in psychology
 - What if I hadn't been speeding?
 - What if I had not been wearing seat belts?
- It may seem pointless, but naturally arise in several settings

Natural Counterfactual Examples

- Legal liability case

- “Did the driver’s intoxicated state cause the accident?”
 - i.e., would accident have happened had driver not been drunk?

- Treatment and Diagnosis

- “We are faced with a car that does not start, but the lights work. Will replacing the battery make the car start?”
- “Patient did not take the prescribed treatment and did not get well and we wish to determine the probability that the patient would have gotten well had he complied with the treatment.”

Types of Queries

- Questions in English:

Who: who is this person? DL

What: what is this flower? DL

When: when did America become independent? DB

Where: where was the declaration of independence signed? DB

Why: why did Americans want independence? AI

How: how did America become independent? AI

- Causal Hierarchy:

Level (Symbol)	Typical Activity	Typical Questions	Examples
1. Association $P(y x)$	Seeing	What is? How would seeing X change my belief in Y?	What does a symptom tell me about a disease? What does a survey tell us about the election results?
2. Intervention $P(y do(x), z)$	Doing Intervening	What if? What if I do X?	What if I take aspirin, will my headache be cured? What if we ban cigarettes?
3. Counterfactuals $P(y_x x', y')$	Imagining, Retrospection	Why? Was it X that caused Y? What if I had acted differently?	Was it the aspirin that stopped my headache? Would Kennedy be alive had Oswald not shot him? What if I had not been smoking the past 2 years?

Counterfactual Queries

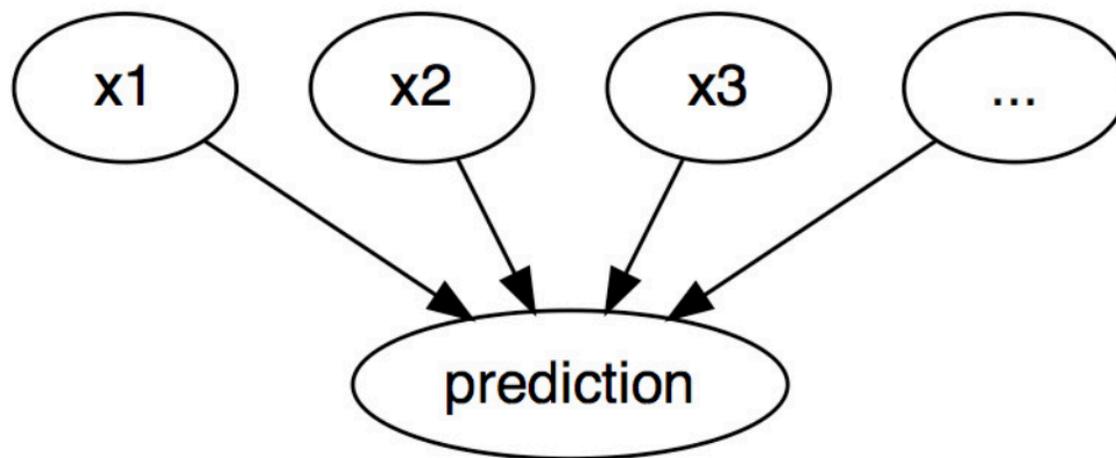
- Causal Query
 - Causal queries discover causality by intervention
 - Thus ruling out mere correlation
- Counterfactual Query
 - A different type of causal query
 - We already have some information about the true state of the world and want to inquire about the state of the world had we been able to intervene
 - And set the values of certain variables
 - E.g., “Would the US have joined World War II had it not been for the attack on Pearl Harbor?”

Counterfactual Explanation in AI

- Explanations are useful for various reasons
 - If life is impacted by a decision
 - Challenge a decision
 - Feedback to fix input
 - Fix bugs
 - Comply with laws
- Counterfactual explanations
 - Users prefer counterfactual explanations to case-based reasoning (which is example based)
 - “Credit score was too low” is less helpful than “Increase income by \$10K”

Counterfactual Explanation in ML

- A counterfactual explanation of a prediction describes the smallest change to the feature values that changes the prediction to a predefined output



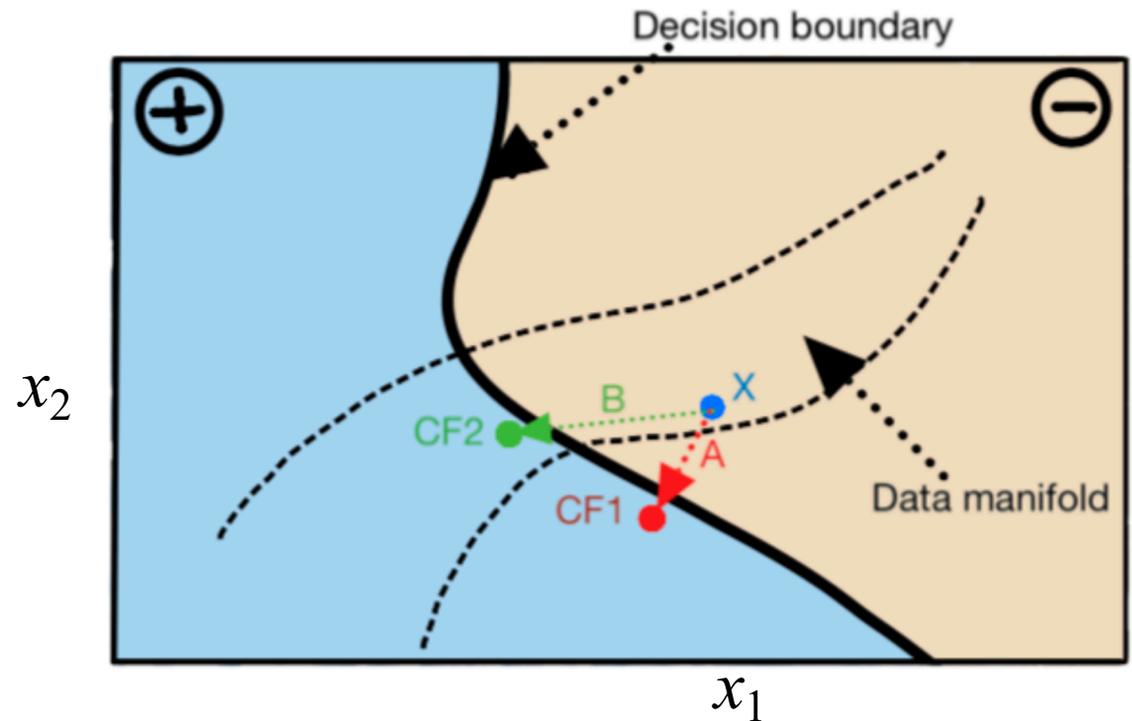
Example of Counterfactual Explanation

- Classifier rejects a loan application based on feature vector
{Income, CreditScore, Education, Age}
- Explanations
 - (1) why was loan denied?
 - (2) what can be different so loan will be approved?
- Latter is basis of a counterfactual explanation
 - What changes could be made to feature vector to end up on other side of classifier decision boundary
 - Increase Income by \$10K or get new MS degree or both

Counterfactuals in Feature space

Two possible paths for a datapoint (shown in blue) originally classified by $f(\mathbf{x})$ to be in the negative class, to cross decision boundary

End points of both paths (shown in red and green) are valid counterfactuals for the original point.



Red path is the shortest, whereas the green path adheres closely to the manifold of the training data, but is longer

Counterfactual explanation as Optimization

- Minimize distance between counterfactual datapoint x' and original datapoint x subject to constraint that classifier output on the counterfactual is the desired label ($y' \in \mathcal{Y}$)

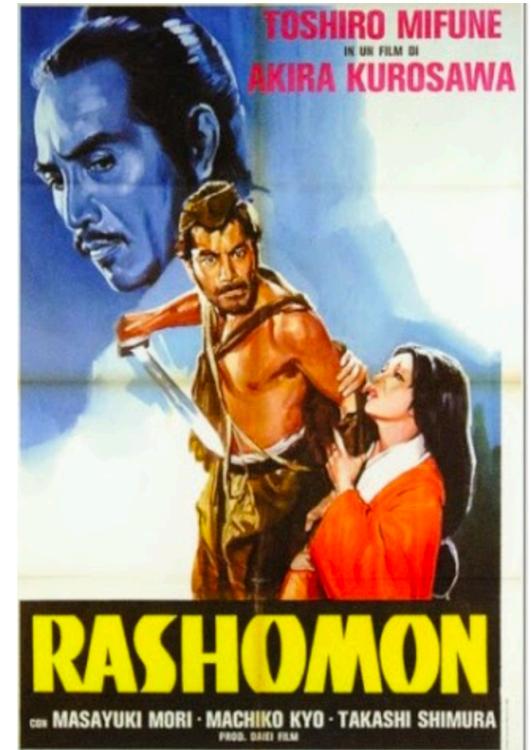
$$\arg \min_{x'} d(x, x') \text{ subject to } f(x') = y' \quad (1)$$

$$\arg \min_{x'} \max_{\lambda} \lambda (f(x') - y')^2 + d(x, x') \quad (2)$$

- First term encourages classifier output on the counterfactual to be close to the desired class
- Second term forces counterfactual to be close to the original datapoint

Rashomon Effect

- Counterfactuals are human friendly explanations
 - Because they contrast instance to focus on few feature changes
- But suffer from 'Rashomon effect'
 - Murder of Samurai told by different people
 - Each story explains outcome equally well, but stories contradict each other
 - Same can happen with counterfactuals, due to multiple counterfactual explanations



Causal Queries

Correlation: $P(y | x)$

Probability of event $Y=y$ given that we observed event $X=x$

Causation: $P(y | do(x), z)$

Probability of event y given that we intervene and set the value of X to x and subsequently observe event $Z=z$

Counterfactuals: $P(y'_{x'} | x, y)$

Probability that Joe's salary would be y' had he finished college given that his actual salary is y and he had only two years of college

Causal Mechanisms

- Causal model is a graph of causal progression
 - The value of each variable is selected, in order, based on parent values, via a causal mechanism
 - Details of this causal mechanism are implicit
 - Assumed that variable X , induces a distribution $P(X | \text{pa}X)$
- Causal query is answered as a probabilistic one
 - This may not always be possible and we may have to reason explicitly about governing mechanism
- Second reason for understanding mechanism:
 - To answer additional types of queries:
 - Counterfactual queries, and certain diagnostic queries.

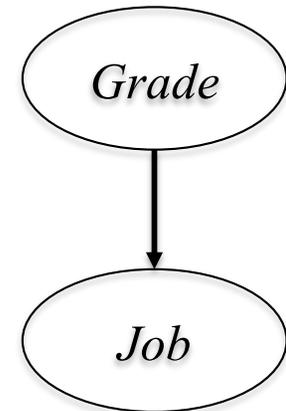
Intuition for a Counterfactual query

- We want to preserve as much as possible between the real and counterfactual worlds
- Mechanism governing each variable, provides reliable inferences about the true world
 - So as to preserve as much as possible when reasoning about the counterfactual events

Example of need for mechanism

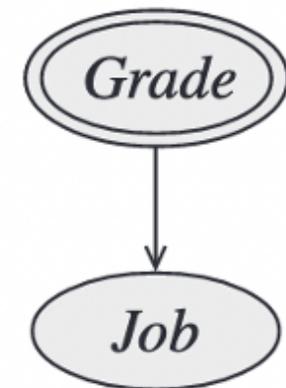
- *Grades* are binary-valued: *High, Low*

$J \setminus G$	g^0	g^1
j^0	0.5	0.5
j^1	0.5	0.5



- Counterfactual Query:

- *Gump* received a Low grade
- He applied to *Acme Consulting*, and he was not hired
- Would *Acme* have hired him had he managed to hack into the registrar's computer system and change his *Grade* to *High*?



Two models, different answers

- Gump applied to *Acme*, and he was not hired
 - Two models
 - Both with two equally likely populations of companies $P(J|G)=(0.5,0.5)$
- 1. First Model
 1. *Not hiring mode*: always reject Gump
 2. *Desperate to hire*: hire anyone regardless of grade
 - Gump not hired, so *Acme* is not hiring. Changing *Grade* is useless
- 2. Second Model
 1. *Highly selective*: hire only if grade is *High*
 2. *Hire underdogs*: hire only if grade is Low (& pay Low salary)
 - Gump not hired, so changing *Grade* makes *Acme* hire *Gump*
- Different answers to counterfactual queries depending on model of the world

Troubleshooting a Broken device

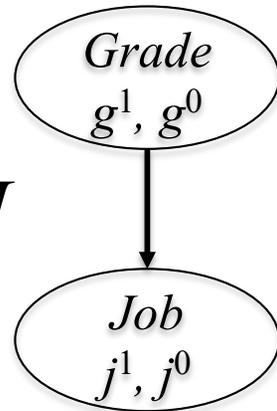
- Observed failure is due to fault in a component (among many)
- Assume we replace one component
 - If not the cause, then replacement has no effect:
 - Probability that device remains broken should be 1
- Mechanism of combinations of faults is complex
 - Uncertain about which one is currently the case
 - However, we wish to have this mechanism persist between the pre-repair and post-repair state
- Examples illustrate need for detailed specification of mechanism of variables

Exogenous vs Endogenous Variables

- Exogenous variables (are external)
 - Unmodeled latent variables
 - Coin Toss depends on *exogenous* variables
 - Coin rotation, forces applied, air movement, etc
 - Given all these variables, coin toss is deterministic
 - Gump's hiring depends on
 - Company goals, funding level, no. applicants, Gump's tie
- Endogenous variables
 - What we choose to include in the model
- We abstract away from exogenous variables and focus on their effect on endogenous variables

Response Variable Example

- U : entire set of exogenous variables
 - affecting outcome of J in causal model $G \rightarrow J$
- J is a deterministic function of G, U
 - for each $\mathbf{u} \in U$ and $g \in G$, J takes the value $f_J(g, \mathbf{u})$
- Mappings between G and J
 - \mathbf{u} induces four possible mappings enumerated as:
 - $\mu_{g^1 \rightarrow j^1, g^0 \rightarrow j^1}^J$: $f_J(\mathbf{u}, g^1) = j^1, f_J(\mathbf{u}, g^0) = j^1$ always hire
 - $\mu_{g^1 \rightarrow j^1, g^0 \rightarrow j^0}^J$: $f_J(\mathbf{u}, g^1) = j^1, f_J(\mathbf{u}, g^0) = j^0$ rational recruiting
 - $\mu_{g^1 \rightarrow j^0, g^0 \rightarrow j^1}^J$: $f_J(\mathbf{u}, g^1) = j^0, f_J(\mathbf{u}, g^0) = j^1$ underdog cases
 - $\mu_{g^1 \rightarrow j^0, g^0 \rightarrow j^0}^J$: $f_J(\mathbf{u}, g^1) = j^0, f_J(\mathbf{u}, g^0) = j^0$ never hire
 - Response variable U^J takes one of the four values



Response Variable Definition

- Let variable X have set of parents Y

$$k=|Val(X)|, m=|Val(Y)|$$

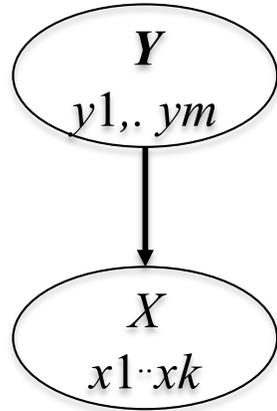
- A response variable for X given Y is a variable U^X

– whose domain is the set of all possible functions $\mu(Y)$ from $Val(Y)$ to $Val(X)$

- i.e., let y_1, \dots, y_m be an enumeration of $Val(Y)$

– $\mu_{y_1, \dots, y_m \rightarrow x_1, \dots, x_k}^J$ denote function that assigns y_i to x_j

– The domain of U^X contains one such function for each such tuple, giving a total of k^m functions of this form



A functional causal model

- A functional causal model C over a set of endogenous variables χ is a causal model
 - Defined over two sets of variables
 - Variables χ , and a set of response variables $\mathcal{U} = \{U^X: X \in \chi\}$
 - Each variable $X \in \chi$ has a set of parents $\text{Pa}_X \subset \chi$, and a response variable parent U^X for X given Pa_X .
 - The model for each variable $X \in \chi$ is deterministic: When $U^X = \mu$ and $\text{Pa}_X = \mathbf{y}$, then $X = \mu(\mathbf{y})$ with probability 1.
- Model C also defines a joint probability distribution over response variables, defined as a BN over \mathcal{U}
 - Thus, each response variable U has a set of parents $\text{Pa}_U \subset \mathcal{U}$, and a CPD $P(U|\text{Pa}_U)$.

Functional causal model for clinical trial

Patients randomly assigned (A): medication and placebo

Assigned Treatment A does not influence the response O directly but only through treatment T

Sample CPDs

	always taker	complier	defier	never taker
	0.05	0.7	0.05	0.2
	$P(U^T)$			
U^T	always well	helped	hurt	never well
always taker	0.2	0.5	0.1	0.2
complier	0.5	0.4	0.05	0.05
defier	0.3	0.4	0.1	0.2
never taker	0.6	0.3	0.05	0.05
	$P(U^O U^T)$			

Resulting joint distribution for response variables U^T and U^O

U^T	U^O	probability
always taker	always well	0.01
always taker	helped	0.025
always taker	hurt	0.005
always taker	never well	0.01
complier	always well	0.35
complier	helped	0.28
complier	hurt	0.035
complier	never well	0.035
defier	always well	0.015
defier	helped	0.02
defier	hurt	0.005
defier	never well	0.02
never taker	always well	0.12
never taker	helped	0.06
never taker	hurt	0.01
never taker	never well	0.01
	$P(U^T, U^O)$	

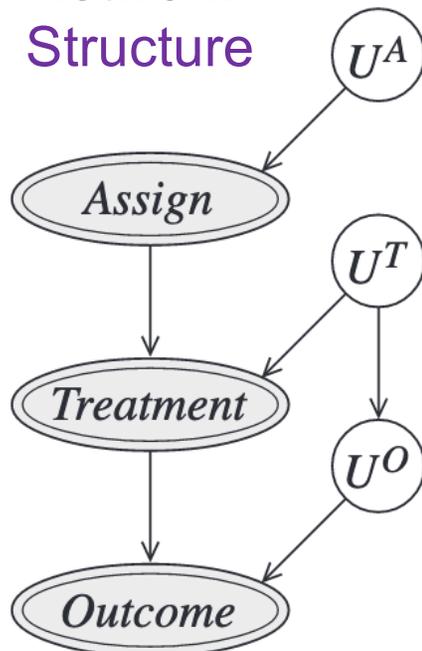
U^T has four functions:

- $\mu_{1 \rightarrow 1, 0 \rightarrow 1}^T$ — *always taker*;
- $\mu_{1 \rightarrow 1, 0 \rightarrow 0}^T$ — *complier (takes medicine if and only if prescribed)*;
- $\mu_{1 \rightarrow 0, 0 \rightarrow 1}^T$ — *defier (takes medicine if and only if not prescribed)*;
- $\mu_{1 \rightarrow 0, 0 \rightarrow 0}^T$ — *never taker*.

U^O has four functions:

- $\mu_{1 \rightarrow 1, 0 \rightarrow 1}^O$ — *always well*;
- $\mu_{1 \rightarrow 1, 0 \rightarrow 0}^O$ — *helped (recovers if and only if takes medicine)*;
- $\mu_{1 \rightarrow 0, 0 \rightarrow 1}^O$ — *hurt (recovers if and only if does not take medicine)*;
- $\mu_{1 \rightarrow 0, 0 \rightarrow 0}^O$ — *never well*.

Network Structure

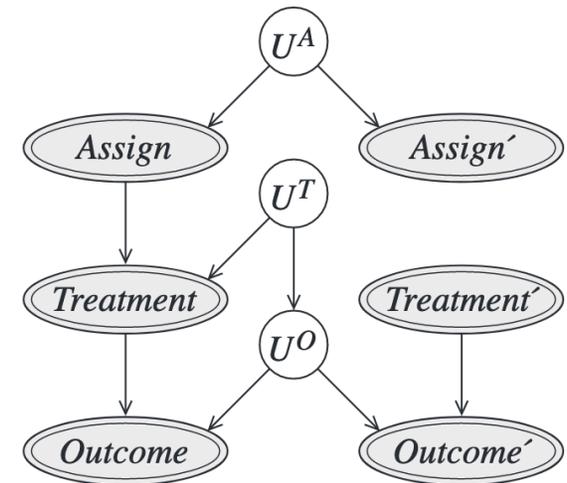


Counterfactual Queries

- Functional causal models can be used to address counterfactual queries
- Arises in a diagnostic setting where we wish to reason about the probabilities of events after taking a repair action

Twinned Networks

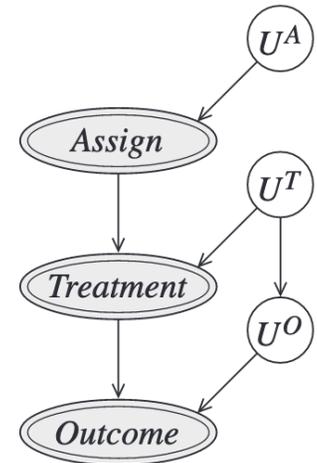
- We are reasoning in parallel about two different worlds: the real one, and the counterfactual one
- A variable may take on value X in the true world, and X' in the counterfactual world
- Both are governed by the same causal model except one involves an intervention



Ex: Counterfactual Query

- Critical property of the desired output for a counterfactual query is that it involves minimal modification to the true world
 - A patient in study took prescribed treatment and did not get well
 - Would he have gotten well had we forced him to comply with the treatment?
 - Formally, this query is

$$P(O '=o^1 \mid T = t^0, O = o^0, do(T ':= t^1))$$



Need for Response Variable

- Formally, this query is $P(O' = o^1 \mid T = t^1, O = o^0, do(T' := t^1))$
 - Intuition says: if we change nothing in the counterfactual world, outcome should be same
 - But even this obvious assertion has significant consequences
 - Clearly, we cannot simply generate a new random assignment to the variables in counterfactual world
 - E.g., we want patient's prescribed course of treatment to be same in both worlds
 - We also want his response to the treatment to be same
- Response variable makes this intuition very precise