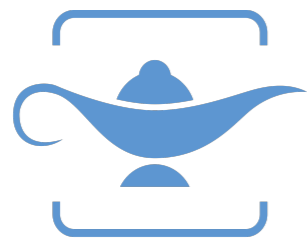


Counterfactual Data Generation using VAEs

Ayan Majumdar, Preethi Lahoti, Junaid Ali, Till Speicher,
Isabel Valera, Krishna Gummadi



MAX PLANCK INSTITUTE
FOR SOFTWARE SYSTEMS



max planck institut
informatik



UNIVERSITÄT
DES
SAARLANDES

Fairness and causation

Is the law school admission process fair?

Jacob is a black male law school applicant. He scored 55 in LSAT and had UGPA 3.3. He was **rejected**.

- Had Jacob been white instead, would he had been **accepted**?
 - *counterfactual*
- Did Jacob's race **cause** him to get negative outcome?
 - *counterfactual fairness (Kusner et al. 2017)*¹

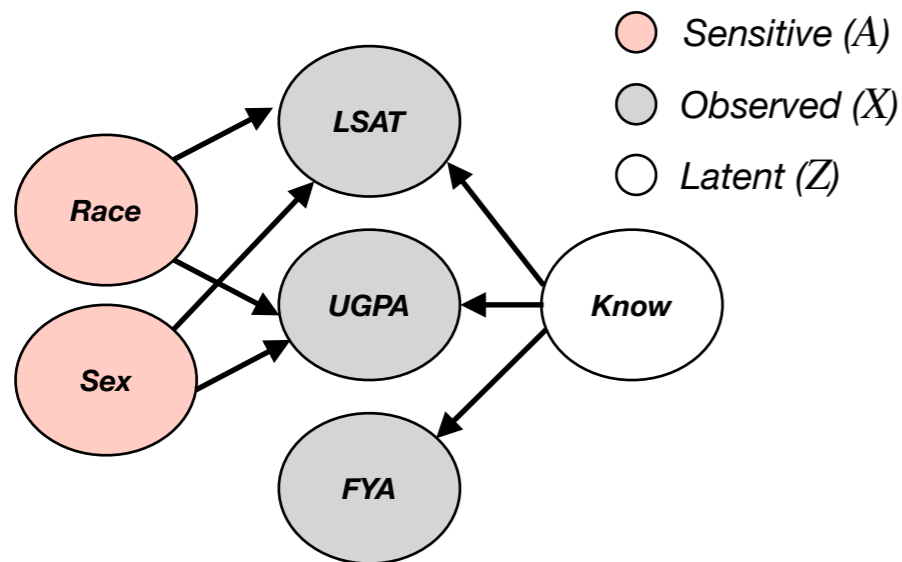
Such questions of fairness need **counterfactual data**!

How do we generate them?



¹Matt J Kusner et al. "Counterfactual Fairness". In *Advances in Neural Information Processing Systems* 30.

Counterfactuals



Causal graph

+

$$\text{LSAT} \sim \mathcal{N}(\exp(b_L + w_L^K K + w_L^R R + w_L^S S), \sigma_L)$$

$$\text{UGPA} \sim \mathcal{N}(b_G + w_G^K K + w_G^R R + w_G^S S, \sigma_G)$$

$$\text{FYA} \sim \mathcal{N}(w_F^K K, 1)$$

$$\text{Know} \sim \mathcal{N}(0, 1)$$

Structural equations

Generating counterfactual (Pearl et al. 2009)² —

1. Abduction: Given observed X , $A = a$ **estimate** Z
2. Action: **Intervene** on A by setting it to a'
3. Prediction: **Re-compute** X using Z under intervention $do(A = a')$

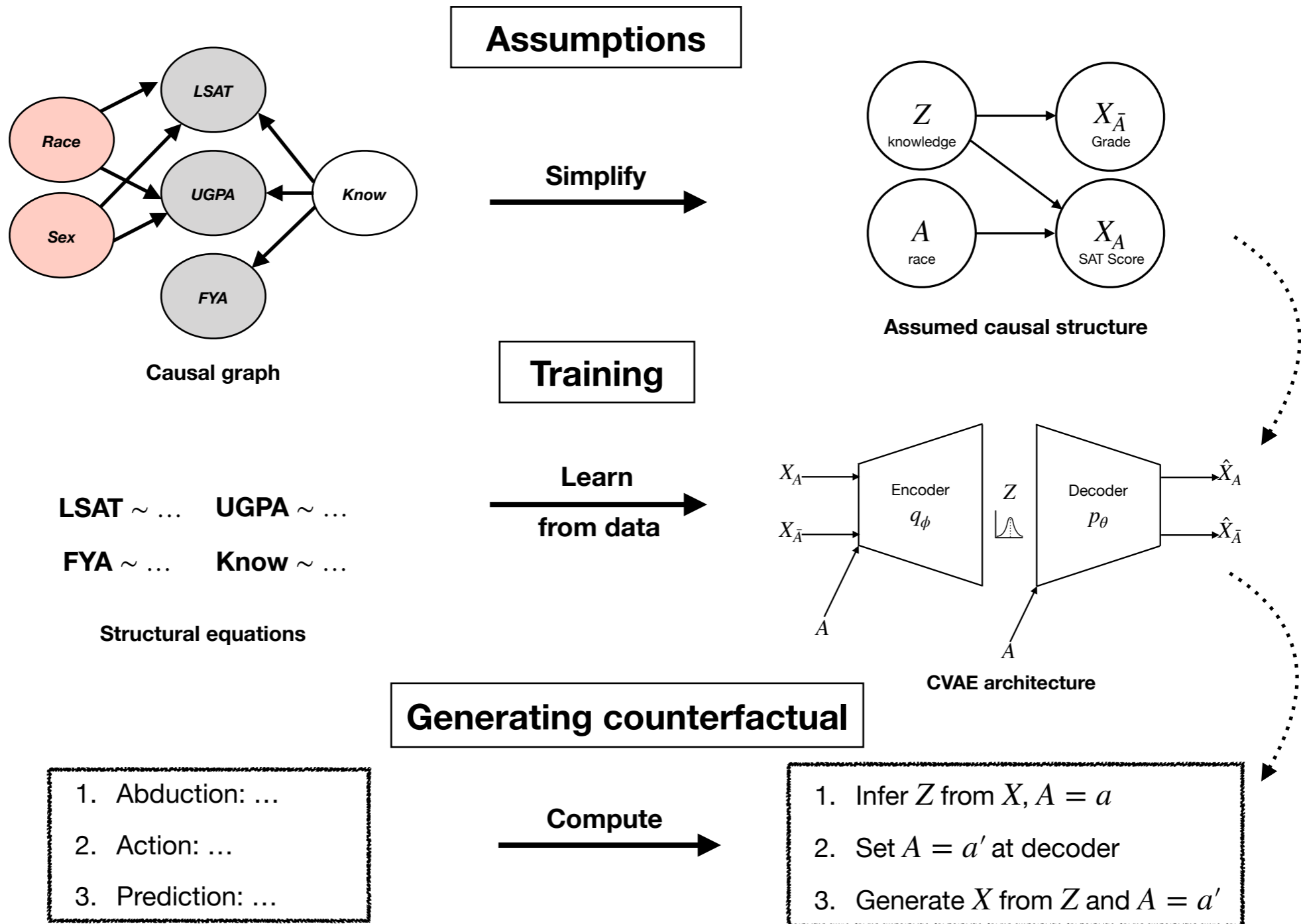
Need complete access to causal model!

Infeasible in real settings.

*Can we generate counterfactual data for fairness
in the absence of the whole causal model?*

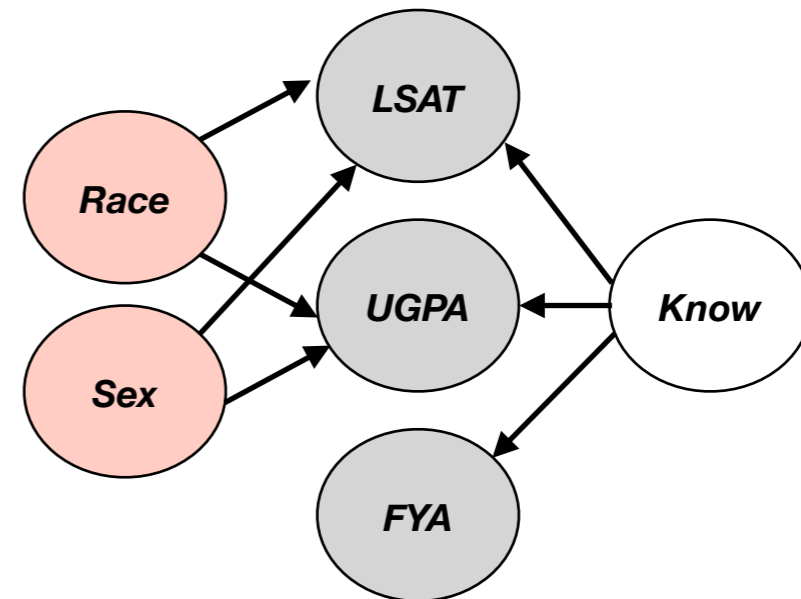
²Pearl, J. (2009). *Causality: Models, reasoning, and inference*, (2nd ed.). New York: Cambridge University Press.

Approach



Can CVAE generate counterfactuals?

- **Train** CVAE on synthetic **generated data**.
- Condition on A (**race, sex**).
- **Metric:** Mean absolute error (MAE), cosine similarity b/w CVAE & causal counterfactuals.



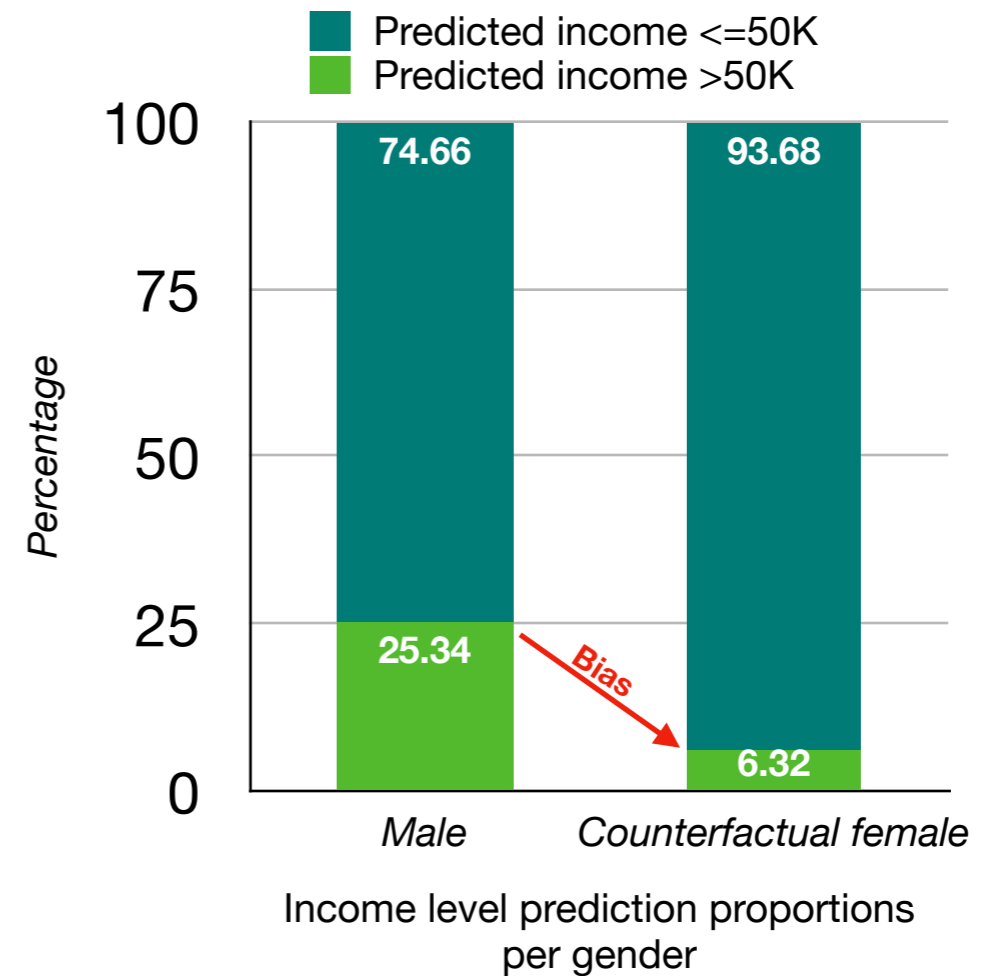
LSAT MAE	UGPA MAE	FYA MAE	Cosine Sim.
0.013	0.028	0.005	0.9997

Counterfactual generation quality (Race: White to Black).

CVAE can generate faithful counterfactuals!

Auditing counterfactual fairness

- Dataset: UCI Adult income
- **Trained** classification model
 - Predict income level ($\leq 50K$; $> 50K$)
 - Sensitive feature (Gender: Male-Female)
- Audit counterfactual fairness:
 - *Male* individual was predicted to have high income.
 - If individual was *female* instead, would the **prediction change**?



High-income males → **low-income** counterfactual females!

Model biased negatively towards females!