

# Fairness in Machine Learning models using Causality



## Supervision:

Patrick Forré,  
Christos Louizos,  
Tamas Erkelens,  
Barteld Braaksma,

University of Amsterdam  
University of Amsterdam  
Municipality of Amsterdam  
Statistics Netherlands

# Today



Problem description



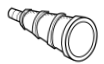
Why Causality?



FairTrade method



Risk profiles in unlawful social welfare



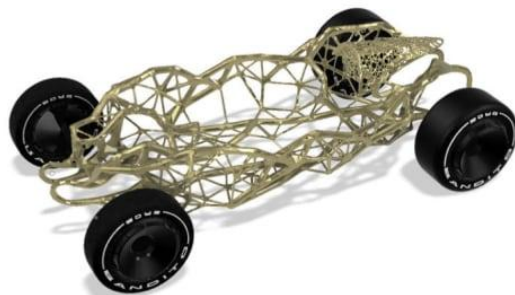
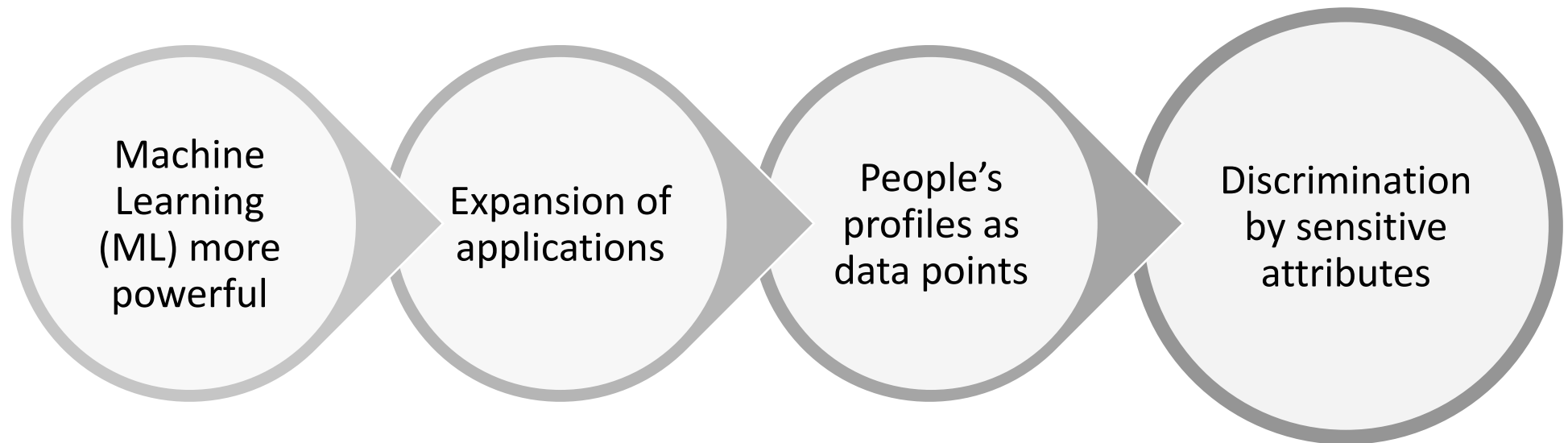
Limitations and Future work



Conclusion



# Problem description





# Problem description

Sensitive attributes: Personal attributes deemed **unfair** to use for prediction models

Discrimination by **sensitive attributes**

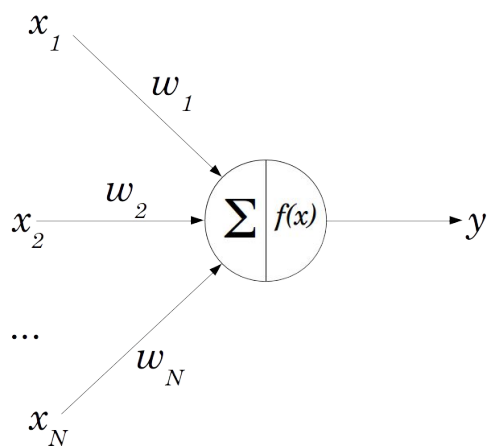
What is fairness?



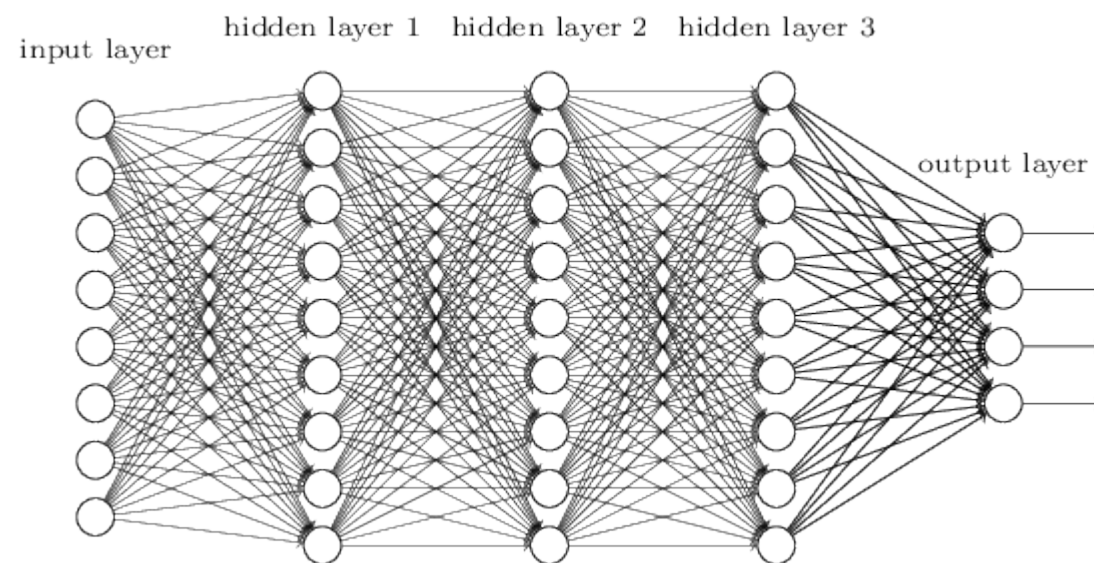


# Problem description

No longer possible to use 'classical' correction due to increased model complexity



Regression model



Neural Network





# Why Causality?

---



# Experiment 1

---

## Fairness by Unawareness experiment

(leaving sensitive variable out of the model)

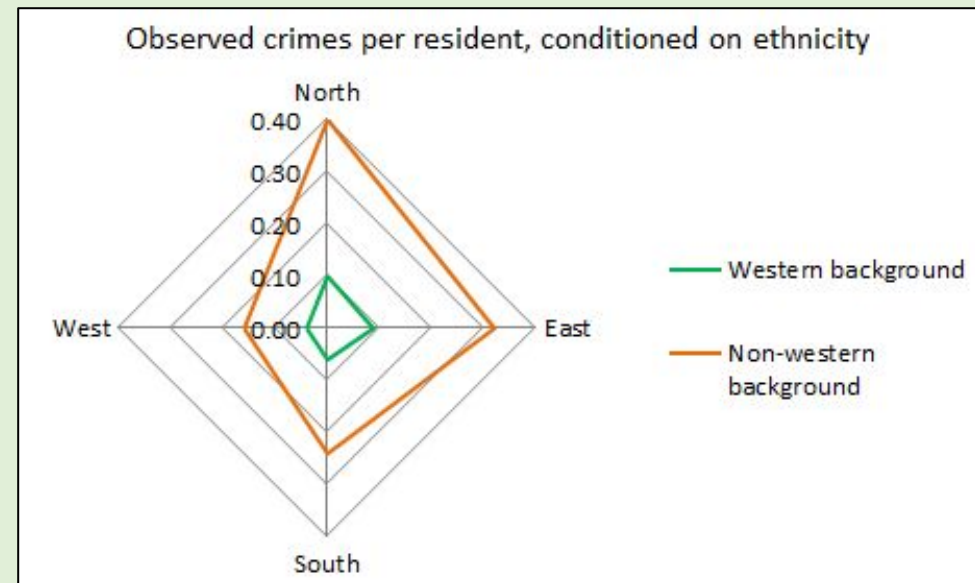
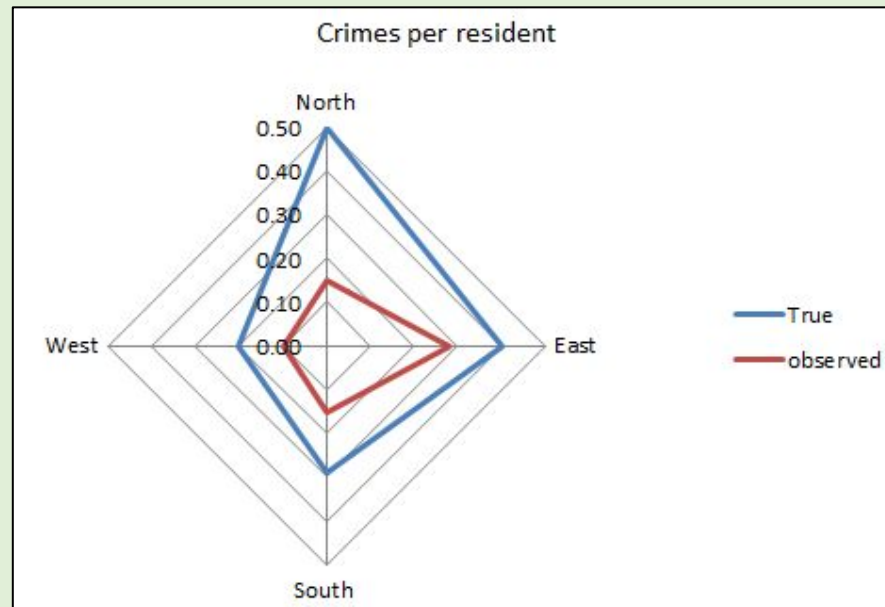


## Simulation experiment:

- Police wants to decide between checking neighbourhood North, East, South or West.
- Data is obtained with a **bias**, people with a non-western background have higher probability of being caught after a crime. Therefore the police wants to **exclude ethnicity** from the model.
- The most crimes occur in **North**.
- Most people with non-western background live in **East**.

# Experiment 1

## Fairness by Unawareness experiment (leaving sensitive variable out of the model)



Simpsons Paradox

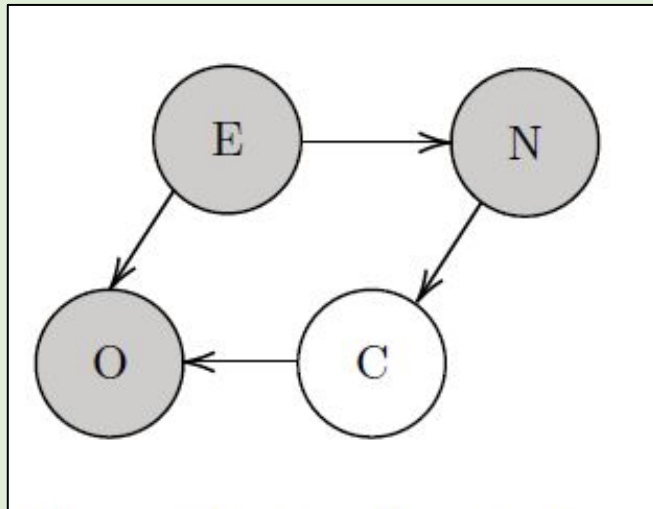


# Experiment 1

---

## Fairness by Unawareness experiment

(leaving sensitive variable out of the model)



Variable	Meaning	Values
E	Ethnicity	Western, Non-western
N	Neighbourhood	North, East, South, West
C	Crime committed	Yes, No
O	Observed Crime	Yes, No



# Why Causality?

---

Solve fundamental problems of observational based fairness metrics by:

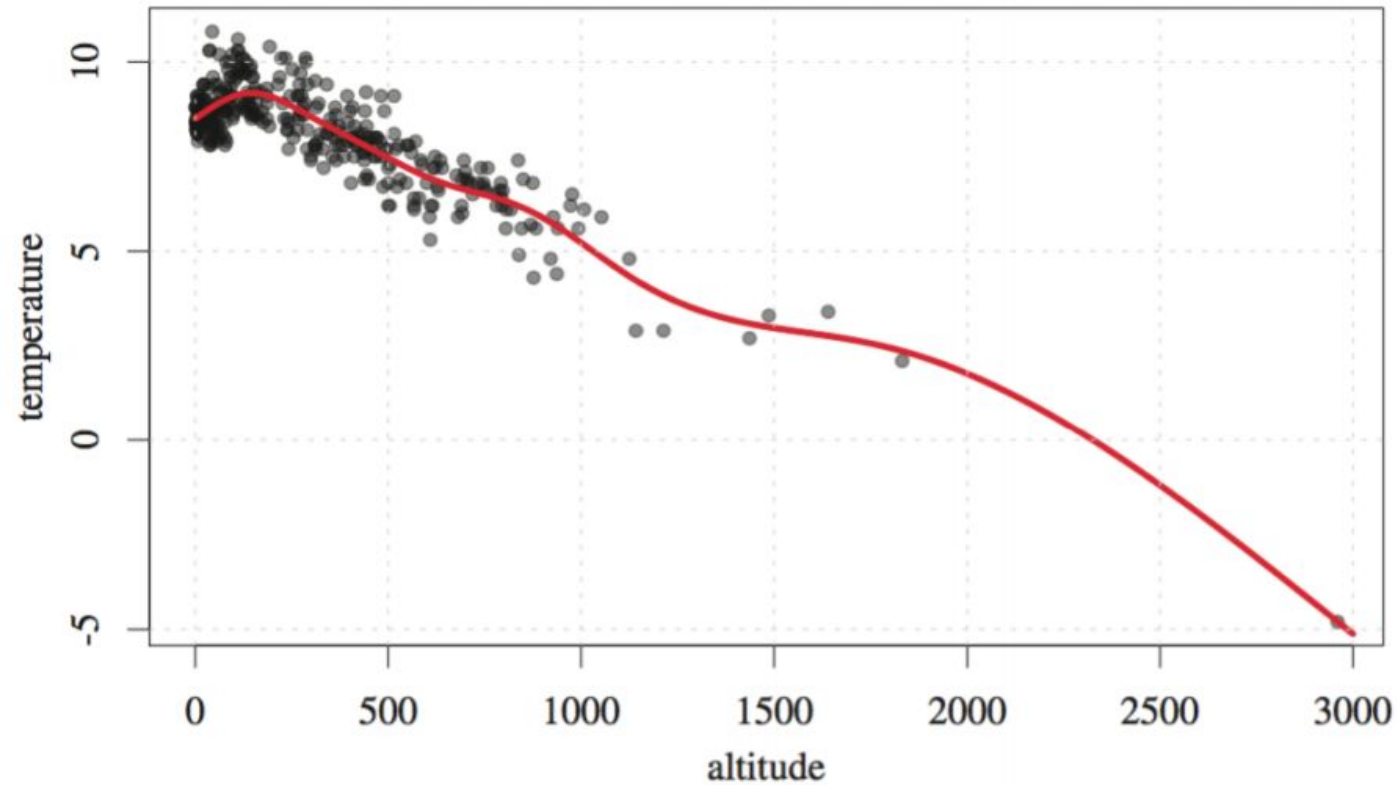
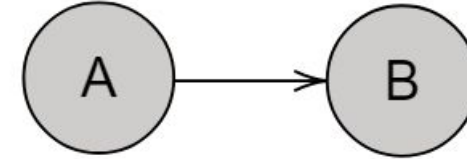
- > Understanding data
- > Interpretation of model effects
- > Control of explicit fairness demands





# Why Causality?

Causality theory considers *causal effects*:





# Why Causality?

Causal solution to *what is fair?*:

Intuition: Intervening on the sensitive attribute should not influence the outcomes

For *Counterfactual Fair* models holds:

$$P(\hat{Y}(a, U) = y | X = x, A = a) = P(\hat{Y}(a', U) = y | X = x, A = a)$$



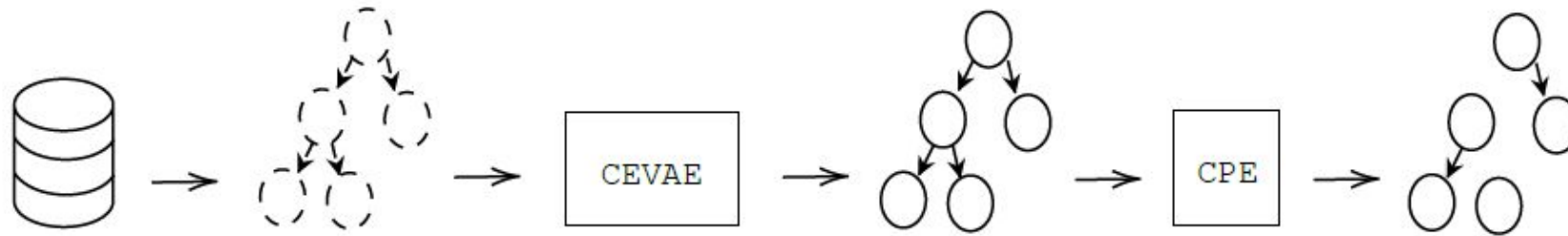
Symbol	Meaning
$\hat{Y}$	Predictor
$A$	Sensitive attribute
$X$	Covariates
$U$	Individual background



# FairTrade method

---

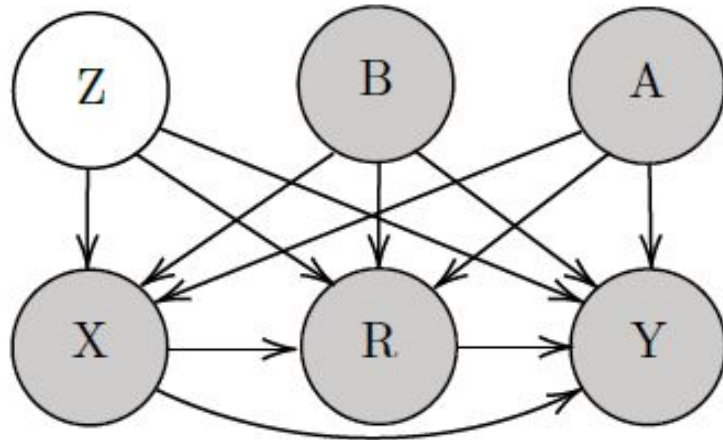
- I. **Assume** Causal graph
- II. **Infer** Causal relations
- III. **Fair** prediction





# FairTrade method

- I. **Assume** Causal graph
- II. Infer Causal relations
- III. Fair prediction



Symbol	Meaning
Y	Label
A	Sensitive attribute
Z	Unobserved confounder
B	Base variables
X	Other variables
R	Resolving variables





# FairTrade method

I. Assume Causal graph

II. Infer Causal relations

III. Fair prediction

Causal Effect Variational Autoencoder (CEVAE):

1. Inference Step: Recover unobserved confounder

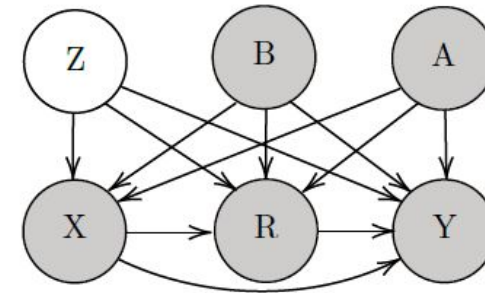
$$q(z | a, b, x, r)$$

2. Generative Step: Reconstruct observed variables from parents

$$p(x | z, a, b)$$

$$p(r | z, a, b, x)$$

$$p(y | z, a, b, x, r)$$





# FairTrade method

---

I. Assume Causal graph

II. Infer Causal relations

III. Fair prediction

Causal Path Enabler (CPE):

Train *auxiliary* model which only has *fair* information as *input*, possible input:

- Non-descendants of the sensitive variable
- Background variables independent of the sensitive variable
- Resolving variables





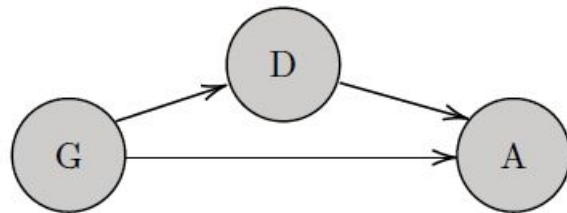


# FairTrade method

- I. Assume Causal graph
- II. Infer Causal relations
- III. Fair prediction

## Resolving variables

Variables deemed fair to use despite influence from sensitive variable



Variable	Meaning
G	Gender
D	Department choice
A	Admission rate

Berkley admission problem, is it fair if admission rate depends on gender via department choice?





# FairTrade method – Possible improvements

---

- I. **Assume** Causal graph
  - > Sensitivity analysis on assumption mistakes
- II. **Infer** Causal relations
  - > Research on recovering of true effects
- III. **Fair** prediction
  - > Formalisation step in input criteria CPE to enable PSE
  - > Evaluation counterfactual distributions





# Risk profiles in unlawful social welfare

---

## Situation Amsterdam:

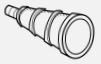
- Around 40.000 people receive social welfare, including an unknown number of fraudulent cases
- Municipality wants to decrease unlawful social welfare, and also has committed to only fair algorithms in the city
- Set of fraud labels is suspected to be biased due to the passive 'signal' approach over the last years

## Goal:

Create a classification model for risk profiles in social welfare, which is counterfactually fair with respect to *ethnicity*.



# Risk profiles in unlawful social welfare



# Experiment 3 - Risk profiles in social welfare

---

## Goal:

Create a classification model for risk profiles in social welfare, which is counterfactually fair with respect to *ethnicity*.

## Method:

FairTrade method

- I. **Assume** Causal graph
- II. **Infer** Causal relations
- III. **Fair** prediction

# Experiment 3 - Risk profiles in social welfare

---

## Data

Proof of concept experiment on CBS data:

- Safe workspace with computing power
- Data from more municipalities to reduce biases
- More attributes to infer personal background

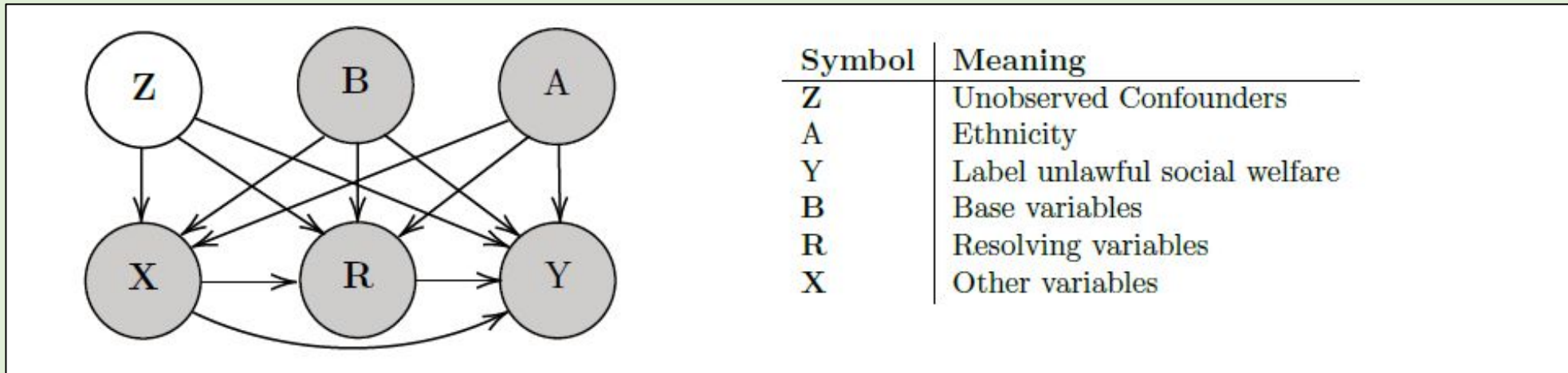
11.230 profiles with balanced fraud labels:

- |             |                         |
|-------------|-------------------------|
| ○ Age       | ○ Crime history         |
| ○ Education | ○ Debt                  |
| ○ Income    | ○ Partner               |
| ○ Housing   | ○ Household             |
| ○ Jobs      | ○ Other social benefits |
| ○ Property  |                         |



# Experiment 3 - Risk profiles in social welfare

- I. **Assume** Causal graph
- II. Infer Causal relations
- III. Fair prediction



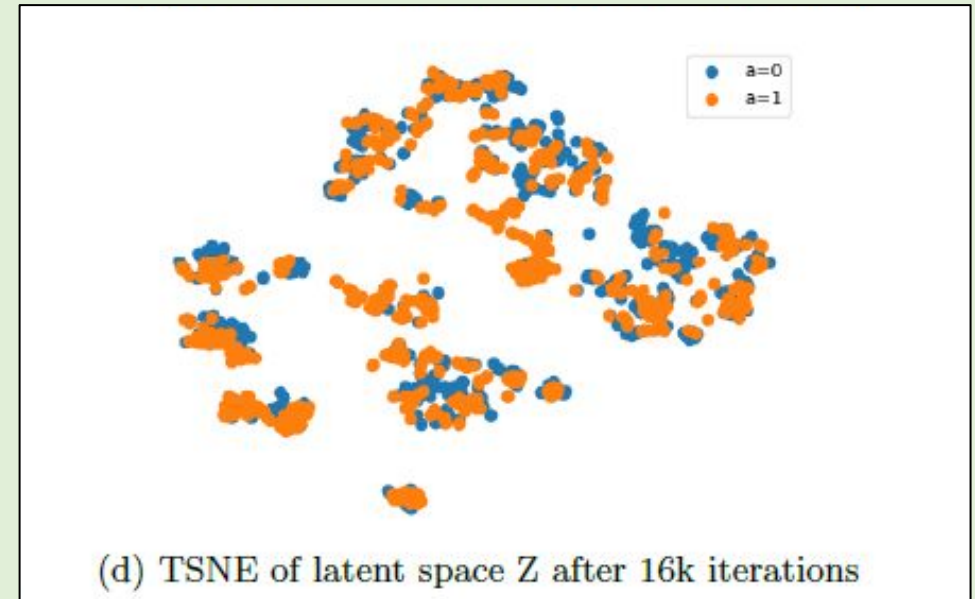
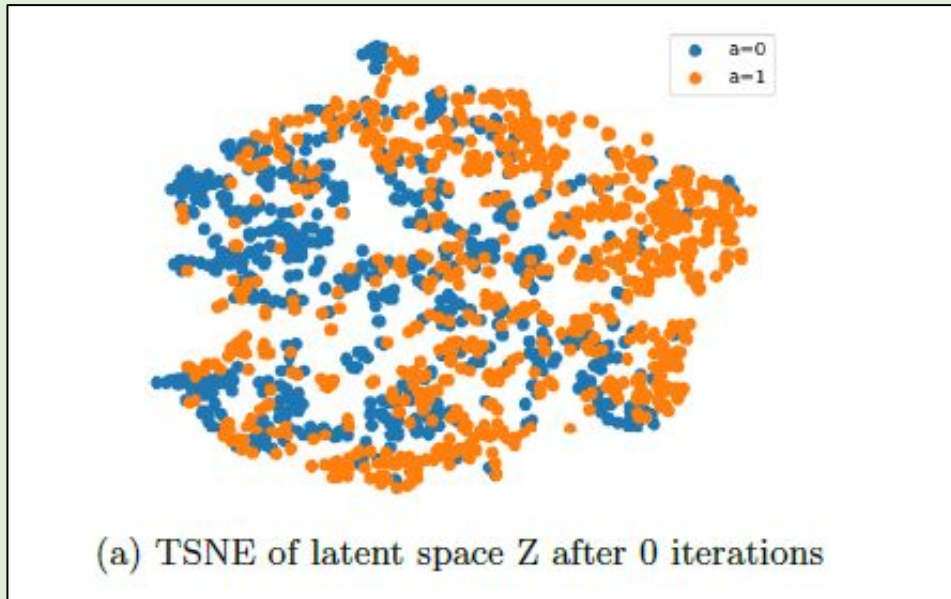
B: {Income mutation, Gender}

R: {Involved in crime, Partner with debt, Recidivism}

X: All other attributes

# Experiment 3 - Risk profiles in social welfare

- I. Assume Causal graph
- II. Infer Causal relations
- III. Fair prediction



Checking implied independencies: background independent of ethnicity?



# Experiment 3 - Risk profiles in social welfare

- I. Assume Causal graph
- II. Infer Causal relations
- III. **Fair prediction**

Statistical Parity: Group level fairness  $[0,1]$ , equals 1 under counterfactual fairness.

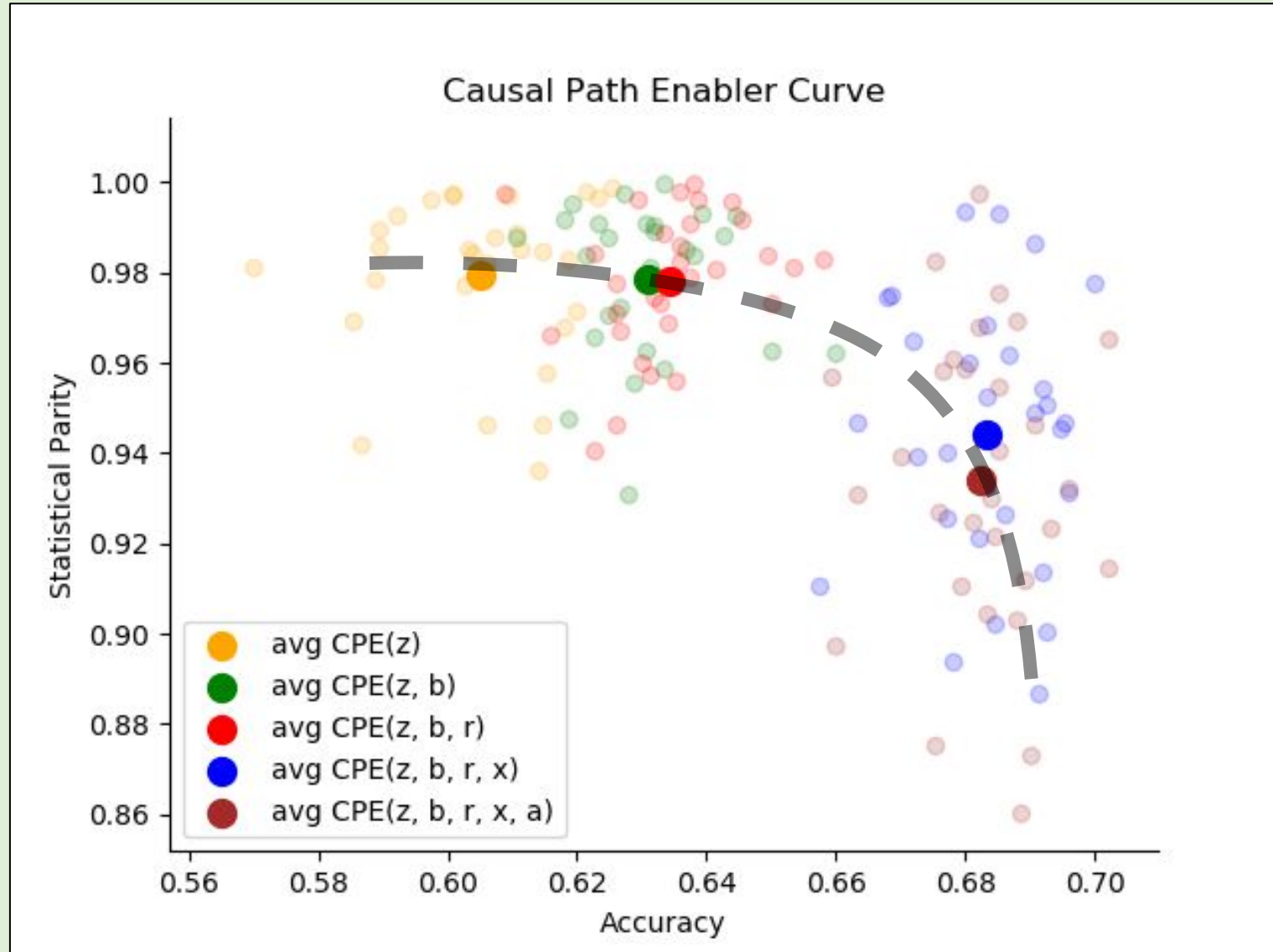
Baseline accuracies

Model	Accuracy (std)
Random Forest	0.6880 (0.01)
MLP	0.6785 (0.02)
Logistic Regression	0.6781 (0.02)

FairTrade model outcomes

	Accuracy (std)	Statistical Parity (std)
CPE(z)	0.605 (0.013)	<b>0.979</b> (0.017)
CPE(z,b)	0.631 (0.010)	0.978 (0.016)
CPE(z,b,r)	0.634 (0.011)	0.978 (0.015)
CPE(z,b,r,x)	0.683 (0.010)	0.944 (0.029)
CPE(z,b,r,x,a)	<b>0.682</b> (0.010)	0.934 (0.033)

# Experiment 3 - Risk profiles in social welfare

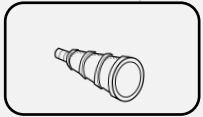


# Experiment 3 - Risk profiles in social welfare

---

## Experiment outcomes:

- Fairness and accuracy show a trade-off curve for different levels of constraints
- Using counterfactual fair information, an accuracy of 63% can be obtained on a balanced data set, compared to a maximum of 68% for models without fairness constraints



# Limitations and Future work

---

## 1. Evaluation

1. Causal assumptions
2. Counterfactual Fairness
3. Approximate inference
4. Path Specific Effects

## 2. Public debate on formal definitions of fairness





# Conclusion

---

- Machine Learning has increasing impact on people's live
- Causality helps to formalise fairness
- The FairTrade method makes it possible to approximate fair models in practical applications
- A trade-off curve between fairness and accuracy is obtained for neural network based classification of exceptionally detailed real data profiles



# Thank you for listening!

