

A

Silvia Chiappa csilvia@deepmind.com

Greek Stochastics 23-26/08/2022



## Outline

Part I: Causal Inference (CI) via Causal Bayesian Networks (CBNs) Part II: CBNs and CI for ML Fairness



## **Machine Learning**

#### **ML Prerequisites**

Linear Algebra Optimization Calculus Probability Statistics

#### **ML Approaches and Methods**

Linear Models Regression, Classification, Maximum Likelihood, Bayesian Inference

Neural Networks Backpropagation, Automatic Differentiation, Gradient Optimization, Deep Learning

Graphical Models Directed, Undirected, Factor Graphs, Conditional Independence, Message Passing

Sequential Decision Making Reinforcement Learning

...

Approximate Inference Variational Inference, Monte-Carlo Methods

#### **Statistical Causality**

**Machine Learning** 

#### ML Prerequisites Basic Understanding of Statistical Causality,

...

#### Core Ingredients of ML Development/Evaluation/Deployment Causally-based Methods/Considerations,

...

#### **Obstacles**

- 1. Knowledge barrier
- 2. Difficulties in identification (black-box)
- 3. Difficulties in estimation
- 4. Not a bottleneck



Part I Causal Inference (CI) via Causal Bayesian Networks (CBNs)



## **Causal Inference**

Separation of Statistical Dependence into Causal Influence and Spurious Correlation

#### Drug Example

 $\rightarrow$ 

- Doctor gives drug more often to male individuals  $G \rightarrow D$
- Recovery is influenced by drug  $D \rightarrow R$  and gender  $G \rightarrow R$  (male individuals recover more often)



<u>Causality: Models, Reasoning, and Inference.</u> J. Pearl. Goal: Learn from data whether giving the drug helps recovery

#### Statistical Dependence (ML Approach)

Conditional **observational distribution** p(R | D)Measures statistical dependence between D and Rcontaining both causal influence  $(D \rightarrow R)$  and spurious correlation  $(D \leftarrow G \rightarrow R)$ 

Causal Influence (Causal Inference Approach)

Conditional interventional distribution

 $p_{\rightarrow D}(R|D) = p(R|\mathrm{do}(D))$ 

Measures only causal influence  $(D \rightarrow R)$ 

#### **Causal Inference via CBNs**

**Identification Problem**: Answer whether it is possible / how to express a causal quantity as a **functional of the observational distribution** *p*, using knowledge of the causal graph

Estimation Problem: How to estimate the functional



## (Causal) Bayesian Networks

Probabilistic Graphical Models: Principles and Techniques. D. Koller, N. Friedman.

Bayesian Reasoning and Machine Learning. <u>D. Barber.</u>

Pattern Recognition and Machine Learning. <u>C. Bishop.</u>

- <u>Causal Inference in Statistics: A Primer.</u> J. Pearl, M. Glymour, N. P. Jewell.
- <u>Elements of Causal Inference</u> Foundations and Learning <u>Algorithms. J. Peters, D. Janzing, B. Schölkopf.</u>
- Causality: Models, Reasoning, and Inference. J. Pearl.



Direct link between two nodes represents statistical relationship

Direct link between two nodes represents causal relationship





## **Bayesian Networks**

Directed acyclic graph  ${\cal G}$  in which each node is associated with conditional distribution given its parents  $p(X_i|{\rm pa}(X_i))$ 

Joint distribution of all nodes is given by the product of all conditional distributions

$$p(X_1,\ldots,X_I) = \prod_{i=1}^{I} p(X_i | \operatorname{pa}(X_i))$$

## **Causal Bayesian Networks**

Bayesian network with special semantic: Each conditional distribution represents a causal mechanism

*X* is a cause of (influences) *Y* if there exists a causal path from *X* to *Y* 

Causal path := directed path

- $\rightarrow X$  is a cause of Y if X is an ancestor of Y
- $\rightarrow X$  is a cause of Y if Y is a descendant of X



- A is a cause of Y
- A is not a cause of Q

p(A,Q,D,Y) = p(Y|A,Q,D)p(D|A)p(A)p(Q)



## Statistical Independence in (Causal) Bayesian Networks (d-separation)

#### d-separation

X is d-separated from Y by  $Z(X \perp _{\mathcal{G}} Y \mid Z)$  if all paths from any element of X to any element of Y are closed (or blocked) given Z

A path is closed if at least one of the following conditions is satisfied

- 1. There is a collider (node with 2 parents) on the path such that neither the collider nor any of its descendants belong to the conditioning set Z
- 2. There is a non-collider on the path which belongs to the conditioning set Z

#### d-separation implies statistical independence

If X is d-separated from Y by Z then X and Y are statistically independent given  $Z(X \perp p Y | Z)$ 

 $X \perp\!\!\!\perp_{\mathcal{G}} Y \,|\, Z \Rightarrow X \perp\!\!\!\perp_p Y |Z$ 





## Statistical Independence in (Causal) Bayesian Networks (d-separation)

#### d-separation

X is d-separated from Y by  $Z(X \perp _{\mathcal{G}} Y \mid Z)$  if all paths from any element of X to any element of Y are closed (or blocked) given Z

A path is closed if at least one of the following conditions is satisfied

- 1. There is a collider (node with 2 parents) on the path such that neither the collider nor any of its descendants belong to the conditioning set Z
- 2. There is a non-collider on the path which belongs to the conditioning set Z

d-separation implies statistical independence

If X is d-separated from Y by Z then X and Y are statistically independent given  $Z(X \perp p Y | Z)$ 

 $X \perp\!\!\!\perp_{\mathcal{G}} Y \mid Z \Rightarrow X \perp\!\!\!\perp_p Y \mid Z$ 



$$A \leftarrow C \leftarrow X \rightarrow Y \text{ is open}$$

## Statistical Independence in (Causal) Bayesian Networks (d-separation)

#### d-separation

X is d-separated from Y by  $Z(X \perp _{\mathcal{G}} Y \mid Z)$  if all paths from any element of X to any element of Y are closed (or blocked) given Z

A path is closed if at least one of the following conditions is satisfied

- 1. There is a collider (node with 2 parents) on the path such that neither the collider nor any of its descendants belong to the conditioning set Z
- 2. There is a non-collider on the path which belongs to the conditioning set Z

#### d-separation implies statistical independence

If X is d-separated from Y by Z then X and Y are statistically independent given  $Z(X \perp p Y | Z)$ 

 $X \perp\!\!\!\perp_{\mathcal{G}} Y \mid Z \Rightarrow X \perp\!\!\!\perp_p Y \mid Z$ 



#### $A \leftarrow E \rightarrow C \leftarrow X \rightarrow Y \, \text{ is closed}$

#### Conditioning on C opens the path





## Identification of Causal Quantities

Probabilistic Graphical Models: Principles and Techniques. D. Koller, N. Friedman.

Bayesian Reasoning and Machine Learning. D. Barber.

Pattern Recognition and Machine Learning. <u>C. Bishop.</u>

- <u>Causal Inference in Statistics: A Primer.</u> J. Pearl, M. Glymour, N. P. Jewell.
- <u>Elements of Causal Inference</u> Foundations and Learning <u>Algorithms. J. Peters, D. Janzing, B. Schölkopf.</u>
- Causality: Models, Reasoning, and Inference. J. Pearl.

## Statistical Dependence between A and Y versus Causal Influence of A on Y



- Causal paths  $\begin{array}{c} A \to Y \\ A \to D \to Y \end{array}$
- Non-causal path (back-door path)

 $A \leftarrow C \to Y$ 



Causal influence  $p_{\rightarrow A}(Y|A)$ 

Conditional distribution of Y given Awithout the information from A traveling through  $A \leftarrow C \rightarrow Y$ 



## **Actual Hard Intervention**



1. Get new data from interventional distribution  $p_{\to A}$  where A i set to value a (does not depend on C)  $p_{\to A}(A) = \delta_{A=a}$ 



## **Causal Influence**

2. Compute conditional interventional distribution  $p_{\rightarrow A}(Y|A)$  using the interventional data

$$\begin{split} p_{\rightarrow A}(Y|A) &= \sum_{C} p_{\rightarrow A}(Y|A,C) p_{\rightarrow A}(C|A) \\ &= \sum_{C} p_{\rightarrow A}(Y|A,C) p_{\rightarrow A}(C) \\ &= \sum_{C} p(Y|A,C) p(C) \end{split}$$



## **Hypothetical Hard Intervention**



- 1. Change p(A|C) into  $p_{\rightarrow A}(A) = \delta_{A=a}$
- 2. Leave other conditional distributions as in p



## **Causal Influence**

3. Express conditional interventional distribution as a functional of the observational distribution

$$\begin{split} p_{\rightarrow A}(Y|A) &= \sum_{C} p_{\rightarrow A}(Y|A,C) p_{\rightarrow A}(C|A) \\ &= \sum_{C} p_{\rightarrow A}(Y|A,C) p_{\rightarrow A}(C) \\ &= \sum_{C} p(Y|A,C) p(C) \end{split}$$



## **Statistical Dependence**



$$p(Y|A) = \sum_{C} p(Y|A, C) p(C|A)$$

## **Causal Influence**



$$p_{\rightarrow A}(Y|A) = \sum_{C} p(Y|A,C)p(C)$$



## **Statistical Dependence**

# $\begin{array}{c|c} \mathcal{G},p \\ \hline \mathcal{C} & p(C) \\ \hline p(A|C) & & & & & & \\ \end{array} \\ p(A|C) & & & & & & & & \\ \end{array} \\ \begin{array}{c} \mathcal{G},p \\ \mathcal{V} & p(Y|A,C) \\ \end{array} \end{array}$

$$p(Y|A) = \sum_{C} p(Y|A, C) p(C|A)$$

## **Causal Influence**



NOTE 
$$p_{\rightarrow A}(Y) = p_{\rightarrow A}(Y|A = a)$$
  
=  $\sum_{C} p(Y|A = a, C)p(C)$ 



#### **Drug Example** Separation of Causal Influence from Spurious Correlation

#### **Data Generation Mechanism**

- Doctor gives drug more often to male individuals  $G \rightarrow D$
- Recovery is influenced by drug  $D \rightarrow R$  and gender  $G \rightarrow R$



•	Causality: Models, Reasoning, and Inference.
	J. Pearl.

Females	R = 0	R = 1		Recovery Rate
No Drug $(D=0)$	21	9	30	30%
Drug $(D=1)$	8	2	10	20%
	29	11	40	
Males	R = 0	R = 1		Recovery Rate
$\frac{\text{Males}}{\text{No Drug } (D=0)}$	$\frac{R=0}{3}$	$\frac{R=1}{7}$	10	Recovery Rate 70%
$\begin{array}{c} \text{Males} \\ \hline \text{No Drug} (D=0) \\ \text{Drug} (D=1) \end{array}$	R = 0 $3$ $12$	R = 1 $7$ $18$	10 30	Recovery Rate 70% 60%
$\begin{array}{c} \text{Males} \\ \text{No Drug} \ (D=0) \\ \text{Drug} \ (D=1) \end{array}$	R = 0 $3$ $12$ $15$	R = 1 $7$ $18$ $25$	10 30 40	Recovery Rate 70% 60%

	R = 0	R = 1		Recovery Rate
No Drug $(D=0)$	24	16	40	40%
Drug $(D=1)$	20	20	40	50%
	44	36	80	

#### **Drug Example** Separation of Causal Influence from Spurious Correlation

#### Example of Simpson's Paradox

Statistical association that holds in several different groups of data is reversed when groups are combined

- According to the female data and to the male data the drug decreases recovery
- According to the combined male and female data the drug increases recovery

We believe this is a paradox as we wrongly believe that conditional distribution gives causal effect

Females	R = 0	R = 1		Recovery Rate
No Drug $(D=0)$	21	9	30	30%
Drug $(D=1)$	8	2	10	20%
	29	11	40	
Males	R = 0	R = 1		Recovery Rate
$\frac{\text{Males}}{\text{No Drug } (D=0)}$	$\frac{R=0}{3}$	$\frac{R=1}{7}$	10	Recovery Rate 70%
$\frac{\text{Males}}{\text{No Drug } (D=0)}$ Drug $(D=1)$	$\frac{R=0}{3}$ 12	$\frac{R=1}{7}$ 18	10 30	Recovery Rate 70% 60%
Males No Drug $(D = 0)$ Drug $(D = 1)$	R = 0 $3$ $12$ $15$	R = 1 $7$ $18$ $25$	10 30 40	Recovery Rate 70% 60%
Males   No Drug (D = 0)   Drug (D = 1)	R = 0 $3$ $12$ $15$	R = 1 $7$ $18$ $25$	10 30 40	Recovery Rate 70% 60%
Males   No Drug (D = 0)   Drug (D = 1)	R = 0 $3$ $12$ $15$	$\frac{R=1}{7}$ $\frac{18}{25}$	10 30 40	Recovery Rate 70% 60%

	R = 0	R = 1		Recovery Rate
No Drug $(D = 0)$	24	16	40	40%
Drug $(D=1)$	20	20	40	50%
	44	36	80	



## Simpson's Paradox

Females	R = 0	R=1		Recovery Rate
No Drug $(D=0)$	21	9	30	30%
$Drug \ (D=1)$	8	2	10	20%
	29	11	40	

Males	R = 0	R = 1		Recovery Rate
No Drug $(D=0)$	3	7	10	70%
Drug $(D=1)$	12	18	30	60%
	15	25	40	

	R = 0	R = 1		Recovery Rate
No Drug $(D = 0)$	24	16	40	40%
Drug $(D=1)$	20	20	40	50%
	44	36	80	

## **Conditional Distribution**

$$p(G = F|D = 0) = \frac{30}{40}, p(G = M|D = 0) = \frac{10}{40}$$
$$p(G = F|D = 1) = \frac{10}{40}, p(G = M|D = 1) = \frac{30}{40}$$

$$p(R = 1|D = 0) = \sum_{G} p(R = 1|D = 0, G)p(G|D = 0)$$
  
=  $\frac{9}{30} \cdot \frac{30}{40} + \frac{7}{10} \cdot \frac{10}{40} = \frac{16}{40} = 0.4$   
$$p(R = 1|D = 1) = \sum_{G} p(R = 0|D = 1, G)p(G|D = 1)$$
  
=  $\frac{2}{10} \cdot \frac{10}{40} + \frac{18}{30} \cdot \frac{30}{40} = \frac{20}{40} = 0.5$ 



## Simpson's Paradox

Females	R = 0	R = 1		Recovery Rate
No Drug $(D = 0)$	21	9	30	30%
Drug $(D=1)$	8	2	10	20%
	29	11	40	

Males	R = 0	R = 1		Recovery Rate
No Drug $(D = 0)$	3	7	10	70%
$Drug \ (D=1)$	12	18	30	60%
	15	25	40	

	R = 0	R = 1		Recovery Rate
No Drug $(D = 0)$	24	16	40	40%
Drug $(D=1)$	20	20	40	50%
	44	36	80	

## **Causal Influence**



$$p(G=F) = \frac{40}{80}, p(G=M) = \frac{40}{80}$$

$$p_{\to D=0}(R=1|D=0) = \sum_{G} p(R=1|D=0,G)p(G)$$
$$= \frac{9}{30} \cdot \frac{40}{80} + \frac{7}{10} \cdot \frac{40}{80} = \frac{9+21}{30} \cdot \frac{1}{2} = 0.5$$
$$p_{\to D=1}(R=1|D=1) = \sum_{G} p(R=1|D=1,G)p(G)$$
$$= \frac{2}{10} \cdot \frac{40}{80} + \frac{18}{30} \cdot \frac{40}{80} = \frac{6+18}{30} \cdot \frac{1}{2} = 0.4$$

## **Identification via Do-calculus**

Three graphical rules that, combined, might allow to express a causal quantity as a functional of the observational distribution

Might: In the presence of latent variables, for certain graph structures identification is not possible





## **Rules of Do-calculus (Simpler Version)**

Rule 2 - Action/Observation Exchange

 $p_{\rightarrow A}(Y|A) = p(Y|A)$ 

 $\text{if } Y \perp \!\!\!\perp_{\mathcal{G}_{A \to}} A$ 



#### Intuition

The graph  $\mathcal{G}_{A\rightarrow}$  has all links emerging from A removed. If Y is independent of A in this graph it means that all backdoor paths (non-causal paths) are closed, therefore the interventional and conditional distributions are equivalent. The only open paths are causal paths! 1. Backdoor path from A to  $Y: A \leftarrow \dots Y$ 

Backdoor paths must be non-causal as otherwise we would create a cycle (A is a cause of Y which is a causal of A)

#### 2. Frontdoor path from A to $Y: A \rightarrow \dots Y$

Only causal frontdoor paths are open, as to be non-causal a frontdoor path must contain a collider  $A \rightarrow \dots \rightarrow Z \leftarrow \dots Y$ 

## **Rules of Do-calculus (Simpler Version)**

Rule 2 – Action/Observation Exchange  $p_{\to A}(Y|A) = p(Y|A)$  if  $Y \perp\!\!\!\!\perp_{\mathcal{G}_{A\to}} A$ 



#### Intuition

The graph  $\mathcal{G}_{A\rightarrow}$  has all links emerging from A removed. If Y is independent of A in this graph it means that all backdoor paths (non-causal paths) are closed, therefore the interventional and conditional distributions are equivalent. The only open paths are causal paths! Rule 3 – Insertion/Deletion of Action $p_{\rightarrow A}(Y|A) = p(Y)$  if  $Y \perp\!\!\!\perp_{\mathcal{G}_{\rightarrow A}} A$ 



#### Intuition

The graph  $\mathcal{G}_{\rightarrow A}$  has all links pointing to A removed. If Y is independent of A in this graph it means that there are no causal paths from A to Y, i.e. A is not a cause of Y, therefore intervening on A or ignoring A is the same. There is no influence reaching Y from A. The only possibly open paths are non-causal paths!



## **Rules of Do-calculus**

#### **Backdoor Criterion**

$$p_{\rightarrow A}(Y|A) = \sum_{Z} p(Y|A, Z) p(Z)$$



Called Adjustment Criterion for a set Z satisfying different conditions (called adjustment set)

Rule 2 Action / Observation exchange

Rule 3 Insertion / Deletion of Action

#### 1. Z is in a backdoor path from A to

Y: Removing links pointing to A blocks the subpaths  $A \leftarrow \dots Z$ A frontdoor path from A to Z must contain a collider otherwise Y would be a cause of A.

2. Z is in a frontdoor path from A

to Y: The subpath  $A \rightarrow \dots Z$  cannot be causal, as otherwise Z would be a descendant of A, and therefore is closed (must contain a collider).

## **Rules of Do-calculus**

**Frontdoor Criterion** 

$$p_{\rightarrow A=a}(Y|A=a) = \sum_{Z} p(Z|A=a) \sum_{A} p(Y|Z,A)p(A)$$

1. If all backdoor paths from Z to Y are blocked by A  $(Y \perp \!\!\!\perp_{\mathcal{G}_{Z \rightarrow}} Z \mid A)$ 2. There is no open backdoor path from A to Z 3. Z intercepts all causal paths from A to Y



 $Y \perp\!\!\!\perp_{\mathcal{G}_{\to A, Z \to}} Z$ 

Rule 2 Action / Observation exchange

 $Z \perp\!\!\!\perp_{\mathcal{G}_{A \to}} A$ 

**Backdoor Criterion** 

#### Rule 2 Action / Observation exchange

Backdoor paths from A to Z: are closed. Causal paths from from A to Z: are cut in the graph with emerging links from A removed.

$$p_{\rightarrow A=a}(Y|A = a) = \sum_{Z} p_{\rightarrow A=a}(Y|A = a, Z)p_{\rightarrow A=a}(Z|A = a)$$

$$= \sum_{Z} p_{\rightarrow A=a, \rightarrow Z}(Y|A = a, Z)p(Z|A = a)$$

$$= \sum_{Z} p_{\rightarrow Z}(Y|Z)p(Z|A = a)$$

$$= \sum_{Z} p(Z|A = a)\sum_{A} p(Y|Z, A)p(A)$$

$$Y \perp \mathcal{G}_{\rightarrow A, Z \rightarrow} A \qquad Y \perp \mathcal{G}_{Z \rightarrow} Z \mid A$$

Rule 3 Insertion / Deletion of Action

Backdoor paths from *A* to *Y*: are cut in the graph with incoming links into *A* removed.

Causal paths from from A to Y: Z intercepts these paths which are cut in the graph with emerging links from Z removed.





## **Estimation of Causal Quantities**

Using knowledge of the CBN, do-calculus enables us (when possible) to obtain identification formulas, i.e. formulas that express a causal quantity as a functional of the observational distribution

- 1. Obtain an estimate of an identification formula with desirable properties
- 2. Choose between several identification formulas

Asymptotically Best Causal Effect Identification with Multi-Armed Bandits. A. Malek, S. Chiappa, 2021.

 $\rightarrow$ 



#### **Selection of Identification Formulas**



## Identification Formulas for $\tau = \mathbb{E}_{p_{\rightarrow A}(Y|A)}[Y]$

- Adjustment criterion using Z<sub>1</sub>
- Frontdoor criterion using  $Z_2$

• Statistical Considerations: Asymptotic variance (each formula associated with asymptotically linear estimator of  $\tau \rightarrow$  convergence rate  $O(\sqrt{n})$ )

Practical Considerations: Costs in observing covariates

#### **Asymptotic Variance**

#### Graphical Criteria for Selection of Adjustment Sets

 [1] Graphical Criteria for Efficient Total Effect Estimation via Adjustment in Causal Linear Models.
 L. Henckel, E. Perković, M. H. Maathuis, 2019.

[2] Efficient Adjustment Sets for Population Average Causal Treatment Effect Estimation in Graphical Models.A. Rotnitzky and E. Smucler, 2020.

[3] **On efficient adjustment in causal graphs.** J. Witte, L. Henckel, M. H. Maathuis, V. Didelez, 2020.

[4] Efficient Adjustment Sets in Causal Graphical Models with Hidden Variables. E. Smucler, F. Sapienza, A. Rotnitzky, 2021.

#### **Graphical Criteria for Selection of Adjustment Sets**

 $B \cap G = \emptyset$ 

#### **Addition of Precision Set**

Adding to the adjustment set *B* a precision set *G*, i.e. variables that d-separated from *A* by *B*  $(A \perp _{\mathcal{G}} G \mid B)$  gives an adjustment set *B* U *G* with smaller asymptotic variance

#### 2.

#### **Removal of Overadjustment Set**

Removing from the adjustment set  $B \cup G$  the *overadjustment* set B, i.e. variables that are d-separated from Y by G and A ( $Y \perp \!\!\!\perp_{\mathcal{G}} B \mid G \cup A$ ) gives an adjustment set G with smaller asymptotic variance



#### Adjustment sets

$$\begin{split} & \emptyset, \{Z_1\}, \{Z_2\}, \{Z_4\} \\ & \{Z_1, Z_2\}, \{Z_1, Z_4\} \\ & \{Z_2, Z_3\}, \{Z_2, Z_4\}, \{Z_3, Z_4\} \\ & \{Z_1, Z_2, Z_3\}, \{Z_1, Z_2, Z_4\} \\ & \{Z_1, Z_3, Z_4\}, \{Z_2, Z_3, Z_4\} \\ & \{Z_1, Z_2, Z_3, Z_4\} \end{split}$$

**2.** Adjustment set  $\{Z_1, Z_2, Z_3, Z_4\}$  $Y \perp _{\mathcal{G}} \{Z_1, Z_2, Z_3\} | Z_4 \cup A \implies \{Z_1, Z_2, Z_3\}$  overadjustment set

2. Adjustment set  $Z_4$  $Y \perp \downarrow \mathcal{G} Z_4 \mid A \implies Z_4$  not an overadjustment set

#### 1. Adjustment set 🖉

 $A \perp\!\!\!\perp_{\mathcal{G}} Z_4 \Rightarrow Z_4$  precision set

2. Adjustment set  $\{Z_1, Z_2\}$  $Y \perp \mathcal{L}_{\mathcal{G}} \{Z_1, Z_2\} \mid A \implies \{Z_1, Z_2\}, Z_1, Z_2$  overadjustment sets

2. Adjustment set  $\{Z_1, Z_2, Z_3\}$  $Y \perp\!\!\!\perp_{\mathcal{G}} Z_1 \mid \{Z_2, Z_3\} \cup A \implies Z_1$  overadjustment set



#### **Graphical Criteria for Selection of Adjustment Sets**

#### **Comparison of Adjustment Sets**

If *B* and *G* are adjustment sets with  $A \perp \!\!\!\perp_{\mathcal{G}} G \setminus B \mid B$ and  $Y \perp \!\!\!\perp_{\mathcal{G}} B \setminus G \mid G \cup A$  then the variance of *G* is smaller than the variance of *B*  Causal nodes cn(A, Y, G): Nodes on causal paths from A to Y, excluding A

Forbidden nodes  $forb(A, Y, \mathcal{G}) = de(cn(A, Y, \mathcal{G}), \mathcal{G}) \cup A$ 

#### **Optimal Adjustment Set**

The variance of the adjustment set  $O(A,Y,\mathcal{G}) = \operatorname{pa}(\operatorname{cn}(A,Y,\mathcal{G}),\mathcal{G}) \setminus \operatorname{forb}(A,Y,\mathcal{G})$  is smaller than for any other adjustment set

 $Z_2 \longrightarrow Z_3 \longrightarrow Z_4$   $\downarrow$   $Z_1 \longrightarrow A \longrightarrow Y$ 

 $cn(A, Y, \mathcal{G}) = \{Y\}$ forb $(A, Y, \mathcal{G}) = \{A, Y\}$ pa $(cn(A, Y, \mathcal{G}), \mathcal{G}) = \{Z_4\}$  $O(A, Y, \mathcal{G}) = \{Z_4\}$ 

Variance of  $Z_4$  smaller than variance of empty set

 $B = \emptyset, G = \{Z_4\} \quad A \perp \!\!\!\perp_{\mathcal{G}} G \setminus B \mid B$  $B = \{Z_4\}, G = \emptyset \quad Y \perp \!\!\!\perp_{\mathcal{G}} B \setminus G \mid G \cup A$ 

Cannot compare  $\{Z_2, Z_3\}$  and empty set

 $B = \emptyset, G = \{Z_2, Z_3\} \quad A \amalg_{\mathcal{G}} G \setminus B \mid B$  $B = \{Z_2, Z_3\}, G = \emptyset \quad Y \amalg_{\mathcal{G}} B \setminus G \mid G \cup A$ 

#### $Z_4$ latent variable

 $\begin{array}{l} \emptyset, \{Z_1\}, \{Z_2\}, \{Z_4\} \\ \{Z_1, Z_2\}, \{Z_1, Z_4\} \\ \{Z_2, Z_3\}, \{Z_2, Z_4\}, \{Z_3, Z_4\} \\ \{Z_1, Z_2, Z_3\}, \{Z_1, Z_2, Z_4\} \\ \{Z_1, Z_3, Z_4\}, \{Z_2, Z_3, Z_4\} \\ \{Z_1, Z_2, Z_3, Z_4\} \end{array}$ 



## Asymptotically Best Causal Effect Identification with Multi-Armed Bandit

A selection method that is applicable to arbitrary identification formulas and accounts for both statistical performance and practical considerations must make use of data

#### Setting

- The investigator collects observational data
- Sequential strategy in which the investigator decides, observation-by-observation (rounds), which subset of variables to observe with the goal of identifying the best formula with the fewest observations

#### Multi-Armed Bandits Formalism

Cast this setting into the best-arm-identification bandit framework, by considering each formula as one arm and by replacing the goal of learning the arm with the best mean with the goal of learning the formula with the best cost-adjusted asymptotic variance

Each arm k corresponds to an estimator  $\hat{\tau}_k$  of  $\tau$  with asymptotic variance  $\sigma_k^2$  and cost  $c_k$ 

Goal:  $k^* = \arg\min_k c_k \sigma_k^2$ 

Leverage algorithms from the bandit literature

Adapt LUCB and Successive Elimination (SE) by introducing finite-sample confidence sequences on  $\hat{\sigma}_k^2(\mathcal{D}_n) - \sigma_k^2$  (confidence bounds that hold for all formulas and all rounds simultaneously)



## Estimators of au

Asymptotically linear estimator  $\hat{\tau}_k$  with nuisance function  $\eta_k$  (e.g. propensity score for IPW)

 $\begin{array}{l} \text{Influence function } \phi_k \text{satisfying} \\ & \mathbb{E}_p[\phi_k(W_k,\eta_k,\tau)] = 0 \quad \mathbb{E}_p[\phi_k^2(W_k,\eta_k,\tau)] < \infty \quad W_k \text{ covariates for formula k, A, Y} \\ & \sqrt{n}(\hat{\tau}_k(\mathcal{D}_n) - \tau) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \phi_k(w_k^i,\eta_k,\tau) + o_p(1) \end{array}$ 

Asymptotically normal with variance  $\sigma_k^2 = \mathbb{E}_p[\phi_k^2(W_k,\eta_k,\tau)]$ 

Can assume an uncentered influence function  $\phi_k(W_k,\eta_k,\tau)=\psi_k(W_k,\eta_k)-\tau$  $\sigma_k^2=\mathrm{Var}_p[\psi_k(W_k,\eta_k)]$ 

When need to estimate  $\eta_k$ , converge of  $\hat{\eta}_k$  at rate slower than  $O(n^{-1/2})$  could cause loss of asymptotically linearity (can happen if e.g.  $\hat{\eta}_k$  is modelled with with neural networks)

Asymptotically linear estimators for any identification formula [Jung, Tian, Bareinboim 2021] using double machine learning

## **Estimating the Asymptotic Variance**

- Our goal: estimate  $\sigma_k^2 = \mathbb{E}_p[\psi_k(W_k, \eta_k) \tau)^2]$  and derive finite-sample confidence sequences
- Our estimator  $\hat{\sigma}_k^2$  for data  $\mathcal{D}_n$  (inspired by [Chernozhukov et al., 2016])
  - $\circ$  Randomly split  $\mathcal{D}_n$  into two folds  $\mathcal{D}_n^\eta, \mathcal{D}_n^\sigma$
  - Fit the nuisance function  $\hat{\eta}_k(\mathcal{D}_n^\eta)$
  - Fit  $\hat{\sigma}_k^2$  with an empirical variance:  $\hat{\sigma}_k^2(\mathcal{D}_n) = \operatorname{Var}_{\mathcal{D}_n^\sigma} \left[ \psi_k(W_k, \hat{\eta}_k(\mathcal{D}_n^\eta)) \right]$

#### **Main Theorem**

 $O(\sqrt{n})$  estimation of the asymptotic variance is possible whenever  $O(\sqrt{n})$  estimation of  $\tau$  is possible



## **Confidence Sequence LUCB (CS-LUCB)**

At each round t

lower bound

- Sample arms  $l_t \leftarrow \arg \min_{k \in [K]} c_k \hat{\sigma}_k^2$ Sample arms  $u_t \leftarrow \arg \min_{k \neq l_t} c_k (\hat{\sigma}_k^2)$
- •
- If  $c_{l_t}\left(\hat{\sigma}_{l_t}^2 + \beta_{l_t}\right) \le c_{u_t}\left(\hat{\sigma}_{u_t}^2 \beta_{u_t}\right) \epsilon \operatorname{return} \hat{k} = l_t$ • otherwise collect new data and update  $l_t$ and  $u_t$

 $\hat{k}$  is at most  $\epsilon$ -suboptimal with probability at least  $1-\delta$ 

$$\mathbb{P}\left(c_{\hat{k}}\sigma_{\hat{k}}^2 \ge \min_k c_k \sigma_k^2 + \epsilon\right) \le \delta$$



Create gap between upper bounds of best arms and lower bounds of remaining arms

## **Linear Experiments**



- 2<sup>M=3</sup> =8 formulas from adjustment criterion
- One formula from frontdoor criterion  $Z_M$



## **Non-Linear Experiments**



Number of Formulas

## Part II CBNs and CI for ML Fairness



## **ML Fairness**

Machine Learning deployed to make decisions that affect people lives can treat individuals unfairly on the basis of sensitive attributes such as race, gender, disabilities, etc.





+



## College Admission Example Characterising Patterns of Unfairness





#### Berkeley's alleged sex bias (1974)

Female applicants were rejected more often than male applicants



## College Admission Example Characterising Patterns of Unfairness







Women spontaneously apply to departements with lower admission rates



## **Different Possible Unfairness Scenarios**



Women spontaneously apply to departements with lower admission rates Women apply to departments with lower admission rates due to due systemic historical or cultural pressures College lower admission rates of departments chosen spontaneously more often by women





## **Path-Specific Influence**

Probabilistic Graphical Models: Principles and Techniques. D. Koller, N. Friedman.

Bayesian Reasoning and Machine Learning. <u>D. Barber.</u>

Pattern Recognition and Machine Learning. <u>C. Bishop.</u>

- <u>Causal Inference in Statistics: A Primer.</u> J. Pearl, M. Glymour, N. P. Jewell.
- <u>Elements of Causal Inference</u> Foundations and Learning <u>Algorithms. J. Peters, D. Janzing, B. Schölkopf.</u>
- Causality: Models, Reasoning, and Inference. J. Pearl.





## **Causal Influence**



Potential Outcome  $Y_a$ Random variable with distribution  $p_{\rm PO}(Y_a):=p_{\rightarrow A=a}(Y)$ 

$$p_{\text{PO}}(Y_a) = p_{\rightarrow A=a}(Y)$$
  
=  $\sum_{C,D} p_{\rightarrow A=a}(C, D, Y)$   
=  $\sum_{C,D} p(Y|A = a, C, D)p(D|A = a)p(C)$ 

## Path-Specific Influence



Path-Specific Potential Outcome  $Y_{\bar{a}}(D_a)$ Random variable with distribution  $p_{\mathrm{PO}}(Y_{\bar{a}}(D_a)) := p_{\rightarrow A, A = \bar{a} \rightarrow Y, A = a \rightarrow D}(Y)$ 

$$p_{\text{PO}}(Y_{\bar{a}}(D_a)) = \sum_{C,D} p(Y|A = \bar{a}, C, D) p(D|A = a) p(C)$$

## **Average Treatment Effect**



Average Treatment Effect  $\mathrm{ATE}_{a\bar{a}} = \mathbb{E}_{p(Y_{\bar{a}})}[Y_{\bar{a}}] - \mathbb{E}_{p(Y_{a})}[Y_{a}]$ 

Measure difference of causal influence when A is set to  $A = \overline{a}$  and to A = a

## Path-Specific Effect



Measure difference of causal influence when A is set to  $A = \overline{a}$  and to A = a along  $A \rightarrow Y$ , whilst keeping A = a along  $A \rightarrow D \rightarrow Y$ 



## **Path-Specific Effect**

Different hard interventions on emerging links



$$\begin{split} p_{\mathrm{PO}}(Y_{\bar{a}}(M_{\bar{a}},L_a(M_{\bar{a}}))) \\ = \int_{C,M,L} p(Y|A=\bar{a},C,M,L) p(L|A=a,C,M) p(M|A=\bar{a},C) p(C) \end{split}$$

#### Different hard interventions on different causal paths



#### Need Structural Causal Models to Identify PSE



## **Path-Specific Effect**

Different hard interventions on different causal paths

#### **Structural Causal Models**

$$\begin{split} & \langle \mathcal{E}, V, F, p(\mathcal{E}) \rangle \\ & \mathcal{E} = \{ \mathcal{E}_{V_1}, \dots, \mathcal{E}_{V_K} \}, p(\mathcal{E}) = \prod_{k=1}^K p(\mathcal{E}_{V_k}) \text{ Exogenous (latent)} \\ & F = \{ f_{V_1}, \dots, f_{V_K} \}, V = \{ V_1, \dots, V_K \} \text{ Endogenous} \\ & V_k = f_{V_k}(\text{pa}(V_k), \mathcal{E}_{V_k}) \end{split}$$



 $L = f_L(A, C, M, \mathcal{E}_L)$  Intervention on L = replace the function with value l

$$A \to M \to L \to Y \qquad A = \bar{a}$$

$$A \to M \to Y \qquad A = \bar{a}$$

$$A \to M \to Y \qquad A = \bar{a}$$

$$A \to Y \qquad A = \bar{a}$$

$$A \to Y \qquad A = \bar{a}$$

$$A \sim \text{Bern}(\pi), C = \epsilon_c$$

$$M = \theta^m + \theta^m_a A + \theta^m_c C + \epsilon_m \qquad \epsilon_c, \epsilon_m, \epsilon_l, \epsilon_y \sim \mathcal{N}(0, 1)$$

$$L = \theta^l + \theta^l_a A + \theta^l_c C + \theta^l_m M + \epsilon_l$$

$$Y = \theta^y + \theta^y_a A + \theta^y_c C + \theta^y_m M + \theta^y_l L + \epsilon_y$$

$$A \to Y \qquad A \to M \to Y \qquad A \to M \to L \to Y A = \bar{a}$$

$$A \to L \to Y \qquad A = \bar{a}$$

$$\begin{split} M_{a} &= \theta^{m} + \theta^{m}_{a}a + \theta^{m}_{c}C + \epsilon_{m} \\ M_{\bar{a}} &= \theta^{m} + \theta^{m}_{a}\bar{a} + \theta^{m}_{c}C + \epsilon_{m} \\ L_{a}(M_{\bar{a}}) &= \theta^{l} + \theta^{l}_{a}a + \theta^{l}_{c}C + \theta^{l}_{m}M_{\bar{a}} + \epsilon_{l} \\ Y_{\bar{a}}(M_{a}, L_{a}(M_{\bar{a}})) &= \theta^{y} + \theta^{y}_{a}\bar{a} + \theta^{y}_{c}C + \theta^{y}_{m}M_{a} + \theta^{y}_{l}L_{a}(M_{\bar{a}}) + \epsilon_{y} \end{split}$$



## **Path-Specific Effect**

Different hard interventions on different causal paths

 $M_{a} = \theta^{m} + \theta^{m}_{a}a + \theta^{m}_{c}C + \epsilon_{m}$   $M_{\bar{a}} = \theta^{m} + \theta^{m}_{a}\bar{a} + \theta^{m}_{c}C + \epsilon_{m}$   $L_{a}(M_{\bar{a}}) = \theta^{l} + \theta^{l}_{a}a + \theta^{l}_{c}C + \theta^{l}_{m}M_{\bar{a}} + \epsilon_{l}$   $Y_{\bar{a}}(M_{a}, L_{a}(M_{\bar{a}})) = \theta^{y} + \theta^{y}_{a}\bar{a} + \theta^{y}_{c}C + \theta^{y}_{m}M_{a} + \theta^{y}_{l}L_{a}(M_{\bar{a}}) + \epsilon_{y}$ 

$$\begin{split} & \mathbb{E}_{p_{\text{PO}}(Y_{\bar{a}}(M_{a},L_{a}(M_{\bar{a}})))}[Y_{\bar{a}}(M_{a},L_{a}(M_{\bar{a}}))] \\ &= \theta^{y} + \theta^{y}_{a}\bar{a} + \theta^{y}_{m}(\theta^{m} + \theta^{m}_{a}a) + \theta^{y}_{l}(\theta^{l} + \theta^{l}_{a}a + \theta^{l}_{m}(\theta^{m} + \theta^{m}_{a}\bar{a})) \\ & \mathbb{E}_{p_{\text{PO}}(Y_{a})}[Y_{a}] \\ &= \theta^{y} + \theta^{y}_{a}a + \theta^{y}_{m}(\theta^{m} + \theta^{m}_{a}a) + \theta^{y}_{l}(\theta^{l} + \theta^{l}_{a}a + \theta^{l}_{m}(\theta^{m} + \theta^{m}_{a}a)) \end{split}$$

$$\begin{split} \mathrm{PSE}_{a\bar{a}} &= \mathbb{E}_{p_{\mathrm{PO}}(Y_{\bar{a}}(M_a, L_a(M_{\bar{a}})))}[Y_{\bar{a}}(M_a, L_a(M_{\bar{a}}))] - \mathbb{E}_{p_{\mathrm{PO}}(Y_a)}[Y_a] \\ &= (\theta_a^y + \theta_l^y \theta_m^l \theta_a^m)(\bar{a} - a) \end{split}$$



- Causal effect along a path = product of coefficients
- Causal effect along a set of paths = sum of causal effects along all paths

## Path-Specific Counterfactual PO



 $\begin{aligned} Y_{\bar{a}}(M_a, L_a(M_{\bar{a}})) \\ p_{\text{PO}}(Y_{\bar{a}}(M_a, L_a(M_{\bar{a}}))) \end{aligned}$ 

Observation for individual/unit n  $o^n = \{a^n = a, c^n, m^n, l^n, y^n\}$  $p_{\text{PO}}(Y_{\bar{a}}(M_a, L_a(M_{\bar{a}})|o^n)$   $\epsilon^n = \{\epsilon^n_c, \epsilon^n_m, \epsilon^n_l, \epsilon^n_y\}$ Specific realization of exogenous variables for individual/unit *n* 

$$\begin{split} c^n &= \epsilon^n_c \\ m^n &= \theta^m + \theta^m_a a + \theta^m_c c^n + \epsilon^n_m \\ l^n &= \theta^l + \theta^l_a a + \theta^l_c c^n + \theta^l_m m^n + \epsilon^n_l \\ y^m &= \theta^y + \theta^y_a a + \theta^y_c c^n + \theta^y_m m^n + \theta^y_l l^n + \epsilon^n_y \end{split}$$

Abduction Infer randomness from observation

$$\begin{split} & \epsilon_c^n = c^n \\ & \epsilon_m^n = m^n - (\theta^m + \theta_a^m a + \theta_c^m c^n) \\ & \epsilon_l^n = l^n - (\theta^l + \theta_a^l a + \theta_c^l c^n + \theta_m^l m^n) \\ & \epsilon_y^n = (y^m - (\theta^y + \theta_a^y a + \theta_c^y c^n + \theta_m^y m^n + \theta_l^y l^n)) \end{split}$$

$$p_{\text{PO}}(Y_{\bar{a}}(M_a, L_a(M_{\bar{a}})|o^n) = p_{\text{PO}}(Y_{\bar{a}}(M_a, L_a(M_{\bar{a}})|a, \epsilon^n)$$

$$\begin{split} M_a|\{a,\epsilon^n\} &= \theta^m + \theta^m_a a + \theta^m_c \epsilon^n_c + \epsilon^n_m = m^n \\ M_{\bar{a}}|\{a,\epsilon^n\} &= \theta^m + \theta^m_a \bar{a} + \theta^m_c \epsilon^n_c + \epsilon^n_m = m^n + \theta^m_a (\bar{a} - a) \\ L_a(M_{\bar{a}})|\{a,\epsilon^n\} &= \theta^l + \theta^l_a a + \theta^l_c \epsilon^n_c + \theta^l_m (m^n + \theta^m_a (\bar{a} - a)) + \epsilon^n_l = l^n + \theta^l_m \theta^m_a (\bar{a} - a) \\ Y_{\bar{a}}(M_a, L_a(M_{\bar{a}}))|\{a,\epsilon^n\} &= \theta^y + \theta^y_a \bar{a} + \theta^y_c \epsilon^n_c + \theta^y_m m^n + \theta^y_l (l^n + \theta^l_m \theta^m_a (\bar{a} - a)) + \epsilon^n_y \\ &= y^n + (\theta^y_a + \theta^y_l \theta^l_m \theta^m_a)(\bar{a} - a) \end{split}$$





## **CBNs and CI for Measuring Unfairness Underlying a Dataset**





## College Admission Example Measuring Unfairness









## Path-Specific Fairness: Measure Influence of A on Y along $A \rightarrow Y$ and $A \rightarrow D \rightarrow Y$ at the Population Level



A=a indicates female applicant  $A=\overline{a}$  indicates male applicant

Path-specific potential outcome  $Y_{\bar{a}}(Q_a, D_{\bar{a}})$ 

 $PSE_{a\bar{a}} = \mathbb{E}_{p_{PO}(Y_{\bar{a}}(Q_a, D_{\bar{a}}))}[Y_{\bar{a}}(Q_a, D_{\bar{a}})] - \mathbb{E}_{p_{PO}(Y_a)}[Y_a]$ 

Measures difference in expectation when setting the value of A to male along  $A \rightarrow Y$  and  $A \rightarrow D \rightarrow Y$  and to female along  $A \rightarrow Q \rightarrow Y$  wrt to setting it to female along all causal paths



## Path-Specific Counterfactual Fairness: Measure Influence of A on Y along $A \rightarrow Y$ and $A \rightarrow D \rightarrow Y$ for a Specific Individual



A=a indicates female applicant  $A=\overline{a}$  indicates male applicant

By conditioning  $Y_{\bar{a}}(Q_a, D_{\bar{a}})$  on observation

from a female individual  $\{a^n = a, q^n, d^n, y^n\}$  we answer the counterfactual question of whether the individual would have been admitted had she being male only along  $A \rightarrow Y$  and  $A \rightarrow D \rightarrow Y$ 

$$p_{\rm PO}(Y_{\bar{a}}(Q_a, D_{\bar{a}})|a^n = a, q^n, d^n, y^n)$$



## Path-Specific Counterfactual Fairness: Measure Influence of A on Y along $A \rightarrow Y$ and $A \rightarrow D \rightarrow Y$ for a Specific Individual



A=a indicates female applicant  $A=\overline{a}$  indicates male applicant

$$\begin{split} A &\sim \operatorname{Bern}(\pi) \\ Q &= \theta^q + \theta^q_a A + \epsilon_q \\ D &= \theta^d + \theta^d_a A + \epsilon_d \\ Y &= \theta^y + \theta^y_a A + \theta^y_q Q + \theta^y_d D + \epsilon_y \end{split}$$

$$Q_a|o^n = \theta^q + \theta^q_a a + \epsilon^n_q$$
  

$$D_{\bar{a}}|o^n = \theta^d + \theta^d_a \bar{a} + \epsilon^n_d$$
  

$$Y_{\bar{a}}(Q_a, D_{\bar{a}})|o^n = \theta^y + \theta^y_a \bar{a} + \theta^y_q q^n + \theta^y_d (d^n + \theta^d_a (\bar{a} - a)) + \epsilon^n_y$$

Useful viewpoint Modification of  $y^n$  obtained by modifying part of the observation for which the value of A is different from the one observed



Gender

## CBNs and CI for Developing Fair Prediction Models

Path-Specific Counterfactual Fairness. S. Chiappa, 2019.

<u>A General Approach to Fairness with Optimal</u> <u>Transport. S. Chiappa, R. Jiang, T. Stepleton, A.</u> <u>Pacchiano, H. Jiang, J. Aslanides, 2020.</u> Graphical Conditions for Introduced Unfairness: Why Fair Labels Can Yield Unfair Predictions. C. Ashurst, R. Carey, S. Chiappa, T. Everitt, 2022



## **Difference with Measuring Unfairness in Datasets**

## 1

2

Need to not absorb bias in data in the model

# Need to maintain accuracy

## 3

Training and deployment settings are different



## Path-Specific Counterfactually Fair Predictor



- Training: Objective function that estimates *p*
- **Deployment:** Observe  $\{a^n = a, q^n, d^n\}$  and assign outcome as mean of

$$\hat{p}_{\text{PO}}(Y_{\bar{a}}(Q_a, D_{\bar{a}})|a^n = a, q^n, d^n)$$

Need to enforce independence in the latent space

## **Path-Specific Counterfactually Fair Predictor**



$$\begin{split} & A \sim \operatorname{Bern}(\pi) \\ & Q \sim p_{\theta}^q(\cdot|A) \\ & H_d \sim p_{\theta}^h(\cdot), \, D \sim p_{\theta}^d(\cdot|A, H_d) \\ & Y \sim p_{\theta}^y(\cdot|A, Q, D) \end{split}$$

Represent conditionals using neural networks

$$\begin{split} &Q_a \sim p_{\theta}^q(\cdot | A = a) \\ &H_d \sim p_{\theta}^h(\cdot), \, D_{\bar{a}} \sim p_{\theta}^d(\cdot | A = \bar{a}, H_d) \\ &Y_{\bar{a}}(Q_a, D_{\bar{a}}) \sim p_{\theta}^y(\cdot | A = \bar{a}, Q_a, D_{\bar{a}}) \end{split}$$

#### Deployment



$$\{a^n=a,q^n,d^n\}$$
  
 $h^{n,i}_d \; i=1,\ldots,I \, ext{ from } p_ heta(H_d|A,Q,D) \, ext{ Abduction}$ 

$$\begin{split} & d_{\bar{a}}^{n,i} \sim p_{\theta}^{d}(\cdot | \bar{a}, h_{d}^{n,i}) & \text{Action-Prediction} \\ & y_{\bar{a}}^{n} = \frac{1}{I} \sum_{i=1}^{I} \mathbb{E}_{p_{\theta}^{y}(Y_{\bar{a}}(Q_{a}, D_{\bar{a}}) | \bar{a}, q^{n}, d_{\bar{a}}^{n,i})} [Y_{\bar{a}}(Q_{a}, D_{\bar{a}})] \end{split}$$



## A-Independent Representation $H_d$ with Variational Autoencoder

Intractable  $p_{\theta}(H_d|A,Q,D)$ 

Variational approximation

 $\epsilon$ 

$$\begin{split} p_{\theta}(H_d|A,Q,D) &\approx q_{\phi}(H_d|A,Q,D) = \mathcal{N}(\mu_{\phi},\sigma_{\phi}) \\ H_d &= \mu_{\phi} + \sigma_{\phi}\epsilon, \ q_{\epsilon} = \mathcal{N}(0,1) \end{split}$$

$$h_{d}^{n,i} = \mu_{\phi}^{n} + \sigma_{\phi}^{n} \epsilon^{i}$$

$$\mu_{\phi}^{n} \quad \sigma_{\phi}^{n}$$
Encoder
$$a^{n}, q^{n}, d^{n}$$

Maximise lower bound on the marginal log-likelihood with a penalty term for enforcing independence on *A* 

$$\mathcal{F}_{\theta,\phi} = \mathbb{E}_{q_{\phi}(H_d|A,Q,D)}[\log p_{\theta}(A,Q,D,Y,H_d)] \\ - \mathbb{E}_{q_{\phi}(H_d|A,Q,D)}[\log q_{\phi}(H_d|A,Q,D)] \\ - \beta \text{MMD}(a,\bar{a})$$

weighting factor that determines the degree of independence

Empirical distributions of  $h_d^{n,i}$  for male and female individuals



CBNs and CI for ML Fairness Formalization Unified Framework Complex Unfairness Scenarios

Unfairness underlying Data

- CBNs for describing different possible unfairness scenarios underlying data
- CBNs for measuring unfairness underlying data

#### Fair Prediction Models

- CBNs for developing fair prediction models
- Causal influence diagrams for auditing fair prediction models



## CBNs and CI for ML Fairness

#### Causal Decision Making

#### **Beyond Prediction Models**

• Design optimized policies under fairness constraint

Pragmatic Fairness: Developing Policies with Outcome Disparity Control. L. Gultchin, S. Guo, A. Malek, S. Chiappa, R. Silva, 2022.





#### Goal

Design a new decision making system that specifies how to select actions *D* that maximize downstream outcome *Y* subject to some fairness constraints

#### Extended view wrt designing fair prediction systems

Distinguish between our choice of policy or action allocation (D), and the downstream effect we *truly* aim for (Y)

#### Granting Loans/Offering College Admission Examples

We would like to be fair with respect to desirable downstream outcomes, such as repayment or academic and economic success, and not just the very allocation of the action

We are not aiming to match a distribution of historical outcomes but to optimize a future expected utility



#### Setting

- A and X are allowed to be associated and potentially directly influence Y
- *D* can only indirectly control this influence

This setting formalizes a situation in which control for association of *A* and *Y* can only be achieved to some extent through a predefined set of available actions, i.e., a given action space. The level to which we can minimize such unfair impact therefore depends on the choice of the action space

#### **Outreach Campaign Example**

Company looking to change its outreach campaign to mitigate imbalances in the demographics A and potentially associated level of experience X of job applicants Y

The company cannot control factors such as cultural preference among applicants for industry sectors, but can induce modifications such as focusing recruiting efforts in events organized by minority groups in relevant conferences, thus optimizing for the things we can control



#### Setting

Observational-data regime: learn optimized policy based on historical data collected using a baseline policy  $\sigma_{\emptyset}$ , representing the action allocation that was in place during the collection of the data, i.e. the "status-quo" pre-optimization.

Reflects considerations such as the cost, ethical constraints, or difficulty of collecting interventional data.

$$\begin{split} \mathcal{D} &= \{a^i, x^i, d^i, y^i\}_{i=1}^N \\ (a^i, x^i, d^i, y^i) \sim p(A, X, D, Y; \sigma_{\emptyset}) \end{split}$$





#### Formalization

 $\mu^{\cdot}$ 

Find policy  $\mu^D_{\sigma_D}(a, x) = \mathbb{E}[D|a, x; \sigma_D]$  that maximizes utility  $\mathbb{E}[Y_{\sigma_D}]$  while controlling for disparity

#### **Moderation Breaking Constraint**

Remove influence of A on  $\mathbb{E}[Y_{\sigma_D}]$  as much as possible what we can control: the allocation of D as determined by  $p(D|A, X; \sigma_D)$ .

Express a possible modulation of the effect of the action on the outcome by the sensitive attribute via the decomposition

$$\begin{aligned} F(a, x, d) &:= \mathbb{E}[Y|a, x, d] \\ &= f(a, x) + g(a, x, d) + h(x, d) \end{aligned}$$



$$\begin{split} \mu_{\sigma_D}^Y(a,x) &:= \mathbb{E}[Y_{\sigma_D}|a,x] = \int_d \mu^Y(a,x,d) p(d|a,x;\sigma_D) \\ &= f(a,x) + g_{\sigma_D}(a,x) + h_{\sigma_D}(x) \end{split}$$

#### Objective

$$\begin{split} &\arg\max_{\sigma_D} \mathbb{E}[Y_{\sigma_D}] \quad \text{s.t.} \\ &(\mathbb{E}[g_{\sigma_D}(a,X) \mid A=a] - \mathbb{E}[g_{\sigma_D}(\bar{a},X) \mid A=\bar{a}])^2 \leq \epsilon, \forall a, \bar{a} \end{split}$$

Model  $\mu^{Y}(a, x, d)$  using a structured neural network that separates into the three components reflecting the decomposition



#### Phase I

Learn parameters of NN<sup>Y</sup>using

#### **Phase II**

 $\begin{array}{l} \text{Model} \ \mu^D_{\sigma_D}(a,x) = \mathbb{E}[D|a,x;\sigma_D] \\ \text{using a MLP neural network} \end{array}$ 



Teal blocks: parameter layers Light green blocks: fixed parameters. Purple diamonds: additive gates

Learn parameters of  $MLP_{\sigma_D}^D$  to learn the objective with

$$\mathbb{E}[Y_{\sigma_D}] \approx \frac{1}{N} \sum_{i=1}^{N} \mu_{\sigma_D}^Y(a^i, x^i)$$



DeepMind

# The end and thank you