Causality for Fairness and Explainability

Karima Makhlouf Ruta Binkyte-Sadauskiene Carlos Pinzon Catuscia Palamidessi **Sami Zhioua**









Fairness: Is the output fair with respect to individuals or subpopulations ? **Explainability:** How the output can be explained in terms of the input features ?

Statistical (Observational) notions of fairness



Purely associational explainability



Purely associational explainability





* Ribeiro, M. T., Singh, S., & Guestrin, C. (2016, August). "Why should i trust you?" Explaining the predictions of any classifier. 22nd ACM SIGKDD.

Purely associational explainability (Counterfactual)*

A *counterfactual* is a generated data point that is as close to the input data point as possible for which the model gives a different outcome.



 $\min_{c} d(\mathbf{x}, \mathbf{c})$
s.t. $f(\mathbf{c}) \neq f(\mathbf{x})$

* Sharma, S., Henderson, J., & Ghosh, J. (2019). Certifai: Counterfactual explanations for robustness, transparency, interpretability, and fairness of artificial intelligence models.

Statistical (Observational) notions of fairness



$$(\hat{Y} \mid A = 0) = P(\hat{Y} \mid A = 1)$$

Statistical Parity

E[S | Y = 1, A = 0] = E[S | Y = 1, A = 1]

$$P(Y = 1 | \hat{Y} = 1, A = 0) = P(Y = 1 | \hat{Y} = 1, A = 1)$$

Predictive Parity

 $P(Y = 1 \mid S = s, A = 0) = P(Y = 1 \mid S = s, A = 1) \quad \forall s \in [0, 1]$ Calibration

How strong is the effect of A on Y?



Why not P(Y|A)? ______ Selection bias TV = P(Y=1|A=1) - P(Y=1|A=0)

The illusion of correlation

"The correlation we observe is an illusion. An illusion we brought upon ourselves by choosing which events to include in our dataset and which to ignore."

Example 2:

Example 1:

Flip two coins 100 times, and write down the results <u>only when</u> at least one of

them comes up head

Notice the dependence: every time coin1 lands tail, coin2 lands head !

Coin 1	Coin 2
Head	head
Tail	head
head	tail
Tail	head
Head	head

Did you notice that among the people you date, the attractive ones are more likely to be jerks ?

You are dating from these:	Attractive Attractive	Jerk Nice
	_Not attractive	Nice
	Not attractive	Jerk

IUDEA PEARL

Why observable association is not reliable to establish the effect of a variable on another variable ?



Simpson's Paradox



Discrimation in favor of women

A = 0	Man
A = 1	Woman

T = 0	Flexible time job
T = 1	Non-flexible time job

_		
1	C	

Y=0

Y=1

Not hired

Hired

Statistical parity = 7/15 – 8/15 = **-1/15**

Discrimination against women

Causality

"The ability to learn causality is considered as a significant component of human-level intelligence and can serve as the foundation of AI" [Pearl 2018]

"The discovery of causal relationships from purely observational data is a fundamental problem in science" [Mooij 2016]

"Almost all of science is about identifying causal relations and the laws or regularities that govern them" [Glymour 2019]

[Pearl 2018] Judea Pearl. 2018. Theoretical impediments to machine learning with seven sparks from the causal revolution

[Mooij 2016] Mooij, J. M., Peters, J., Janzing, D., Zscheischler, J., & Schölkopf, B. (2016). Distinguishing cause from effect using observational data: methods and benchmarks..

[Glymour 2019] Glymour, Clark, Kun Zhang, and Peter Spirtes. "Review of causal discovery methods based on graphical models." (2019).



























In medical studies: select half of individuals randomly, and give them the treatment

In fairness problems: select half of candidates and *set* their gender to protected group (female).



 $P(y_a)$

TE = ACE = P(Y=1|do(A=1)) - P(Y=1|do(A=0))



Estimating P(Y|do(A=a)) from observed data

Definition 3.3.1 (The Backdoor Criterion) *Given an ordered pair of variables* (X, Y) *in a directed acyclic graph G, a set of variables Z satisfies the* backdoor criterion *relative to* (X, Y) *if no node in Z is a descendant of X, and Z blocks every path between X and Y that contains an arrow into X.*

If a set of variables Z satisfies the backdoor criterion for X and Y, then the causal effect of X on Y is given by the formula

$$P(Y = y | do(X = x)) = \sum_{z} P(Y = y | X = x, Z = z) P(Z = z)$$



Estimating P(Y|do(A=a)) from observed data

Definition 3.3.1 (The Backdoor Criterion) Given an ordered pair of variables (X, Y) in a directed acyclic graph G, a set of variables Z satisfies the backdoor criterion relative to (X, Y) if no node in Z is a descendant of X, and Z blocks every path between X and Y that contains

If a set of variables Z satisfies the backdoor criterion for X and Y, then the causal effect of X on Y is given by the formula

$$P(Y = y | do(X = x)) = \sum_{z} P(Y = y | X = x, Z = z) P(Z = z)$$



Estimating P(Y|do(A=a)) from observed data

Definition 3.4.1 (Front-Door) A set of variables Z is said to satisfy the front-door criterion relative to an ordered pair of variables (X, Y) if

- 1. Z intercepts all directed paths from X to Y.
- 2. There is no unblocked path from X to Z.
- 3. All backdoor paths from Z to Y are blocked by X.

Theorem 3.4.1 (Front-Door Adjustment) If *Z* satisfies the front-door criterion relative to (X, Y) and if P(x, z) > 0, then the causal effect of *X* on *Y* is identifiable and is given by the formula

$$P(y|do(x)) = \sum_{z} P(z|x) \sum_{x'} P(y|x', z) P(x')$$
(3.16)

Estimating P(Y|do(A=a)) from observed data















 $P(Y=y|do(A=a)) = P(y_a)$

Mediation Analysis



Discrimination ? It depends on Z

Causality and out-of-distribution (OOD) Learning

* Schölkopf, B., Locatello, F., Bauer, S., Ke, N. R., Kalchbrenner, N., Goyal, A., & Bengio, Y. (2021). Towards causal representation learning.

Current and future work

Causal discovery (structure learning) algorithms

- Does the shape of data (type, binarization, pre-processing, etc.) has an impact on the generated causal model/graph ?
- Do causal discovery algorithms provide confidence levels on the edges?
- What is the impact of confidence levels on the causal effects ?

Mediation analysis for explainability

Experiments on Causal graph generation Using Tetrad

1- Compas dataset

Compas

Before binarization (age, prior_counts)

After binarization (age, prior_counts)

Purely associational explainability

Based only on the CBN, how to tell if two variables are dependent/independent/conditionally-independent?

If they are conditionally independent, on which variables we should condition on ?

Definition 2.4.1 (*d*-separation) A path p is blocked by a set of nodes Z if and only if

- 1. *p* contains a chain of nodes $A \rightarrow B \rightarrow C$ or a fork $A \leftarrow B \rightarrow C$ such that the middle node B is in Z (i.e., B is conditioned on), or
- 2. *p* contains a collider $A \rightarrow B \leftarrow C$ such that the collision node B is not in Z, and no descendant of B is in Z.

If Z blocks every path between two nodes X and Y, then X and Y are d-separated, conditional on Z, and thus are independent conditional on Z.

d-separation and variables independence

Definition 2.4.1 (*d*-separation)^{*} A path p is blocked by a set of nodes Z if and only if

- 1. *p* contains a chain of nodes $A \rightarrow B \rightarrow C$ or a fork $A \leftarrow B \rightarrow C$ such that the middle node B is in Z (i.e., B is conditioned on), or
- 2. *p* contains a collider $A \rightarrow B \leftarrow C$ such that the collision node B is not in Z, and no descendant of B is in Z.
- If Z blocks every path between two nodes X and Y, then X and Y are d-separated, conditional on Z, and thus are independent conditional on Z.

Total (causal) Effect:

TE = ACE = P(Y=1|do(A=1)) - P(Y=1|do(A=0))

 $P(y_{A\leftarrow a})$ $P(y_{a\leftarrow a})$ Two random variables X and Y are called independent, if for each values of X and Y, x and y,

-
$$P(X = x, Y = y) = P(X = x) \cdot P(Y = y)$$
 or

$$- P(X = x | Y = y) = P(X = x) \text{ or } P(Y = y | X = x) = P(Y = y)$$

- Denoted by $X \perp Y$
- Two random variables X and Y are called conditionally independent given Z, if for each values of (X, Y, Z), (x, y, z),
 - $P(X = x, Y = y | Z = z) = P(X = x | Z = z) \cdot P(Y = y | Z = z)$ or

-
$$P(X = x | Y = y, Z = z) = P(X = x | Z = z)$$
 or

- P(Y = y | X = x, Z = z) = P(Y = y | Z = z)
- Denoted by $X \perp Y | Z$
- Note: conditional independence neither implies nor is implied by independence.

Random variables A and Y are independent

 $\mathsf{P}(\mathsf{Y} | \mathsf{A}) = \mathsf{P}(\mathsf{Y})$

Our belief Y remains unchanged upon learning A.

Random variables A and Y are conditionally independent given C

P(Y | A,C) = P(Y | C)

Once we know C, our belief Y remains unchanged upon learning A.

A and Y are independent in the new dataset created by filtering on C.

Causal Bayesian Network (CBN)

represents a set of variables and their conditional dependencies via a directed acyclic graph (DAG).

Suppose we have a distribution P defined on n discrete variables which we may order arbitrarily as X_1, X_2, ... X_n. The chain rule allows to decompose the joint distribution P as:

$$P(x_1, \dots, x_n) = \prod_j P(x_j \mid x_1, \dots, x_{j-1}) = \prod_i P(x_i \mid pa_i).$$

Definition 1.2.1 (Markovian Parents) Let $V = \{X_1, ..., X_n\}$ be an ordered set of variables, and let P(v) be the joint probability distribution on these variables. A set of variables PA_j is said to be Markovian parents of X_j if PA_j is a minimal set of predecessors of X_j that renders X_j independent of all its other predecessors. In other words, PA_j is any subset of $\{X_1, ..., X_{j-1}\}$ satisfying

 $P(x_j \mid pa_j) = P(x_j \mid x_1, \dots, x_{j-1}) \qquad P(x_1, x_2, x_3, x_4, x_5) = P(x_1)P(x_2 \mid x_1)P(x_3 \mid x_1)P(x_4 \mid x_2, x_3)P(x_5 \mid x_4).$

SEASON

 X_2

WET

SLIPPERY

X4

RAIN

SPRINKLER (X_3)

and such that no proper subset of PA_j satisfies (1.32).⁵

Definition 1.2.2 (Markov Compatibility)

If a probability function P admits the factorization of (1.33) relative to DAG G, we say that G represents P, that G and P are compatible, or that P is Markov relative to G^{6} .

Based only on the CBN, how to tell if two variables are dependent/independent/conditionally-independent ?

If they are conditionally independent, on which variables we should condition on ?

Definition 2.4.1 (*d*-separation) A path p is blocked by a set of nodes Z if and only if

- 1. *p* contains a chain of nodes $A \rightarrow B \rightarrow C$ or a fork $A \leftarrow B \rightarrow C$ such that the middle node B is in Z (i.e., B is conditioned on), or
- 2. *p* contains a collider $A \rightarrow B \leftarrow C$ such that the collision node B is not in Z, and no descendant of B is in Z.

If Z blocks every path between two nodes X and Y, then X and Y are d-separated, conditional on Z, and thus are independent conditional on Z.

• Example (blocking of paths)

$$X \longrightarrow Z \longrightarrow U \longrightarrow Y$$

- Path from X to Y is blocked by conditioning on $\{U\}$ or $\{Z\}$ or both $\{U, Z\}$
- Example (unblocking of paths)

- Path from X to Y is blocked by \emptyset or $\{U\}$
- Unblocked by conditioning on {Z} or {W} or both {Z, W}

• Example (*d*-separation)

- We have following *d*-separation relations
 - $(X \perp Y | Z)_G, (X \perp Y | U)_G, (X \perp Y | ZU)_G$
 - $(X \perp Y | ZW)_G, (X \perp Y | UW)_G, (X \perp Y | ZUW)_G$
 - $-(X \perp Y | VZUW)_G$
- However we do NOT have
 - $(X \perp Y | VZU)_G$

Based only on the CBN, how to tell if two variables are dependent/independent/conditionally-independent?

If they are conditionally independent, on which variables we should condition on ?

Definition 2.4.1 (*d*-separation) A path p is blocked by a set of nodes Z if and only if

- 1. *p* contains a chain of nodes $A \rightarrow B \rightarrow C$ or a fork $A \leftarrow B \rightarrow C$ such that the middle node B is in Z (i.e., B is conditioned on), or
- 2. *p* contains a collider $A \rightarrow B \leftarrow C$ such that the collision node B is not in Z, and no descendant of B is in Z.

If Z blocks every path between two nodes X and Y, then X and Y are d-separated, conditional on Z, and thus are independent conditional on Z.

Structural Causal Model (SCM)

A causal model is triple $\mathcal{M} = \langle U, V, F \rangle$, where

- U is a set of exogenous (hidden) variables whose values are determined by factors outside the model;
- V = {X₁, ..., X_i, ...} is a set of endogenous (observed) variables whose values are determined by factors within the model;
- $F = \{f_1, \dots, f_i, \dots\}$ is a set of deterministic functions where each f_i is a mapping from $U \times (V \setminus X_i)$ to X_i . Symbolically, f_i can be written as

$$x_i = f_i(\boldsymbol{p}\boldsymbol{a}_i, \boldsymbol{u}_i)$$

where pa_i is a realization of X_i 's parents in V, i.e., $Pa_i \subseteq V$, and u_i is a realization of X_i 's parents in U, i.e., $U_i \subseteq U$.

Structural Causal Model (SCM)

Example:

Assume U_I, U_H, U_W, U_E are mutually independent.

How strong is the causal dependence of Y on A (causal effect of A on Y)?

Example 1:

Flip two coins 100 times, and write down the results only when at least one of

them comes up head

Notice the dependence:
every time coin1 lands
tail, coin2 lands head !

Coin 1	Coin 2
Head	head
Tail	head
head	tail
Tail	head
Head	head

The illusion of correlation

"The correlation we observe is an illusion. An illusion we brought upon ourselves by choosing which events to include in our dataset and which to ignore."

Example 2:

Did you notice that among the people you date, the attractive ones are more likely to be jerks ?

You are dating from these:	Attractive Attractive	Jerk Nice
	Not attractive	Nice
	Not attractive	Jerk

JUDEA PEARL WINNER OF THE TURING AWARD AND DANA MACKENZIE