

Survey on Causal-based Machine Learning Fairness Notions

Karima Makhlouf
Université du Québec à Montréal
Montréal, Canada
makhlouf.karima@courrier.uqam.ca

Sami Zhioua
Higher Colleges of Technology
Dubai, UAE
szhioua@hct.ac.ae

Catuscia Palamidessi
Inria, École Polytechnique, IPP
Paris, France
catuscia@lix.polytechnique.fr

ABSTRACT

Addressing the problem of fairness is crucial to safely use machine learning algorithms to support decisions with a critical impact on people’s lives such as job hiring, child maltreatment, disease diagnosis, loan granting, etc. Several notions of fairness have been defined and examined in the past decade, such as, statistical parity and equalized odds. The most recent fairness notions, however, are causal-based and reflect the now widely accepted idea that using causality is necessary to appropriately address the problem of fairness. This paper examines an exhaustive list of causal-based fairness notions, in particular their applicability in real-world scenarios. As the majority of causal-based fairness notions are defined in terms of non-observable quantities (e.g. interventions and counterfactuals), their applicability depends heavily on the identifiability of those quantities from observational data. In this paper, we compile the most relevant identifiability criteria for the problem of fairness from the extensive literature on identifiability theory. These criteria are then used to decide about the applicability of causal-based fairness notions in concrete discrimination scenarios.

KEYWORDS

Fairness, machine learning, causal-based, identifiability.

1 INTRODUCTION

Fairness is getting an increasing attention when designing automated decision-making systems as such systems can lead to discrimination against individuals or sub-populations. For the affected individuals, it is indeed very unpleasant, and sometimes devastating, to experience outcomes perceived as unfair and caused by factors outside their control. Hence, when choosing policies and designing systems that will impact people’s lives (e.g. employment, college admission, credit, insurance, etc.), decision should be made independent from factors outside an individual’s control, such as gender, ethnicity, or place of birth.

The first endeavor of the research community to achieve fairness consisted in developing correlation/association-based notions, including statistical parity [6], equalized odds [9], predictive parity [5], and calibration [5] which primarily focus on discovering the discrepancy of statistical metrics between individuals or sub-populations. The problem of these notions is that they lack in reasoning about cause-effect relations between attributes. Moreover, the research community is realizing that fairness cannot be well assessed based only on mere correlation or association [9, 10, 13, 15, 37, 38]. A canonical example is Simpson’s paradox [21], where the statistical conclusions drawn from the sub-populations differ from that from the whole population. On the other hand, discrimination claims usually require plaintiffs to demonstrate a causal

connection between the challenged decision (e.g. hiring, firing, admission) and the sensitive feature (e.g. gender, race). It is then necessary to investigate the causal relationship between the sensitive attribute and the decision rather than the associated relationship. Various causal-based fairness notions have been recently proposed to tackle the problem of automated discrimination through causal inference lenses.

Causal-based fairness notions differ from the previous statistical fairness approaches in that they are not totally based on data but consider additional knowledge about the structure of the world, in the form of a causal model. This additional knowledge helps us understand how data is generated in the first place and how changes in variables propagate in a system. Most of these fairness notions are defined in terms of non-observable quantities such as interventions (to simulate random experiments) and counterfactuals (which consider other hypothetical worlds, in addition to the actual world). Such quantities cannot be always uniquely computed from observational data. This problem is known as identifiability and hinders significantly the usefulness of causal-based notions in practical scenarios.

There is a large body of work in the literature that addresses the problem of identifiability, typically through complex graphical criteria [1, 8, 11, 17, 20, 21, 28–30, 33, 35].

In this paper, we first provide an exhaustive list of causal-based fairness notions. Then, we compile the most relevant identifiability criteria for the specific problem of discrimination discovery. These criteria are grouped into intervention, counterfactual, direct and indirect effect, and path-specific effect identifiability results. By placing the fairness notions in the three rungs causation ladder of Pearl [22] and using the identifiability criteria, we propose a guideline for applying causal-based fairness notions in real-scenarios.

2 THE NEED FOR CAUSALITY: AN EXAMPLE

Consider the hypothetical example¹ of an automated system for deciding whether to fire a teacher at the end of the academic year. Deployed teacher evaluation systems have been suspected of bias in the past. For example, IMPACT is a teacher evaluation system used in the city of Washington DC and have been found to be unfair against teachers from minority groups [19, 23, 24]. Assume that the system takes as input two features, namely, the location of the school where the teacher is working (C) and the initial² average level of the students in her class (A). The outcome is whether to fire the teacher (Y). Assume also that all 3 variables are binary with the following values: if the school is located in a high-income neighborhood, $C = 1$, otherwise (the school is located in a low-income neighborhood), $C = 0$. If the initial average score for the

¹Inspired by the prior convictions example in [18].

²At the beginning of the academic year.

students assigned to the teacher is high, $A = 1$, otherwise (initial level is low), $A = 0$. Firing a teacher corresponds to $Y = 1$, while retaining her corresponds to $Y = 0$. The level of students in a given class can be influenced by several variables, but in this example, assume that it is only influenced by the location of the school; students in high-income neighborhoods are more advantaged and typically perform better in school.

Assume now that the automated decision system is suspected to be biased by the initial level of students assigned to the teacher. That is, it is claimed that the system will more likely fire teachers who have been assigned classes with low level students at the beginning of the academic year which is clearly unfair. The sensitive attribute in this case is the initial level of students assigned to the teacher (A). For concreteness, consider the prediction system that yields the following conditional probabilities:

$$\begin{aligned} P(Y = 1 \mid A = 1, C = 0) &= 0.02 & P(A = 1 \mid C = 0) &= 0.2 \\ P(Y = 1 \mid A = 1, C = 1) &= 0.0675 & P(A = 1 \mid C = 1) &= 0.8 \\ P(Y = 1 \mid A = 0, C = 0) &= 0.01 & P(A = 0 \mid C = 0) &= 0.8 \\ P(Y = 1 \mid A = 0, C = 1) &= 0.25 & P(A = 0 \mid C = 1) &= 0.2 \end{aligned}$$

and that the dataset is collected from a population where schools are located with equal proportions in high-income and low-income neighborhood, that is, $P(C = 1) = P(C = 0) = 0.5$. Assume also that the proportion of classes with a low initial average level of students is the same as the one with high average initial level of students, that is, $P(A = 1) = P(A = 0) = 0.5$. To keep the scenario simple, assume that the level of students A does not depend on any other feature except C and that the firing decision Y depends only on A and C .

A simple approach to check the fairness of the firing decision Y with respect to the sensitive attribute A is to contrast the conditional probabilities: $P(Y = 1 \mid A = 0)$ and $P(Y = 1 \mid A = 1)$ which quantify, respectively, the likelihood of firing a teacher given that she is assigned students with an initial low level and the likelihood of firing a teacher given that she is assigned students with an initial high level class. Such probabilities can be computed as follows:

$$\begin{aligned} P(Y = 1 \mid A = a) &= \sum_{c \in \{0,1\}} P(Y = 1 \mid A = a, C = c) \\ &\quad \times P(C = c \mid A = a) \end{aligned}$$

Hence,

$$P(Y = 1 \mid A = 1) = 0.02 \times 0.2 + 0.0675 \times 0.8 = 0.058$$

$$P(Y = 1 \mid A = 0) = 0.01 \times 0.8 + 0.25 \times 0.2 = 0.058$$

As the values are equal, the rates of firing between teachers who were assigned low level students and high level students appear to be equal and hence no discrimination is detected³. This conclusion is flawed because it doesn't consider the mechanism by which the data is generated. In particular, the location of the school in which the teacher is working influences both the initial level of students assigned to her as well as the decision to fire or retain her. In such situations, if the data is collected in a certain way, the above conditional probabilities may end up being equal despite the presence of bias. The above example is one of such cases since the conditional probability $P(A = 0 \mid C = 0) = 0.8$ is quite extreme.

³This corresponds to statistical parity.

Notice that teachers in low-income neighborhoods are more likely to be assigned classes where the average level of students is low, but not to that extent (80% of classes in low-income neighborhood have an average low level). In general, any statistical fairness notion which relies solely on correlation between variables, will fail to detect such bias.

To avoid such misleading conclusions, the causal relationships between variables should be considered. Figure 1 illustrates the causal relations between the three variables of the above example where the location of the school C is a confounder. Based on such causal graph, a firing decision system is fair if is as likely to fire teachers in the following two hypothetical cases: (1) when *all teachers in the population are assigned students of low level on average*, and (2) when all teachers in the population are assigned students of high level on average. This is achieved using intervention ($do()$ operator)⁴ and allows to break the problematic dependence between A and C . The probabilities of firing a teacher in these two hypothetical cases are expressed as $P(Y_{A=0} = 1) = P(Y = 1 \mid do(A = 0))$ and $P(Y_{A=1} = 1) = P(Y = 1 \mid do(A = 1))$ respectively. In this simple graph, and assuming no other variable is used in the prediction, these probabilities can be computed as follows:

$$P(Y_{A=a} = 1) = \sum_{c \in \{0,1\}} P(Y = 1 \mid A = a, C = c) \times P(C = c)$$

Hence,

$$P(Y_{A=1} = 1) = 0.02 \times 0.5 + 0.0675 \times 0.5 = 0.0437$$

$$P(Y_{A=0} = 1) = 0.01 \times 0.5 + 0.25 \times 0.5 = 0.13$$

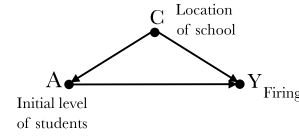


Figure 1: Causal graph of the firing example.

The values confirm the existence of a bias against teachers which are assigned students of low level on average. These teachers are more likely to belong to minority groups.

Another example that is used extensively in the literature to illustrate the importance of going beyond mere correlation to appropriately assess fairness is the gender bias in Berkeley admission [4, 16].

3 PRELIMINARIES AND NOTATION

Variables are denoted by capital letters. In particular, A is used for the sensitive variable (e.g., gender, race, age) and Y is used for the outcome of the automated decision system (e.g., hiring, admission, releasing on parole). Small letters denote specific values of variables (e.g., $A = a'$, $W = w$). Bold capital and small letters denote a set of variables and a set of values, respectively.

A structural causal model [21] is a tuple $M = \langle U, V, F, P(U) \rangle$ where:

⁴Intervention and the $do()$ operator will be explained further in Section 5.1.

- U is a set of exogenous variables which cannot be observed or experimented on but constitute the background knowledge behind the model.
- V is a set of observable variables which can be experimented on.
- F is a set of structural functions where each f_i is mapping $U \cup V \rightarrow V \setminus \{V_i\}$ which represents the process by which variable V_i changes in response to other variables in $U \cup V$.
- $P(u)$ is a probability distribution over the unobservable variables U .

Causal assumptions between variables are captured by a causal diagram G which is a directed acyclic graph (DAG) where vertices represent variables and directed edges represent functional relationships between the variables. Directed edges can have two interpretations. A probabilistic interpretation where the edge represents a dependency among the variables such that the direction of the edge is irrelevant. A causal interpretation where the edge represents a causal influence between the corresponding variables such that the direction of the edge matters. Unobserved variables U , which are typically not represented in the causal diagram, can be either mutually independent (Markovian model) or dependent from each others. In case the unobserved variables can be dependent and each $U_i \in U$ is used in at most two functions in F , the model is called semi-Markovian. In causal diagrams of semi-Markovian models, dependent unobservable variables (unobserved confounders) are represented by a dotted bi-directed edge between observable variables. Figure 2 shows causal graphs of Markovian and semi-Markovian models, respectively. An intervention, noted $do(V =$

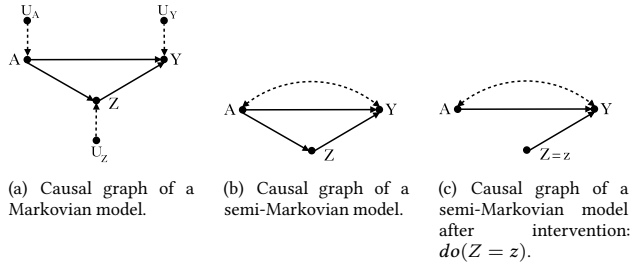


Figure 2

v), is a manipulation of the model that consists in fixing the value of a variable (or a set of variables) to a specific value regardless of the corresponding function f_v . Graphically, it consists in discarding all edges incident to the vertex corresponding to variable V . Figure 2(c) shows the causal diagram of the manipulated model after intervention $do(Z = z)$ denoted $M_{Z=z}$ or M_z for short. The intervention $do(V = v)$ induces a different distribution on the other variables. For example, in Figure 2(c), $do(Z = z)$ results in a different distribution on Y , namely, $P(Y|do(Z = z))$. Intuitively, while $P(Y|Z = z)$ reflects the population distribution of Y among individuals whose Z value is z , $P(Y|do(Z = z))$ reflects the population distribution of Y if *everyone in the population* had their Z value fixed at z . The obtained distribution $P(Y|do(Z = z))$ can be considered as a *counterfactual* distribution since the intervention forces Z to take a value different from the one it would take in the actual world. Such counterfactual variable is noted $Y_{Z=z}$ or Y_z

for short⁵. $P(Y = y|do(Z = z)) = P(Y_{Z=z} = y) = P(Y_z = y) = P(y_z)$ is used to define the causal effect of z on Y . The term counterfactual quantity is used for expressions that involve explicitly multiple worlds. In Figure 2(b), consider the expression $P(y_{a'}|Y = y, A = a) = P(y_{a'}|y, a)$. Such expression involves two worlds: an observed world where $A = a$ and $Y = y$ and a counterfactual world where $Y = y$ and $A = a'$ and it reads “the probability of $Y = y$ had A been a' given that we observed $Y = y$ and $A = a$ ”. In the common example of job hiring, if A denotes race (a :white, a' :non-white) and Y denotes the hiring decision (y :hired, y' :not hired), $P(y_{a'}|y, a)$ reads “given that a white applicant has been hired, what is the probability that the same applicant is still being hired had he been non-white”. Nesting counterfactuals can produce complex expressions. For example, in the relatively simple model of Figure 2(b), $P(y_{a,z,a'}|y_{a'}) = P(y(a, z(a'))|y'(a'))$ reads the probability of $Y = y$ had (1) A been a' and (2) Z been z when A is a' , given that an intervention $A = a'$ produced y' . This expression involves three worlds: a world where $A = a$, a world where $Z = z$, and a world where $A = a'$. Such complex expressions are used to characterize direct, indirect, and path-specific effects.

Causal-based discrimination discovery aims at telling if the outcome of an automated decision making is fair or discriminative. Several causal-based fairness notions are defined in the literature (Section 4) and expressed in terms of joint, conditional, interventional, and counterfactual probabilities. The application of a fairness notion requires as input a dataset D and a causal graph G . While joint probabilities (e.g. $P(X = x, Y = y, Z = z)$) and conditional probabilities (e.g. $P(Y = y|X = x)$) can be trivially estimated from the dataset D , probabilities involving interventions or counterfactuals cannot always be estimated from D and G . When a probability can be estimated from observable data (D), it is said to be *identifiable*. Otherwise it is *unidentifiable*. More formally, let M_1 and M_2 be two causal models sharing the same causal graph (not including the unobservable variable U) and the same set of probability distributions ψ , a quantity Q (e.g. intervention or counterfactual) is identifiable using ψ (noted ψ -identifiable), if the value of Q is unique and computable from ψ in any models M_1 and M_2 . In other words, if there exists two models M_1 and M_2 sharing the same graph structure and the same probability distributions, but yielding different Q values, then Q is unidentifiable. Typically, the identifiability of interventional and counterfactual quantities depends on the structure of the graph, in particular, the location of unobserved confounding variables. Identifiability criteria are summarized in Section 5.

4 CAUSALITY-BASED FAIRNESS NOTIONS

Without loss of generality, assume the sensitive attribute and outcome A and Y as binary variables where $A = a_0$ denotes the privileged group (e.g. male), typically considered as the reference in characterizing discrimination, and $A = a_1$ the disadvantaged group (e.g. female).

The most common non-causal fairness notion is total variance (TV), known as statistical parity, demographic parity, or risk difference. The total variance of $A = a_1$ on the outcome $Y = y$ with

⁵The notations $Y_{Z=z}$ and $Y(z)$ are used in the literature as well.

reference $A = a_0$ is defined using conditional probabilities as follows:

$$TV_{a_1, a_0}(y) = P(y | a_1) - P(y | a_0) \quad (1)$$

Intuitively, $TV_{a_1, a_0}(y)$ measures the difference between the conditional distributions of Y when we (passively) observe A changing from a_0 to a_1 . The main limitation of TV is purely statistical nature which makes it unable to reflect the causal relationship between A and Y , that is, it is insensitive to the mechanism by which data is generated. Total effect (TE) [21] is the causal version of TV and is defined in terms of experimental probabilities as follows:

$$TE_{a_1, a_0}(y) = P(y_{a_1}) - P(y_{a_0}) \quad (2)$$

TE measures the effect of the change of A from a_1 to a_0 on $Y = y$ along all the causal paths from A to Y . TE is normally estimated using randomized controlled trials (RCT) [7]. However, given the impracticality of controlling on the sensitive attribute in a satisfactorily randomized way, intervention (y_a) is used instead. Intuitively, while TV reflects the difference in proportions of $Y = y$ in the current cohort, TE reflects the difference in proportions of $Y = y$ in the entire population. A more involved causal-based fairness notion considers the effect of a change in the sensitive attribute value (e.g. gender) on the outcome (e.g. probability of admission) given that we already observed the outcome for that individual. This typically involves an impossible situation which requires to go back in the past and change the sensitive attribute value. Mathematically, this can be formalized using counterfactual quantities. The simplest fairness notion using counterfactuals is the effect of treatment on the treated (ETT) [21] which is defined as:

$$ETT_{a_1, a_0}(y) = P(y_{a_1} | a_0) - P(y | a_0) \quad (3)$$

$P(y_{a_1} | a_0)$ reads the probability of $Y = y$ had A been a_1 , given A had been observed to be a_0 . Such probability involves two worlds: an actual world where $A = a_0$ and a counterfactual world where for the same individual $A = a_1$. Notice that $P(y_{a_0} | a_0) = P(y | a_0)$, a property called consistency [21].

TV , TE , and ETT fall into the framework of disparate impact [2] which aims at ensuring the equality of outcomes among all groups (protected and unprotected). An alternative framework is the disparate treatment [2] which seeks equality of treatment achievable through prohibiting the use of the sensitive attribute in the decision process.

Common fairness notions from the disparate treatment framework include direct effect, indirect effect, and path-specific effect. An effect can be deemed fair or unfair by an expert of the scenario at hand. Unfair effect is called discrimination. Direct discrimination is assessed using causal effect along direct edge from A to Y while indirect discrimination is measured using the causal effect along causal paths that pass through redlining/proxy attributes⁶. Figure 3 presents another causal graph of the firing example (Section 1) involving a redlining variable R (e.g. the number of disruptive incidents reported in the classrooms, which is clearly dependent on the initial level of students in the class, and hence a proxy for that sensitive attribute) and an explaining variable E (e.g. the difference between the final level (at the end of the academic year)

⁶Redlining/proxy attributes are attributes that cannot be objectively justified if used in the decision making process.

and the initial level, which is a legitimate feature to consider in the decision making.

Direct discrimination is then transmitted along the path $A \rightarrow Y$, and indirect discrimination is transmitted along the path $A \rightarrow R \rightarrow Y$. Assume that the use of the difference between a student initial and final average level E can be objectively justified since it is reasonable to fire a teacher given the dissuasive performance of her students. Path $A \rightarrow E \rightarrow Y$ is an explainable effect of A on Y since it is objectively justifiable to fire a teacher based on the poor improvement of the level of her students. Natural direct effect

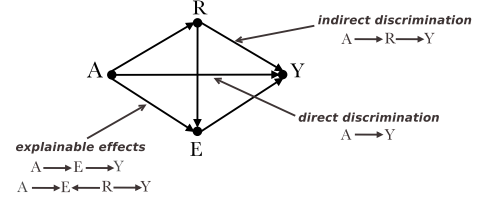


Figure 3: Teacher firing example where A : initial level of students, Y : firing decision, R : number of reported disruptive incidents in class, E : difference between final level and initial levels of students. $A \rightarrow Y$: direct discrimination. $A \rightarrow R \rightarrow Y$: indirect discrimination through the redlining variable R . $A \rightarrow E \rightarrow Y$: explainable effect of A on Y due to E .

(NDE) [20] is a common notion that measures the direct discrimination and is defined as:

$$NDE_{a_1, a_0}(y) = P(y_{a_1, Z_{a_0}}) - P(y_{a_0}) \quad (4)$$

Where Z is the set of mediator variables and $P(y_{a_1, Z_{a_0}})$ is the probability of $Y = y$ had A been a_1 and had Z been the value it would naturally take if $A = a_0$. That is, A is set to a_1 in the single direct path $A \rightarrow Y$ and is set to a_0 in all other indirect paths ($A \rightarrow R \rightarrow Y$ and $A \rightarrow E \rightarrow Y$).

On the other hand, natural indirect effect (NIE) [20] measures the indirect effect of A on Y and is defined as:

$$NIE_{a_1, a_0}(y) = P(y_{a_0, Z_{a_1}}) - P(y_{a_0}) \quad (5)$$

The problem of NIE is that it does not distinguish between the fair (explainable) and unfair (indirect discrimination) effects. Path-specific effect [21] is a more nuanced measure that characterizes the causal effect in terms of specific paths.

Given a path set π , the π -specific effect is defined as:

$$PSE_{a_1, a_0}^{\pi}(y) = P(y_{a_1 | \pi, a_0 | \bar{\pi}}) - P(y_{a_0}) \quad (6)$$

where $P(y_{a_1 | \pi, a_0 | \bar{\pi}})$ is the probability of $Y = y$ in the counterfactual situation where the effect of A on Y with the intervention (a_1) is transmitted along π , while the effect of A on Y without the intervention (a_0) is transmitted along paths not in π (denoted by: $\bar{\pi}$).

4.1 No unresolved discrimination

No unresolved discrimination [14] is a fairness notion that falls into the disparate treatment framework and focuses on the indirect causal effects from A to Y . No unresolved discrimination is satisfied when no directed path from A to Y is allowed, except via a resolving variable E . A resolving variable is any variable in a causal

graph that is influenced by the sensitive attribute in a manner that it is accepted as nondiscriminatory. Figure 4 presents two alternative causal graphs for the teacher firing example. The graph at the left exhibits unresolved discrimination along the heavy paths: $A \rightarrow R \rightarrow Y$ and $A \rightarrow Y$. By contrast, the graph at the right does not exhibit any unresolved discrimination as the effect of A on Y is justified by the resolved variable E : $A \rightarrow E \rightarrow Y$.



Figure 4: Two alternative graphs for a teacher firing example. Y exhibits unresolved discrimination in the left graph (along the heavy paths), but not in the right one.

4.2 No proxy discrimination

Similarly to no unresolved discrimination, no proxy discrimination [14] focuses on indirect discrimination. A causal graph exhibits potential proxy discrimination if there exists a path from the protected attribute A to the outcome Y that is blocked by a proxy/redlining variable R . It is called proxy because it is used to decide about the outcome Y while it is a descendent of A which is significantly correlated with it in such a way that using the proxy in the decision has almost the same impact as using A directly. An outcome variable Y exhibits no proxy discrimination if the equality:

$$P(Y | R_r) = P(Y | R_{r'}) \quad \forall r, r' \in \text{dom}(R) \quad (7)$$

holds for any potential proxy R .

Figure 5 shows two similar causal graphs for the same firing example. The causal graph at the left presents a potential proxy discrimination via the path: $A \rightarrow R \rightarrow Y$. However, the graph at the right is free of proxy discrimination as the edge between A and its proxy R has been removed.

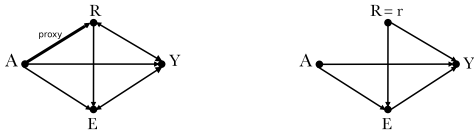


Figure 5: The graph at the left exhibits a potential proxy discrimination (along the heavy edge between A and R) but not in the right one.

4.3 Fair on Average Causal Effect (FACE) and Fair on Average Causal Effect on the Treated (FACT)

Khademi et al. [13] defined two group-based fairness notions based on the potential outcome framework [12] called FACE and FACT. FACE considers the average causal effect of the sensitive attribute A on the outcome Y at a population level while FACT focuses on the same effect but on the sub-population/group level. Let $Y_i^{(a)}$ be

the potential outcome of a data point i had A been a . Using FACE, an outcome Y is said to be fair, on average over all individuals (data points) in the population, with respect to A , iff:

$$E[Y_i^{(a_1)} - Y_i^{(a_0)}] = 0 \quad (8)$$

where $E[\cdot]$ is the expectation of a random variable over all data inputs. This is equivalent to the expected value of the $TE_{a_1, a_0}(Y)$ (Equation (2)) in the causal model of Pearl [21]. Using FACT, an outcome Y is said to be fair with respect to the sensitive attribute A , on average over individuals of the group $A = a_0$, iff:

$$E[Y_i^{(a_1)} - Y_i^{(a_0)} | A^i = a_0] = 0 \quad (9)$$

This is equivalent to the expected value of $ETT_{a_1, a_0}(Y)$ (Equation (3)). Being defined in terms of potential outcomes, these two notions can be estimated using tools from the potential outcome framework. Khademi et al. estimated FACE using Inverse Probability Weighting (IPW) [25] and FACT using matching methods [31].

4.4 Counterfactual fairness

Counterfactual fairness [15] is a fine-grained variant of ETT conditioned on all attributes. That is, an outcome Y is counterfactually fair if under any assignment of values $X = x$,

$$P(y_{a_1} | X = x, A = a_0) = P(y_{a_0} | X = x, A = a_0) \quad (10)$$

where $X = V \setminus \{A, Y\}$ is the set of all remaining variables. Since conditioning is done on all remaining variable X , counterfactual fairness is an individual notion. According to Equation (10), counterfactual fairness is satisfied if the probability distribution of the outcome Y is the same in the actual and counterfactual worlds, for every possible individual. Kusner et al. [15] tested counterfactual fairness by generating, for every individual in the population, another sample with counterfactual sensitive value. Then, they compared the density functions of the actual samples with the counterfactual samples. To be fair, a predictor should produce outcome values where actual and counterfactual density plots are identical.

4.5 Counterfactual Effects

By conditioning on the sensitive attribute $A = a$, Zhang and Bareinboim [38] defined two variants of NDE (Equation (4) and NIE (Equation (5)) which focus on the direct and indirect effect for a specific group. In addition, they characterize a third type of effect, spurious, which considers the back-door paths between A and Y , that is, paths with an arrow into A . For instance, in Figure 6, the observed disparities between the protected and unprotected groups can be decomposed into direct ($A \rightarrow Y$), indirect ($A \rightarrow R \rightarrow Y$ and $A \rightarrow E \rightarrow Y$) and spurious ($A \leftarrow C \rightarrow Y$) effects.

The three effects are defined as follows:

$$DE_{a_1, a_0}(y|a) = P(y_{a_1, Z_{a_0}} | a) - P(y_{a_0} | a) \quad (11)$$

$$IE_{a_1, a_0}(y|a) = P(y_{a_0, Z_{a_1}} | a) - P(y_{a_0} | a) \quad (12)$$

$$SE_{a_1, a_0}(y) = P(y_{a_0} | a_1) - P(y_{a_0} | a_0) \quad (13)$$

where in Equations (11) and (12), a can be a_0 or a_1 . Considering the simple job hiring example and focusing on the unprotected group ($A = a_1$ e.g. females), $DE_{a_1, a_0}(y|a_1)$ measures the change in the probability of Y (e.g. hiring) had A been a_1 (e.g. female), while mediators Z are kept at the level they would take had A been a_0

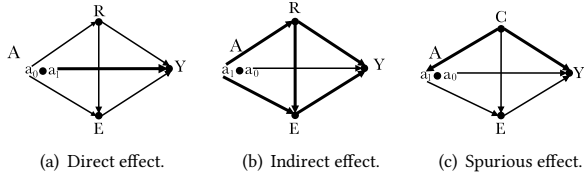


Figure 6: Figure 6(a) illustrates how counterfactual direct effect might be applied through the direct path: $A \rightarrow Y$, Figure 6(b) examines the indirect discriminative effect of A on Y through the path: $A \rightarrow R \rightarrow Y$ and Figure 6(c) reveals the presence of spurious effect via the back-door path: $A \leftarrow C \rightarrow Y$.

(e.g. male), in particular for the individuals $A = a_1$ (e.g. female). $SE_{a_1, a_0}(y)$ reads the change in the probability of hiring Y had A been a_0 (e.g. female) for the individuals that would naturally possess $A = a_0$ versus a_1 . In Figure 6(c), the heavy path represents the spurious effect between A and Y through the confounder C . If there is no such back-door path, $SE_{a_1, a_0}(y)$ would be zero.

Interestingly, by considering these fine grained variants of NDE and NIE , it is possible to decompose $TV_{a_1, a_0}(y)$ (Equation (1)) and $ETT_{a_1, a_0}(y)$ (Equation (3)) as follows:

$$TV_{a_1, a_0}(y) = SE_{a_1, a_0}(y) + IE_{a_1, a_0}(y|a_1) - DE_{a_0, a_1}(y|a_1) \quad (14)$$

$$ETT_{a_1, a_0}(y) = DE_{a_1, a_0}(y|a_0) - IE_{a_0, a_1}(y|a_0) \quad (15)$$

In linear models, the three types of effects sum up to the total variation:

$$TV_{a_1, a_0}(Y) = SE_{a_1, a_0}(Y) + IE_{a_1, a_0}(Y|a_1) + DE_{a_1, a_0}(Y|a) \quad (16)$$

4.6 Counterfactual Error Rates

Equalized odds is an important statistical fairness notion which requires equality of error rates (TPR and FPR) across sub-populations, that is,

$$ER_{a_1, a_0}(\hat{y}|y) = P(\hat{y} | a_1, y) - P(\hat{y} | a_0, y) = 0 \quad (17)$$

where \hat{y} denotes the prediction while y denotes the true outcome. The problem of this statistical notion is the difficulty to identify the causes behind the discrimination if any. Zhang and Bareinboim [37] decompose equalized odds (Equation (17)) using three counterfactual measures corresponding to the direct, indirect and spurious effects of A on \hat{Y} . The three measures are counterfactual direct error rate, counterfactual indirect error rate, and counterfactual spurious error rate. Let $\hat{y} = f(\hat{\mathbf{p}}_A)$ be a classifier where $\hat{\mathbf{p}}_A$ is the set of input features (parent variables of \hat{Y}) for the classifier. The counterfactual error rates for a sub-population a, y (with prediction $\hat{y} \neq y$) are defined as:

$$ER_{a_1, a_0}^d(\hat{y} | a, y) = P(\hat{y}_{a_1, y, (\hat{\mathbf{p}}_A \setminus A)_{a_0, y}} | a, y) - P(\hat{y}_{a_0, y} | a, y) \quad (18)$$

$$ER_{a_1, a_0}^i(\hat{y} | a, y) = P(\hat{y}_{a_0, y, (\hat{\mathbf{p}}_A \setminus A)_{a_1, y}} | a, y) - P(\hat{y}_{a_0, y} | a, y) \quad (19)$$

$$ER_{a_1, a_0}^s(\hat{y} | y) = P(\hat{y}_{a_0, y} | a_1, y) - P(\hat{y}_{a_0, y} | a_0, y) \quad (20)$$

For example, the counterfactual direct error rate (Equation (18)) measures the error rate (disparity between the true and the predicted outcome) in terms of the direct effects of the sensitive attribute A on the prediction \hat{Y} . In the simple job hiring example, considering the rejected females sub-population ($a = a_1$ and $y = \text{rejected}$), it reads: for a rejected female candidate, how would the prediction \hat{Y} change had the candidate been a female (A been a_1), while keeping all the other features $\hat{\mathbf{p}}_A \setminus A$ at the level that they would attain had “she was male”, compared to the prediction \hat{Y} she would receive had “she was male” and being rejected? Figure 7 illustrates

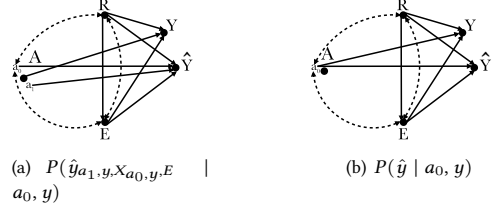


Figure 7: Illustration of the counterfactual direct error rate in the job hiring example. There are unobserved confounders between A, R and E because the definition is conditioning on a collider (Y).

illustrates how counterfactual direct error rate is applied. Thus, the input of A to the direct path $A \rightarrow \hat{Y}$ is changed from a_0 (baseline) to a_1 , while keeping the value of A to the other variables (R and E) fixed at the baseline level (a_0, y). Since the direct path $A \rightarrow \hat{Y}$ is the only difference between models of Figure 7 (a-b), the change in \hat{Y} measures the direct effect of A on \hat{Y} . The variable $E_{a_0, y} = E$ since E is a non-descendant node of A and Y .

Interestingly, the statistical equalized odd error rate (Equation (17)) can be decomposed in terms of the three above causal-based error rates:

$$ER_{a_1, a_0}(\hat{y} | y) = ER_{a_1, a_0}^d(\hat{y} | a_0, y) - ER_{a_0, a_1}^i(\hat{y} | a_0, y) - ER_{a_0, a_1}^s(\hat{y} | y) \quad (21)$$

4.7 Individual direct discrimination

Individual direct discrimination [41] aims to discover the direct discrimination at the individual level. It is based on situation testing [3], a legally grounded technique for analyzing the discrimination at an individual level. It consists in comparing the individual with similar individuals from both groups (protected and unprotected). That is, for an individual in question i , find the k other individuals which are the most similar to i in the group $A = a_0$ and k similar individuals from the group $A = a_1$. The first set is denoted S^+ while the second S^- . The target individual is considered as discriminated if the difference observed between the rate of positive decisions in S^- and S^+ is higher than a predefined threshold τ (typically 5%).

Causal inference is used to define the distance function $d(i, i')$ required to select the elements of S^- and S^+ . First, only attributes that are direct causes of the outcome should be considered in the computation of the distance. That is, based on the causal graph, $Q = Pa(Y) \setminus \{A\}$ denotes the set of variables that should be used in the distance function. Second, the causal effect of each of the

selected attributes ($Q_k \in Q$) on the the outcome should be considered in the function definition. In particular, for each variable Q_k , $CE(q_k, q'_k)$ measures the causal effect on the outcome when the value of Q_k changes from q_k to q'_k and is defined as:

$$CE(q_k, q'_k) = P(y_q) - P(y_{q'_k, q \setminus \{q_k\}}) \quad (22)$$

where $(P(y_q))$ is the effect of the intervention that forces the set Q to take the set of values q , and $(P(y_{q'_k, q \setminus \{q_k\}}))$ is the effect of the intervention that forces Q_k to take value q'_k and other attributes in Q to take the same values as q .

4.8 Non-Discrimination Criterion

Non-discrimination criterion [42] is a group fairness notion that aims to discover and to quantify direct discrimination. This notion is based on the identification of meaningful partitions of attributes that can be used to provide quantitative assessment for discrimination. A partition includes a subset of non-protected attributes (called *block set*) which blocks the causal effect from A to Y while a group is specified by a value assignment to these attributes. Given a block set Q , the discrimination between protected and unprotected groups is assessed by computing the risk difference [26]:

$$|\Delta P|_q = |P(y | a_1, q) - P(y | a_0, q)| \quad (23)$$

where q is a value assignment for the block set Q and the absolute value to consider both positive and negative discriminations. If the risk difference is less than τ for all combination of values of all block sets, no direct discrimination is reported. Equation (23) holds for each value assignment q of each block set $Q(Q = \text{Par}(Y) \setminus \{A\})$. This notion is similar to the individual direct discrimination except that instead of using the set Q as an input to measure the similarity between individuals, it is used to define meaningful partitions to assess discrimination among groups.

4.9 PC-Fairness

PC-Fairness [36] is a general fairness formalization that covers various causal-based fairness notions. That is, by differently tuning its parameters, it is possible to match several of the causal-based fairness notions mentioned above. Given a factual condition $O = o$ where $O \subseteq \{A, V, Y\}$ and a causal path set π , an outcome Y achieves the PC-fairness iff:

$$PCE_{a_1, a_0}^\pi(y | o) = 0 \quad (24)$$

where $PCE_{a_1, a_0}^\pi(y | o) = P(y_{a_1 | \pi, a_0 | \bar{\pi}} | o) - P(y_{a_0} | o)$

Intuitively, $PCE_{a_1, a_0}^\pi(y | o)$ represents the path-specific counterfactual effect of the value change of A from a_0 to a_1 on Y through the specific causal path set π (with reference a_0) and given the factual observation o .

Most of the aforementioned causal-based fairness notions can be expressed as special cases of PC-Fairness. For instance, if the set π includes all possible paths and $O = V \setminus \{Y\}$, PC-fairness corresponds to counterfactual fairness (Equation (10)). If the set π includes all possible paths and $O = \emptyset$, PC-fairness corresponds to total effect (Equation (2)).

4.10 Equality of Effort

Equality of effort [10] fairness notion identifies discrimination by assessing how much effort is needed by the disadvantaged individual/group to reach a certain level of outcome. A treatment variable T (considered as a legitimate variable) is selected and used to address the question: "to what extent this treatment variable T should change to make the individual (or a group of individuals) achieve a certain outcome level?". Hence, this notion focuses on whether the effort to reach a certain outcome level is the same for the protected and unprotected groups. Considering the simple job hiring example, the education level E is a good choice for the treatment variable. Two equality of effort notions are defined based on the potential outcome framework [12], individual γ -Equal effort and system γ -Equal effort. Let $Y_i^{(t)}$ be the potential outcome for individual i had T been t and $E[Y^{(t)}]$ be the expected outcome for individual i . As $Y_i^{(t)}$ is not observable, situation testing is used to estimate it in a similar way as individual direct discrimination (Section 4.7). Let S^+ and S^- be the two sets of similar individuals with $A = a_0$ and $A = a_1$, respectively, and $E[Y_{S^+}^{(t)}]$ be the expected outcome under treatment t for the subgroup S^+ . The minimal effort needed to achieve γ level of outcome variable within the subgroup S^+ is defined as:

$$\Psi_{S^+}(\gamma) = \underset{t \in T}{\operatorname{argmin}} \{E[Y_{S^+}^{(t)}] \geq \gamma\} \quad (25)$$

Individual γ -Equal effort is satisfied for individual i if:

$$\Psi_{S^+}(\gamma) = \Psi_{S^-}(\gamma) \quad (26)$$

System γ -Equal effort is satisfied for a sub-population (e.g. $A = a_1$) if:

$$\Psi_{D^+}(\gamma) = \Psi_{D^-}(\gamma) \quad (27)$$

where D^+ and D^- are the subsets of the entire dataset with sensitive attributes a_0 and a_1 , respectively.

5 IDENTIFIABILITY

The identifiability of causal quantities has been extensively studied in the literature: causal effect (intervention) identification [8, 11, 21, 29, 30, 32–34], counterfactual identification [27–29, 35], direct/indirect effects [20] and path-specific effect identification [1, 17, 27, 39, 40]. This section summarizes the main identifiability conditions as they relate to the specific problem of discrimination discovery.

5.1 Identifiability of causal effect (intervention)

The causal effect of a cause variable X on an effect variable Y is captured by $P(Y_x) = P(Y | do(X = x))$. In discrimination setup, the cause is typically the sensitive attribute A . So the causal effect of changing the value of the sensitive attribute on the outcome variable Y is $P(Y | do(A = a))$. A basic case where identifiability can be avoided altogether is when it is possible to perform experiments by intervening on the sensitive attribute A . When this is possible, randomized controlled trial (RCT) [7] can be used to estimate the causal effect. RCT consists in randomly allocating subjects to two or more groups according to the sensitive attribute (e.g. selecting randomly men and women). Then, comparing the outcome Y

of all groups. The keyword in this process is “randomized” which is very difficult to guarantee in practice. For instance, to estimate the causal effect of race on loan granting, one can consider all applicants to loans in a specific city or in a specific year and select randomly applicants from every ethnicity. The experimenter may truly use randomization in the selection, but the result is not guaranteed to be randomized for reasons such as: individuals from a certain race are more likely to apply for loans in that city, low-income individuals from a given race are less likely to apply for loans in that specific year of the study, etc.

In Markovian models (no unobserved confounding), the causal effect is always identifiable (Corollary 3.2.6 in [21]). The simplest case is when there is no confounding between A and Y (Figure 8(a)). In that case, the causal effect matches the conditional probability regardless of any mediator:

$$P(y_a) = P(y|do(a)) = P(y|a) \quad (28)$$

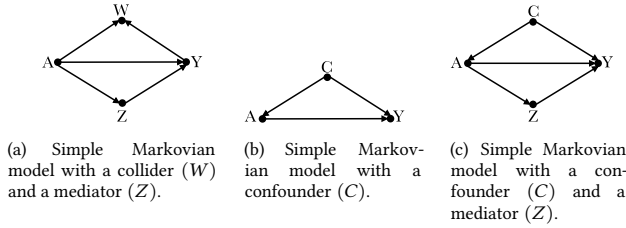


Figure 8

Note that marginalizing over the variable W invalidates the equality in (28) as W is a collider [21]. In presence of an observable confounder (Figure 8(b)), $P(y_a)$ is identifiable by adjusting on the confounder:

$$P(y_a) = \sum_C P(y|a, c) P(c) \quad (29)$$

where the summation is on values c in the domain (sample space) of C denoted $dom(C)$. Equation (29) is called the back-door formula⁷. In Markovian models, the causal effect $P(y_a)$ can also be computed using the truncation formula:

$$P(y_a) = \sum_{Y=y} \prod_{V \in V \setminus \{A, Y\}} P(v|pa_V) \quad (30)$$

where pa_V denotes the parent variables of V .

Using the back-door formula (29) or the truncated formula (30) on Figure 8(c) produces the same result:

$$\begin{aligned} P(y_a) &= \sum_C \sum_Z P(y|a, z, c) P(z|a) P(c) \\ &= \sum_C P(y|a, c) P(c) \end{aligned}$$

while the joint probability

$$P(y, a, c, z) = P(y|a, c, z) P(z|a) P(a|c) P(c)$$

and the conditional probability

$$P(y|a) = \sum_C \sum_Z P(y|a, c, z) P(c, z|a)$$

⁷Called also adjustment formula or stratification.

For semi-Markovian models, identifiability of $P(y_a)$ is not guaranteed. The simplest graph in which the causal effect between A and Y is not identifiable is the “bow” graph (Figure 9(a)). This simple unidentifiability criterion can be generalized to a more complex graphs called c-tree. A c-tree is a graph that is at the same time a tree⁸ and a c-component. A c-component is a set of vertices (variables) such that every pair of vertices are connected by a bi-directed path (composed only of unobservable confounding edges). Figure 9(b) shows an example of c-tree. If the causal graph is a c-tree rooted in the outcome variable Y , $P(y_a)$ is unidentifiable [29].

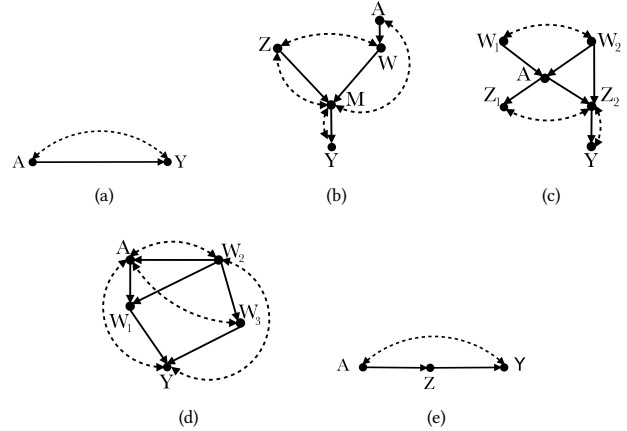


Figure 9: Figure 9(a) presents the “bow” graph, Figure 9(b) illustrates the structure of a c-tree, Figure 9(c) shows a semi-Markovian model where $P(y_a)$ is observable, Figure 9(c) presents a semi-Markovian model where $P(y_a)$ is identifiable and Figure 9(e) illustrates a simple example of a front-door criterion.

The simplest case where $P(y_a)$ is identifiable in a semi-Markovian model is when the sensitive attribute A is not involved in any confounding. That is, there is no bi-directional edge connected to A . This matches Theorem 3.2.5 in [21] which states that if all parents of a cause variable A are observable, the causal effect of that variable is identifiable. Figure 9(c) shows a graph diagram of a semi-Markovian model where $P(y_a)$ is identifiable. If such criterion is satisfied, the causal effect can be computed as follows [33]:

$$P(y_a) = P(y|a, pa_a) P(pa_a) \quad (31)$$

where pa_a is the set of values of the parents of A . A more complex criterion of identifiability of $P(y_a)$ is when the sensitive attribute A is involved in confounding, but no child of A is involved in confounding. Figure 9(d)⁹ shows a graph satisfying this criterion, and hence, allowing the identification of $P(y_a)$. In such case [33],

$$P(y_a) = \left(\prod_{c \in cha} P(c|pa_c) \right) \sum_{a' \in dom(A)} \frac{P(y, a')}{\prod_{c \in cha} P(c|pa_c)} \quad (32)$$

⁸Notice that, in this paper, the direction of the arrows between nodes is reversed compared to the usual tree structure.

⁹The same example as in [21] page 92.

For example, in Figure 9(d),

$$\begin{aligned} P(y_a) &= P(w_1|w_2, a) \sum_{a' \in \text{dom}(A)} \frac{P(y, a')}{P(w_1|w_2, a')} \\ &= \sum_{w_1} \sum_{w_2} \sum_{a'} P(y|w_1, w_2, a') P(a'|w_2) \times P(w_1|w_2, a) P(w_2) \end{aligned}$$

Pearl [21] obtains the same result using the do-calculus (to be covered later in the section).

All the above criteria can be generalized to the case where the sensitive attribute is not connected to any of its children through a confounding path. The main idea of this result is that since the causal graph is not composed of a single c-component (in which case, it is not identifiable as mentioned earlier), variables (vertices) can be partitioned into a disjoint set of c-components. For example, in Figure 9(c), there are three c-components: $\{W_1, W_2\}$, $\{A\}$, $\{Z_1, Z_2, Y\}$. Hence, as long as there is no confounding path connecting A to any of its direct children, $P(y_a)$ is identifiable and can be computed as [33]:

$$P(y_a) = \frac{P(y)}{Q^A} \sum_{a'} Q^A \quad (33)$$

where Q^A is the c-factor of the c-component containing A (S^A) computed as follows:

$$Q^A = \prod_{v \in S^A} P(v|v^{-1}) \quad (34)$$

where v^{-1} is the set of values of all previous variables to V , assuming a topological order $V_1 < V_2 < \dots < V_n$ over all variables. For instance, in Figure 9(d), $W_2 < A < W_1 < W_3 < Y$ is a valid topological order. This criterion can be slightly generalized to be: $P(y_a)$ is identifiable if there is no confounding path connecting A to any of its children in $G_{An(Y)}$ which is the subgraph of G composed only of ancestors of the outcome variable Y .

In case there is unobservable confounding between the sensitive attribute A and the outcome Y , all above criteria will fail. However, $P(y_a)$ can still be identifiable using the front-door criterion. This criterion is satisfied in Figure 9(e)) and consists in having a mediator variable (Z) such that:

- there are no backdoor paths from A to Z
- all backdoor paths from Z to Y are blocked by A .

A backdoor path from A to Z is any path starting at A with a backward edge \leftarrow into A (e.g. $A \leftarrow \dots Z$). If such criterion is satisfied, $P(y_a)$ can be computed as follows:

$$\begin{aligned} P(y_a) &= \sum_Z P(y|do(z)) P(z|do(a)) \\ &= \sum_Z P(y|z, a) P(a) P(z|a) \end{aligned} \quad (35)$$

In addition to these graphical criteria, Pearl [21] proposed a do-calculus composed of three rules allowing to turn interventional probabilities to observational ones:

- (1) $P(y_a|z, w) = P(y_a|z)$ provided that the set of variables Z blocks all backdoor paths from W to Y after all arrows leading to A have been deleted.
- (2) $P(y_a|z) = P(y|a, z)$ provided that the set of variables Z blocks all backdoor paths from A to Y .

- (3) $P(y_a) = P(y)$ provided that there are no causal paths between A and Y .

In [21] Section 3.5, Pearl gives examples of the above criteria. Shpitser and Pearl [29] proved that all the unidentifiable cases of the causal effect $P(y_a)$ boil down to a general graphical structure called the hedge criterion. Based on this criterion, they designed a complete identifiability algorithm called *ID* which outputs the expression of $P(y_a)$ if it is identifiable, or the reason of the unidentifiability otherwise.

5.2 Identifiability of counterfactuals

Most of causality-based fairness notions listed in Section 4 are defined in terms of counterfactual quantities. Hence, the applicability of those notions depends heavily on the identifiability of the counterfactuals composing them. If a fairness notion is relying on a counterfactual quantity which is proven to be unidentifiable in the causal model at hand, the notion cannot be used to compute the level of discrimination¹⁰. In Markovian, as well as semi-Markovian models, if all parameters of the causal model are known (including $P(u)$), any counterfactual is identifiable and can be computed using the three steps abduction, action, and prediction (Theorem 7.1.7 in [21]). Let $P_* = \{P_x | X \subseteq V, x \text{ a value assignment of } X\}$ be the set of all interventional distributions in a given causal model. While the identifiability of interventional probabilities $P(y_a)$ is characterized based on observational probabilities $P(v)$, in this section, the identifiability of counterfactuals is characterized in terms of interventional probabilities P_* . Then, combining these results with the criteria of the previous section, a counterfactual can, in turn, be identified using observational probabilities $P(v)$.

Given a causal graph G of a Markovian model and a counterfactual expression $\gamma = v_x|e$ with e some arbitrary set of evidence, identifying and computing $P(\gamma)$ requires to construct a counterfactual graph which combines parallel worlds. Every world is represented by a model M_x corresponding to each subscript in the counterfactual expression. For example, given the causal graph in Figure 10(a) and the counterfactual expression $y_{a'}|a, y$, the resulting counterfactual graph is shown in Figure 10(b). The vertex $A = a'$ represents the intervention $do(a')$ and therefore the absence of any incoming arrow edge into $A = a'$ vertex. The unobservable confounding variable U_Y highlights the interaction between both worlds (the actual and the counterfactual). Notice that, starting from the Markovian model of Figure 10(a), representing the counterfactual expression resulted in the semi-Markovian model of Figure 10(b). The counterfactual graph should be “reduced” by merging together vertices that share the same causal mechanism (**make-cg** algorithm in [29] automates this procedure). The resulting counterfactual graph can be considered as a typical causal graph for a larger causal model. Consequently, all the graphical criteria listed in Section 5.1 for the identifiability of causal effect apply on the counterfactual graph to identify counterfactual quantities, in particular, the c-component factorization of the counterfactual graph [28].

The simplest unidentifiable counterfactual quantity is $P(y_{a'}, y'_a)$ which is called the probability of necessity and sufficiency. The corresponding counterfactual graph is the W-graph that has the same

¹⁰But it can be used to find bounds on the actual value.

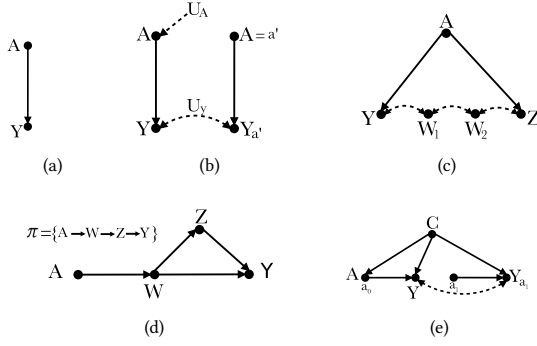


Figure 10: Causal graphs.

structure as to Figure 10(b). This simple criterion can be generalized to the zig-zag graph (Figure 10(c)) where the counterfactual $P(y_a, w_1, w_2, z_x)$ is not identifiable.

Pearl [21] proves two results about the identifiability of counterfactuals. First, for linear causal models (i.e. the functions F are linear), any counterfactual is experimentally (using P_*) identifiable whenever the model parameters are identified. Second, in linear causal models, if some of the model parameters are unknown, any counterfactual of the form $E(Y_a|e)$ where e is some arbitrary set of evidence, is identifiable provided that $E(y_a)$ is identifiable. Finally, there is no single necessary and sufficient criterion for the identifiability of counterfactuals in semi-Markovian models [1]. ID^* and IDC^* algorithms [29] automate the identifiability and computation of counterfactuals based on all above criteria.

To illustrate the computation of a counterfactual probability, consider the teacher firing example of Figure 1 and the counterfactual probability $P(y_{a_1}|a_0)$ which reads the probability of firing a teacher who is assigned a class with a high initial level of students (a_0) had she been assigned a class with a low initial level of students (a_1). Applying **make-cg** algorithm based on this counterfactual quantity produces the counterfactual graph in Figure 10(e) which combines two worlds: the actual world where the teacher has normally $A = a_0$ and the counterfactual world where *the same* teacher is assigned $A = a_1$. Both variables C are reduced to a single variable and Y and Y_{a_1} are connected by an unobservable confounder. The counterfactual graph is composed of three c-components $\{C\}, \{A\}, \{Y, Y_{a_1}\}$. Applying algorithm IDC^* [29] results in:

$$P(y_{a_1}|a_0) = \frac{\sum_{y,c} Q(c) Q(a_0) Q(y, y_{a_1})}{P(a_0)} \quad (36)$$

where $Q(v) = P(v|pa(V))$ in the counterfactual graph. Hence,

$$\begin{aligned} P(y_{a_1}|a_0) &= \frac{\sum_{y,c} P(c) P(a_0|c) P(y, y_{a_1}|c)}{P(a_0)} \\ &= \frac{\sum_c P(c) P(a_0|c) P(y_{a_1}|c)}{P(a_0)} \end{aligned} \quad (37)$$

$$\begin{aligned} &= \frac{\sum_c P(c) P(a_0|c) P(y|a_1, c)}{P(a_0)} \quad (38) \\ &= \frac{0.5 \times 0.8 \times 0.25 + 0.5 \times 0.2 \times 0.01}{0.5} \\ &= 0.202 \end{aligned}$$

y in Equation (37) is cancelled by summation while $P(y_{a_1}|c)$ in the same equation is transformed into $P(y|a_1, c)$ in Equation (38) using rule (2) of the *do*-calculus (Section 5.1).

5.3 Identifiability of direct and indirect effects

In Markovian models, the average natural direct effect NDE and the average natural indirect effect NIE are always identifiable (from observational data) and can be computed as follows [20]:

$$NDE_{a_1, a_0}(Y) = \sum_s \sum_z \left(E[Y|a_1, z] - E[Y|a_0, z] \right) P(z|a_0, s) P(s) \quad (39)$$

$$NIE_{a_1, a_0}(Y) = \sum_s \sum_z E[Y|a_0, z] \left(P(z|a_1, s) - P(z|a_0, s) \right) P(s) \quad (40)$$

where Z is a set of mediator variables and S is any set of variables satisfying the back-door criterion between the sensitive variable A and the mediator variables Z , that is, (i) no variable in S is a descendant of A and (ii) S blocks all back-door paths between A and Z . A simpler formulation can be used in case there is no confounding between A and Z , where the need for S is dropped altogether:

$$NDE_{a_1, a_0}(Y) = \sum_z \left(E[Y|a_1, z] - E[Y|a_0, z] \right) P(z|a_0) \quad (41)$$

$$NIE_{a_1, a_0}(Y) = \sum_z E[Y|a_0, z] \left(P(z|a_1) - P(z|a_0) \right) \quad (42)$$

In semi-Markovian models, NDE and NIE are not generally identifiable, even if we have the luxury to perform any experiment using *RCT*, because of the nested counterfactual $P(Y_{a_1}, Z_{a_0})$ in Equation (4). Nevertheless, these quantities are identifiable *from experimental data* provided that there is a set of variables W which are parents of the outcome variable Y but non-descendants of A and Z such that $Y_{a_0, z} \perp\!\!\!\perp Z_{a_0} | W$ (reads: $Y_{a_0, z}$ and Z_{a_0} are independent conditional of W). This condition can be easily checked from the causal graph as follows: W d-separates Y and Z in the graph formed by deleting all arrows emanating from A and Z , denoted simply as $(Y \perp\!\!\!\perp Z | W)_{G_{AZ}}$.

If such graphical condition is satisfied, NDE and NIE can be computed from experimental quantities as follows:

$$NDE_{a_1, a_0}(Y) = \sum_{z, w} \left(E[Y_{a_1, z}|w] - E[Y_{a_0, z}|w] \right) P(Z_{a_0} = z|w) P(w) \quad (43)$$

$$NIE_{a_1, a_0}(Y) = \sum_{z, w} E[Y_{a_0, z}|w] \left(P(Z_{a_1} = z|w) - P(Z_{a_0} = z|w) \right) P(w) \quad (44)$$

5.4 Identifiability of path-specific effects

The identifiability of $PSE_{\pi}(a_1, a_0)$ in Markovian models depends on whether $P(y|do(a_1|_{\pi}, a_0|_{\bar{\pi}}))$ is identifiable. Avin et al. [1] gave a single necessary and sufficient criterion for the identifiability of $P(y|do(a_1|_{\pi}, a_0|_{\bar{\pi}}))$ in Markovian models called recanting witness criterion. This criterion holds when there is a vertex W along the causal path π that is connected to Y through another causal path not in π . For instance, Figure 10(d) satisfies the recanting witness criterion when $\pi = A \rightarrow W \rightarrow Z \rightarrow Y$ with W as witness. The corresponding graph structure is called “kite” graph. When this criterion is satisfied, $P(y|do(a_1|_{\pi}, a_0|_{\bar{\pi}}))$ is not identifiable, and consequently, $PSE_{\pi}(a_1, a_0)$ is not identifiable. Shpitser [27] generalizes this criterion to semi-Markovian models known as recanting district criterion.

6 APPLICABILITY

In his book, *The Book of Why* [22], Pearl describes a causation ladder with three rungs: statistical observations (seeing), intervention (doing), and counterfactual (imagining). In this section, all causal-based fairness notions (Section 4) are placed in the causation ladder which will help us address the problem of their applicability in real-scenarios. The causation ladder is structured in such a way that quantity at a certain rung can be identified in terms of quantities at the rung just below it. As a consequence, the higher the rung, the more challenging the problem of identifiability, and hence the less applicable a fairness notion defined at that rung.

The diagram in Figure 11 shows the causation ladder and indicates at which rung every causal-based fairness notion stands. *TV* which is the only non-causal fairness notion covered in this paper is at rung 1. It is always applicable provided that a set of observations (dataset) is available. No unresolved and non-discrimination criteria are placed midway between rungs 1 and 2 as they are applicable provided that the causal graph is available along the dataset. Non-discrimination criterion, however, requires the Markov property to be applicable because causal dependence through unobservable paths cannot be blocked. It also has an exponential complexity since it considers all combination of values of the parent variables of the outcome Y . A relaxation is described by the authors [42] but the notion remains computationally intractable.

Fairness notions at rung 2 (*TE*, *FACE*, No-proxy discrimination, and individual direct discrimination) are applicable in any scenario where either experiments are possible (RCT) or hypothetical interventions are identifiable. As mentioned in Section 5.1, in Markovian models any intervention probability is identifiable from observational data. Hence, these fairness notions are always applicable in Markovian models, except for individual direct discrimination which assumes in addition that the sensitive attribute A has no parent [41]. In semi-Markovian models, the applicability of these rung 2 notions depend on the identifiability of the intervention terms used in their respective definitions. For instance, for individual direct discrimination, the term in question is $CE(q_k, q'_k)$ in Equation (22). The identifiability of intervention quantities in semi-Markovian models is discussed in Section 5.1. For the particular case of *FACE* notion, even if the intervention quantity $E[Y_i^{(a_1)}]$ in Equation (8) is not identifiable, it can be still estimated empirically since *FACE* is defined in the potential outcome framework [12].

The bulk of causal-based fairness notions are defined in terms of counterfactual quantities and hence are placed in rung 3 of the causation ladder. In Figure 11, the counterfactual notions are ranked from top to bottom according to their degree of applicability. For instance, counterfactual effects are placed on top of equality of effort and counterfactual fairness to indicate that the former is applicable in more scenarios than the latter. In Markovian models, the top 5 notions (*ETT*, *FACT*, *NDE*, *NIE*, and counterfactual effects) are always identifiable and hence applicable. That is, specific formula are already available to compute each counterfactual term used in their definitions. For instance, given a Markovian model, the three counterfactual effects (Section 4.5) can be computed from observational data as follows:

$$DE_{a_1, a_0}(y|a) = \sum_{Z, W} (P(y | a_1, z, w) - P(y | a_0, z, w))P(z | a_0, w)P(w | a)$$

$$IE_{a_1, a_0}(y|a) = \sum_{Z, W} P(y | a_0, z, w)P(z | a_1, w) - P(z | a_0, w))P(w | a)$$

$$SE_{a_1, a_0}(y) = \sum_{Z, W} P(y | a_0, z, w)P(z | a_0, w)(P(w | a_1) - P(w | a_0))$$

where Z and W are sets of mediator and confounder variables, respectively.

In Markovian models, the identifiability of counterfactual fairness and individual equality of effort depends on the identifiability of the term $P(y_{a_1}|X = x, A = a_0)$ which is only identifiable if X does not contain any variable which is at the same time descendant of A and ancestor of Y , that is, $X \cap B = \emptyset$ where $B = An(Y) \cap De(A)$ [35]. Path-specific counterfactual fairness is applicable provided that the model is Markovian and the recanting witness criterion is not satisfied (Section 5.4). In semi-Markovian models, unless all model parameters are known (including $P(u)$)¹¹, the identifiability of rung 3 fairness notions depends on the criteria discussed in Section 5.2, which rarely hold in practice. For *FACT* and equality of effort, as they are defined in the potential outcome framework, they come with a detailed procedure for empirical estimation [10, 13].

Finally, counterfactual error rate is a special case of rung 3 fairness notions as it is the only notion that conditions on the true outcome Y to assess the fairness of the prediction \hat{Y} (Equations (18), (19), and (20)). Such conditioning has an important implication on identifiability since Y is a collider, and conditioning on a collider creates a dependence between the previous variables [21]. This leads to unobservable confounding between the causes of the Y (Figure 7). Hence, even if the causal model is Markovian, applying counterfactual error rate turns it into a semi-Markovian model. Zhang and Bareinboim [37] define an identifiability criterion for counterfactual error rate in Markovian models called explanation criterion.

7 CONCLUSION

Applying causal-based fairness notions in practical scenarios is significantly hindered by the problem of identifiability. Based on an extensive body of work on identifiability theory, we summarized the most relevant identifiability criteria to the problem of discrimination discovery. These criteria are then used to characterize the

¹¹In that case, it is possible to use the three steps abduction, action, and prediction [21].

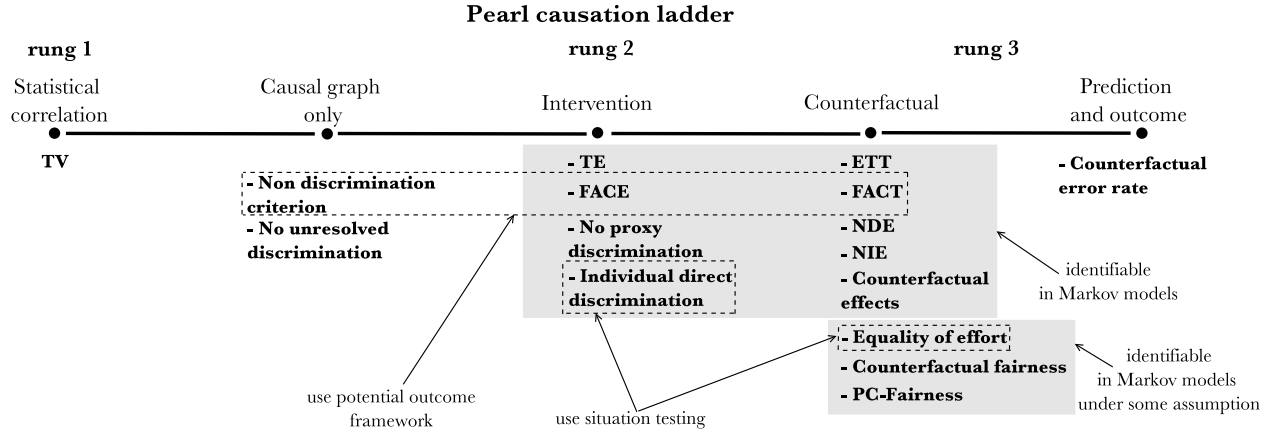


Figure 11

situations where causal-based fairness notions can be used in practical scenarios.

The main objective of this paper is to bridge the gap between the practical scenarios of automated (and generally unintentional) discrimination discovery and the mostly theoretical tackling of the problem in the literature. This effort can be of particular interest to civil right activists, civil right associations, anti-discrimination law enforcement agencies, and practitioners in fields where automated decision making systems are increasingly used.

8 ACKNOWLEDGEMENTS

The work of Catuscia Palamidessi was supported by the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme. Grant agreement № 835294.

REFERENCES

- [1] Chen Avin, Ilya Shpitser, and Judea Pearl. 2005. Identifiability of path-specific effects. (2005).
- [2] Solon Barocas and Andrew D Selbst. 2016. Big data’s disparate impact. *Calif. L. Rev.* 104 (2016), 671.
- [3] Marc Bendick. 2007. Situation testing for employment discrimination in the United States of America. *Horizons stratégiques* 3 (2007), 17–39.
- [4] Peter J Bickel, Eugene A Hammel, and J William O’Connell. 1975. Sex bias in graduate admissions: Data from Berkeley. *Science* 187, 4175 (1975), 398–404.
- [5] Alexandra Chouldechova. 2017. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data* 5, 2 (2017), 153–163.
- [6] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*. 214–226.
- [7] Ronald Aylmer Fisher. 1992. Statistical methods for research workers. In *Breakthroughs in statistics*. Springer, 66–70.
- [8] David Galles and Judea Pearl. 2013. Testing identifiability of causal effects. *arXiv preprint arXiv:1302.4948* (2013).
- [9] Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of opportunity in supervised learning. In *Advances in neural information processing systems*. 3315–3323.
- [10] Wen Huan, Yongkai Wu, Lu Zhang, and Xintao Wu. 2020. Fairness through Equality of Effort. In *Companion Proceedings of the Web Conference 2020*. 743–751.
- [11] Yimin Huang and Marco Valtorta. 2006. Identifiability in causal bayesian networks: A sound and complete algorithm. In *Proceedings of the national conference on artificial intelligence*, Vol. 21. Menlo Park, CA; Cambridge, MA; London: AAAI Press; MIT Press; 1999, 1149.
- [12] Guido W Imbens and Donald B Rubin. 2015. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press.
- [13] Aria Khademi, Sanghack Lee, David Foley, and Vasant Honavar. 2019. Fairness in algorithmic decision making: An excursion through the lens of causality. In *The World Wide Web Conference*. 2907–2914.
- [14] Niki Kilbertus, Mateo Rojas Carulla, Giambattista Parascandolo, Moritz Hardt, Dominik Janzing, and Bernhard Schölkopf. 2017. Avoiding discrimination through causal reasoning. In *Advances in Neural Information Processing Systems*. 656–666.
- [15] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual fairness. In *Advances in neural information processing systems*. 4066–4076.
- [16] Joshua R Loftus, Chris Russell, Matt J Kusner, and Ricardo Silva. 2018. Causal reasoning for algorithmic fairness. *arXiv preprint arXiv:1805.05859* (2018).
- [17] Daniel Malinsky, Ilya Shpitser, and Thomas Richardson. 2019. A potential outcomes calculus for identifying conditional path-specific effects. *Proceedings of machine learning research* 89 (2019), 3080.
- [18] Razieh Nabi and Ilya Shpitser. 2018. Fair inference on outcomes. In *Proceedings of the... AAAI Conference on Artificial Intelligence*. AAAI Conference on Artificial Intelligence, Vol. 2018. NIH Public Access, 1931.
- [19] Catherine O’Neill. 2016. Weapons of math destruction. *How Big Data Increases Inequality and Threatens Democracy* (2016).
- [20] Judea Pearl. 2001. Direct and indirect effects. In *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*. 411–420.
- [21] Judea Pearl. 2009. *Causality*. Cambridge university press.
- [22] Judea Pearl and Dana Mackenzie. 2018. *The book of why: the new science of cause and effect*. Basic Books.
- [23] Kimberly Quick. 2015. The Unfair Effects of IMPACT on Teachers with the Toughest Jobs. *The Century Foundation* (2015).
- [24] Michelle Rhee. 2019. IMPACT: The DCPS Evaluation and Feedback System for School-Based Personnel.
- [25] James M Robins, Miguel Angel Hernán, and Babette Brumback. 2000. Marginal Structural Models and Causal Inference in Epidemiology. *Epidemiology* 11, 5 (2000), 551.
- [26] Andrea Romei and Salvatore Ruggieri. 2011. A multidisciplinary survey on discrimination analysis.
- [27] Ilya Shpitser. 2013. Counterfactual graphical models for longitudinal mediation analysis with unobserved confounding. *Cognitive science* 37, 6 (2013), 1011–1035.
- [28] Ilya Shpitser and Judea Pearl. 2007. What counterfactuals can be tested. In *23rd Conference on Uncertainty in Artificial Intelligence, UAI 2007*. 352–359.
- [29] Ilya Shpitser and Judea Pearl. 2008. Complete identification methods for the causal hierarchy. *Journal of Machine Learning Research* 9, Sep (2008), 1941–1979.
- [30] Ilya Shpitser and Judea Pearl. 2012. Identification of conditional interventional distributions. *arXiv preprint arXiv:1206.6876* (2012).
- [31] Elizabeth A Stuart. 2010. Matching methods for causal inference: A review and a look forward. *Statistical science: a review journal of the Institute of Mathematical Statistics* 25, 1 (2010), 1.
- [32] Jin Tian. 2004. Identifying linear causal effects. In *AAAI*. 104–111.
- [33] Jin Tian and Judea Pearl. 2002. A general identification condition for causal effects. In *Aaai/iaai*. 567–573.
- [34] Jin Tian and Ilya Shpitser. 2003. On the identification of causal effects. (2003).
- [35] Yongkai Wu, Lu Zhang, and Xintao Wu. 2019. Counterfactual Fairness: Unidentification, Bound and Algorithm. In *IJCAI*. 1438–1444.
- [36] Yongkai Wu, Lu Zhang, Xintao Wu, and Hanghang Tong. 2019. Pc-fairness: A unified framework for measuring causality-based fairness. In *Advances in Neural Information Processing Systems*. 3404–3414.

- [37] Junzhe Zhang and Elias Bareinboim. 2018. Equality of opportunity in classification: A causal approach. In *Advances in Neural Information Processing Systems*. 3671–3681.
- [38] Junzhe Zhang and Elias Bareinboim. 2018. Fairness in decision-making—the causal explanation formula. In *Proceedings of the... AAAI Conference on Artificial Intelligence*.
- [39] Lu Zhang and Xintao Wu. 2017. Anti-discrimination learning: a causal modeling-based framework. *International Journal of Data Science and Analytics* 4, 1 (2017), 1–16.
- [40] Lu Zhang, Yongkai Wu, and Xintao Wu. 2016. A causal framework for discovering and removing direct and indirect discrimination. *arXiv preprint arXiv:1611.07509* (2016).
- [41] Lu Zhang, Yongkai Wu, and Xintao Wu. 2016. Situation Testing-Based Discrimination Discovery: A Causal Inference Approach. In *IJCAI*, Vol. 16. 2718–2724.
- [42] Lu Zhang, Yongkai Wu, and Xintao Wu. 2017. Achieving non-discrimination in data release. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 1335–1344.