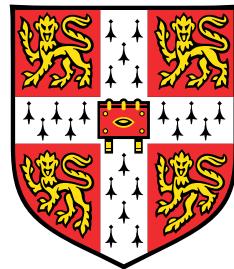# Beyond traditional assumptions in fair machine learning

**Niki Kilbertus**

Supervisor: Prof. Dr. Carl E. Rasmussen

Advisor: Dr. Adrian Weller

Department of Engineering
University of Cambridge

This thesis is submitted for the degree of
*Doctor of Philosophy*

Pembroke College                                             October 2020

# Acknowledgments

# Abstract

This thesis scrutinizes common assumptions underlying traditional machine learning approaches to fairness in consequential decision making. After challenging the validity of these assumptions in real-world applications, we propose ways to move forward when they are violated. First, we show that group fairness criteria purely based on statistical properties of observed data are fundamentally limited. Revisiting this limitation from a causal viewpoint we develop a more versatile conceptual framework, causal fairness criteria, and first algorithms to achieve them. We also provide tools to analyze how sensitive a believed-to-be causally fair algorithm is to misspecifications of the causal graph. Second, we overcome the assumption that sensitive data is readily available in practice. To this end we devise protocols based on secure multi-party computation to train, validate, and contest fair decision algorithms without requiring users to disclose their sensitive data or decision makers to disclose their models. Finally, we also accommodate the fact that outcome labels are often only observed when a certain decision has been made. We suggest a paradigm shift away from training predictive models towards directly learning decisions to relax the traditional assumption that labels can always be recorded. The main contribution of this thesis is the development of theoretically substantiated and practically feasible methods to move research on fair machine learning closer to real-world applications.

# Table of contents

# List of figures

# List of tables

<div align="right">

# 1

</div>

# Introduction and overview

## 1.1 Motivation

As machine learning penetrates all aspects of our everyday lives and the push for automation moves beyond industrial settings to decisions about human beings, we are facing a multitude of new challenges. While the performance goals for narrowly confined industrial automation tasks are often easy to express, automated decisions affecting people's livelihood and well-being must meet more complex requirements. We want to verify that these systems make ethically agreeable decisions and that they do so for the right reasons. These considerations include a broad range of concerns about fairness, trust, accountability, transparency, and privacy. Similar issues are also faced by human decision makers and have been the subject of various fields of research long before the advent of computers and large scale data analytics. However, software is a human creation and thereby, in principle, we expect to have full control over its actions. Hence, we may also hope to exert full control over the consequences of its deployment. However, modern data-driven machine learning systems can be notoriously hard to interpret and their downstream impact unpredictable. Moreover, because of their scalability, the potential for harm due to misspecifications of an automated system often by far exceeds what any individual human decision maker can incur. Unsurprisingly, undesirable effects of machine learning systems have already been observed in practice.

What do we mean by "undesirable effects"? There is an ongoing debate whether morality is innate or a socially constructed ideal among ethicists and philosophers. Arguably, the nuances of our present-day understanding still depend on subjective human judgment both on the level of individuals as well as societies and cultures. Therefore, throughout this thesis, we will avoid attempts of an objective definition or even opinion of "right or wrong" and "fair or unfair". This is not to say that there will not be definitions of fairness. In fact, we will encounter whole families of formally defined fairness criteria. However, we explicitly consider those to be mere candidates that may possibly cover some reasonable dimensions

of what is considered fair by some people within a given context. Even though our language will not always be as cautious as in the previous sentence, we ask the reader to consider what follows from this viewpoint.

Now, to illustrate the dangers of deploying machine learning systems in the social context, we briefly highlight a small collection of recent examples of harmful bias. In many cases we do not know the specifics of the underlying systems, in particular, whether they were powered by machine learning. They are meant to serve as clarifying examples that could have plausibly been generated by machine learning algorithms and share a vast scale of impact enabled by digital technologies. The wide deployment of such systems as well as their opaqueness have also been identified by Cathy O'Neil as key characteristics of "weapons of math destruction" (O'Neil, 2016). For a popular introduction containing plentiful examples of how algorithms can scale and accelerate inequality and injustice, we also point the reader to Noble (2018); Benjamin (2019); Broussard (2018); Eubanks (2018); Kearns & Roth (2019). For further reading on the risks and dangers of data-driven decision support systems and the sources of bias or discrimination from a more technical perspective, we refer the reader to the excellent online book by Barocas et al. (2019) as well as reports by Barocas & Selbst (2016a); Muñoz et al. (2016).

**Criminal justice.**   A decision support tool called *Correctional Offender Management Profiling for Alternative Sanction* (COMPAS) was used in U.S. courts to estimate the recidivism risk of defendants for decisions on pre-trial detention. Reporters at ProPublica claimed that COMPAS is biased against blacks, because among convicts that did not go on to re-offend within two years of their release, blacks had received systematically worse scores than whites (Angwin et al., 2016).

**Bias in search and recommendation.**   In web searches for names, Google's advertisement delivery algorithm AdSense was found to deliver ads suggestive of arrest more often for names primarily assigned to black people than for those assigned to white people (Sweeney, 2013). Furthermore, stereotype exaggeration and systematic underrepresentation of women were found in image search results for a variety of occupations (Kay et al., 2015).

**Exam grading.**   In a randomized study of discrimination in grading in India, it was found that teachers assigned worse scores to exams seemingly written by children from lower castes consistent with statistical discrimination (Hanna & Linden, 2012).

**Prime free same-day delivery.**   When Amazon rolled out free same-day delivery, the region of availability excluded predominantly black ZIP codes in some cities (Ingold & Soper, 2016).

Figure 1.1 These translation results on `https://translate.google.com` were obtained on August 11, 2018 and show stereotypical gender association of certain occupations.

**Image recognition.**   After a public media outcry over Google Photos labeling pictures of black people as gorillas in 2015 (Dougherty, 2015), Amazon received bad publicity in 2018, because their commercial facial recognition system "Rekognition" falsely matched 28 members of congress with mugshots. Black members were disproportionately affected (Snow, 2018). An academic study of commercially available gender classification software from images confirms substantial disparities in accuracy for different demographic groups formed by gender and skin tone with darker-skinned females being the most misclassified group (Buolamwini & Gebru, 2018).

**Bias in semantics derived from human language.**   Automatically deriving features and semantics from large corpora of human language is an important tool in natural language processing. Such systems have been shown to contain human-like biases, often reflecting problematic associations regarding race and gender (Caliskan et al., 2017). For example, popular word embeddings have been shown to establish a similar relation between "man" and "computer programmer" as between "woman" and "homemaker" among similar examples (Bolukbasi et al., 2016). These imprints of historical bias can also be found in machine learning powered translation software such as Google translate, see Figure 1.1. As yet another example, a chatbot named "Tay" launched by Microsoft to engage in conversations with people on Twitter in 2016 soon posted wildly racist tweets and reprehensible images as it learned from its conversations (Lee, 2016).

**Focus of this thesis.**   These examples underline the breadth of the potential for algorithmic discrimination. In this thesis we will focus on scenarios in which machine learning systems are involved in taking concrete decisions (typically one out of a finite set of possible actions) that directly affect individuals well-being and livelihoods. The criminal justice and exam grading examples fit well into this category. Other applications could be lending, hiring, or insurance decisions and treatment choices in health-care. Furthermore, we typically

assume the algorithm to directly trigger decisions deterministically. In practice, the more likely scenario may be one in which the algorithm only supports or informs human decision makers in various ways. Ultimately, the core issues we discuss remain regardless of whether the algorithm takes, supports, or informs decisions.

## 1.2   Outline and contributions

Albeit not immediately relevant for the technical contributions of this thesis, we will briefly dive a little bit deeper into fundamental questions of fairness and justice from a philosophical viewpoint in Chapter 2. Readers only interested in the methodological contributions of the thesis may skip this chapter.

Chapter 3 contains an introduction to existing approaches to fair machine learning for decision making. After an overview of the general setting and some notation, we focus on notions of fairness for classification or decision tasks with a finite set of actions. We specifically emphasize the assumptions required for the respective approaches and the resulting limitations. The remaining chapters each propose solutions that allow us to relaxing some of these restrictive assumptions.

In Chapter 4, we take a critical look at common statistical fairness criteria. Through the lens of causality, we crisply articulate why and when specific group fairness notions fail. These insights suggest to shift attention from the elusive question "What is the right fairness criterion?" to "What do we want to assume about our model of the causal data generating process?" We propose the language of causality as a useful technical tool to scrutinize and talk about such assumptions. Within that framework, we further expose previously ignored subtleties fundamental to the problem, e.g., how to interpret protected attributes in a causal sense. For example, in most of the real-world examples at the end of Section 1.1, we use terms such as "perceived" or "seemingly" in conjunction with group membership indicators such as "black" or "disadvantaged". We propose two complementary viewpoints for interpreting causal influences that enable practitioners to better navigate the conceptual fuzziness around these terms. Finally, we put forward natural causal non-discrimination criteria and develop predictive algorithms that satisfy them.

Subsequently, in Chapter 5, we also scrutinize one of the major assumptions of causal reasoning, namely that the true causal graph is known. Potential misspecifications of the assumed causal model can invalidate conclusions about the fairness of learned models. One common way for misspecification to occur is via *unmeasured confounding*: the true causal effect between variables is partially due to unobserved quantities. Causally fair classifiers critically rely on assumptions about the causal structure. In case the real world does not satisfy these assumptions, such a proposed classifier will actually be unfair during

deployment. Chapter 5 develops computationally efficient tools to assess the sensitivity of fairness measures to unobserved confounding. These allow decision makers to measure how unfair their proposed classifier becomes in case the real world deviates from the causal assumptions within quantitatively measurable limits. Hence, the main contribution of this chapter is to bring conceptually useful causal fairness notions closer to application by quantifying how they break, when the required assumptions do not hold perfectly and instead may be slightly violated. We demonstrate how to perform sensitivity analysis with our tools on two real-world datasets.

Chapter 6 focuses on the pressing question of how training of fair machine learning systems may compromise user privacy. Typically, sensitive data is required to train fair models. However, users may be wary of providing this data to decision makers. We suggest to tackle this predicament via cryptographic tools that allow us to train fair models with access only to encrypted sensitive data. The main contribution of this chapter is to extend existing *secure multi-party computation* protocols to also incorporate linear constraints in the optimization of machine learning models. Thereby, we bring fair model training into the realm of computing on encrypted inputs—in our case encrypting sensitive features for the training data. Moreover, we introduce a notion of accountability by allowing external entities to check the fairness properties of decision models and verify their outputs without revealing sensitive user data or the model itself, which may be considered secretive intellectual property.

Finally, in Chapter 7, we highlight that most traditional approaches to fair machine learning are based on predictive models learned from static datasets in a supervised fashion, implicitly assuming training data to be an i.i.d. sample from the distribution expected during deployment. However, as we derive decisions from these predictions, the data distribution observed at test time may depend on these decisions. Specifically, we consider a setting where labels only come into existence when a positive decision is made—if a loan is denied, there is not even an option for the individual to pay it back. Our first contribution is to formally show that in such a setting, if an imperfect predictive model is used for data collection, the observed data will not suffice to train an optimal predictive model via constrained risk minimization—the most common technique in existing works on statistical fairness. Instead we propose a paradigm shift away from minimizing a predictive loss with observed labels towards maximizing a utility that directly measures the quality of decisions in terms of accuracy and fairness. We show that this requires non-deterministic decision rules that "explore over all inputs", i.e., give positive decisions with non-zero probability for any input. Finally, we develop an algorithm to solve the utility maximization and demonstrate its efficacy in terms of accuracy and fairness on synthetic and real-world data. Beyond enabling fair model training under more realistic circumstances, namely when labels only exist for certain

decisions, this chapter also formalizes why it is crucial to distinguish between predictions and decisions in a social context.

Chapter 8 first summarizes the contributions of the thesis. We then draw our conclusions and highlight directions for future work before eventually circling back to the broader context of injustice and discrimination.

## 1.3 Publications

This thesis is based on the following publications.

Chapter 4 is based on Kilbertus et al. (2017):

---
AVOIDING DISCRIMINATION THROUGH CAUSAL REASONING

Niki Kilbertus, Mateo Rojas-Carulla, Giambattista Parascandolo, Moritz Hardt, Dominik Janzing, Bernhard Schölkopf.

*Neural Information Processing Systems (NeurIPS), 2017*

---

Chapter 5 is based on Kilbertus et al. (2019):

---
THE SENSITIVITY OF COUNTERFACTUAL FAIRNESS TO UNMEASURED CONFOUNDING

Niki Kilbertus, Philip J. Ball, Matt Kusner, Adrian Weller, Ricardo Silva.

[https://github.com/nikikilbertus/cf-fairness-sensitivity]

*Uncertainty in Artificial Intelligence (UAI), 2019*

---

Chapter 6 is based on Kilbertus et al. (2018a):

---
BLIND JUSTICE: FAIRNESS WITH ENCRYPTED SENSITIVE ATTRIBUTES

Niki Kilbertus, Adrià Gascón, Matt Kusner, Michael Veale, Krishna P. Gummadi, Adrian Weller.

[https://github.com/nikikilbertus/blind-justice]

*International Conference on Machine Learning (ICML), 2018*

---

Chapter 7 is based on Kilbertus et al. (2020b):

---
FAIR DECISIONS DESPITE IMPERFECT PREDICTIONS

Niki Kilbertus, Manuel Gomez-Rodriguez, Bernhard Schölkopf, Krikamol Muandet, Isabel Valera.

[https://github.com/nikikilbertus/fair-decisions]

*International Conference on Artificial Intelligence and Statistics (AISTATS), 2020*

---

During my PhD I also contributed to the following papers, which will not be described in this thesis (Hron et al., 2020; Kilbertus et al., 2020a; Träuble et al., 2020; Gebhard et al., 2019; Kilbertus et al., 2018b; Parascandolo et al., 2018; Gebhard et al., 2017):

---

EXPLORATION IN TWO-STAGE RECOMMENDER SYSTEMS

Jiri Hron*, Karl Krauth*, Michael I. Jordan, Niki Kilbertus.

*ACM RecSys* workshop *REVEAL: Bandit and Reinforcement Learning from User Interactions, 2020*

---

A CLASS OF ALGORITHMS FOR GENERAL INSTRUMENTAL VARIABLE MODELS

Niki Kilbertus, Matt J. Kusner, Ricardo Silva.

*Neural Information Processing Systems (NeurIPS), 2020*

---

IS INDEPENDENCE ALL YOU NEED? ON THE GENERALIZATION OF REPRESENTATIONS LEARNED FROM CORRELATED DATA

Frederik Träuble, Elliot Creager, Niki Kilbertus, Anirudh Goyal, Francesco Locatello, Bernhard Schölkopf, Stefan Bauer.

*under review*

---

CONVOLUTIONAL NEURAL NETS: A MAGIC BULLET FOR GRAVITATIONAL-WAVE DETECTION?

Timothy Gebhard*, Niki Kilbertus*, Ian Harry, Bernhard Schölkopf.

[https://github.com/timothygebhard/magic-bullet]

*Physical Review D, 2019*

---

GENERALIZATION IN ANTI-CAUSAL LEARNING

Niki Kilbertus*, Giambattista Parascandolo*, Bernhard Schölkopf.

*Neural Information Processing Systems (NeurIPS)* workshop on *Critiquing and Correcting Trends in Machine Learning, 2018*

---

LEARNING INDEPENDENT CAUSAL MECHANISMS

Giambattista Parascandolo, Niki Kilbertus, Mateo Rojas-Carulla, Bernhard Schölkopf.

*International Conference on Machine Learning (ICML), 2018*

---

SEARCHING FOR GRAVITATIONAL WAVES WITH FULLY CONVOLUTIONAL NEURAL NETS

Timothy Gebhard*, Niki Kilbertus*, Giambattista Parascandolo, Ian Harry, Bernhard Schölkopf.

[https://github.com/nikikilbertus/convwave]

*Neural Information Processing Systems (NeurIPS)* workshop on *Deep Learning for Physical Sciences, 2017*

---

* equal contribution

# 2

# Background

In this chapter we provide some background about the fundamental challenges around fairness and justice from various viewpoints. We take a top down approach, starting from abstract philosophical notions of justice, to concrete social challenges related to discrimination, to the specifics of machine learning. While we continuously narrow down the scope, we hope that the big picture context provided in this chapter helps the reader not to lose sight of the core issues once we dive into technical details in the following chapters. Readers only interested in the technical contributions can skip this chapter. It starts with terminology and connections to other disciplines in Section 2.1, illuminating how terminology is used throughout the thesis and providing broader context for how to think about fairness in machine learning. This discussion naturally leads to the question, whether we can reasonably expect to be able to practically mitigate unfairness in machine learning systems. In Section 2.2, we explore on a high level which challenges we are facing when trying to achieve lasting positive impact on society via automated decisions. There, we will also introduce the main theme of the thesis: *scrutinizing common assumptions underlying traditional approaches to fairness in machine learning and proposing solutions for when these assumptions are violated.* After a brief overview of where and how bias can enter the machine learning pipeline in Section 2.3, we then summarize and conclude this introductory chapter in Section 2.4.

## 2.1   Terminology and insights from other disciplines

While some sort of immoral bias is present in all the examples in Section 1.1, it can be hard to pinpoint—let alone quantify—what exactly is "wrong". There are multiple dimensions along which one can characterize these difficulties to break down the field of algorithmic bias into more specific subproblems.

In this section we take a step back and assay the nature of such problems from a philosophical and ethical viewpoint. Among other disciplines, scholars in ethics, philosophy, and law have been thinking about issues of justice, fairness, discrimination, and bias from different

perspectives for much longer than the computer science community. If we want to mitigate ethically objectionable impact of machine learning models, we need to acquire at least a basic understanding of why and what is considered "unfair" and to what extent we can hope this understanding to be universal. We cannot provide a faithful account of centuries of research in philosophy, ethics, and law. Thus this section is neither exhaustive nor unbiased. It serves to illustrate some relevant viewpoints and provide useful terminology. As a starting point for further reading on fairness from the perspective of political philosophy for computer scientists, we refer to Binns (2018), who elucidates some connections between recent machine learning research and philosophy. For a purely philosophical treatise, we suggest readers start with a modern overview of philosophical theories of justice by Sandel (2010) and then dive deeper into more focused works (Bentham, 1780; Young, 1995; Roemer, 1998; Moulin, 2004; Roemer, 2009; Rawls, 2009).

**Terminology and definitions.**　　Within the machine learning community, *fairness* (as in *fairness in machine learning*, *fair learning*, or *algorithmic fairness*) developed as the prevailing umbrella term for a range of concerns about unethical behavior or impact of automated systems. The word "fairness" is multi-faceted and can have different technical meanings across disciplines. Just to name a few examples, *fair value* in economics often refers to a rational estimate of the potential market price of a good. In other fields of computer science, fairness is discussed both as a property of concurrency in the context of unbounded nondeterminism, as well as in the context of network engineering. There, it describes how resources are shared and allocated among applications. *Fair division* in a game theoretic context typically describes the division of goods among players with subjective valuations subject to various constraints. Within economics, the term *fairness* was used to describe a market anomaly in which firms do not strictly maximize profits, because it may be considered unfair to raise prices to exploit shifts in demand (Kahneman et al., 1986). This notion of fairness was an important concept in the emergence of behavioral economics. Within each of these domains, fairness describes a rather narrow concept and there is typically no universally accepted definition. Occasional confusion is expected when narrow mathematical constructs are named after real life problems.

Within machine learning, fairness has primarily discussed distributive aspects of social justice. Since empirical evidence suggests that people do not reach a consensus in defining fairness intuitively (Yaari & Bar-Hillel, 1984), distributive justice in the social context typically refers to *subjective fairness* rather than *objective fairness*. Yaari & Bar-Hillel (1984) conclude that

> "[. . . ] a satisfactory theory of distributive justice would have to be endowed with considerable detail and finesse. Sweeping solutions and world-embracing theories are not

*likely to be adequate for dealing with the intricacies inherent in the problem of How to*
*Distribute."*

Note that the problem of *how to distribute* considers a wide variety of goods and services, including decisions and outcomes in applications such as criminal justice, hiring, insurance, or lending. This conclusion should remind us that—as a relatively young community—we may be prone to misunderstandings. Currently, we do not expect a single technical definition to emerge as the "right" one. Therefore, we must be transparent about our assumptions, be specific about context, and refrain from relying on an intuitive understanding of terminology used in our definitions.

We have already used the words *fairness* and *justice* almost synonymously without further explanation. Together with *discrimination* and *bias* they are ubiquitous in the literature on fairness in machine learning. There is no general agreement on their meaning, neither in a technical nor necessarily an intuitive sense. One usually has to rely on the context for proper interpretation. In the following, we provide clarifying remarks on the differences and connections between those terms, starting with popular dictionary definitions of fairness and justice.

**Fairness**  is the quality of making judgments that are free from discrimination.

**Justice**  is an action that is morally right and fair.

From this perspective, fairness may be considered an idealized concept, whereas justice is about taking the right actions (potentially when unfairness has already occurred).

John Rawls even argues that fairness is more fundamental than justice (Rawls, 1958):

*"It might seem at first sight that the concepts of justice and fairness are the same, and that there is no reason to distinguish them, or to say that one is more fundamental than the other. I think that this impression is mistaken. In this paper I wish to show that the fundamental idea in the concept of justice is fairness;"*

He further considers fairness to imply mutual agreement among parties that no one is taken advantage of, or forced to give in to claims that anyone of them considers illegitimate. Justice on the other hand, does not necessarily require such mutual agreement, as is usually the case in law enforcement.

From a more applied perspective, Friedman & Nissenbaum (1996) define *bias* in one of the earliest works on fairness in computer systems as follows:

*"In its most general sense, the term bias means simply "slant." Given this undifferentiated usage, at times the term is applied with relatively neutral content. [...] At other times, the term bias is applied with significant moral meaning. [...] we use the term*

*bias to refer to computer systems that systematically and unfairly discriminate against certain individuals or groups of individuals in favor of others. A system discriminates unfairly if it denies an opportunity or a good or if it assigns an undesirable outcome to an individual or group of individuals on grounds that are unreasonable or inappropriate."*

By now, the reader has certainly noted the inherent circularity or ambiguity of these definitions. While there may be an innate mutual understanding of certain moral values (Turiel, 2002; Blair, 1995), other aspects of morality are inherently subjective. This is reflected for example by considerably different legislation across nations. Hence, the circularity of the provided definitions as well as their reliance on similarly fuzzy terms such as "right", "unfairly", "unreasonable", and "inappropriate" merely reflect the elusive nature of the subject.

As mentioned above, we will consider systems that take, trigger, or inform actions in the physical world. Following the above dictionary definitions, we are less concerned with fairness (as an idealized concept), than with justice (taking the right action). Even though *non-discriminatory* or *just* may often be more appropriate terms than *fair*, the community has converged to the term *fairness*. It will be instructive for categorizing existing work on fairness, to borrow some more insights from philosophers' viewpoints on justice. Instead of describing general frameworks such as consequentialist, deontological, contractualist, or virtue ethics, we will right away dive into distributive justice following the focus area of the thesis.

**Philosophical theories of justice.** Philosophical theories of distributive justice need to specify which goods are to be distributed among which entities, and what constitutes a proper distribution. They do not answer questions about whether a particular distribution is objectively correct, or who should have the right to choose and enforce a certain distribution. Opposed to this *distributive* approach—defining a favorable distribution in a consequentialist fashion—we can also follow a certain procedure to assign outcomes to individuals, called *procedural approach* to fairness, which is closer to a deontological approach. By adhering to such a (defined to be) fair procedure, we consequently accept the arising distribution as fair.

We can further divide fairness notions along another dimension into *normative* (or *prescriptive*) and *descriptive* (or *comparative*) approaches. Crudely, the former start from a fixed distribution or procedure, which is assumed to be fair by general agreement and thus ought to be reached or followed. Contrary, *descriptive* approaches seek to empirically crowd source society's opinion on what is perceived—and thus accepted to be—fair.

Most democracies today deploy a somewhat procedural approach by accepting the distribution that arises naturally from adhering to the law and social norms. Of course, there is substantial feedback in that the procedure, i.e., the legislation, is heavily influenced by

critically reflecting on and reacting to the resulting distribution. Similarly, while the legal system has a predominantly normative character, the prescription is largely determined by a comparative assessment of society's opinion at large—at least in most democracies. Consequently, we argue that our reality of justice hinges on a procedural approach that is normative in its execution, but heavily influenced by a comparative democratic process and regularized by theoretic distributive ideals. A large body of work on fairness in machine learning seeks to derive technical fairness criteria from normative, distributive theories. While this can be an instructive endeavor, one should not lose sight of the current procedural legal practices and descriptive influences (Green & Hu, 2018).

## 2.2   Fundamental challenges of fair decision making

Let us now dissect specific challenges of fair decision making and point out important concepts to factor into the process.

**The role of historicity.**   In a fictitious world where our moral understanding (and thus our legal system) had been stationary over time, and all humans had always acted morally benign (lawfully), it is conceivable that nobody would voice concerns about unfairness or discrimination. This fictitious world is the ideal we strive for when fighting discrimination. The thought experiment also highlights that virtually all situations, in which fairness arises as an issue, are invariably tied to historicity. Unfairness occurs when past or current practices conflict with today's moral or legal standards. Considering these historical events and practices to leave behind unfavorable disparities as an imprint on our society, we have to consider multiple goals for fair machine learning systems.

One may choose to accept existing disparities to be an unalterable truth about the state of our world that can fairly be exploited in a decision making process. Such a viewpoint could be advocated for by not wanting to take blame for historic wrongdoing. Following this idea of "the data are what the data are", or "what you see is what you get" (Friedler et al., 2016), the goal could be to not amplify, perpetuate, or create new harmful biases. Decision makers may be drawn to such arguments when adhering to stricter correctional fairness measures would negatively impact their utility. On the other hand, we may seek to correct for past misdoing by taking affirmative action to dynamically transition the current, unfair situation to a favorable state that is better aligned with today's moral standards. This often requires an understanding of how society evolves dynamically as a consequence of decisions, i.e., how such interventions feed back into the state of society and social dynamics.

Let us give two examples showing that both modes of reasoning are not entirely clear-cut. Is it fair if a warehouse owner only employs individuals that can lift 40 pounds even though this may be heavily correlated with gender and age? Or should they be required to work

towards equalizing such differences, perhaps by providing free access to strength training? Another example in which such questions become particularly relevant is retributive or corrective justice: to what extent does punishment contribute to a better society? Is it fair to imprison people with a higher potential to commit more severe crimes even if that correlates with race? Or should we instead focus more resources on supporting these communities to fight the root cause of their criminal activity?

While we have seen in the previous section how theories of justice may describe desirable states or actions, there are yet more choices to be made when setting out to operationalize these criteria. Beyond the generic issue that there is no principled way to measure, evaluate, or compare different fairness notions—especially for normative approaches—we also need to discuss appropriate fairness goals as well as their time horizon in any given context.

**Context sensitivity and long- versus short-term interventions.** Binns (2018) provides three compelling examples to highlight the context dependence of fairness judgments.

1. When it comes to **voting**, we tend to favor an egalitarian outcome distribution: every person has one and only one vote regardless of skill or effort.

2. However, the same does not hold true for **job applications**. While we still generally seek to provide equal opportunity, it is morally acceptable to condition on skill and effort. In this scenario, we may find that existing disparities in the distribution of skill and effort can be explained by systematic historical discrimination and its perpetuation, leading to a debate whether we should take affirmative action or accept the status quo as a property of society. After all, the distributions simply differ and it may be hard to determine whether the root cause is past discrimination, or a legitimate alternative explanation. Rawls (2009) argues that

   > *"The natural distribution is neither just nor unjust; nor is it unjust that persons are born into society at some particular position. These are simply natural facts. What is just and unjust is the way that institutions deal with these facts."*

   There is some ambiguity in whether the *natural distribution* is the one observed today (potentially reflecting historic bias), or some idealistic, unobservable distribution. Rawls' statement represents the viewpoint that justice is about the actions we take in a given situation. However, it does not answer how precisely to deal with disparities in the natural or observed distribution.

   This begs naive questions such as "How much of the past should we aim to correct for?" or "How far into the future should we aim to reach a desirable state?" Partially motivated by these questions, Hu & Chen (2018) study a simplified dynamic model of the labor market, and argue that *"A dynamic model recognizes the powerful ripple effect of*

*the past and calls for a fairness intervention that carries momentum into the future.*" Similarly, Liu et al. (2018) analyze a one-step feedback model of how decisions according to some fairness criteria affect the well-being of different groups. They show that enforcing common fairness criteria (see Chapter 3) may actually cause harm in the long run. This indicates that there is a trade-off between optimizing for short-term (doing the right thing now) and long-term (achieving the right state in the future) fairness goals even when measured with the same criterion.

3. In his third example, Binns (2018) considers **airport security screenings**. One may find it appropriate, perhaps due to social solidarity, to distribute scrutiny equally among all individuals in an egalitarian fashion. However, the data indicate real differences in the base rates, i.e., some visually identifiable subgroups of the population are statistically at higher risk of attempting a terroristic attack. Is it fair to use a predictive system to subject this group to more rigorous examination? After all, statistically speaking, focusing more attention to these groups can prevent unnecessary deaths.

He concludes that "*We therefore can't assume that fairness metrics which are appropriate in one context will be appropriate in another.*" Thus, the appropriateness of different fairness criteria can only be judged in a narrow context and with a good understanding of the specific historical, sociological and legal perspective.

**Subjectivity and non-stationarity.** We have argued before that fairness is not only context dependent, but also subjective. Moral values vary greatly across nations and cultures, which is also reflected in the respective legal systems. Even within specific cultures, individuals often hold wildly different opinions on questions of justice and morality. Again, this contributes to the difficulty of nailing down precise notions of fairness, even after fixing the domain and specific goals. Arguably, what matters to the individual is the comparative hardship suffered by those at risk of being discriminated against and the mechanisms or circumstances leading to such disparate hardship. Hence, beyond domain knowledge, a deep empathic understanding of experiences and emotions is crucial to devise effective fairness interventions.

Looking back, our understanding of fairness and egalitarianism has changed drastically over time. Today, one may find the late onset and slow progress of worldwide emancipation for gender equality, continuing racial discrimination, or the idea of capital punishment shocking—if not repelling. We speculate that past societies on average felt roughly as assertive about their contemporary moral values as we do now. They may similarly have found preceding practices repugnant. This indicates that even if we could agree on specific fair machine learning methods today, following generations may well despair over our current moral wrongdoing and struggle to deal with the biased imprints we have left behind.

**Individual perception and the comparative nature of fairness.**   Specific personal complaints about unfairness are typically based on comparative arguments. For example, when being denied a loan, a black person may complain that their white neighbor, who has an otherwise similar socio-economic background, received a loan. On an individual level we tend to care less about being treated according to some abstract fairness maxim than about how we are treated relative to others. In such comparisons, the similarity with respect to all relevant aspects except for group membership is crucial. We regularly even ponder counterfactual questions like "Would I have gotten the loan if I were white?" or "Would I get paid more if I were a man?". Such counterfactuals can be thought of as comparisons with a version of oneself, where only the group membership in question has been changed. The contrast between this comparative nature and the distributive theories discussed in the previous sections leads to the distinction between *group* and *individual fairness* common in the machine learning literature.

Most of our discussion so far was concerned with fairness on a group level. In contrast, individual fairness demands that, colloquially, "similar people are treated similarly". This notion goes back to Aristotle's principle that, paraphrased, "equals should be treated equally". A formalization of this idea requires a measure of similarity for both individuals and outcomes (Dwork et al., 2012). For now, we note that any concrete distributive ideal of what would constitute a perfectly fair state of society may still inevitably lead to perceived unfairness on an individual level, especially when trying to move from a non-ideal state to the ideal one. We will come back to technical definitions of both group and individual fairness notions in Chapter 3.

**Moving to machine learning.**   From the picture we have drawn so far, it becomes apparent that fair decision making can only be successfully tackled on a systemic level, taking into account all factors from historic and current societal circumstances to design decisions by individual engineers. Narrowing down the scope further, we now describe the specific context for the machine learning systems that we will later analyze in isolation.

We assume a fixed set of demographic groups, at least one of which face disadvantages when it comes to access to certain opportunities, goods, or services. Such groups are often declared by the legal code. For example, the Equality Act (2010) in the United Kingdom specifies the following *protected characteristics*[1]: age, gender reassignment, being married or in a civil partnership, being pregnant or on maternity leave, disability, race including color, nationality, ethnic or national origin, religion or belief, sex, and sexual orientation. Discrimination with respect to these characteristics is legally prohibited at work, in education, as a consumer, when using public services, when buying or renting property, and as a

---

[1]https://www.gov.uk/discrimination-your-rights; http://www.legislation.gov.uk/ukpga/2010/15/contents

member or guest of a private club or association. In Germany, according to the "Allgemeines Gleichbehandlungsgesetz",[2] it is illegal to discriminate based on age, sex, sexual identity, disability, race or ethnic origin, and religion or belief within a similar scope. As a final example, in the United States legally recognized protected classes include race, color, sex, religion, national origin, citizenship, age, pregnancy, familial status, disability status, veteran status, and even genetic information by a number of legal texts such as the Civil Rights Act (1964, 1968), the Equal Pay Act (1963), the Immigration Reform and Control Act, the Age Discrimination in Employment Act (1967), the Pregnancy Discrimination Act, the Rehabilitation Act (1973), the Americans with Disabilities Act (1990), the Vietnam Era Veterans' Readjustment Assistance Act (1974), or the Genetic Information Nondiscrimination Act.[3]

Against the backdrop of these legal regulations, we commonly refer to *protected* or *sensitive attributes* in the following chapters without further explanation. Moreover, we assume that individuals can be unambiguously assigned to one of a finite set of such socially salient groups relevant to the given decision scenario. We highlight that this is a crude simplification, which can by itself introduce new issues. These challenges around group membership on an algorithmic level are touched upon in the next section and again briefly in Chapter 4. We note that legislation does not provide a moral framework, let alone concrete non-discrimination definitions that could be formalized mathematically or operationalized algorithmically.

Moreover, we restrict ourselves to scenarios in which machine learning systems are used to either fully take over or otherwise inform and support decisions that directly affect the well-being and livelihood of individuals. We often refer to such situations as making *consequential decisions*. In this context, the goal of group fair machine learning is to make consequential decisions such that no protected group is disadvantaged or discriminated against. Such discrimination could manifest itself in different ways. For example, different groups could be treated differently—there are different decision processes in place for different groups—which may be unethical. On the other hand, some group may experience undesirable downstream impacts even though—or even because—inputs are processed in the exact same way for everyone. We will return to ways of quantifying specific dimensions of disadvantage or discrimination in Chapter 3.

Since the focus of this thesis is on data-driven systems, data itself plays a crucial role. What we have previously called "the state of society" or "status quo", to a machine learning system is typically represented by the data, which is also its main input. The first essential step towards fairness, even before analyzing the algorithms themselves, is to understand

---

[2]https://www.bmas.de/DE/Service/Gesetze/allgemeines-gleichbehandlungsgesetz.html
[3]We thank Moritz Hardt for summarizing these legally protected classes and their corresponding legal texts in public talks.

whether data adequately represents the "status quo" in sufficient detail with respect to the given task and context.

Even more broadly, one must take into account the entire socioalgorithmic system they are embedded in, including data collection and downstream impact of decisions, which may alter the data we get to collect in the future. In the next section, we will highlight some ways in which bias can creep into the different stages of what Barocas et al. (2019) call the *machine learning loop*. For more details and examples we refer the reader to the work of Barocas et al. (2019, Chapter 1) as well as popular books on the topic (O'Neil, 2016; Noble, 2018; Broussard, 2018; Eubanks, 2018; Benjamin, 2019; Kearns & Roth, 2019).

## 2.3 From data to decisions and back

**Measurement.** Machine learning algorithms see the world through data. While architecture choices, learning algorithms and inductive biases also play an important role, final decisions are most affected by the training data. Hence, we must not assume that data are an unemotional and unopinionated mirror of the "true state of the world" that we can blindly trust in. Instead they must be scrutinized, distrusted, and continuously analyzed as part of a larger data-driven system. Unfortunately, machine learning engineers often seem to believe that their job only starts after the measurement phase. Data are readily available and rarely questioned, but instead treated like mere facts about the state of the world. However, measuring data about humans is a different process from, say, a physical measurement using a well-understood and calibrated measurement device. Often, we seek to measure social constructs, such as the race or gender a person identifies with, for which we do not have well-defined scales. Even the available options frequently change. Which properties to measure as well as the scales on which to measure them are chosen by a small set of individuals. In addition, the measurements themselves may be self-reported, introducing a new host of potential difficulties.

The lack of objective ground truth about gender or race assignments poses great conceptual challenges when we discuss causal notions of fairness. Broadly speaking, in causal modeling we care about how certain quantities causally affect each other, i.e., for which there are invariable mechanisms that ensure that changing one variable will inevitably change another in a specific way. In fairness, naturally arising questions are how a protected attribute, say race or gender, causally influenced a consequential decision. For such statements to make sense, we need to clearly define what any given node, especially nodes marked as protected, references. Most work on causal fairness notions consider protected attributes such as race or sex to take on a fixed value at birth. Thus all other measurable features come later and are therefore causal descendants of the protected attribute. We will challenge this viewpoint to some extent in Chapter 4. However, there are even more fundamental ontological and

epistemic questions about the validity of causal models, which are described with great clarity in Barocas et al. (2019, Chapter 4). Recently, scholars started to challenge the stability of the typically proposed ontologies of causal statements Kohler-Hausmann (2018); Hu & Kohler-Hausmann (2020); Dembroff et al. (2020). We believe this debate to become of crucial importance not only for causal modeling of fairness in the machine learning realm, but for a deeper understanding of quantifying discrimination and injustice in general.

Even when we do not collect data about humans, the measurement process may capture social disparities. For example, there has been an effort to map out potholes in the city of Boston with a smartphone app that automatically reports the location of potholes recognized by its sensors. Crawford (2013) points out that the collected data reflect different levels of smartphone penetration for different demographics. As a consequence, disproportionately fewer potholes have been recorded in lower-income areas and regions with predominantly elderly inhabitants. Barocas et al. (2019) describe the inherent messiness of measurement as a "*manifestation of the limitations of data-driven techniques*".

While data are important, they are far from the only issue. In particular, we cannot expect to fix all concerns about discrimination by only working on datasets. While one may hope to carefully design and collect an "unbiased dataset" that precludes any systematic bias downstream the machine learning pipeline, such efforts are futile.

**Existing disparities.**   Most machine learning techniques in the social context still cling to the traditional objective of maximizing accuracy (or some related performance measure) with respect to some *target variable* or *label*. We will see in Chapter 3 that many attempts at fair machine learning merely add constraints to this eager optimization paradigm to incorporate fairness. As the addition of constraints can at most reduce the achieved accuracy, this gives rise to a tension between fairness and accuracy. In the literature this has often been referred to as "the cost of fair classification" or the "fairness accuracy trade-off".

However, the constrained optimization approach to fairness appears to be inconsistent in the implicit assumptions about the recorded target labels. On the one hand, suggesting fairness constraints acknowledges perpetuated stereotypes, a history of explicit discrimination, and other factors that lead to statistical disparities in today's society, and thus in the recorded target labels. On the other hand, accuracy with respect to these recorded target labels is still considered a meaningful goal to optimize for as the main objective. Issues surrounding the validity of the recorded labels as an optimization target are referred to as forms of *label bias*. Recently Wick et al. (2019) describe this dualism well in that "*[. . . ] phenomena such as label bias are appropriately acknowledged as a source of unfairness when designing fair models, only to be tacitly abandoned when evaluating them.*" While making assumptions and especially also ethical judgments about the available data (and the distribution it comes from) explicit can help on a conceptual level (Friedler et al., 2016), there still exists no widely accepted

way of incorporating and dealing with such assumptions algorithmically. In other words, despite several attempts, ethical frameworks do not directly allow for operationalization in algorithmic terms.

As a consequence, the definition of the target variable plays a prominent role and its choice is heavily influenced by the parties involved in a given decision scenario. Moreover, the chosen measurable target is often a proxy for the "true goal", an elusive, abstract (social) construct. For instance, we use credit score as a proxy for creditworthiness, likelihood of recidivism for the influence of incarceration on character and behavior, or perhaps even generated revenue for the overall value of an employee. Similar considerations hold true for the integrity of the features used as inputs to a machine learning pipeline. Together with the messiness of measurement and the subjective choice of (proxies for) relevant variables, existing disparities become a complex source of bias that is hard to model and account for in observed datasets.

**Dynamics, adaptivity and feedback.**   Just like the machine learning loop begins before the modeling part, namely with data collection, it does not end after a prediction has been made either. Consequential actions affect individuals and thus the communities they are embedded in. Again, this is in stark contrast with natural sciences, where the laws of physics do not react to our actions and we can justifiably consider ourselves passive observers of the outcomes of our experiments (barring measurements in quantum mechanics). This interpretation breaks for algorithmic predictive systems, as soon as their predictions are revealed and influence actions. In consequential decision settings, traditional assessment of machine learning on static datasets is rarely a good measure. This observation is rooted in the difficulty of controlling deployment (or test-time) conditions for consequential decisions. In particular, there is a host of plausible mechanisms through which the deployment conditions may depend on the choice of the decision policy.

For example, a traffic congestion prediction may cause drivers to congest routes shown as clear and clear up routes predicted to be congested (Barocas et al., 2019). In general, widely accessible predictions by automated systems may influence people's actions, thereby breaking the assumptions that went into the prediction, which consequently invalidates itself. This is related to self fulfilling prophecies, where a prediction of a raise in stock prices can cause growing demand, which in turn actually leads to an increase in price. Similar feedback loops have been observed in predictive policing (Lum & Isaac, 2016; Ensign et al., 2018b), as well as in the effects of pretrial detention on conviction, future crime and employment (Dobbie et al., 2018).

In another form of adaptivity, individuals may strategically and deliberately invest effort to change their features (or disclosure thereof) adaptively to receive favorable decisions (Hardt et al., 2016a; Dong et al., 2018). These efforts can either be an attempt to *game the system* or to

*legitimately self-improve*, warranting interpretations as either a nuisance to be countered by robust classification, or an opportunity for mechanism design (Miller et al., 2020). Recently, it has been shown that the different costs of strategic behavior for different groups can give rise to further disparities (Milli et al., 2019). Similarly, the mere act of forming a prediction or an immediate consequence thereof can influence the distribution of the prediction target, which is referred to as *performativity* (Perdomo et al., 2020). Finally, the data we get to observe to train a decision system may depend on the decisions taken. For example, in a *selective labels* setting we only get to observe true outcomes when we take a positive decision (we only get to observe repayment of a loan if one was granted in the first place). We will return to this specific setting in Chapter 7. Differences in the ability to adapt and the environmental dynamics for members of different demographic groups can result in perpetuated or even amplified disparities. Since anticipating and modeling such complex societal interactions is generically error prone, they pose a serious challenge for fair machine learning.

**Closing the loop.**   We have highlighted some intricately interrelated challenges when moving discussions about justice and fairness from the philosophical to the technical realm of machine learning algorithms. Our moral valuation of the current state of affairs is hard to formalize mathematically, measurement and data collection is a messy, error-prone process, and our decision systems affect both in non-obvious ways which are hard to model. In short, our decision algorithms are part of an evolving socioalgorithmic system with numerous feedback loops and interactions. Hence, their immediate and long-term consequences are hard to predict, let alone control.

## 2.4   The role of algorithms and their developers

To conclude this chapter, we raise the question to what extent machine learning is the right tool to tackle issues of fairness and discrimination. One of the key goals is to automate decisions at scale. At the same time, we argued in this chapter that fairness requires a case-by-case analysis of the context and goals. Can machine learning still be a suitable tool to increase justice and diversity? Which goals are we equipped and entitled to take on as machine learning researchers and engineers? Should all algorithmically supported decisions be closely monitored by humans and judged by our current moral values, or should we even entrust them with shaping our society on their own terms, guiding the way to a more equitable future? Let us exemplify the importance of these questions.[4]

Pretrial risk assessment is usually framed as a predictive task in trying to estimate the risk of a person not appearing to court or recidivating in the meantime. A high risk of not appearing to court or of recidivism is then converted into a decision to detain the defendant. Once we

---

[4]We thank Moritz Hardt for pointing us to this example.

accept this viewpoint, machine learning indeed has a lot to offer in terms of making accurate predictions and thus seemingly contributing to the larger social good. However, in this example it may be more productive to challenge the basic assumption rather than specifics of the predictive algorithms and to recognize that people fail to appear because they do not have access to child care or transportation, cannot afford not to go to work, or even have other overlapping court appointments. Hence, appropriate countermeasures would include enabling functional two-way communication and provide child care and transportation support to defendants. Indeed some of these measures have been part of a Harris County lawsuit settlement in 2019 (Despart, 2019). However, whenever algorithms are used within a larger system, they may also be at fault. To still make progress in those settings, we will typically require strict simplifying assumptions about which effects we consider in a given problem formulation. We will clearly state these assumptions near the beginning of each of the Chapters 4-7.

We will not engage in a broader discussion of unequal power and access to machine learning and the resulting responsibilities. We only note that while there are democratic mechanisms in place to prevent individual policy makers from instilling non-representative, subjective opinions into our laws, no such mechanisms currently exist for machine learning engineers, software developers, and data collectors in both industry and government. To date, the demographic makeup of these professions considerably misrepresents society as a whole, which leaves us worrying about a lack of diversity of viewpoints when designing and building potentially impactful algorithms. Harms in deploying algorithmic systems can only be anticipated and detected if a diverse team with a broad set of viewpoints, experiences and conceptual frameworks dedicates conscious effort into it. While diversity and inclusion are perhaps more important for designing equitable systems than algorithmic considerations, a thorough discussion thereof goes beyond the scope of this thesis. We refer the interested reader to recent work by Mohamed et al. (2020), who explore connections between society and machine learning from the viewpoint of critical science and decolonial theories.

After this rather pessimistic perspective of fair machine learning, we must not disguise that we are already facing very real concrete issues with automated decisions. Even though we may not hope to solve them entirely, some are understood well enough to be readily improved. Since large scale machine learning applications are still commonly deployed without any consideration for fairness throughout the process, even imperfect methods can have a substantial positive impact. Moreover, algorithmic methods may also help to uncover and quantify existing discrimination in the first place (Kleinberg et al., 2018b; Abebe et al., 2020). At the least, machine learning and statistics provide valuable tools to quantify biases, potentially providing the smoking gun evidence required to trigger introspection and positive change.

# 3
# Existing approaches

In this section we provide a survey of fair machine learning for consequential decision making. We categorize the existing body of work by characteristics of the employed fairness criteria. Other possible segmentations would be *group* versus *individual* fairness (see Section 3.4), or according to implementation aspects such as the distinction of *pre-processing*, *in-processing*, and *post-processing* methods. We will adhere to "fairness notion" as the primary partition and mention such orthogonal properties for relevant works as we discuss them. As a general warning, because fairness is a fairly young field within the machine learning community and for reasons we discussed in Section 2.1, there is no agreement on the naming of fairness criteria. Some of them have been developed concurrently by different teams of researchers who gave them different names. We try to use common terminology and clearly reference the corresponding works.

As becomes apparent from the examples in the motivation in Section 1.1, there are plenty applications of machine learning outside of automated decision making for which fairness can be a major concern. There is also an extensive literature on fairness in settings such as ranking, recommender systems, computer vision, natural language processing, bandit learning, clustering, dimensionality reduction, as well as many others which we will not discuss in this thesis. Chouldechova & Roth (2020) provide a concise summary of the current research frontier including pointers to the relevant literature.

## 3.1  Notation

Most works on fairness in machine learning for consequential decisions have focused on what we call *outcome-based* notions. Outcome-based fairness encompasses scenarios in which a system determines (or supports the decision) whether to grant a loan to an applicant based on their financial history, to release a defendant for parole based on a questionnaire and their criminal history, or to invite a job applicant to an interview based on their written

application. The mentioned inputs are only examples among a great variety of possible features to take into account.

Throughout this thesis, we use the following notation. Typically, we denote random variables as well as matrices by upper case letters and their domains (and other sets) as the corresponding calligraphic upper case letters. The elements of such sets, i.e., also specific values of random variables, are denoted by lower case letters. We sometimes highlight that values are multi-dimensional (i.e., vectors) using bold font. Distributions of random variables are denoted by $\mathrm{P}(\cdot)$ and we overload notation by also using $\mathrm{P}(\cdot)$ for probabilities. For example, if $X$ is a continuous real-valued random variable and $Y$ is a binary random variable taking values in $\{0, 1\}$, we commonly denote their joint and marginal distributions by $\mathrm{P}(X, Y)$, and $\mathrm{P}(X), \mathrm{P}(Y)$ respectively (instead of, e.g., $\mathrm{P}_{X,Y}, \mathrm{P}_X, \mathrm{P}_Y$). Similarly, we use $\mathrm{P}(Y = 1 \mid X = x)$ and $\mathrm{P}(Y = 1 \mid x)$ equivalently for the probability of $Y = 1$ given $X = x$. Expectations are denoted by $\mathbb{E}[\cdot]$ where the source of randomness is either clear from the context, or explicitly added as a subscript of $\mathbb{E}$. Unless explicitly stated otherwise, we assume densities exist for all continuous distributions and thus sometimes use $\mathrm{P}$ even for densities whenever there is no ambiguity. When explicitly highlighting that we are referring to densities, we use $\mathrm{p}(\cdot)$. As an example, for a random variable $X$, we write $\mathrm{p}(x)$ for the density, dropping the subscript as in $\mathrm{p}_X(x)$ for simplicity. We denote statistical independence between two random variables by $\cdot \perp\!\!\!\perp \cdot$ and conditional independence by $\cdot \perp\!\!\!\perp \cdot \mid \cdot$ accordingly. Our notation, especially overloading the symbol $\mathrm{P}$, glances over measure theoretic subtleties. Since those will not be relevant for the content of this thesis, we opted for a notation that is easy to parse and understand. We now describe common concepts and variable names that we will use throughout the thesis.

$Z$ is (are) the **sensitive** or **protected attribute(s)** we want to protect for, e.g., gender, race, age, disability, see Section 2.2.

$X$ are the **non-sensitive features**, e.g., salary, SAT scores, etc. We will often refer to $X$ just as **features**, implying that they are non-sensitive.

$Y$ is the **true outcome**, i.e., the observed quantity of interest—for example whether a person defaulted on a loan or re-offended when let out on parole. Note that $Y$ is typically only available in historical data, i.e., a potential training set for supervised machine learning algorithms. Moreover, the term "true" does not imply that this outcome is fair or agreeable, merely that it has been observed or measured, see Section 2.3 for a more detailed discussion.

$\hat{Y}$ is the **predicted outcome** for $Y$ from an automated (machine learning) system. This notation implies that predicting the chosen target $Y$ has been accepted to be a worthwhile optimization goal. In a slight abuse of notation, we will consider $\hat{Y}$ as a function mapping $\hat{Y} : \mathcal{X} \times \mathcal{Z} \to \mathcal{Y}$, where $\hat{Y}(X, Z)$ can still be considered a random variable.

The randomness can come from both, the randomness of $X$, $Z$ and because $\hat{Y}$ may be a randomized mapping.

Note that predicted outcomes may or may not be directly translated into decisions. While most of the literature does not explicitly distinguish between predictions and decisions, we elaborate on the importance of this difference in Chapter 7. Until then, we assume that predictions are closely linked to decisions and even use prediction, outcome, and decisions synonymously.

Outcome-based fairness is concerned with how we can define, measure, and mitigate (un)fairness of the predictions $\hat{Y}$ (taking values in $\mathcal{Y}$) with respect to the features $X$, protected attributes $Z$, and the true outcome $Y$ (taking values in given domains $\mathcal{X}$, $\mathcal{Y}$, and $\mathcal{Z}$ respectively). The most common and highly relevant scenario considers $p \in \mathbb{N}$ categorical (mostly binary) sensitive attributes $\mathcal{Z}$, binary outcomes $\mathcal{Y} = \{0,1\}$, and $d$-dimensional features (real-valued or categorical), $\mathcal{X} \subset \mathbb{R}^d$. We typically assume $\hat{Y} = 1$ to be the more desirable decision or prediction and $Y = 1$ to be the more desirable true outcome or label. Also, in settings with binary sensitive attributes we consider $Z = 1$ to indicate membership in a minority or disadvantaged group.

As described in Section 1.1, this framework only captures a fraction of situations in which fairness may be a concern. However, it entails common settings in which automated systems arguably have the most decisive direct impact on people's lives. Indeed, our excursion into philosophical aspects of fairness in Section 2.1 almost exclusively dealt with notions of justice that are best described in the outcome-based framework. In addition, it allows for the rather simple setup outlined above, thus facilitating in-depth technical discussion of concrete challenges, both theoretical and empirical.

## 3.2 Fairness through unawareness

The first idea to ensure fairness that comes to mind in the outcome-based framework is to simply omit the sensitive attributes $Z$ from the inputs to the algorithm. This approach is called **fairness through unawareness** and can be formally written as $\hat{Y} : \mathcal{X} \to \mathcal{Y}$, i.e., we omit $Z$ as an input for $\hat{Y}$. Due to possible correlations between sensitive attributes and other features, which may readily be exploited by a machine learning algorithm, this approach cannot reliably exclude systematically biased outcomes. For example, while the name of a person is typically not considered a protected attribute, it can correlate strongly with gender or ethnicity. Thus an algorithm could internally "recover" the gender of an individual from their name and exploit existing correlations in the training data in its decisions (Dwork et al., 2012). While fairness through unawareness is generally considered insufficient in most

applications, in certain scenarios there have also been arguments in favor of this approach (Grgić-Hlača et al., 2016).

Fairness through unawareness is often said to avoid **disparate treatment**, a US-centric legal term denoting decisions that are explicitly based on protected attributes, as well as intentional discrimination. To first order, disparate treatment occurs whenever a model uses membership in a protected group as input and differentiates based on it. Fore more details about the legal interpretation of disparate treatment in the context of data-driven decisions we refer the reader to work by Barocas & Selbst (2016b).

## 3.3   Observational group matching criteria

**Observational criteria** are fairness measures that only depend on the joint distribution $P(X, Y, Z, \hat{Y})$. All probabilities and distributions in this section are with respect to this joint distribution and we only consider notions of group fairness, see Section 2.2. For ease of notation, in the examples here we will assume $\mathcal{Z} = \{0, 1\}$. Most criteria easily extend to the case of multiple, categorical protected attributes.

**Demographic parity and disparate impact.**   Proactively adapting to the shortcomings of fairness through unawareness, **demographic parity** (DP) (also called **statistical parity**) suggests to ensure statistical independence of the protected attribute and the predicted outcome: $\hat{Y} \perp\!\!\!\perp Z$. We can write this as

$$P(\hat{Y} \mid Z) = P(\hat{Y}) \,,$$

which for the binary setting is equivalent to

$$P(\hat{Y} = 1 \mid Z = 0) = P(\hat{Y} = 1 \mid Z = 1) \,.$$

In words, the overall probability for a member of group $Z = 0$ to receive outcome $\hat{Y} = 1$ (or $\hat{Y} = 0$), is the same as for a member of group $Z = 1$. While this equality is intended to describe a state with zero unfairness, in practice we often wish to work with a measure of various degrees of unfairness. Two possibilities immediately come to mind: the difference (Calders & Verwer, 2010)

$$P(\hat{Y} = 1 \mid Z = 0) - P(\hat{Y} = 1 \mid Z = 1) \in [-1, 1] \tag{3.1}$$

(or its absolute value), and the ratio (Feldman et al., 2015; Zafar et al., 2017b)

$$\frac{P(\hat{Y} = 1 \mid Z = 0)}{P(\hat{Y} = 1 \mid Z = 1)} \in [0, \infty) \,. \tag{3.2}$$

In eq. (3.1) demographic parity is achieved for a value of 0, whereas in eq. (3.2) for the value 1. The ratio definition in eq. (3.2) is inspired by the legal notion of **disparate impact** (DI) (Barocas & Selbst, 2016a). Therefore, the term is also sometimes used for eq. (3.2) and due to the similarity even for eq. (3.1). Note that due to different interpretations of the legal term *disparate impact*, it is also used as an umbrella term for various kinds of outcome disparity that arise despite the decision not being actively based on the protected attribute. In Chapter 6, we will use disparate impact as such an umbrella term. In Zafar et al. (2017a), disparate impact is used to mimic the p%-rule introduced in the legal literature on employment (Biddle, 2006) by formalizing it as

$$\min \left\{ \frac{P(\hat{Y} = 1 \mid Z = 1)}{P(\hat{Y} = 1 \mid Z = 0)}, \frac{P(\hat{Y} = 1 \mid Z = 0)}{P(\hat{Y} = 1 \mid Z = 1)} \right\} \geq \frac{p}{100}. \tag{3.3}$$

Demographic parity has been criticized for two specific reasons in particular (Hardt et al., 2016b). Both of them are due to the fact that it is based on too little information, caring only about group membership and the outcome. First, it allows $\hat{Y}$ to choose equal fractions differently in the two groups. In the hiring example, it could select the least qualified 10% of applicants in the disadvantaged group and the top qualified 10% of applicants in the privileged group for a phone interview based on the written application. As long as the same fraction of each group is given the opportunity, demographic parity is satisfied. Because it formally allows for such unfair practice, demographic parity can be too weak.

On the other hand, it completely disregards potentially legitimate correlation between $Z$ and $\hat{Y}$. This limitation can be exemplified in the Berkeley college admission case (Bickel et al., 1975). Bickel had shown that a lower college-wide admission rate for women than for men could be explained by the fact that women applied for more competitive departments. When adjusted for department choice, women experienced a slightly higher acceptance rate compared with men in each individual department. In this case, the department choice correlated with the applicant's gender as well as with the chance of admission due to differences in competitiveness. One may believe to understand the nature of both mechanisms reasonably well and tend to conclude that the subsequent correlation between gender (protected attribute) and the admission decision is legitimate and need not be corrected for by the admission committee. In this case ($\mathrm{cov}(Y, A) \neq 0$), demographic parity does not allow the optimal and arguably fair predictor $\hat{Y} = Y$. In this sense demographic parity is too strong. We will revisit the Berkeley college admission example from a causal perspective in Section 3.5.

After reading Chapter 1, we hope our interpretation of the situation made the reader feel uncomfortable or at least a bit skeptical. The observed differences may be rooted much deeper. Perhaps certain departments are known for a hostile and poisonous environment

towards women or existing gender imbalance and a lack of (visibility of) role models further discourage women from applying to certain departments? Again, real social change will require larger-scale systemic reforms and are unlikely to emerge from the specific choice of the admission decision system ignoring other aspects. However, we will continue to try to isolate the effects of algorithms within such a system to try to better understand their role in it.

**Equal odds and equal opportunity.** We now describe attempts to remedy the shortcomings of demographic parity. Ignoring the true outcome $Y$ seemed to leave too little information for a meaningful fairness criterion. **Equalized odds** therefore requires the predicted outcome to be independent of the protected attribute conditioned on the true outcome: $\hat{Y} \perp\!\!\!\perp Z \mid Y$ (Hardt et al., 2016b). The intuition behind the criterion can be roughly described as: *The predicted outcome should not be informed by group membership except for information that comes from the true outcome and thus legitimizes differences.* Under our assumptions, it can be written equivalently as

$$P(\hat{Y} = 1 \mid Y = y, Z = 0) = P(\hat{Y} = 1 \mid Y = y, Z = 1) \quad \text{for } y \in \{0,1\} \,.$$

Note that this already implies an analogous statement for $\hat{Y} = 0$. This definition is also sometimes interpreted as: *Equalized odds seeks to balance false positive rates (FPR) and false negative rates (FNR) between different demographic groups.* The intuitive idea is that we should not make wrong predictions at different rates for the two groups.

As a more specific variant of equalized odds, **equal opportunity** recognizes that in many scenarios we care more about not falsely denying a desirable opportunity, than about providing an opportunity undeservedly (Hardt et al., 2016b). If the desired outcome is encoded by $\hat{Y} = 1$, equal opportunity amounts to balancing the FNR

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

whereas if the desired outcome is encoded by $\hat{Y} = 0$, it corresponds to balancing the FPR

$$P(\hat{Y} = 1 \mid Y = 0, Z = 0) = P(\hat{Y} = 1 \mid Y = 0, Z = 1) \,.$$

A similar criterion was concurrently proposed by Zafar et al. (2017a), who refer to it as avoiding **disparate mistreatment**.

**Calibration and predictive parity.** Calibration is a concept that was not originally associated with fairness. It describes how well reality follows the predictions, i.e., an algorithm is well-calibrated if an event that it predicts to occur with some probability actually happens with this probability. In other words, calibration is sometimes described as *the predictions*

Table 3.1 Confusion tables for two different demographic groups.

| group $Z = 0$ | $\hat{Y} = 0$ | $\hat{Y} = 1$ |
|---|---|---|
| $Y = 0$ | $a$ | $b$ |
| $Y = 1$ | $c$ | $d$ |

| group $Z = 1$ | $\hat{Y} = 0$ | $\hat{Y} = 1$ |
|---|---|---|
| $Y = 0$ | $a'$ | $b'$ |
| $Y = 1$ | $c'$ | $d'$ |

*mean what the are supposed to mean.* In the fairness literature, **calibration** amounts to $Y \perp\!\!\!\perp Z \mid \hat{Y}$ and can be written as (Chouldechova, 2017; Kleinberg et al., 2017b)

$$P(Y = \hat{Y} \mid \hat{Y} = \hat{y}, Z = 0) = P(Y = \hat{Y} \mid \hat{Y} = \hat{y}, Z = 1) \quad \text{for } \hat{y} \in \mathcal{Y}.$$

In words, a given prediction should actually come true with the same probability for members in both groups. We note that this definition differs slightly from non-fairness related notions of calibrated score functions. Perhaps more appropriately, it has also been called **predictive parity** (PP), and further specialized to **positive predictive parity** (PPP) and **negative predictive parity** (NPP) for

$$P(Y = 1 \mid \hat{Y} = 1, Z = 0) = P(Y = 1 \mid \hat{Y} = 1, Z = 1) \text{ and}$$
$$P(Y = 0 \mid \hat{Y} = 0, Z = 0) = P(Y = 0 \mid \hat{Y} = 0, Z = 1)$$

respectively. Closely related criteria matching so called **false omission rates** (FOR) and **false discovery rates** (FDR) read (Zafar et al., 2017a)

$$P(Y = 0 \mid \hat{Y} = 1, Z = 0) = P(Y = 0 \mid \hat{Y} = 1, Z = 1) \text{ and}$$
$$P(Y = 1 \mid \hat{Y} = 0, Z = 0) = P(Y = 1 \mid \hat{Y} = 0, Z = 1).$$

**Generic matching of conditional probabilities.** At this point a clear pattern emerges. Let us draw the confusion tables of true outcomes $Y$ and predictors $\hat{Y}$ for two groups denoted by blue ($Z = 0$) and orange ($Z = 1$) as in Table 3.1. One can then formulate a multitude of observational **matching criteria** by equating different quantities made up of $a$, $b$, $c$, and $d$ across protected groups (Berk et al., 2017; Zafar et al., 2017a). We provide a structured overview of some criteria and the corresponding expressions to be matched in Table 3.2 and Table 3.3. Note that in the binary case, there is some redundancy. For example, the matching criteria for calibration and predictive parity are identical.

Table 3.2 Expressions that may be required to be matched in different fairness criteria in terms of the entries of the confusion tables in Table 3.1.

| expression | name |
|---|---|
| $\frac{b+d}{a+b+c+d}$ | acceptance rate (AR) |
| $\frac{a+d}{a+b+c+d}$ | accuracy (AC) |
| $\frac{b+c}{a+b+c+d}$ | error rate (ER) |
| $\frac{a}{a+b}$ | true negative rate (TNR) |
| $\frac{d}{c+d}$ | true positive rate (TPR) |
| $\frac{b}{a+b}$ | false positive rate (FPR) |
| $\frac{c}{c+d}$ | false negative rate (FNR) |
| $\frac{a}{a+c}$ | negative predictive value (NPV) |
| $\frac{d}{b+d}$ | positive predictive value (PPV) |
| $\frac{c}{a+c}$ | false discovery rate (FDR) |
| $\frac{b}{b+d}$ | false omission rate (FOR) |

Table 3.3 The requirements of some common observational matching criteria in terms of the expressions in Table 3.2.

| name(s) | matching criteria | references |
|---|---|---|
| demographic parity | $AR = AR'$ | (Calders & Verwer, 2010) (Kamishima et al., 2012) (Zemel et al., 2013) (Edwards & Storkey, 2016) (Feldman et al., 2015) (Zafar et al., 2017b) |
| balanced classification rate | $\frac{TPR+FPR}{2} = \frac{TPR'+FPR'}{2}$ | (Friedler et al., 2019) |
| equalized odds | $FPR = FPR'$ and $FNR = FNR'$ | (Hardt et al., 2016b) |
| equality of opportunity | $FNR = FNR'$ (if $\hat{Y} = 1$ is desirable) | (Hardt et al., 2016b) |
| calibration | $FDR = FDR'$ and $FOR = FOR'$ | (Chouldechova, 2017) (Kleinberg et al., 2017b) |
| predictive parity | $PPV = PPV'$ and $NPV = NPV'$ | (Zafar et al., 2017a) |

The fairness criteria outlined in Table 3.3 can be categorized into three generic groups described by (conditional) independences between $Y, Z, \hat{Y}$. Barocas et al. (2019) refer to

$$
\begin{aligned}
\textbf{independence:} \quad & \hat{Y} \perp\!\!\!\perp Z \,, \\
\textbf{separation:} \quad & \hat{Y} \perp\!\!\!\perp Z \,|\, Y \,, \\
\textbf{sufficiency:} \quad & Y \perp\!\!\!\perp Z \,|\, \hat{Y} \,.
\end{aligned}
$$

All observational group matching criteria can be assigned in spirit to one of these three categories.

We highlight at this point, that all the criteria mentioned so far are observational, i.e., they can be formulated with reference only to the joint distribution $P(X, Y, \hat{Y}, Z)$ (where we have not yet made use of $X$). Given a dataset for features, protected attributes, true outcomes, and predicted outcomes from some system, no further assumptions are needed to directly verify each notion on the empirical distribution. Moreover, for each of them there is at least one scenario in which it can be framed as a desirable property in terms of fairness. A principled comparison of their respective utility for fair decision making is virtually impossible without context and a specific application. Unfortunately, one also can not have them all.

**Impossibility results.** Discussions about fairness often arise in situations, where the observed **base rates** differ across protected groups. The base rate is the fraction of people with positive true outcome $Y = 1$. In our notation, this means $p := {}^{a+b}/_{a+b+c+d} \neq {}^{a'+b'}/_{a'+b'+c'+d'} =: p'$, where we denote the base rates for the two groups by $p$ and $p'$. From these expressions, one can directly verify that within each protected group (Chouldechova, 2017)

$$
\mathrm{FPR} = \frac{p}{1-p} \frac{1 - \mathrm{PPV}}{\mathrm{PPV}} \mathrm{FNR} \quad \text{and} \quad \mathrm{FPR}' = \frac{p'}{1-p'} \frac{1 - \mathrm{PPV}'}{\mathrm{PPV}'} \mathrm{FNR}' \,. \tag{3.4}
$$

Therefore, if we wish to achieve equalized odds ($\mathrm{FPR} = \mathrm{FPR}'$, $\mathrm{FNR} = \mathrm{FNR}'$) and calibration (in particular $\mathrm{PPV} = \mathrm{PPV}'$) at the same time, eq. (3.4) implies that $p = p'$ or $\mathrm{FPR} = \mathrm{FPR}' = \mathrm{FNR} = \mathrm{FNR}' = 0$.

As a consequence, except for the degenerate cases of equal base rates or a perfect predictor, no system can satisfy calibration and equalized odds at the same time (Kleinberg et al., 2017b; Chouldechova, 2017). In general, unless the confusion matrices in Table 3.1 are scalar multiples of each other, or all off-diagonal entries are 0, many pairs of observational fairness criteria (e.g., Table 3.3) cannot be satisfied jointly. For example, Pleiss et al. (2017) studied which fairness criteria are compatible with calibration. Similarly, impossibility results can be phrased elegantly in terms of separation and sufficiency (Barocas et al., 2019). Kim et al. (2020) provide a model agnostic characterization of the (in)compatibility of many combinations of criteria listed in Table 3.3. Comparing different outcome-based (un)fair

predictors to one another remains a challenge that has only recently seen some first advances (Speicher et al., 2018). The impossibility of "playing it safe", by satisfying all potentially desirable group fairness criteria simultaneously, highlights the necessity of comparing and discussing their usage in the context of a specific application.

Even before the first mention of these results (Kleinberg et al., 2017b; Chouldechova, 2017), a closely related debate was carried out in the COMPAS example which we encountered in Section 1.1. Angwin et al. (2016) from ProPublica argued that COMPAS does not satisfy equality of opportunity, because among the people who did not go on to re-offend ($Y = 1$), black defendants ($Z = 1$) got worse scores than white defendants ($Z = 0$), i.e., $P(\hat{Y} = 1 \mid Y = 1, Z = 0) \neq P(\hat{Y} = 1 \mid Y = 1, Z = 1)$. In a rebuttal to the ProPublica article, the company behind COMPAS (at the time it was called *Northpointe*, in the meantime they have been renamed to *Equivant*) argued for calibration as a more appropriate fairness criterion in their application and demonstrated that COMPAS indeed satisfies calibration (Dieterich et al., 2016). This specific real-world example of incompatibility and the resulting trade-offs has also be studied in more detail from an algorithmic and ethical perspective (Corbett-Davies et al., 2016; Flores et al., 2016; Corbett-Davies et al., 2017).

## 3.4 Sub-group and individual fairness

An interesting non-observational fairness definition is **individual fairness** (Dwork et al., 2012), which assumes the existence of a similarity measure for individuals, and requires that any two similar individuals should receive a similar distribution over outcomes. Hence, we here consider a mapping $\hat{Y} : \mathcal{X} \to \Delta(\mathcal{Y})$, where $\Delta(\mathcal{Y})$ denotes the set of all probability distributions on $\mathcal{Y}$. Then the resulting criterion for individual fairness takes the form

$$D(\hat{Y}(x), \hat{Y}(x')) \leq d(x, x') \quad \text{for } x, x' \in \mathcal{X},$$ (3.5)

where $D : \Delta(\mathcal{Y}) \times \Delta(\mathcal{Y}) \to [0, \infty)$ and $d : \mathcal{X} \times \mathcal{X} \to [0, \infty)$ are metrics. More recent work lends additional conceptual and theoretical support to such a definition (Friedler et al., 2016).

However, Dwork et al. (2012) also acknowledge the main difficulty of the approach, namely choosing metrics $D$ and $d$ in practice. This issue has been approached in an online learning setting, where a regulator "knows unfairness when they see it" to circumvent a normative definition of similarity (Gillen et al., 2018). Along similar lines, Kim et al. (2018) assume that instead of having access to the metric in analytical form, we can query it a bounded number of times, where the queries also may be answered by experts. Their fairness notion aligns with ideas from previous work aiming to relax individual fairness by attempting to interpolate between group and individual fairness (Hébert-Johnson et al., 2018; Kearns

et al., 2018). Lahoti et al. (2019) model side-information, such as fairness judgments from a variety of sources (including human judgments) as a fairness graph and then learn a unified presentation capturing pairwise fairness to tackle the issue of learning a similarity metric. Learning such a metric from human annotated data has also been explored in detail specifically for the criminal recidivism prediction task on COMPAS data (Wang et al., 2019).

These attempts rely on ensuring parity not between group statistics, but between many (possibly overlapping) subgroups of the population with identical (or similar) features. Besides operationalizing individual fairness as proposed by Dwork et al. (2012), a key motivation for such approaches to subgroup fairness is to avoid *fairness gerrymandering* (Kearns et al., 2018; Hébert-Johnson et al., 2018). In concurrent work, Yona & Rothblum (2018) arrive at the conclusion that ensuring the strict definition eq. (3.5) exactly is generally computationally intractable. However, these intractabilities can be overcome for a relaxed notion they call *approximate metric-fairness*. For the remainder of this thesis we will focus on group fairness notions.

## 3.5 Causal fairness criteria

In this section, we motivate the study of fairness through the lens of causality and describe some existing causal notions of fairness. An in-depth discussion of our own contributions in going beyond observational criteria follows in Chapter 4.

**Why causality?**  To motivate the importance of understanding cause-effect relations in fairness, let us come back to the Berkeley college admission example (Bickel et al., 1975). The potential allegation in this study was based on the observation that the overall acceptance rate for females was lower. If one does not believe this to be a coincidence, the conclusion that the acceptance decisions are based on gender stands to reason. After all, it appears unlikely that the admission decision influenced the gender. However, this seemingly plausible explanation neglects possible alternatives due to confounding. When presented with two correlated quantities, we all too commonly conclude that either the first influences the second or vice versa. Thereby, we overlook the possibility of a third variable influencing both quantities in question. In the Berkeley college admission case, only a closer analysis revealed the underlying Simpson's paradox and opened up the possibility that gender may not have directly influenced the admission decision. As soon as one conditions on department choice, the effect reverses and females experience higher acceptance rates in each department.

Crucially, when having access only to the joint distribution of admission decision $\hat{Y}$, gender $Z$, and some features $X$, e.g., SAT scores, we cannot tell, whether $Z$ had a direct influence

on $\hat{Y}$. If this were the case, it would raise serious concerns about blatantly direct gender discrimination. However, the same joint distribution could arise under a different generative model, where gender has a direct causal effect only on the department choice, which in turn affects the outcome. In short, different causal structures can entail the same joint distribution over all variables. As a consequence, all observational criteria—based only on the joint distribution—cannot distinguish between causally different scenarios. Since the example shows that whether a decision is perceived as fair or not depends heavily on the causal structure of the process, this is bad news for observational criteria.

The importance of *unidentifiability* of causal structure from observational data for fairness can be traced back to Pearl (2009, Section 6). Hardt et al. (2016b) explicitly mention that observational criteria are too weak to distinguish two intuitively different scenarios. However, the work does not provide a formal mechanism to articulate why and how these scenarios should be considered different. Proposing a causality-based approach to cope with this issue is one of our contributions described in Chapter 4.

To further motivate a causal framework to analyze fairness, recall the intuition behind some of the observational criteria. For example, demographic parity asks the outcome to be statistically independent of the protected attribute. However, we conjecture that the way the general public intuitively thinks about demographic parity is that *membership in a protected group should not have any effect on the predicted outcome*. Similarly, equalized odds informally asks that the predicted outcome should not be informed by group membership except for information that comes from the true outcome. This is likely understood as *the predicted outcome should only be affected by the group membership to the extent to which this is justified or mediated by the true outcome*. Both informal explanations hinge on causal statements rather than mere statistical dependence. Thus, a causal framework may be more appropriate to formalize them.

The arguments in favor of causality so far were based on issues of group fairness criteria. However, causality offers a similar conceptual benefit to dealing with individual fairness. Initial work on individual fairness focused on distance measures of the features $X$ of individuals in the context of a given application to assess similarity. As outlined in Section 2.2, the most immediate comparison conceptually is with a different version of ourselves, where everything is equal except for the protected attribute. Causality allows us to formally reason about such counterfactuals, rendering it an appealing tool to investigate individual notions of fairness. The idea of a similarity measure of individuals in observed data can then be addressed by the method of *matching* used in counterfactual reasoning (Rosenbaum & Rubin, 1983; Qureshi et al., 2019). That is, evaluating approximate counterfactuals by comparing individuals with similar values of covariates $X$ but different protected attribute $Z$.

**Structural equation models and causal graphs.**    We assume a basic understanding of the key ideas of causality, in particular *interventions* and *counterfactuals* and predominantly work with the graphical approach to causality and the *do-calculus*, see Pearl (2009); Peters et al. (2017). We also borrow notation from these books. Let us now recall the definition of structural equation models and provide a subjective opinion on how to interpret them.

For our purposes, a **causal graph** is a directed, acyclic graph whose nodes reference random variables. The directed edges represent causal influences between those variables. In that way, causal graphs are a convenient tool to organize assumptions about the data generating process. More formally, we work with Structural Equation Models. For a rigorous introduction of structural equation models with more attention to measure-theoretic considerations, we refer the reader to Bongers et al. (2016). A **Structural Equation Model (SEM)** is a set of equations of the form

$$V_i = f_i(\mathrm{pa}(V_i), N_i) \quad \text{for } i \in \{1, \ldots, n\},$$

where

- the $V_i$ are real-valued random variables,

- $\mathrm{pa}(V_i) \subset \{V_1, \ldots, V_n\}$ are called the **parents** of $V_i$,

- the real-valued random variables $N_i$ are independently distributed according to $\mathrm{P}(N_i)$,

- the **causal graph** induced by the structural equations $f_i$, i.e., the directed graph with nodes $\{V_i\}_{i=1}^n$ and edges $V_i \to V_j$ if and only if $V_i \in \mathrm{pa}(V_j)$, is acyclic.

We typically denote the resulting directed acyclic graph (DAG) by $\mathcal{G}$ and call $\mathrm{pa}(V_i)$ the **direct causes** of $V_i$. The random variables $N_i$ are also referred to as **noise variables** or **exogenous variables**.

The structural equations of an SEM are interpreted as assignments. Because we assume acyclicity, following the nodes of the graph in topological order, we can recursively compute the values of all variables, given the noise variables. In most applications, we assume the $\{f_i\}_{i=1}^n$ to be deterministic, i.e., all the randomness comes from the noise variables. Hence, given specific values for the $N_i$, the values of all $V_i$ are fixed. A structural equation model entails a unique joint distribution over all variables, leading us to view structural equation models as **data generating models**. As mentioned above, the same joint distribution can usually be entailed by multiple structural equation models with distinct causal graphs.

Some graph-theoretical terminology will come in handy in the following two chapters. A **path** in a DAG $\mathcal{G}$ is a sequence of distinct nodes $V_1, \ldots, V_k$, for $k \geq 2$, such that $V_i \to V_{i+1}$ or $V_{i+1} \to V_i$ for all $i \in \{1, \ldots, k-1\}$. A path is **directed**, if $V_i \to V_{i+1}$ for all $i \in \{1, \ldots, k-1\}$, i.e., all arrows "point in the same direction". We say a path between $V_1$ and $V_k$ in a DAG $\mathcal{G}$

is **blocked by a set of nodes** $S$, where $V_1, V_k \notin S$ when one of the following two conditions is met: (a) There exists a $V_i \in S$ such that $V_{i-1} \to V_i \to V_{i+1}$ or $V_{i+1} \to V_i \to V_{i-1}$. (b) Or there exists a $V_i$ such that $V_{i-1} \to V_i \leftarrow V_{i+1}$ and neither $V_i$ nor any of its descendants is in $S$. The notion of blocked paths finally allows us to define **d-separation**: Two disjoint sets of nodes $V$ and $W$ in a DAG $\mathcal{G}$ are **d-separated** by a third pairwise disjoint set of nodes $S$ if every path between nodes in $V$ and $W$ is blocked by $S$.

Any structural equation model by definition entails a DAG $\mathcal{G}$ and a unique distribution over all its variables. This raises a general question about the compatibility of graphs with joint distributions over the same variables in terms of d-separation on the one hand and conditional independences on the other hand. In particular, the following definitions will be relevant for Chapter 4. Let $\mathcal{G}$ be a DAG and P a joint distribution over all variables in $\mathcal{G}$ that allows for a density with respect to the product measure. Then we say P **is Markov** with respect to $\mathcal{G}$ if for all disjoint sets of nodes $V, W$ and $S$ we have that

$$V \text{ and } W \text{ are d-separated by } S \Rightarrow V \perp\!\!\!\perp W \mid S.$$

Two DAGs $\mathcal{G}_1$ and $\mathcal{G}_2$ are called **Markov equivalent** if the sets of all distributions (with a density) that are Markov with respect to a $\mathcal{G}_1$ and $\mathcal{G}_2$ respectively, are equal. Finally, from this follows the **Markov equivalence class of a DAG** $\mathcal{G}$, which we define as the set of all DAGs that are Markov equivalent to $\mathcal{G}$.

Since two graphs are Markov equivalent if and only if they satisfy the same set of d-separations, by the Markov property, we can infer that they satisfy the exact same set of conditional independence relations. As we have seen, the Markov property talks about when graphical d-separation criteria imply the corresponding conditional independence relation. There is an analogous definition for the other direction: A joint distribution P is **faithful to the DAG** $\mathcal{G}$ if for all disjoint sets of nodes $V, W$ and $S$ we have that

$$V \perp\!\!\!\perp W \mid SV \text{ and } \Rightarrow W \text{ are d-separated by } S.$$

Faithfulness can be violated when different paths in a given causal graph cancel each other out such that there is an independence relation satisfied by the joint distribution, for which the corresponding d-separation is violated. In Chapter 4 we will argue that such a situation may indeed occur in practice.

**Interpretation.** We consider the causal influences captured by the $f_i$ to be mechanisms that are intrinsic to how the universe works. Prototypical examples would be laws of physics, however, we will also assume that there exist invariant mechanisms between fuzzier concepts that can still be captured in a functional form. For example, in causal inference for health-care, we may assume that one can meaningfully model concepts such as "smoker"

and "has lung cancer" as random variables as well as the causal influence between the two as a mathematical function. Similarly, economists may seek to compute the causal effect of the "health of institutions" on the gross domestic product. None of these concepts are unambiguously defined properties of the physical world and—as somewhat constructed notions—are themselves not "causing" anything in the physical sense. However, we still believe that there exist causal mechanisms relating them and that we can meaningfully represent those by mathematical equations.

In these examples, we essentially claim that the ontology of "smoker", "lung cancer", "health of institutions", and "gross domestic product" as well as the assumed interactions between them are stable enough. Intuitively, this means that we can pinpoint and agree upon what these terms reference with a sufficient degree of certainty and precision. The mass or electric charge of a physical object have a more stable ontology than the health of institutions. However, we may still agree on what is referenced by the health of institutions and can agree on the methods that led us to believe in that ontology. As we have discussed in Section 2.3, moving to concepts such as race and gender, ontological stability becomes a serious concern. How to practically address these challenges is the subject of ongoing research (Kohler-Hausmann, 2018; Hu & Kohler-Hausmann, 2020).

In the context of this thesis, the predictor $\hat{Y}$ maps inputs, e.g., the features $X$ and the protected attribute $Z$, to a predicted outcome. Hence, we model $\hat{Y}$ as a childless node, whose parents are its input variables. This predictor node has a special role, because its structural equation does not correspond to universally true mechanism of the universe, but can be chosen by the decision-maker, for example by training a predictive model. While we commonly display $\hat{Y}$ as part of a causal graph like all other variables, it is important to keep in mind that its incoming edges have a different meaning from the others.

We note that a key difference between causal SEMs and probabilistic graphical models is that in SEMs we only care about causal factorizations of the joint probability. This means that by writing down a causal DAG, we explicitly formulate the assumption that the parents of any given variable are precisely the variables that actually have a direct causal influence on that variable. In contrast, probabilistic graphical models allow for any kind of factorization of the joint distribution and thus typically do not require parents to be causes of their descendants. Arguably the two most important tools the language of causal SEMs offers beyond the standard statistical toolkit of probabilistic graphical models are formal notions of interventions and counterfactuals. Hence, it is little surprising that most existing causal fairness criteria are based on one or the other. We now provide a short overview of works at the intersection of causality and fairness in machine learning.

**Causal fairness criteria.** Kusner et al. (2017) were among the first to put forward one possible causal definition, namely the notion of **counterfactual fairness**. It ensures that the

outcome any individual receives in the real world is the same as the one she would receive in a *counterfactual world* in which only the protected attribute of the individual has changed, everything else remaining equal. This can be formalized as

$$P(\hat{Y}_{Z=z} = 1 \mid X = x, Z = z) = P(\hat{Y}_{Z=z'} = 1 \mid X = x, Z = z) \quad \text{for } x \in \mathcal{X}, \, z, z' \in \mathcal{Z}. \quad (3.6)$$

We interpret the right hand side of this equation as *the probability that $\hat{Y}$ predicts $1$ for a given individual for which we observe features $x$ and protected attribute $z$, had the protected attribute been $z'$ instead of $z$.* An analogous interpretation of the left hand side of eq. (3.6) asserts that it is equal to $P(\hat{Y} = 1 \mid X = x, Z = z)$.

Note that this definition requires modeling counterfactuals on a per individual level, which is a delicate task. Even determining the effect of *race* at the group level is difficult as discussed by VanderWeele & Robinson (2014); Hu & Kohler-Hausmann (2020). Interventions on (often ill-defined) protected attributes such as gender or race are generically hard to conceive. For example, imagine a pregnant woman's job application gets rejected. How should we conceive of the world in which she has been a man? Should we think of her as being born male, or being perceived as male during the hiring procedure? Is she a pregnant man now? The counterfactual where somebody was "born a different person", may lead to the comparison of vastly different (fictitious) individuals, undermining the key motivation for counterfactual fairness. Such challenges are aggravated when moving beyond inappropriately simplified binary gender assignments. In Chapter 4 we propose a proxy-based approach to formalize a continuum of possible interventions that may be harder or easier to conceive and—perhaps even more important—to perform in practice.

In contrast to counterfactuals, for a pure intervention there is no abduction step, i.e., we do not update the distribution of the exogenous variables conditioned on the observations. Instead, we simply replace the structural equation for the variable we intervene on by setting it to a constant (or a fixed distribution). Intuitively, this notion is most often compared to randomized trials. For example, before having recruiters screen written applications, we could intervene by randomizing the names—perhaps one of the strongest salient indicators for gender—of the applicants. Note that in this case we only intervene on a proxy for gender. The intervention does not account for potential gender discrimination applicants have experienced before this specific application. However, it can easily be carried out in practice. Indeed, such a randomized study has shown that applications with typical White-sounding names receive considerably higher callback rates for interviews than African-American-sounding names (Bertrand & Mullainathan, 2004).

Formally, **interventional fairness** can be stated in the form

$$P(\hat{Y} = 1 \mid do(Z = 0), X = x) = P(\hat{Y} = 1 \mid do(Z = 1), X = x),$$

where we may or may not condition on $x \in \mathcal{X}$. Note that even when we condition on the features, we still match probabilities within groups of potentially different individuals who happen to have the same features. This notion is different from the counterfactual statement in eq. (3.6), where we condition on having observed a specific individual to update the exogenous variables before intervening on the protected attribute. We further remark that we still used the letter $Z$, even though in the example we spoke of a proxy of the protected attribute. This distinction will be made clear in Chapter 4.

**Related work on causality-based fairness.** More generally, causality has already been employed for the discovery of discrimination in existing datasets (Bonchi et al., 2017; Qureshi et al., 2019). Causal graphical conditions to identify *meaningful partitions* have been proposed for the discovery and prevention of certain types of discrimination by preprocessing the data (Zhang & Wu, 2017). These conditions rely on the evaluation of *path specific effects*, which relates closely to earlier work by Pearl (2009, Section 4.5.3). Nabi & Shpitser (2018) picked up this notion and generalized Pearl's approach by a constraint based prevention of discriminatory path specific effects arising from counterfactual reasoning. The idea of path specific effects to distinguish fair and unfair causal pathways, and further to suppress influences along the unfair ones, have also been combined with deep learning and variational inference approaches to render it applicable in complex, non-linear scenarios (Chiappa & Gillam, 2018). Further, Zhang & Bareinboim (2018) propose three fine grained types of causal influence, based on which they attempt to aid the practitioner to jointly model the data generating mechanisms and an appropriate fairness criterion. For more details, we refer the reader to a survey on causal reasoning for algorithmic fairness by Loftus et al. (2018).

**Technical challenges of causal reasoning for fairness.** At the beginning of this section, we argued that causality is conceptually useful to tackle some of the subtle issues when it comes to fairness criteria. However, there is a price to pay. First, most works start from the usual assumption that the causal graph is known. Since causal discovery is generically difficult (Peters et al., 2017), this is not an assumption to make lightly. Russell et al. (2017) present some ideas on how to achieve approximate counterfactual fairness for multiple competing causal models simultaneously to account for some degree of uncertainty about the correct one. Furthermore, causal inference often comes with a set of hard-to-verify or even untestable standard assumptions such as identifiability, no unobserved confounding, or faithfulness (Pearl, 2009). Finally, in particular for notions based on counterfactuals, there is usually no data available for training and verification, because by their very definition, we never observe the counterfactual outcomes. We develop tools to deal with some of these issues. In Chapter 4, we address potential misspecifications of the causal graph. The issue of labels missing depending on our decisions are discussed in Chapter 7.

## 3.6    Fair representations

Supervised representation learning offers a natural approach to demographic parity. The basic idea is to transform the input features via a mapping $t : \mathcal{X} \to \widetilde{\mathcal{X}}, x \mapsto \widetilde{x}$ into a new space $\widetilde{\mathcal{X}}$ such that the following conditions hold:

1. We can train a predictor $\hat{Y} : \widetilde{\mathcal{X}} \to \mathcal{Y}$ with high accuracy. Informally, the new representation $\widetilde{x}$ of the features $x$ should still contain all relevant information about the true outcome $y$.

2. It is hard to learn a predictor $\hat{Z} : \widetilde{\mathcal{X}} \to \mathcal{Z}$, i.e., the new representation $\widetilde{x}$ of the features $x$ contains no information about the corresponding protected attribute $z$.

The aspiration of learned representations is usually their usefulness for potentially unknown downstream tasks. In the fairness context, the hope is that once we have learned such a fair transformation $t$, one can freely train any unconstrained predictor $\hat{Y} : \widetilde{\mathcal{X}} \to \mathcal{Y}$ without having to worry about introducing unfairness.

Zemel et al. (2013); Calmon et al. (2017) formulate this pre-processing task as an optimization problem and directly learn $t$. A large body of work in this direction is based on *adversarial training*, i.e., learning $t$ jointly with $\hat{Y}$ and $\hat{Z}$ by optimizing for the two objectives—minimizing the loss for $\hat{Y}$ and maximizing the loss for $\hat{Z}$—simultaneously by alternating gradient descent (Edwards & Storkey, 2016; Sokolic et al., 2017; Madras et al., 2018). Another approach exploits variational auto-encoders (Kingma & Welling, 2014) with an additional loss terms to render the latent representation a bad predictor for the protected attribute (Louizos et al., 2016).

What we described here as *fair representation learning* is still narrowly concerned with outcome-based fairness. Other issues with biased representations, some of which were mentioned in Section 1.1, are beyond the scope of this thesis.

## 3.7    Procedural approaches and human perception

Almost all existing work discussed in this chapter falls into the broad category of distributive fairness. In a series of intriguing papers, Grgić-Hlača et al. (2016; 2018b;a; 2020) study procedural approaches, or *process fairness*, based on the collective moral judgment of humans. For example, the surveys ask which features people find fair to use, e.g., in the COMPAS recidivism risk assessment setting, and further analyze how perceptions of fairness relate to demographics and personal experiences of individuals. Such empirical approaches may help describe the procedural component of fairness and justice. For example, in a related setting, the landmark study on moral decisions in trolley problems related to autonomous

driving called *the moral machine experiment* provided insights that may inform regulatory decisions regarding autonomous vehicles (Awad et al., 2018).

# 4
# Causality and fairness

In this chapter, we go beyond the assumption that only observational data without any additional information is available. By framing the problem of discrimination based on protected attributes in the language of causal reasoning, we can circumvent the limitations of observational criteria. This viewpoint shifts attention from "What is the right fairness criterion?" to "What do we want to assume about our model of the causal data generating process?" Through the lens of causality, we make several contributions. First, we crisply articulate why and when observational criteria fail, thus formalizing what was before a matter of opinion. Second, our approach exposes previously ignored subtleties and why they are fundamental to the problem. For example, we do not assume we can evaluate individual causal effects and meaningfully conceive interventions on protected attributes. Finally, we put forward natural causal non-discrimination criteria and develop algorithms that satisfy them.

The main content of this chapter has been published in the following paper:

> AVOIDING DISCRIMINATION THROUGH CAUSAL REASONING
> Niki Kilbertus, Mateo Rojas-Carulla, Giambattista Parascandolo, Moritz Hardt, Dominik Janzing, Bernhard Schölkopf.
> *Neural Information Processing Systems (NeurIPS), 2017*

## 4.1  Introduction

We start from the fact that observational criteria are insufficient to conclusively capture fairness simply because there exist scenarios with *intuitively* different social interpretations that admit identical joint distributions over $(\hat{Y}, Z, Y, X)$ (Hardt et al., 2016b). Thus, no observational criterion can distinguish them. Inspired by Pearl's causal interpretation of Simpson's paradox (Pearl, 2009, Section 6), we propose causality as a way of coping with this

unidentifiability result. Carefully using the language of causal reasoning supports several contributions:

- Revisiting the two scenarios proposed by Hardt et al. (2016b) that cannot be differentiated by observational criteria, we articulate a natural causal criterion that formally distinguishes them.

- We point out subtleties in fair decision making that arise naturally from a causal perspective, but have gone widely overlooked in the past. Specifically, we formally argue for the need to distinguish between the underlying concept behind a protected attribute, such as race or gender, and its *proxies* available to the algorithm, such as visual features or name.

- We introduce and discuss two natural causal criteria centered around the notion of *interventions* (relative to a causal graph) to formally describe specific forms of discrimination.

- Finally, we initiate the study of algorithms that avoid these forms of discrimination. Under certain linearity assumptions about the underlying causal model generating the data, an algorithm to remove a specific kind of discrimination leads to a simple and natural heuristic.

At a higher level, our work proposes a shift from trying to find a single statistical fairness criterion to arguing about properties of the data and which assumptions about the generating process are justified. Causality provides a flexible framework for organizing such assumptions. In particular, we will introduce two complementary thought frameworks, which we call the *benevolent* and the *skeptic* viewpoint. The benevolent viewpoint is characterized by generally not assuming causal influences of the sensitive attribute on the decision to be unfair. Modelers may then mark specific variables along such paths that they deem should be causally irrelevant for the decision as *proxy variables*. In this scenario our goal is to make decisions that are not causally influenced by such proxy variables. In the skeptic viewpoint, we take the stance that all causal influences from the sensitive attribute on the decision are by default considered inappropriate. However, as modelers and stakeholder we may declare some variables along such paths as *resolving variables* if we believe that they can fairly be used for making decisions despite being influenced by the sensitive attribute. For example, while an employer may be prohibited from discriminating applicants based on their country of origin, it may still be appropriate to demand certain language skills, which are often influenced by the country of origin. This process of closely analysing an assumed causal model underlying the observed data and making deliberate judgment calls about the fairness of causal influences, i.e., marking variables as proxies or resolving variables, is at the heart of our proposal and calls for a more situational and context dependent assessment of fairness.

Figure 4.1 The admission decision $\hat{Y}$ does not only directly depend on gender $Z$, but also on department choice $X$, which in turn is also affected by gender $Z$.

In addition to our introduction of structural equation models in Section 3.5, one more concept that will be useful in our exposition is that of *terminal ancestors*. We will refer to the *terminal ancestors of a node $V$ in a causal graph $\mathcal{G}$*, denoted by $ta^{\mathcal{G}}(V)$, which are those ancestors of $V$ that are also root nodes of $\mathcal{G}$.

## 4.2    Unresolved discrimination and limitations of observational criteria

To bear out the limitations of observational criteria, we again turn to the commentary on claimed gender discrimination in Berkeley college admissions in Section 4.5.3 of Pearl (2009). To recap, Bickel et al. (1975) had shown earlier that a lower college-wide admission rate for women than for men was explained by the fact that women applied in more competitive departments. When adjusted for department choice, women experienced a slightly higher acceptance rate compared with men. From the causal point of view, arguably what matters is the *direct effect* of the protected attribute (here, gender $Z$) on the decision (here, college admission $\hat{Y}$) that cannot be ascribed to a mediator that has been judged to be admissible in the decision process without fairness concerns. We shall use the term *resolving variable* for any such variable in the causal graph that is influenced by $Z$ in a manner that we accept as non-discriminatory, such as department choice $X$ in our example, see Figure 4.1.[1] With this convention, the criterion can be stated as follows.

**Definition 1** (Unresolved discrimination). A variable $V$ in a causal graph exhibits *unresolved discrimination* if there exists a directed path from $Z$ to $V$ that is not blocked by a resolving variable and $V$ itself is non-resolving.

Pearl's commentary is consistent with what we call the *skeptic viewpoint*. All paths from the protected attribute $Z$ to $\hat{Y}$ are problematic, unless they are justified by a resolving variable. The presence of unresolved discrimination in the predictor $\hat{Y}$ is worrisome and demands further scrutiny. In practice, $\hat{Y}$ is not a priori part of a given graph. Instead it is our objective to construct it as a function of the features $X$, some of which might be resolving. Hence we should first look for unresolved discrimination in the features. A canonical way to

---

[1]We remark again that disarming any discrimination concerns by only considering direct effects to be discriminatory is a strong normative statement. The root cause may be much closer to different departments following discriminatory policies or nurturing a hostile environment towards women, see Section 3.5.

Figure 4.2 Two graphs that may generate the same joint distribution for the Bayes optimal unconstrained predictor $\hat{Y}^*$. If $X_1$ is a resolving variable, $\hat{Y}^*$ exhibits unresolved discrimination in the right graph (along the red paths), but not in the left one.

avoid unresolved discrimination in $\hat{Y}$ is to only input the set of features that do not exhibit unresolved discrimination. However, the remaining features might be affected by non-resolving *and* resolving variables. In Section 4.4 we investigate whether one can exclusively remove unresolved discrimination from such features. A related notion of "explanatory features" in a non-causal setting was introduced by Kamiran et al. (2013).

The definition of unresolved discrimination in a predictor has some interesting special cases worth highlighting. If we take the set of resolving variables to be empty, we intuitively get a causal analog of demographic parity. No directed paths from $Z$ to $\hat{Y}$ are allowed, but $Z$ and $\hat{Y}$ can still be statistically dependent. Similarly, if we choose the set of resolving variables to be the singleton set $\{Y\}$ containing the true outcome, we obtain a causal analog of equalized odds where strict independence is not necessary. The causal intuition implied by "the protected attribute should not affect the prediction", and "the protected attribute can only affect the prediction when the information comes through the true label", is neglected by (conditional) statistical independences $Z \perp\!\!\!\perp \hat{Y}$, and $Z \perp\!\!\!\perp \hat{Y} \mid Y$, but well captured by only considering dependences mitigated along directed causal paths.

We will next show that observational criteria are fundamentally unable to determine whether a predictor exhibits unresolved discrimination or not. This is true even if the predictor is *Bayes optimal*. In passing, we also note that fairness criteria such as equalized odds may or may not exhibit unresolved discrimination, but this is again something an observational criterion cannot determine.

**Theorem 1.** *Given a joint distribution over the protected attribute $Z$, the true label $Y$, and some features $X_1, \ldots, X_n$, in which we have already specified the resolving variables, no observational criterion can generally determine whether the Bayes optimal unconstrained predictor or the Bayes optimal equal odds predictor exhibit unresolved discrimination.*[2]

---

[2]The Bayes optimal equal odds predictor is the Bayes optimal predictor among all predictors that satisfy the equal odds fairness criterion.

*Proof.* Let us consider the two graphs in Figure 4.2. First, we show that these graphs can generate the same joint distribution $P(Z, Y, X_1, X_2, \hat{Y}^*)$ for the Bayes optimal unconstrained predictor $\hat{Y}^*$.

We choose the following structural equations for the graph on the left[3]: $Z = \text{Ber}(1/2)$, $X_1$ is a mixture of Gaussians $\mathcal{N}(Z + 1, 1)$ with weight $\sigma(2Z)$ and $\mathcal{N}(Z - 1, 1)$ with weight $\sigma(-2Z)$, $Y = \text{Ber}(\sigma(2X_1))$, $X_2 = X_1 - Z$, and $\hat{Y}^* = X_1$. Throughout this proof we assume that the Bernoulli distribution Ber has support $\{-1, 1\}$ instead of $\{0, 1\}$.

For the graph on the right, we define the structural equations $Z = \text{Ber}(1/2)$, $Y = \text{Ber}(\sigma(2Z))$, $X_2 = \mathcal{N}(Y, 1)$, $X_1 = Z + X_2$, and $\hat{Y}^* = X_1$.

First we show that in both scenarios $\hat{Y}^*$ is actually an optimal score. In the first scenario $Y \perp\!\!\!\perp Z \mid X_1$ and $Y \perp\!\!\!\perp X_2 \mid X_1$ thus the optimal predictor is only based on $X_1$. We find

$$P(Y = y \mid X_1 = x_1) = \sigma(2x_1 y), \tag{4.1}$$

which is monotonic in $x_1$. Hence optimal classification is obtained by thresholding a score based only on $\hat{Y}^* = X_1$.

In the second scenario, because $Y \perp\!\!\!\perp X_1 \mid \{Z, X2\}$ the optimal predictor only depends on $Z, X_2$. We compute for the densities (which all exist and are positive)

$$
\begin{aligned}
P(Y \mid X_2, Z) &= \frac{P(Y, X_2, Z)}{P(X_2, Z)} \\
&= \frac{P(X_2, Z \mid Y) \, P(Y)}{P(X_2, Z)} \\
&= \frac{P(X_2 \mid Y) \, P(Z \mid Y) \, P(Y)}{P(X_2, Z)} \\
&= \frac{P(X_2 \mid Y) \frac{P(Y \mid Z) \, P(Z)}{P(Y)} \, P(Y)}{P(X_2, Z)} \\
&= \frac{P(X_2 \mid Y) \, P(Y \mid Z) \, P(Z)}{P(X_2, Z)},
\end{aligned}
$$

where for the third equal sign we use $Z \perp\!\!\!\perp X_2 \mid Y$. In the numerator we have

$$p(x_2 \mid Y = y) \, p(y \mid Z = z) \, p(z) = f_{\mathcal{N}(y,1)}(x_2) f_{\text{Ber}(\sigma(2z))}(y) f_{\text{Ber}(1/2)}(z), \tag{4.3}$$

where $f_D$ is the probability density function of the distribution $D$. The denominator can be computed by summing up eq. (4.3) for $y \in \{-1, 1\}$. Overall this results in

$$P(Y = y \mid X_2 = x_2, Z = z) = \sigma(2y(z + x_2)).$$

---

[3] $\sigma(x) = 1/(1 + e^{-x})$

Since by construction $X_1 = Z + X_2$, the optimal predictor is again $\hat{Y}^* = X_1$. If the joint distribution $P(Z, Y, \hat{Y}^*)$ is identical in the two scenarios, so are the joint distributions $P(Z, Y, X_1, X_2, \hat{Y}^*)$, because of $X_1 = \hat{Y}^*$ and $X_2 = X_1 - Z$. To show that the joint distributions $P(Z, Y, \hat{Y}^*) = P(Y \mid Z, \hat{Y}^*) P(\hat{Y}^* \mid Z) P(Z)$ are the same, we compare the conditional distributions in the factorization. Let us start with $P(Y \mid Z, \hat{Y}^*)$. Since $\hat{Y}^* = X_1$ and in the first graph $Y \perp\!\!\!\perp Z \mid X_1$, we already found the distribution in eq. (4.1). In the right graph, $P(Y \mid \hat{Y}^*, Z) = P(Y \mid X_2 + Z, Z) = P(Y \mid X_2, Z)$ which we have found in eq. (4.2) and coincides with the conditional in the left graph because of $X_1 = Z + X_2$.

Now consider $\hat{Y}^* \mid Z$. In the left graph we have $P(\hat{Y}^* \mid Z) = P(X_1 \mid Z)$ and the distribution $P(X_1 \mid Z)$ is just the mixture of Gaussians defined in the structural equation model. In the right graph $\hat{Y}^* = Z + X_2 = Y + \mathcal{N}(Z, 1)$ and thus $P(\hat{Y}^* \mid Z) = \mathcal{N}(Z \pm 1)$ for $Y = \pm 1$. Because of the definition of $Y$ in the structural equations of the right graph, following a Bernoulli distribution with probability $\sigma(2Z)$, this is the same mixture of Gaussians as the one we found for the left graph. The distribution of $Z$ is identical in both cases. Consequently the joint distributions agree. When $X_1$ is a resolving variable, the optimal predictor in the left graph does not exhibit unresolved discrimination, whereas the graph on the right does.

The proof for the equal odds predictor $\widetilde{Y}$ is immediate once we show $\widetilde{Y} = X_2$. This can be seen from the graph on the right, because here $X_2 \perp\!\!\!\perp Z \mid Y$ and both using $Z$ or $X_1$ would violate the equal odds condition. Because the joint distribution in the left graph is the same, $\widetilde{Y} = X_2$ is also the optimal equal odds score.                                  $\square$

The two graphs in Figure 4.2 are taken from Hardt et al. (2016b), which we here reinterpret in the causal context to prove Theorem 1. We point out that there is an established set of conditions under which unresolved discrimination can, in fact, be determined from observational data. Note that the two graphs are not Markov equivalent. Therefore, to obtain the same joint distribution we must violate faithfulness.[4] If we do assume the Markov condition and faithfulness, then conditional independences determine the graph up to its Markov equivalence class. We later argue that violation of faithfulness is by no means pathological, but emerges naturally when designing predictors $\hat{Y}$. In any case, interpreting conditional dependences can be difficult in practice (Cornia & Mooij, 2014).

## 4.3   Proxy discrimination and interventions

We now turn to an important aspect of our framework. Determining causal effects in general requires modeling interventions. Interventions on deeply rooted individual properties such as *gender* or *race* are notoriously difficult to conceptualize—especially at an individual level, and impossible to perform in a randomized trial. VanderWeele & Robinson (2014) discuss the

---

[4]See Section 3.5 for definitions of *faithfulness*, *Markov property*, and *Markov equivalence class*.

problem comprehensively in an epidemiological setting and Hu & Kohler-Hausmann (2020) analyze it in the context of fair machine learning using the example of what constitutes sex as a sensitive attribute. From a machine learning perspective, it thus makes sense to separate the protected attribute $Z$ from its potential *proxies*, such as name, visual features, languages spoken at home, etc. Intervention based on proxy variables poses a more manageable problem. By deciding on a suitable proxy we can find an adequate mounting point for determining and removing its influence on the prediction. Moreover, in practice we are often limited to imperfect measurements of $Z$ in any case, making the distinction between root concept and proxy prudent.

As was the case with resolving variables, a *proxy* is a priori nothing more than a descendant of $Z$ in the causal graph that we choose to label as a proxy. Nevertheless in reality we envision the proxy to be a clearly defined observable quantity that is significantly correlated with $Z$, yet in our view should not affect the prediction.

**Definition 2** (Potential proxy discrimination). A variable $V$ in a causal graph exhibits *potential proxy discrimination*, if there exists a directed path from $Z$ to $V$ that is blocked by a proxy variable and $V$ itself is not a proxy.

Potential proxy discrimination articulates a causal criterion that is in a sense dual to unresolved discrimination. From the *benevolent viewpoint*, we *allow* any path from $Z$ to $\hat{Y}$ unless it passes through a proxy variable, which we consider worrisome. This viewpoint acknowledges the fact that the influence of $Z$ on other variables may be complex and it can be too restraining to rule out all but a few designated features. In practice, as with unresolved discrimination, we can naively build an unconstrained predictor based only on those features that do not exhibit potential proxy discrimination. Then we must not provide proxy $P$ as input to $\hat{Y}$;[5] unawareness, i.e., excluding $P$ from the inputs of $\hat{Y}$, suffices. However, by granting $\hat{Y}$ access to $P$, we can carefully tune the function $\hat{Y}(P, X)$ to cancel the implicit influence of $P$ on features $X$ that exhibit potential proxy discrimination by the explicit dependence on $P$. Due to this possible cancellation of paths, we called the path based criterion *potential* proxy discrimination. When building predictors that exhibit no *overall proxy discrimination*, we precisely aim for such a cancellation.

Fortunately, this idea can be conveniently expressed by an *intervention* on $P$, which is denoted by $do(P = p)$ (Pearl, 2009). Visually, intervening on $P$ amounts to removing all incoming arrows of $P$ in the graph; algebraically, it consists of replacing the structural equation of $P$ by $P = p$, i.e., we put point mass on the value $p$.

---

[5]We note that the symbols P (used for the distributions of random variables) and $P$ (denoting proxy variables) may be hard to distinguish. However, the respective meaning will usually be clear from the context.

**Definition 3** (Proxy discrimination). A predictor $\hat{Y}$ exhibits no *proxy discrimination* based on a proxy $P$ if for all $p, p'$

$$\mathrm{P}(\hat{Y} \,|\, do(P = p)) = \mathrm{P}(\hat{Y} \,|\, do(P = p')) . \tag{4.4}$$

The interventional characterization of proxy discrimination leads to a simple procedure to remove it in causal graphs that we will turn to in the next section. It also leads to several natural variants of the definition that we discuss in Section 4.4. We remark that Equation (4.4) is an equality of probabilities in the "do-calculus" that cannot in general be inferred by an observational method, because it depends on an underlying causal graph (Pearl, 2009). However, in some cases, we do not need to resort to interventions to avoid proxy discrimination.

**Proposition 1.** *If there is no directed path from a proxy to a feature, unawareness as in Section 3.2 avoids proxy discrimination.*

*Proof.* An unaware predictor $\hat{Y}$ is given by $\hat{Y} = r(X)$ for some function $r$ and features $X$. If there is no directed path from proxies $P$ to $X$, i.e., $P \notin ta^{\mathcal{G}}(X)$, then $\hat{Y} = r(X) = r(ta^{\mathcal{G}}(X)) = r(ta^{\mathcal{G}}(X) \setminus \{P\})$. Thus $\mathrm{P}(\hat{Y} \,|\, do(P = p)) = \mathrm{P}(\hat{Y})$ for all $p$, which avoids proxy discrimination. $\qquad\square$

## 4.4  Procedures for avoiding discrimination

Having motivated the two types of discrimination that we distinguish, we now turn to building predictors that avoid them in a given causal model. First, we remark that a more comprehensive treatment requires individual judgment of not only variables, but the legitimacy of every existing path that ends in $\hat{Y}$, i.e., evaluation of *path-specific effects* (Zhang & Wu, 2017; Nabi & Shpitser, 2018), which is tedious in practice. The natural concept of proxies and resolving variables covers most relevant scenarios and allows for natural removal procedures.

**Avoiding proxy discrimination.**  While presenting the general procedure, we illustrate each step in the example shown in Figure 4.3. A protected attribute $Z$ affects a proxy $P$ as well as a feature $X$. Both $P$ and $X$ have additional unobserved causes $N_P$ and $N_X$, where $N_P, N_X, Z$ are pairwise independent. Finally, the proxy also has an effect on the features $X$ and the predictor $\hat{Y}$ is a function of $P$ and $X$. Given labeled training data, our task is to find a good predictor that exhibits no proxy discrimination within a hypothesis class of functions $\hat{Y}_\theta(P, X)$ parameterized by a real valued vector $\theta$.

Figure 4.3 A template graph $\tilde{\mathcal{G}}$ for proxy discrimination (left) with its intervened version $\mathcal{G}$ (right). While from the benevolent viewpoint we do not generically prohibit any influence from $Z$ on $\hat{Y}$, we want to guarantee that the proxy $P$ has no overall influence on the prediction, by adjusting $P \to \hat{Y}$ to cancel the influence along $P \to X \to \hat{Y}$ in the intervened graph.

Figure 4.4 A template graph $\tilde{\mathcal{G}}$ for unresolved discrimination (left) with its intervened version $\mathcal{G}$ (right). While from the skeptical viewpoint we generically do not want $Z$ to influence $\hat{Y}$, we first intervene on $E$ interrupting all paths through $E$ and only cancel the remaining influence from $Z$ to $\hat{Y}$.

We now work out a formal procedure to solve this task under specific assumptions. While the general procedure in principle works for arbitrarily large graphs and potentially non-linear structural equations, we will simultaneously illustrate the algorithm in a small fully linear example for clarity. We will use gray font for the specific example, in which the structural equations are given by

$$P = \alpha_P Z + N_P, \qquad X = \alpha_X Z + \beta P + N_X, \qquad \hat{Y}_\theta = \lambda_P P + \lambda_X X.$$

Note that we choose linear functions parameterized by $\theta = (\lambda_P, \lambda_X)$ as the hypothesis class for $\hat{Y}_\theta(P, X)$ in this example.

In the procedure we clarify the notion of *expressibility*, which is an assumption about the relation of the given structural equations and the hypothesis class we choose for $\hat{Y}_\theta$.

**Proposition 2.** *If there is a choice of parameters $\theta_0$ such that $\hat{Y}_{\theta_0}(P, X)$ is constant with respect to its first argument and the structural equations are* expressible, *the following procedure returns a predictor from the given hypothesis class that exhibits no proxy discrimination and is non-trivial in the sense that it can make use of features that exhibit potential proxy discrimination.*

1. Intervene on $P$ by removing all incoming arrows and replacing the structural equation for $P$ by $P = p$. For the example in Figure 4.3,

$$P = p, \qquad X = \alpha_X Z + \beta P + N_X, \qquad \hat{Y}_\theta = \lambda_P P + \lambda_X X. \qquad (4.5)$$

2. Iteratively substitute variables in the equation for $\hat{Y}_\theta$ from their structural equations until only root nodes of the intervened graph are left, i.e., write $\hat{Y}_\theta(P, X)$ as $\hat{Y}_\theta(P, g(ta^{\mathcal{G}}(X)))$

for some function $g$. Since the causal graph is acyclic, we can write any variable as a function of only root nodes by iteratively substituting parent nodes with their structural equations. In the example, $ta(X) = \{Z, P, N_X\}$ and

$$\hat{Y}_\theta = g(ta(X)) = (\lambda_P + \lambda_X \beta)p + \lambda_X(\alpha_X Z + N_X). \tag{4.6}$$

3. We now require the distribution of $\hat{Y}_\theta$ in eq. (4.6) to be independent of $p$. In the example, we require for all $p, p'$ that

$$P((\lambda_P + \lambda_X \beta)p + \lambda_X(\alpha_X Z + N_X)) = P((\lambda_P + \lambda_X \beta)p' + \lambda_X(\alpha_X Z + N_X)). \tag{4.7}$$

In general, the goal is to write the predictor as a function of $P$ and all the other roots of $\mathcal{G}$ separately. If our hypothesis class is such that there exists a parameter vector $\tilde{\theta}$ and a function $\tilde{g}$ such that $\hat{Y}_\theta(P, g(ta(X))) = \hat{Y}_{\tilde{\theta}}(P, \tilde{g}(ta(X) \setminus \{P\}))$, we call the structural equation model and hypothesis class specified in eq. (4.5) *expressible*. Equation (4.7) then yields the *non-discrimination constraint* $\tilde{\theta} = \theta_0$. Our example model is expressible with $\tilde{\theta} = (\lambda_P + \lambda_X \beta, \lambda_X)$ and $\tilde{g} = \alpha_X Z + N_X$. A possible $\theta_0$ is $\theta_0 = (0, \lambda_X)$, which simply yields $\lambda_P = -\lambda_X \beta$.

4. Given labeled training data, we can optimize the predictor $\hat{Y}_\theta$ within the hypothesis class as given in eq. (4.5), subject to the non-discrimination constraint. In our linear example, we obtain
$$\hat{Y}_\theta = -\lambda_X \beta P + \lambda_X X = \lambda_X(X - \beta P),$$
with the free parameter $\lambda_X \in \mathbb{R}$ that we can optimize for accuracy.

In general, the non-discrimination constraint eq. (4.7) is by construction just $P(\hat{Y} \mid do(P = p)) = P(\hat{Y} \mid do(P = p'))$, coinciding with Definition 3. Thus Proposition 2 holds by construction of the procedure. The choice of $\theta_0$ strongly influences the non-discrimination constraint. However, as the example shows, it allows $\hat{Y}_\theta$ to exploit features that exhibit potential proxy discrimination.

**Avoiding unresolved discrimination.** We proceed analogously to the previous section using the example graph in Figure 4.4. Instead of the proxy, we consider a resolving variable $E$. For the running example, the causal dependences are equivalent to the ones in Figure 4.3 and we again assume linear structural equations for a running example

$$E = \alpha_E Z + N_E, \qquad X = \alpha_X Z + \beta E + N_X, \qquad \hat{Y}_\theta = \lambda_E E + \lambda_X X.$$

Let us now try to adjust the previous procedure to the context of avoiding unresolved discrimination.

1. Intervene on $E$ by fixing it to a random variable $\eta$ with $P(\eta) = P(E)$, the marginal distribution of $E$ in $\tilde{\mathcal{G}}$, see Figure 4.4. In the example we find

$$E = \eta, \qquad X = \alpha_X Z + \beta E + N_X, \qquad \hat{Y}_\theta = \lambda_E E + \lambda_X X. \tag{4.8}$$

2. By iterative substitution of structural equations, write $\hat{Y}_\theta(E, X)$ as $\hat{Y}_\theta(E, g(ta^{\mathcal{G}}(X)))$ for some function $g$. In the example

$$\hat{Y}_\theta = g(ta^{\mathcal{G}}(X)) = (\lambda_E + \lambda_X \beta)\eta + \lambda_X \alpha_X Z + \lambda_X N_X. \tag{4.9}$$

3. We now demand the distribution of $\hat{Y}_\theta$ in eq. (4.9) be invariant under interventions on $Z$, which coincides with conditioning on $Z$ whenever $Z$ is a root of $\tilde{\mathcal{G}}$. Hence, in the example, for all $z, z'$

$$P((\lambda_E + \lambda_X \beta)\eta + \lambda_X \alpha_X z + \lambda_X N_X) = P((\lambda_E + \lambda_X \beta)\eta + \lambda_X \alpha_X z' + \lambda_X N_X). \tag{4.10}$$

Here, the subtle asymmetry between proxy discrimination and unresolved discrimination becomes apparent. Our approach relies on "bundling" directed paths from the sensitive variable $Z$ to the predictor $\hat{Y}_\theta$ into the ones that pass through $P$ and $E$ respectively. In the benevolent setting, when talking about proxy discrimination, the goal is to correct for the influence of all paths passing through $P$. Since our predictor $\hat{Y}_\theta$ can use the proxy $P$ directly as input, under the expressibility assumption, we can tune it such that it precisely counteracts the aggregate influence $P$ has on $\hat{Y}_\theta$ mediated by all other possible features.

On the other hand, in the skeptic setting, when talking about resolved discrimination, the goal is to only allow for the influences mediated by $E$. Hence we would similarly seek to cancel out all paths from $Z$ on $\hat{Y}_\theta$ that do not pass through $E$. However, in this setting $\hat{Y}_\theta$ is not explicitly a function of $Z$. Therefore, we cannot tune the predictor to cancel implicit influences of $Z$ through $X$. There might still be a $\theta_0$ such that $\hat{Y}_{\theta_0}$ indeed fulfills eq. (4.10), but without being able to use $Z$ as a direct input for the predictor there is no principled way for us to construct it.

In the example, eq. (4.10) suggests the obvious *non-discrimination constraint* $\lambda_X = 0$. We can then proceed as before and, given labeled training data, optimize $\hat{Y}_\theta = \lambda_E E$ by varying $\lambda_E$.

However, by setting $\lambda_X = 0$, we also cancel the path $Z \to E \to X \to \hat{Y}$, even though it is blocked by a resolving variable. In general, if $\hat{Y}_\theta$ does not have access to $Z$, we can not adjust for unresolved discrimination without also removing resolved influences from $Z$ on $\hat{Y}_\theta$. If, however, $\hat{Y}_\theta$ is a function of $Z$, i.e., we add the term $\lambda_Z Z$ to $\hat{Y}_\theta$ in eq. (4.8), the non-discrimination constraint is $\lambda_Z = -\lambda_X \alpha_X$ and we can proceed analogously to the procedure for avoiding proxy discrimination.

We want to end on a final remark about the intuition underlying the expressibility assumption. The key point in removing proxy discrimination is to "separate out" the effect of $P$ on $\hat{Y}_\theta$ via features $X$ such that it can be "neutralized" by using $P$ as a separate input. For example, in a fully linear system, $P$ can only enter as a single additive linear term via features $X$. The influence of $P$ can thus easily be eliminated or neutralized by subtracting precisely that term, which we can do via direct access to $P$. A non-linear example of functional forms with a similarly meaningful notion of "neutralization" is when $P$ enters as an overall multiplicative term that only depends on $P$ and is non-zero. In this case, we can divide the entire expression by that term. The expressibility assumption captures all the scenarios in which the influence of $P$ can be meaningfully eliminated. We note that there are also non-linear scenarios in which this is not possible. For example, if $P$ enters $\hat{Y}_\theta$ through $X$ in a functional form such as $\sin(e^{X_1 \cdot P})$ (where $X_1$ is one of the observed non-sensitive features), it is not clear how one should "neutralize" the effect of $P$ by providing it as a separate argument to $\hat{Y}_\theta$.

**Relating proxy discrimination to other notions of fairness.** Motivated by the algorithm to avoid proxy discrimination, we discuss some natural variants of the notion that connect our interventional approach to individual fairness and other proposed criteria. Given the multitude of proposed fairness criteria outlined in Chapter 3 together with their incompatibilities and inconsistent naming conventions, it is useful to compare criteria and find specific assumptions under which some of them may become mathematically equivalent. In this part, we explore some such connections for a generic graph structure as shown on the left in Figure 4.5. The proxy $P$ and the features $X$ could be multidimensional. The empty circle in the middle represents any number of variables forming a DAG that respects the drawn arrows. This general DAG may also be empty. That means that, e.g., arrows between $X$ and $P$ directly are also allowed. Figure 4.3 is an example thereof. All dashed arrows are optional depending on the specifics of the situation.

**Definition 4.** A predictor $\hat{Y}$ exhibits no *individual proxy discrimination*, if for all $x$ and all $p, p'$

$$P(\hat{Y} \mid do(P = p), X = x) = P(\hat{Y} \mid do(P = p'), X = x).$$

A predictor $\hat{Y}$ exhibits no *proxy discrimination in expectation*, if for all $p, p'$

$$\mathbb{E}[\hat{Y} \mid do(P = p)] = \mathbb{E}[\hat{Y} \mid do(P = p')].$$

Figure 4.5 *Left:* A generic graph $\tilde{\mathcal{G}}$ to describe proxy discrimination. *Right:* The graph corresponding to an intervention on $P$. The circle labeled "DAG" represents any sub-DAG of $\tilde{\mathcal{G}}$ and $\mathcal{G}$ containing an arbitrary number of variables that is compatible with the shown arrows. Dashed arrows can, but do not have to be present in a given scenario.

Individual proxy discrimination aims at comparing examples with the same features $X$, for different values of $P$. Note that this can be individuals with different values for the unobserved non-feature variables. A true individual-level comparison of the form "What would have happened to me, if I had always belonged to another group" is captured by counterfactual fairness (Kusner et al., 2017; Nabi & Shpitser, 2018).

For an analysis of proxy discrimination, we need the structural equations for $P, X, \hat{Y}$ in Figure 4.5

$$P = \hat{f}_P(\mathrm{pa}(P)) ,$$
$$X = \hat{f}_X(\mathrm{pa}(X)) = f_X(P, ta^{\mathcal{G}}(X) \setminus \{P\}) ,$$
$$\hat{Y} = \hat{f}_{\hat{Y}}(P, X) = f_{\hat{Y}}(P, ta^{\mathcal{G}}(\hat{Y}) \setminus \{P\}) .$$

For convenience, we will use the notation $ta_P^{\mathcal{G}}(X) := ta^{\mathcal{G}}(X) \setminus \{P\}$. We can find $f_X, f_{\hat{Y}}$ from $\hat{f}_X, \hat{f}_{\hat{Y}}$ by first rewriting the functions in terms of root nodes of the *intervened graph*, shown on the right side of Figure 4.5, and then assigning the *overall* dependence on $P$ to the first argument.

We now compare proxy discrimination to other existing notions.

**Theorem 2.** *Let the influence of P on X be additive and linear, i.e.,*

$$X = f_X(P, ta_P^{\mathcal{G}}(X)) = g_X(ta_P^{\mathcal{G}}(X)) + \mu_X P$$

*for some function $g_X$ and $\mu_X \in \mathbb{R}$. Then any predictor of the form*

$$\hat{Y} = r(X - \mathbb{E}[X \mid do(P)])$$

*for some function r exhibits no proxy discrimination.*

*Proof.* It suffices to show that the argument of $r$ is constant with respect to $P$, because then $\hat{Y}$ and thus $\mathrm{P}(\hat{Y})$ are invariant under changes of $P$. We compute

$$\mathbb{E}[X \mid do(P)] = \mathbb{E}[g_X(ta_P^{\mathcal{G}}(X)) + \mu_X P \mid do(P)]$$
$$= \underbrace{\mathbb{E}[g_X(ta_P^{\mathcal{G}}(X)) \mid do(P)]}_{=0} + \mathbb{E}[\mu_X P \mid do(P)]$$
$$= \mu_X P .$$

Hence,

$$X - \mathbb{E}[X \mid do(P)] = g_X(ta_P^{\mathcal{G}}(X))$$

is constant with respect to $P$. □

Note that in general $\mathbb{E}[X \mid do(P)] \neq \mathbb{E}[X \mid P]$. Since in practice we only have observational data from $\tilde{\mathcal{G}}$, one cannot simply build a predictor based on the "regressed out features" $\tilde{X} := X - \mathbb{E}[X \mid P]$ to avoid proxy discrimination. In the scenario of Figure 4.3, the direct effect of $P$ on $X$ along the arrow $P \to X$ in the left graph cannot be estimated by $\mathbb{E}[X \mid P]$, because of the common confounder $Z$. The desired interventional expectation $\mathbb{E}[X \mid do(P)]$ coincides with $\mathbb{E}[X \mid P]$ only if one of the arrows $Z \to P$ or $Z \to X$ is not present. Estimating direct causal effects is a hard problem, well studied by the causality community and often involves instrumental variables (Angrist & Krueger, 2001).

In Section 3.6 we provided a short overview of learning fair representation. A common theme in this field is to learn a representation $\tilde{X}$ that is statistically independent of the sensitive data, but still contains information to predict $Y$. A natural idea is to use $X - \mathbb{E}[X \mid Z]$ as input to a downstream classifier instead of $X$ to "regress out" linear correlations. Within our setting, where we try to avoid proxy discrimination, we typically want to remove the effect of proxies and not the protected attribute directly. Therefore, we can ask the question when using the representation $\tilde{X} := X - \mathbb{E}[X \mid P]$ suffices to remove proxy discrimination, i.e., when the difference between $\mathbb{E}[X \mid P]$ and $\mathbb{E}[X \mid do(P)]$ becomes irrelevant for our purposes.

**Corollary 1.** *Under the assumptions of Theorem 2, if all directed paths from any ancestor of $P$ to $X$ in the graph $\mathcal{G}$ are blocked by $P$, then any predictor based on the* adjusted features $\tilde{X} := X - \mathbb{E}[X \mid P]$ *exhibits no proxy discrimination and can be learned from the observational distribution $\mathrm{P}(P, X, Y)$ when target labels $Y$ are available.*

*Proof.* Let $A$ denote the set of ancestors of $P$. Under the given assumptions $A \cap ta^{\mathcal{G}}(X) = \emptyset$, because in $\mathcal{G}$ all arrows into $P$ are removed, which breaks all directed paths from any variable in $A$ to $X$ by assumption. Hence the distribution of $X$ under an intervention on $P$ in $\tilde{\mathcal{G}}$, where the influence of potential ancestors of $P$ on $X$ that does not go through $P$ would not be affected, is the same as simply conditioning on $P$. Therefore $\mathbb{E}[X \mid do(P)] = \mathbb{E}[X \mid P]$, which

can be computed from the joint observational distribution, since we observe $X$ and $P$ as generated by $\tilde{\mathcal{G}}$. □

Our definition of proxy discrimination in expectation in eq. (4) is motivated by a weaker notion proposed by Calders & Verwer (2010). It asks for the expected outcome to be the same across the different populations $\mathbb{E}[\hat{Y} \,|\, P = p] = \mathbb{E}[\hat{Y} \,|\, P = p']$. Again, when talking about proxies, we must be careful to distinguish conditional and interventional expectations, which is captured by the following proposition and its corollary.

**Proposition 3.** *Any predictor of the form $\hat{Y} = \lambda(X - \mathbb{E}[X \,|\, do(P)]) + c$ for $\lambda, c \in \mathbb{R}$ exhibits no proxy discrimination in expectation.*

*Proof.* We directly test the definition of proxy discrimination in expectation using linearity of expectation

$$
\begin{aligned}
\mathbb{E}[\hat{Y} \,|\, do(P = p)] &= \mathbb{E}[\lambda(X - \mathbb{E}[X \,|\, do(P)]) + c \,|\, do(P = p)] \\
&= \lambda(\mathbb{E}[X \,|\, do(P = p)] - \mathbb{E}[X \,|\, do(P = p)]) + c \\
&= c \,.
\end{aligned}
$$

This holds for any $p$, hence proxy discrimination in expectation is achieved. □

From this and the proof of Corollary 1 we conclude the following Corollary.

**Corollary 2.** *If all directed paths from any ancestor of $P$ to $X$ are blocked by $P$, any predictor of the form $\hat{Y} = r(X - \mathbb{E}[X \,|\, P])$ for linear $r$ exhibits no proxy discrimination in expectation and can be learned from the observational distribution $\mathrm{P}(P, X, Y)$ when target labels $Y$ are available.*

Finally, we provide an additional statement that is a first step towards the "opposite direction" of Theorem 2, i.e., whether we can infer information about the structural equations, when we are given a predictor of a special form that does not exhibit proxy discrimination.

**Theorem 3.** *Let the influence of $P$ on $X$ be additive and linear and let the influence of $P$ on the argument of $\hat{Y}$ be additive linear, i.e.,*

$$
\begin{aligned}
f_X(ta^{\mathcal{G}}(X)) &= g_X(ta_P^{\mathcal{G}}(X)) + \mu_X P \\
f_{\hat{Y}}(P, ta_P^{\mathcal{G}}(X)) &= h(g_{\hat{Y}}(ta_P^{\mathcal{G}}(X)) + \mu_{\hat{Y}} P)
\end{aligned}
$$

*for some functions $g_X, g_{\hat{Y}}$, real numbers $\mu_X, \mu_{\hat{Y}} \in \mathbb{R}$ and a smooth, strictly monotonic function $h$. Then any predictor that avoids proxy discrimination is of the form*

$$
\hat{Y} = r(X - \mathbb{E}[X \,|\, do(P)])
$$

*for some function r.*

*Proof.* From the linearity assumptions we conclude that

$$\hat{f}_{\hat{Y}}(P, X) = h(g_X(ta_P^{\mathcal{G}}(X)) + \mu_X P + \hat{\mu}_{\hat{Y}} P),$$

with $\hat{\mu}_{\hat{Y}} = \mu_{\hat{Y}} - \mu_X$ and thus $g_X = g_{\hat{Y}}$. That means that both the dependence of $X$ on $P$ along the path $P \to \cdots \to X$ as well as the direct dependence of $\hat{Y}$ on $P$ along $P \to \hat{Y}$ are additive and linear.

To avoid proxy discrimination, we need

$$\begin{aligned}
P(\hat{Y} \,|\, do(P = p)) &= P(h(g_{\hat{Y}}(ta_P^{\mathcal{G}}(X)) + \mu_{\hat{Y}} p)) \\
&\stackrel{!}{=} P(h(g_{\hat{Y}}(ta_P^{\mathcal{G}}(X)) + \mu_{\hat{Y}} p')) = P(\hat{Y} \,|\, do(P = p')).
\end{aligned}$$

Because $h$ is smooth an strictly monotonic, we can conclude that already the distributions of the argument of $h$ must be equal, otherwise the transformation of random variables could not result in equal distributions, i.e.,

$$P(g_{\hat{Y}}(ta_P^{\mathcal{G}}(X)) + \mu_{\hat{Y}} p) \stackrel{!}{=} P(g_{\hat{Y}}(ta_P^{\mathcal{G}}(X)) + \mu_{\hat{Y}} p').$$

Since, up to an additive constant, we are comparing the distributions of the *same* random variable $g_{\hat{Y}}(ta_P^{\mathcal{G}}(X))$ and not merely identically distributed ones, the following condition is not only sufficient, but also necessary for eq. (4.12)

$$g_{\hat{Y}}(ta_P^{\mathcal{G}}(X)) + \mu_{\hat{Y}} p \stackrel{!}{=} g_{\hat{Y}}(ta_P^{\mathcal{G}}(X)) + \mu_{\hat{Y}} p'.$$

This holds true for all $p, p'$ only if $\mu_{\hat{Y}} = 0$, which is equivalent to $\hat{\mu}_{\hat{Y}} = -\mu_X$.

Because as in the proof of Theorem 2

$$\mathbb{E}[X \,|\, do(P)] = \mu_X P,$$

under the given assumptions any predictor that avoids proxy discrimination is simply

$$\hat{Y} = X + \hat{\mu}_{\hat{Y}} P = X - \mathbb{E}[X \,|\, do(P)].$$

$\square$

## 4.5 Conclusion

The perspective developed in this chapter so far naturally addresses shortcomings of earlier statistical approaches. Causal fairness criteria are suitable whenever we are willing to make assumptions about the (causal) generating process governing the data. Whilst not always feasible, the causal approach naturally creates an incentive to scrutinize the data more closely and work out plausible assumptions to be discussed alongside any conclusions regarding fairness. Key concepts of our conceptual framework are *resolving variables* and *proxy variables* that play a dual role in defining causal discrimination criteria. We develop a practical procedure to remove proxy discrimination given the structural equation model and analyze a similar approach for unresolved discrimination. In the case of proxy discrimination for linear structural equations, the procedure has an intuitive form that is similar to heuristics already used in the regression literature. Our framework is limited by the assumption that we can construct a valid causal graph. The removal of proxy discrimination moreover depends on the functional form of the causal dependencies. This dependence is captured by the notion of expressibility. The causal perspective suggests a number of interesting new directions at the technical, empirical, and conceptual level.

# Sensitivity of causal fairness

In this chapter, we go one step beyond our findings in Chapter 4 and also scrutinize one of the primary assumptions in causal modeling, namely that the causal graph is known. Potential misspecifications of the causal model introduce new opportunities for bias. One common way for misspecification to occur is via *unmeasured confounding*: the true causal effect between variables is partially described by unobserved quantities. We develop tools to assess the sensitivity of fairness measures to this confounding for the popular class of non-linear additive noise models (ANMs). Specifically, we give a procedure for computing the maximum difference between two counterfactually fair predictors, where one has become biased due to confounding. For the case of bivariate confounding our technique can be swiftly computed via a sequence of closed-form updates. For multivariate confounding we give an algorithm that can be efficiently solved via automatic differentiation. We demonstrate our new sensitivity analysis tools in real-world fairness scenarios to assess the bias arising from confounding.

The main content of this chapter has been published in the following paper:

<div style="border:1px solid">

THE SENSITIVITY OF COUNTERFACTUAL FAIRNESS TO UNMEASURED CONFOUNDING

Niki Kilbertus, Philip J. Ball, Matt Kusner, Adrian Weller, Ricardo Silva.

[https://github.com/nikikilbertus/cf-fairness-sensitivity]

*Uncertainty in Artificial Intelligence (UAI), 2019*

</div>

## 5.1 Introduction

In this section, we will focus on *counterfactual fairness* (CF) as introduced by Kusner et al. (2017), an individual-specific criterion aimed at answering the counterfactual question: "What would have been my prediction if—all else held causally equal—I was a member of another protected group?". Despite the utility of such causal criteria, they are often contested, because they are based on strong assumptions that are hard to verify in practice.

First and foremost, all causal fairness criteria proposed in the literature assume that the causal structure of the problem is known. Typically, one relies on domain experts and methods for causal discovery from data to construct a plausible causal graph. While it is often possible with few variables to get the causal graph approximately right, one often needs untestable assumptions to construct the full graph. The most common untestable assumption is that there is no unmeasured confounding between some variables in the causal graph. Because we cannot measure it, this confounding can introduce bias that is unaccounted for by causal fairness criteria.

As a solution to this problem, we will introduce tools to measure the sensitivity of the popular *counterfactual fairness* criterion to unmeasured confounding. Our tools are designed for the commonly used class of non-linear additive noise models (ANMs, Hoyer et al., 2008). Specifically, they describe how counterfactual fairness changes under a given amount of confounding. The core ideas here described can be adapted for sensitivity analysis of other measures of causal effect, such as the average treatment effect (ATE), itself a topic not commonly approached in the context of graphical causal models. Note that counterfactual fairness poses extra challenges compared to the ATE, as it requires the computation of counterfactuals in the sense of Pearl (2009). Concretely, in the remainder of this chapter we will develop the following tools:

- For confounding between two variables, we design a fast procedure for estimating the worst-case change in counterfactual fairness due to confounding. It consists of a series of closed-form updates assuming linear models with non-linear basis functions. This family of models is particularly useful in graphical causal models where any given node has only few parents.

- For more than two variables, we fashion an efficient procedure that leverages automatic differentiation to estimate worst-case counterfactual fairness. In particular, compared to standard sensitivity analysis (typically applied to ATE problems, see e.g. Dorie et al., 2016), we formulate the problem in a multivariate setting as opposed to the typical bivariate case. The presence of other modeling constraints brings new challenges not found in the standard literature.

- We demonstrate that our method allows us to understand how fairness guarantees degrade based on different confounding levels. We also show that even under high levels of confounding, learning counterfactually fair predictors has lower fairness degradation than standard predictors using all features or using all features save for the protected attributes.

## 5.2   Background

We focus on a subclass of SEMs called *additive noise models* (ANMs) (Hoyer et al., 2008). This means that the structural equation for a node $X$ of the causal graph is given by $X = f_X(\text{pa}_{\mathcal{G}}(X)) + \epsilon$ for a non-linear function $f_X$. Here, we use $\epsilon$ for independent noise variables instead of $N$ for notational convenience. To make model fitting efficient, we will consider (a) functions $f_X$ that derive all their non-linearity from an embedding function $\boldsymbol{\phi}$ of their direct parents, and are linear in this embedding; and (b) Gaussian noise (error) $\epsilon$ so that:

$$X = \boldsymbol{\phi}(\text{pa}_{\mathcal{G}}(X))^{\top} \boldsymbol{w}_X + \epsilon \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \sigma_X),$$

where $\boldsymbol{w}_X$ are weights. Later on, we will consider ANMs over observed variables, where the noises may be correlated. Note that this class of ANMs is not closed under marginalization. For a more detailed analysis of the testable implications of the ANM assumption, see (Peters et al., 2017). Neither of our choices (a) and (b) are a fundamental limitation of our framework: the framework can easily be extended to general non-linear, or even non-parametric functions $f_X$, as well as non-Gaussian noises. We make this choice to balance flexibility and computational cost.

For convenience, we restate the definition of counterfactual fairness, where we allow for continuous target values $y \in \mathcal{Y} \subset \mathbb{R}$:

$$P(\hat{Y}_{Z=z'} = y \mid X = x, Z = z) = P(\hat{Y}_{Z=z} = y \mid X = x, Z = z), \quad (5.1)$$

where $\hat{Y}_{Z=z'}$ is the counterfactual prediction, imagining $Z = z'$ (note that, because in reality $Z = z$, we have that $\hat{Y}_{Z=z} = \hat{Y}$), and $x$ is a realization of other variables in the causal system. In ANMs $\hat{Y}_{Z=z'}$ can be computed in four steps:

1. Fit the parameters of the assumed causal model using the observed data: $\mathcal{D} = \{x_i, z_i\}_{i=1}^{n}$;

2. Using the fitted model and data $\mathcal{D}$, estimate all noise variables $\epsilon$;

3. Replace $Z$ with counterfactual value $z'$ in all causal model equations;

4. Using the fitted parameters, estimated noise variables, and $z'$, recompute all variables affected (directly or indirectly) by $Z$, and recompute the prediction $\hat{Y}$.

To learn a CF predictor satisfying eq. (5.1) it is sufficient to use any variables that are non-descendants of $Z$, such as the noise variables $\epsilon$ (Kusner et al., 2017). It may appear as if any predictor using only the noise variables as input is going to have low accuracy. however, we note that in our setting the so called noise variables are not to be thought of as mere measurement noise that contains little information, but summarizes all other influences on

their corresponding variables that have not been explicitly observed. In this sense, the noise variables may be highly predictive of their respective variables and other quantities in the causal graph. Thus it is not unreasonable to build counterfactually fair predictors using only the noise variables as inputs, which are non-descendants of the sensitive attribute by assumption.

One key assumption on which CF relies is that there is no *unmeasured confounding* relationship missing in the causal model. For our purposes, we formalize unmeasured confounding as non-zero correlations between any two noise variables in $\epsilon$ which are assumed to follow a multivariate Gaussian distribution. Without accounting for this, the above counterfactual procedure will compute noise variables that are not guaranteed to be independent of $Z$. Thus any predictor trained on these exogenous variables is not guaranteed to satisfy counterfactual fairness eq. (5.1). This setup captures the idea that often we have a decent understanding of the causal structure, but might overlook confounding effects, here in the form of pairwise correlations of noise variables. At the same time, such confounding is often unidentifiable (save for specific parameterizations). Thus assessing confounding is not a model selection problem but a sensitivity analysis problem. To perform such analysis we propose tools to measure the worst-case deviation in CF due to unmeasured confounding. Before describing these tools, we first place them in the context of the long tradition of sensitivity analysis in causal modeling.

**Traditional sensitivity analysis**   Sensitivity analysis for quantities such as the average treatment effect can be traced back at least to the work by Jerome Cornfield on the General Surgeon study concerning the smoking and lung cancer link (Rosenbaum, 2002). Rosenbaum cast the problem in a more explicit statistical framework, addressing the question on how the ATE would vary if some degree of association between a treatment and a outcome was due to unmeasured confounding. The logic of sensitivity analysis can be described in a simplified way as follows: i) choose a level of "strength" for the contribution of a latent variable to the structural equation(s) of the treatment and/or outcome; ii) by fixing this confounder contribution, estimate the corresponding ATE; iii) vary steps i) and ii) through a range of "confounding effects" to report the level of unmeasured confounding required to make the estimate ATE be statistically indistinguishable from zero; iv) consult an expert to declare whether the level of confounding required for that to happen is too strong to be plausible, and if so, conclude that the effect is real to the best of one's knowledge. This basic idea has led to a large literature, see (Dorie et al., 2016; Robins et al., 2000) for two noteworthy examples.

Note the crucial difference between sensitivity analysis and just fitting a latent variable model: we are not learning a latent variable distribution, as the confounding effect for a single cause-effect pair is *unidentifiable*. By holding the contribution of the confounder

as constant and known, the remaining parameters become identifiable. We can vary the sensitivity parameter without assuming a probability measure on the confounding effect. The hypothesis test mentioned in the example above can be substituted by other criteria of practical significance.

Much of the work in the statistics literature on sensitivity analysis addresses pairs of cause-effects as opposed to a causal system with intermediate outcomes, and focuses on the binary question on when an effect is non-zero. The *grid search* idea of attempting different levels of the confounding level does not necessarily translate well to a full SEM: grid search grows exponentially with the number of pairs of variables. In our problem formulation described in the sequel, we are interested in bounding the maximum magnitude $p_{\max}$ of the noise correlation matrix entries, while maximizing a measure of counterfactual unfairness to understand how it varies by the presence of unmeasured confounding. The solution is not always to set all entries to $p_{\max}$, since among other things we may be interested in keeping a subset of noise correlations to be zero. In this case, a sparse correlation matrix with all off-diagonal values set to either $0$ or $p_{\max}$ is not necessarily positive-definite. A multidimensional search for the entries of the confounding correlation matrix is then necessary, which we will do in Section 5.4 by encoding everything as a fully differentiable and unconstrained optimization problem.

Finally, we note that unmeasured confounding is but one of many possible misspecifications of causal assumptions. the true causal effect between variables is partially described by un-observed quantities. For example, we still assume that all edges between observed variables are correctly identified. Russell et al. (2017) previously addressed how to enforce counterfactual fairness across a small enumeration of different competing—but identifiable—models. Instead, we only focus on misspecification in the form of unmeasured confounding.

## 5.3   Tool #1: grid-based analysis

The notion of sensitivity analysis in an SEM can be complex, particularly when the estimated quantity involves counterfactuals. Therefore, we first describe a tool that estimates the effect of confounding on counterfactual fairness, when the confounding is limited to two variables (i.e., *bivariate confounding*). This procedure is computationally efficient for this setting. For the general setting of confounding between any number of variables (*multivariate confounding*) we will introduce a separate tool in Section 5.4. We now describe our fast two-variable tool using a real-world example.

**A motivating example.**   To motivate our approach, let us revisit the example about law school success analyzed by Kusner et al. (2017). In this task, we want to predict the first year average grade ($Y$) of incoming law school students from their grade-point average ($G$) before

Figure 5.1 Causal models for the law school example. Model A is the guessed model that has no unobserved confounding. Model B includes confounding via the covariance matrix $\Sigma$, which is captured by a bidirected edge using the standard acyclic directed mixed graph notation (ADMG, Richardson, 2003). Our techniques will estimate the worst case difference in the estimation of counterfactual fairness due to such confounding (we will consider a more complicated setup in Section 5.5).

entering law school and their law school admission test scores ($L$). In the original work, the goal was to train a predictor $\hat{Y}$ that is counterfactually fair with respect to race.

To evaluate any causal notion of fairness, we need to first specify the causal graph. Here we assume $G \to L$ with noises $\epsilon_G, \epsilon_L$, where $G$ and $L$ are both influenced by the sensitive attribute $Z$, see **Model A** in Figure 5.1. Given this specification, the standard way to train a counterfactually fair classifier is using $\epsilon_G, \epsilon_L$—the non-descendants of $Z$. To do so, we first learn them from data as the residuals in predicting $G$ and $L$ from their parents.

The validity of causal estimates rely on the assumption that the constructed causal model and its respective graph (here Model A) captures the true data-generating mechanism. While previous work addressed how to enforce counterfactual fairness across a small enumeration of identifiable competing models (Russell et al., 2017), in this work we consider misspecification in the lack of *unidentifiable* unmeasured confounding. In our example, this means violation of the assumed independence of the noise variables $\epsilon_G$ and $\epsilon_L$.

To capture such confounding, we introduce **Model B** in Figure 5.1. Here the noise variables are not independent, they co-vary: $(\epsilon_G, \epsilon_L)^\top \sim \mathcal{N}(\mathbf{0}, \Sigma)$ where,

$$\Sigma = \begin{pmatrix} \sigma_G^2 & p\,\sigma_G\,\sigma_L \\ p\,\sigma_G\,\sigma_L & \sigma_L^2 \end{pmatrix}.$$

Here, $\sigma_\bullet$ is the standard deviation of $\bullet$ and $p \in [-1, 1]$ is the correlation, such that the overall covariance matrix $\Sigma$ is positive semidefinite. Before going into the detailed procedure of our sensitivity analysis, let us give a general description of what we mean by Model A and Model B throughout this work.

**Model A** is the "guessed" causal graph model used to build a counterfactually fair predictor. The assumption we will scrutinize is that this guess at the true underlying causal model is perfectly correct.

**Model B** is a version of Model A that allows for further unobserved confounding between pairs of noise variables not originally featured in A. Model B will play the role of a hypothetical ground truth that simulates "true" counterfactual versions of the predictions made within Model A.

Our tool allows us to answer the following question: how does a predictor that is counterfactually fair under Model A perform in terms of counterfactual unfairness under the confounded Model B? Our goal is to quantify how sensitive counterfactual unfairness is to misspecifications of the causal model, in particular to unobserved confounding. To do so, we will introduce a measure which we will call *counterfactual unfairness* (CFU). Given this, we describe how to compute the worst-case violation of counterfactual fairness within a certain confounding budget, which we characterize by the correlation $-1 \leq p_{\max} \leq 1$ in Model B. By varying the confounding budget, we can assess how robust Model A is to different degrees of model misspecification. Like in classical sensitivity analysis, we can alternatively start from a level of unacceptable CFU, search for the minimum $p_{\max}$ whose worst-case CFU reaches this level, and leave it to domain experts to judge the plausibility of such a degree of unmeasured confounding $p_{\max}$.

**Notation and problem setup.**   For both Model A and B the structural equations are:

$$G = \boldsymbol{\phi}_G(Z)^\top \boldsymbol{w}_G + \epsilon_G \, , \quad L = \boldsymbol{\phi}_L(Z, G)^\top \boldsymbol{w}_L + \epsilon_L \, , \tag{5.2}$$

where $\boldsymbol{\phi}_G : \mathcal{Z} \to \mathbb{R}^{d_G}$ and $\boldsymbol{\phi}_L : \mathcal{Z} \times \mathbb{R} \to \mathbb{R}^{d_L}$ denote *fixed* embedding functions for $Z$ and $Z, G$ respectively, $Z \in \mathcal{Z}$ indicates the membership in a protected group (where $\mathcal{Z}$ is the set of possible groups), and $\boldsymbol{w}_G \in \mathbb{R}^{d_G}$, $\boldsymbol{w}_L \in \mathbb{R}^{d_L}$ are the weights of the model.

In order to simplify notation, for observed data $\{(z_i, g_i, l_i)\}_{i=1}^n$, we define

$$\boldsymbol{x}_i = \begin{pmatrix} g_i \\ l_i \end{pmatrix} \in \mathbb{R}^2 \, , \quad \boldsymbol{w} = \begin{pmatrix} \boldsymbol{w}_G \\ \boldsymbol{w}_L \end{pmatrix} \in \mathbb{R}^{d_G + d_L} \, , \quad \Phi_i = \begin{pmatrix} \boldsymbol{\phi}_{G_i}^\top & \boldsymbol{0}^\top \\ \boldsymbol{0}^\top & \boldsymbol{\phi}_{L_i}^\top \end{pmatrix} \in \mathbb{R}^{2 \times (d_G + d_L)} \, , \tag{5.3}$$

where we write $\boldsymbol{\phi}_{G_i} = \boldsymbol{\phi}_G(z_i)$ and $\boldsymbol{\phi}_{L_i} = \boldsymbol{\phi}_L(z_i, g_i)$ for brevity. In eq. (5.3) as well as the remainder of this chapter, equations and assignments with subscripts $i$ on both sides hold for all $i \in \{1, \ldots, n\}$.[1]

**Model A: fit a counterfactually fair predictor.** First, we build a counterfactually fair predictor with our guessed unconfounded Model A via the following steps.

**1.** Fit Model A via regularized maximum likelihood:

$$\min_{w, \sigma_G, \sigma_L} \sum_{i=1}^{n} (x_i - \Phi_i w)^\top \Sigma^{-1} (x_i - \Phi_i w) + \lambda \|w\|_2^2 + n \log \det(\Sigma), \tag{5.4}$$

where

$$\Sigma = \begin{pmatrix} \sigma_G^2 & 0 \\ 0 & \sigma_L^2 \end{pmatrix}.$$

Note that we can alternately solve for $w$ and $\sigma_G, \sigma_L$ as follows. First fix $\sigma_G = \sigma_L = 1$ and compute

$$\tilde{w}^\dagger = \left( \sum_{i=1}^{n} \Phi_i^\top \Phi_i + \lambda \, I \right)^{-1} \left( \sum_{i=1}^{n} \Phi_i^\top x_i \right).$$

The optimal standard deviations $\sigma_G, \sigma_L$ are then simply given by the empirical standard deviations of the residuals under $\tilde{w}^\dagger$. Thus, the optimum of eq. (5.4) is

$$w^\dagger = \left( \sum_{i=1}^{n} \Phi_i^\top \Sigma^{-1} \Phi_i + \lambda \, I \right)^{-1} \left( \sum_{i=1}^{n} \Phi_i^\top \Sigma^{-1} x_i \right),$$

where $\Sigma = \mathrm{diag}(\sigma_G^2, \sigma_L^2)$.

**2.** Given fitted weights $w^\dagger$, estimate the noises $\epsilon_G, \epsilon_L$,

$$\hat{\epsilon}_i \equiv (\hat{\epsilon}_{g_i}, \hat{\epsilon}_{l_i})^\top \equiv x_i - \Phi_i w^\dagger.$$

**3.** Fit a counterfactually fair predictor $\hat{y}_i \equiv f_\theta(\hat{\epsilon}_i)$ with parameters $\theta$ to predict outcomes $y_i$ via

$$\theta^\dagger = \arg\min_\theta \sum_{i=1}^{n} \mathcal{L}(f_\theta(\hat{\epsilon}_i), y_i),$$

for some loss function $\mathcal{L}$. Note that because we are only using non-descendants of $Z$, namely the noise variables, as input, the predictor is counterfactually fair by design. While virtually any predictive model can be used in the two-variable case, in the general case we require the counterfactually fair predictor to be differentiable, such that it is amenable to gradient-based

---

[1] Note that $Z$ need not be exogenous. Since we would need to include additional—standard but occluding—steps in the algorithm to handle discrete variables, this assumption is solely to simplify the presentation.

optimization. The definition of counterfactual fairness constrains the optimization for any loss function. Here, we use the sufficient condition for counterfactual fairness that the predictor $\hat{Y}$ depends only on the noise terms, which are non-descendants of $Z$ (Kusner et al., 2017).

**Model B: evaluate counterfactual unfairness.** Next, we evaluate how the predictor $f_{\theta^\dagger}$ obtained in the previous section breaks down in the presence of unobserved confounding, i.e., in Model B. To do so, we fit Model B and generate "true" counterfactuals $x'$. If we were handed these counterfactuals and we wanted to make predictions using $f_{\theta^\dagger}$ we would compute their noise terms $\hat{\epsilon}'$ using step 2 above. If Model A was in fact the model that generated the counterfactuals $x'$ then the predictions on the noise terms for the real data and the counterfactuals would be *identical*: $f_{\theta^\dagger}(\hat{\epsilon}) = f_{\theta^\dagger}(\hat{\epsilon}')$.

However, because the counterfactuals were generated by the "true" weights $w^*$ of Model B, not the weights $w^\dagger$ of Model A, there will be a difference between the real data and counterfactual predictions $f_{\theta^\dagger}(\hat{\epsilon}) \neq f_{\theta^\dagger}(\hat{\epsilon}')$. It is this discrepancy we will quantify with our measure of counterfactual unfairness (CFU). Here is how we compute it for a given confounding budget $p_{\max}$.

**1.** Fit Model B via regularized maximum likelihood:

$$\min_{w, \sigma_G, \sigma_L} \sum_{i=1}^{n} (x_i - \Phi_i w)^\top \Sigma^{-1} (x_i - \Phi_i w) + \lambda^\dagger \|w\|_2^2 + n \log \det(\Sigma),$$

where

$$\Sigma \equiv \begin{pmatrix} \sigma_G & 0 \\ 0 & \sigma_L \end{pmatrix} \underbrace{\begin{pmatrix} 1 & p_{\max} \\ p_{\max} & 1 \end{pmatrix}}_{=:P} \begin{pmatrix} \sigma_G & 0 \\ 0 & \sigma_L \end{pmatrix}.$$

As before, we can alternately solve for $w$ (closed-form) and $\sigma_G, \sigma_L$ (via coordinate descent).[2] Let $w^*$ be the final weights after optimization.

**2.** Given weights $w^*$, estimate the noises of Model B,

$$\hat{\delta}_i = (\hat{\delta}_{g_i}, \hat{\delta}_{l_i})^\top = x_i - \Phi_i w^*.$$

**3.** For a fixed counterfactual value $z' \in \mathcal{Z}$, compute the Model B counterfactuals of $G$ and $L$ for all $i \in \{1, \ldots, n\}$,

$$g_i' = \boldsymbol{\phi}_G(z_i')^\top w_G^* + \hat{\delta}_{g_i}, \quad l_i' = \boldsymbol{\phi}_L(z_i', g_i')^\top w_L^* + \hat{\delta}_{l_i},$$

---

[2]In fact we optimize $\log(\sigma_G), \log(\sigma_L)$ to ensure the standard deviations are positive.

where $w^* = (w_G^*, w_L^*)^\top$. If $x_i' \equiv (g_i', l_i')^\top$, we can write the above equation as

$$x_i' = \Phi_i' w^* + \hat{\delta}_i ,$$

and $\Phi_i' \equiv \mathrm{diag}(\phi_G(z_i'), \phi_L(z_i', g_i'))$ is defined in general by sequential propagation of counterfactual values according to the ancestral ordering of the SEM.

**4.** Compute the "incorrect" noise terms of the counterfactuals using the same procedure as in step 2 of Section 5.3, using weights $w^\dagger$ of Model A:

$$\hat{\epsilon}_i' = (\hat{\epsilon}_{g_i'}, \hat{\epsilon}_{l_i'})^\top = x_i' - \Phi_i' w^\dagger .$$

Again, the predictions on the above quantity $f_{\theta^\dagger}(\hat{\epsilon}_i')$ will differ from those made on the real-data noise terms $f_{\theta^\dagger}(\hat{\epsilon}_i)$ (unless the counterfactuals were also generated according to Model A).

**5.** To measure the discrepancy, we propose to quantify counterfactual unfairness as the squared difference between the above two quantities:

$$\mathrm{CFU}_i = (f_{\theta^\dagger}(\hat{\epsilon}_i) - f_{\theta^\dagger}(\hat{\epsilon}_i'))^2 .$$

Ultimately, to summarize the aggregate unfairness, we will compute the average counterfactual unfairness:

$$\mathrm{CFU} = \frac{1}{n} \sum_{i=1}^{n} \mathrm{CFU}_i . \tag{5.5}$$

A quick note: in the two-variable setting, given a confounding budget $p_{\max}$, the worst-case CFU occurs precisely at $p_{\max}$, which need not be the case for multivariate confounding as we show at the end of the current section.

Thus, the above procedure computes the maximum CFU with bivariate confounding budget equal to $p_{\max}$. CFU measures how the counterfactual responses $\hat{Y}(z)$ and $\hat{Y}(z')$, defined using Model A, would differ "in reality", i.e., if Model B were "true". What qualifies as bad CFU is problem dependent and requires interaction with domain experts, who can make judgment calls about the plausibility of the misspecification $p_{\max}$ that is required to reach a breaking point. Here, a breaking point could be the CFU of a predictor that completely ignores the causal graph.

To summarize: we learn $\hat{Y} \equiv f_{\theta^\dagger}$ as a function of $X$ and $Z$, where $X$ and $Z$ are implicit in the expression of the (estimated) noise terms $\epsilon$ that are computed using the assumptions of the working Model A. We assess how "unfair" $\hat{Y}$ is by comparing for each data point the two counterfactual values $\hat{Y}(z) \equiv f_{\theta^\dagger}(\hat{\epsilon}_i)$ and $\hat{Y}(z') \equiv f_{\theta^\dagger}(\hat{\epsilon}_i')$ where the "true" counterfactual is

generated according to the world assumed by Model B. The space of models to which Model B belongs is a continuum indexed by $p_{\max}$, which will allow us to visualize the sensitivity of Model A by a one-dimensional curve. We will do this by finding the best fitting model (in terms of structural equation coefficients and noise variances) at different values of $p_{\max}$, so that the corresponding CFU measure is determined by $p_{\max}$ only (results on the above law school model are shown in Section 5.5). We assume that the free confounding parameter is not identifiable from data (as it would be the case if the model was linear and the edge $Z \to L$ was missing, the standard instrumental variable scenario).

As we have mentioned previously, in the multivariate setting, the worst-case counterfactual unfairness with a confounding budget of $p_{\max}$ is not necessarily obtained when all non-zero entries of the correlation matrix are set to $p_{\max}$. Even though intuitively that would lead to the largest allowed correlation between any pair of variables and thus in some sense to the largest deviation from the assumed Model A, we will now show that this setting does not necessarily lead to the biggest deviation in predictions. In particular, we will show that such a matrix with all non-zero correlation entries set to $p_{\max}$ is not a valid correlation matrix. To this end, it suffices to find a symmetric matrix $P$ with 1s on the diagonal that is not positive semidefinite when all its non-zero off-diagonal entries are set to the same value, which we define to be the considered confounding budget $p_{\max}$. Since each valid correlation matrix must be positive semidefinite, the correlation matrix for the worst-case counterfactual unfairness must be different from $P$ (while maintaining the zero entries). Because all off-diagonal entries are upper bounded by $p_{\max}$, at least one of them must be smaller than the corresponding value in $P$.

For example, consider

$$P = \begin{pmatrix} 1 & p_{\max} & p_{\max} \\ p_{\max} & 1 & 0 \\ p_{\max} & 0 & 1 \end{pmatrix}.$$

Since the eigenvalues of $P$ are $1$, $1 - \sqrt{2}\, p_{\max}$, and $1 + \sqrt{2}\, p_{\max}$, we see that $P$ is not positive semidefinite for $p_{\max} > 1/\sqrt{2}$.

In general, for $m \in \mathbb{N}$ with $m > 2$, the matrix $P \in \mathbb{R}^{m \times m}$ with $P_{ii} = 1$ for $i \in \{1, \ldots, m\}$, $P_{1i} = P_{i1} = p_{\max}$ for $i \in \{2, \ldots, m\}$ and $P_{ij} = 0$ for all remaining entries, has the eigenvalues (without multiplicity) $1$, $1 - \sqrt{m-1}\, p$, and $1 + \sqrt{m-1}\, p$. Therefore, $P$ is not positive semidefinite for $p_{\max} > 1/\sqrt{m-1}$. We conclude that as the dimensionality of the problem increases, we may encounter such situations for ever smaller confounding budget. This leads us to the second tool for sensitivity analysis in the multi-variate setting.

## 5.4 Tool #2: optimization-based analysis

In this section, we generalize the procedure outlined for the two-variable case in Section 5.3 to the general multivariate case.

**Notation and problem setup.** Besides the protected attribute $Z$ and the target variable $Y$, let there be $m$ additional observed feature variables $X_j$ in the causal graph $\mathcal{G}$ each of which comes with an unobserved noise variable $\epsilon_j$.

As before, we express the assignment of the structural equations for a specific realization of observed features $x = (x_1, \dots, x_m)^\top$ and noise terms $\epsilon = (\epsilon_1, \dots, \epsilon_m)^\top$ as the following operation, i.e., $x = \Phi w + \epsilon$. Here $\Phi$ has $m$ rows and $d = \sum_{V \in \text{has-parents}(\mathcal{G})} d_V$ columns, where $d_V$ is the dimensionality of embedding $\phi_V : \mathbb{R}^{|\text{pa}_\mathcal{G}(V)|} \to \mathbb{R}^{d_V}$ for each node $V \in \mathcal{G}$ that has parent nodes. Without loss of generality, we assume the nodes $\{Z\} \cup \{X_j\}_{j=1}^m$ to be topologically sorted with respect to $\mathcal{G}$ with $Z$ always being first. We combine the individual weights as $w = (w_{X_1}, w_{X_2}, \dots, w_{X_m}) \in \mathbb{R}^d$ and represent $\Phi$ once evaluated on a specific sample $z, x$ of the variables $Z, X_1, \dots, X_m$ as

$$
\Phi = \begin{pmatrix} \boldsymbol{\phi}_{X_1}^\top & & \mathbf{0}^\top \\ & \ddots & \\ \mathbf{0}^\top & & \boldsymbol{\phi}_{X_m}^\top \end{pmatrix},
$$

where $\boldsymbol{\phi}_{X_j}$ is based on the parents of $X_j$. The covariance matrix of the noise terms is given by

$$
\Sigma \equiv \text{diag}(\sigma_1, \dots, \sigma_m) \, P \, \text{diag}(\sigma_1, \dots, \sigma_m),
$$

where $\sigma_1, \dots, \sigma_m$ are the standard deviations of each variable and $P$ is a correlation matrix.

**The optimization problem.** In the general case our goal is to find a correlation matrix $P$ that satisfies a "confounding budget" $p_{\max}$. In particular we would like to constrain the correlation $P_{jk}$ between any two different variables $X_j$ and $X_k$ for $j \neq k$, while allowing $P_{jj} = 1$ for all $j$. Additionally, we want to take into account any prior knowledge that certain variable pairs should have no correlation, if available. The most intuitive way to budget the amount of confounding is to limit the absolute size of any correlation by $p_{\max}$ as: $|P_{jk}| \leq p_{\max}$ for all $j \neq k$. This captures a notion of "restricted unobserved confounding" and leads to the

following optimization problem

$$
\begin{aligned}
\underset{P}{\text{maximize}} \quad & \sum_{i=1}^{n} \text{CFU}_i \\
\text{subject to} \quad & P_{jj} = 1 \quad \text{for } j \in \{1, \ldots, m\}\,, \\
& |P_{jk}| \le p_{\max} < 1 \quad \text{for all } (j,k) \in \mathcal{C}\,, \\
& P_{jk} = 0 \quad \text{for all } (j,k) \notin \mathcal{C} \text{ with } j \ne k\,.
\end{aligned}
\tag{5.6}
$$

$\mathcal{C}$ is the set of correlations that are allowed to be non-zero. We remark that the setting where there are zero correlations can be captured using the standard acyclic directed mixed graph notation (ADMG, Richardson, 2003). Specifically, this can be represented with ADMGs by removing bidirected edges between any two noise terms whose correlation is fixed to zero.

As in the bivariate case, CFU is a direct function of $P$ only: all other parameters will be determined given the choice of correlation matrix by maximizing likelihood. Note that eq. (5.6) contains multiple nested optimization problems: The outer optimization over correlation matrices $P$ and the inner optimization for the counterfactually fair model weights $\theta^\dagger$ as well as the weights and standard deviations of Models A and B. To solve it efficiently, we will parameterize $P$ in a way that facilitates optimization via off-the-shelf, unconstrained, automatic differentiation tools.

**Algorithm.**   We use the following approach to accommodate the constraints in eq. (5.6) in a way such that our algorithm does not require a constrained optimization subroutine. Assume first that $P$ has no correlations that should be zero. We compute $LL^\top$ for a matrix $L \in \mathbb{R}^{m \times m}$, whose entries are the parameters we eventually optimize. To constrain the off-diagonals to a given range and ensure that $P$ has 1s on the diagonal, we define $P$ as,

$$
P := \tanh_{p_{\max}}(LL^\top) := I + p_{\max}\,(J - I) \odot \tanh(LL^\top)\,,
$$

where $\odot$ denotes element-wise multiplication of matrices, $J$ is a matrix of all ones, and $I$ is the identity matrix. This way $P$ is symmetric, differentiable with respect to the entries of $L$, has 1s on the diagonal, and its off-diagonal values are squashed to lie within $(-p_{\max}, p_{\max})$. While it is natural to directly mask and clamp the diagonal, there are various ways to squash the off-diagonals to a fixed range in a smooth way, which bears close resemblance to barrier methods in optimization. We choose $\tanh()$ because of its abundance in the machine learning literature and availability in computational frameworks (including gradient implementations), but other forms of $P$ may work better for specific applications. Note that this formulation does not guarantee $P$ to be positive semidefinite.

---

**Algorithm 1** MAXCFU: Maximize counterfactual unfairness under a certain confounding budget constraint.

**Input:** data $\{x_i, y_i, z_i,\}_{i=1}^n$, confounding budget $p_{\max}$, learning rate $\alpha$, minibatch size $B$

1: $\{\hat{\epsilon}_i\}_{i=1}^n, w^\dagger, \theta^\dagger \leftarrow \text{FITMODELA}(\{x_i, y_i, z_i,\}_{i=1}^n)$
2: $\mathcal{D} \leftarrow \{x_i, y_i, z_i, \Phi_i, \hat{\epsilon}_i\}_{i=1}^n$      ▷ full dataset
3: $L \leftarrow \text{INITIALIZEPARAMETERS}()$
4: **for** $t = 1 \dots T$ **do**      ▷ iterations
5:     $\mathcal{D}^{(t)} \leftarrow \text{SAMPLEMINIBATCH}(\mathcal{D}, B)$
6:     $\Delta \leftarrow \nabla_L \text{CFU}(\mathcal{D}^{(t)}, w^\dagger, \frac{B}{n}\lambda^\dagger, \theta^\dagger, L)$      ▷ autodiff
7:     $L \leftarrow L + \alpha\Delta$      ▷ gradient ascent step
8: **return** $\text{CFU}(\mathcal{D}, w^\dagger, \lambda^\dagger, \theta^\dagger, L)$

9: **function** FITMODELA($\{x_i, y_i, z_i,\}_{i=1}^n$)
10:     $w^+ \leftarrow \left(\sum_{i=1}^n \Phi_i^\top \Phi_i + \lambda^\dagger I\right)^{-1} \left(\sum_{i=1}^n \Phi_i^\top x_i\right)$
11:     $\Sigma \leftarrow \text{diag}(\text{var}(\{x_i - \Phi_i w^+\}_{i=1}^n))$
12:     $w^\dagger \leftarrow \left(\sum_{i=1}^n \Phi_i^\top \Sigma^{-1} \Phi_i + \lambda^\dagger I\right)^{-1} \left(\sum_{i=1}^n \Phi_i^\top \Sigma^{-1} x_i\right)$
13:     $\hat{\epsilon}_i \leftarrow x_i - \Phi_i w^\dagger$
14:     $\theta^\dagger \leftarrow \arg\min_\theta \sum_{i=1}^n \mathcal{L}(f_\theta(\hat{\epsilon}_i), y_i)$
15:     **return** $\{\hat{\epsilon}_i\}_{i=1}^n, w^\dagger, \theta^\dagger$

16: **function** CFU($\mathcal{D}, w^\dagger, \lambda^\dagger, \theta^\dagger, L$)
17:     $w^*, \sigma^* \leftarrow \min_{w,\sigma} \sum_{i=1}^n (x_i - \Phi_i w)^\top \Sigma^{-1}(x_i - \Phi_i w) + \lambda^\dagger \|w\|_2^2 + n \log\det(\Sigma)$
     where $\Sigma = \text{diag}(\sigma) \tanh_{p_{\max}}(LL^\top) \text{diag}(\sigma)$
18:     $\hat{\delta}_i \leftarrow x_i - \Phi_i w^*$
19:     $z_i' \leftarrow 1 - z_i$ and $x_i' \leftarrow \Phi_i' w^* + \hat{\delta}_i$
     where $\Phi_i'$ is computed via iterative assignment
20:     $\hat{\epsilon}_i' \leftarrow x_i' - \Phi_i' w^\dagger$
21:     $\text{CFU} \leftarrow \frac{1}{n}\sum_{i=1}^n (f_{\theta^\dagger}(\hat{\epsilon}_i) - f_{\theta^\dagger}(\hat{\epsilon}_i'))^2$
22:     **return** CFU

---

In Algorithm 1, we describe our procedure to maximize counterfactual unfairness given a confounding budget $p_{\max}$ and observational data $\{x_i, y_i, z_i,\}_{i=1}^n \subset \mathbb{R}^m \times \{0,1\}^2$. The algorithm closely follows the procedure described in Section 5.3 for the bivariate case. Since we use automatic differentiation provided by PyTorch (Paszke et al., 2017) to obtain gradients,

we only show the forward pass in Algorithm 1. For the initialization INITIALIZEPARAME-
TERS(), we simply populate $L$ as a lower triangular matrix with small random values for the
off-diagonals and 1s on the diagonal.

The main place where code optimization can take place is step 17 of Algorithm 1. Alterna-
tives to the (local) penalized maximum likelihood taking place there could be suggested
(perhaps using spectral methods). It is hard though to say much in general about Step 20,
as counterfactual fairness allows for a large variety of loss functions usable in supervised
learning. In the case of linear predictors, it is still a non-convex problem due to the complex
structure of the correlation matrix, and for now we leave as an open problem whether
non-gradient based optimization may find better local minima.

If $\mathcal{C}$ indicates some correlations should be zero, we suggest the following standard "clique
parameterization": $L$ is a $m \times c$ matrix where $c$ is the number of cliques in $\mathcal{C}$, with $L_{ik}$
being a non-zero parameter if and only if vertex $i$ is in clique $k$. $L_{ik} \equiv 0$ otherwise. It
follows that such a matrix will have zeros at precisely the locations not in $\mathcal{C}$.[3] See Silva et al.
(2007); Barber (2009) for examples of applications of this idea. For large cliques, further
refinements are necessary to avoid unnecessary constraints, such as creating more than
one row per clique of size four or larger. Our experiments will not make use of sparse
$P$ (note that this parameterization also assumes that the number of cliques is tractable).
Note that individual parameters $L_{ik}$ may not be identifiable. However, identifiability is not
necessary here, all we care about is the objective function: CFU. As a matter of fact, multiple
globally optimal solutions are to be expected even in the space of $P$ transformations. A more
direct parameterization of sparse $P$, with exactly one parameter per non-zero entry of the
upper covariance matrix, is discussed by Drton & Richardson (2004). Computationally, this
minimal parameterization does not easily lead to unconstrained gradient-based algorithms
for optimizing sparse correlation matrices with bounded entries. We suggest the clique
parameterization as a pragmatic alternative. Special cases may be treated with more efficient
specialized approaches. Cinelli et al. (2019) provide a thorough discussion of fully linear
models.

In Section 5.5 we will demonstrate this approach on a real-world 3-variable-confounding sce-
nario to showcase our approach. Before moving to empirical results and the implementation
of our tools, we briefly describe a methodological extension of our approach to path-specific
effects (Shpitser, 2013). Specifically, we describe an example that illustrates how notions of
path-specific effects can be easily pipelined with our sensitivity analysis framework.

Consider Figure 5.2, where the path from $Z \rightarrow U$ is considered unfair and $Z \rightarrow F$ is
considered fair, in the sense that we do not want a non-zero path-specific effect of $Z$ on $\hat{Y}$
that is comprised by a possible path $Z \rightarrow U \rightarrow \hat{Y}$ in the causal graph implied by the chosen

---

[3]Barring unstable parameter cancellations that have measure zero under continuous measures on $\{L_{ik}\}$.

Figure 5.2 A path-specific model where the path from protected attribute $Z$ to feature $U$ is unfair and the path from $Z$ to feature $F$ is fair.

construction of $\hat{Y}$. Then a path-specific counterfactually fair predictor is one that uses $\{\epsilon_U, F\}$ as input. Note that the only difference this makes in our grid-based tool is that we only estimate the noise $\epsilon_U$ for the unfair path in Model A (step 2) and fit a predictor on $\{\epsilon_U, F\}$ (step 3). Additionally, we only compute the incorrect noise terms of the counterfactuals in Model B, using the weights of Model A (step 4). For the optimization-based tool we would change lines 13, 14, and 20 in the same way.

## 5.5 Experiments

We compare the grid-based and the optimization-based tools introduced in Sections 5.3 and 5.4 on two real datasets.

In all experiments our embedding $\phi$ is a polynomial basis up to a fixed degree. The degree is determined via cross validation (5-fold) jointly with the regularization parameter $\lambda^\dagger$. Our counterfactually fair predictor is regularized linear regression on the noise terms $\epsilon$:

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^n (y_i - \boldsymbol{\phi}(\hat{\boldsymbol{\epsilon}}_i)^\top \boldsymbol{\theta})^2 + \lambda \|\boldsymbol{\theta}\|_2^2 \,.$$

For this model, counterfactual unfairness is:

$$\text{CFU}_i = \left( \left( \boldsymbol{\phi}(\hat{\boldsymbol{\epsilon}}_i) - \boldsymbol{\phi}(\hat{\boldsymbol{\epsilon}}_i') \right)^\top \boldsymbol{\theta}^\dagger \right)^2 \,.$$

For comparison, we train two baselines that also use regularized ridge regression (degree and regularization are again selected by 5-fold cross-validation):

**unconstrained:** an unconstrained predictor using all observed variables as input $f_{\text{uc}}$ : $(Z, X_1, \ldots, X_m) \mapsto Y$.

**blind unconstrained:** an unconstrained predictor using all features, but not the protected attribute, as input $f_{\text{buc}} : (X_1, \ldots, X_m) \mapsto Y$.

Analogous to our definition of CFU in eq. (5.5), we compute the unfairness of these baselines as the mean squared difference between their predictions on the observed data and the predictions of the counterfactually fair predictor on the observed data: $\frac{1}{n} \sum_{i=1}^n \left( f_{\boldsymbol{\theta}^\dagger}(\hat{\boldsymbol{\epsilon}}_i) - \right.$

Figure 5.3 Counterfactual unfairness for the law school dataset. See text for details.



**features**

satisfaction with

$O$: organization
$M$: manager
$J$: job

Figure 5.4 The assumed to be true causal graph for the NHS Staff Survey dataset.

$\hat{y}_i^{(b)uc})^2$, where $\hat{y}_i^{uc} = f_{uc}(z_i, x_i)$ and $\hat{y}_i^{buc} = f_{buc}(x_i)$. This choice is motivated by the fact that in practice we care about how much potential predictions deviate from predictions satisfying a fairness measure. For our grid-based approach we repeatedly fix $p_{max} \in (-1, 1)$ to a particular value and then use the procedure in Section 5.3 to compute CFU. For the optimization approach we similarly fix $p_{max} \in [0, 1)$ in the constraint of eq. (5.6). For efficiency we use the previously found correlation matrix $P$ as initialization for the next setting of $p_{max}$.

**Law School data.** As our first experiment we consider the motivating example introduced in Section 5.3 on law school success (recall eq. (5.2) and Figure 5.1 for details on the causal models). Our data comes from law school students in the US (Wightman, 1998). As our causal model investigates confounding between two variables, we will use the grid-based approach introduced in Section 5.3 to calculate the maximum CFU. Recall that for the bivariate approach we fix a confounding level $p = p_{max}$ and then compare predictions between real data based on Model A versus counterfactuals generated from Model B. Figure 5.3 shows the CFU for the grid-based approach (black), alongside the baselines (green/red), as the

Figure 5.5 Counterfactual unfairness as a function of $p_{\max}$ for the multivariate NHS dataset.

correlation $p$ varies. We first note that the confounding is not symmetric around $p = 0$. For the law school data, negative correlations have smaller CFU. In general, this is a data-specific property.

Additionally, we notice that as $p_{\max}$ moves away from 0 (where we expect CFU $= 0$ up to numerical errors) it increases noticeably, then plateaus in roughly $[0.1, 0.9]$ and finally increases again. For large $p_{\max} \geq 0.9$ we cannot exclude numeric instability as the covariance matrix becomes nearly negative definite. Finally, we note that both baseline approaches have higher CFU than found with any grid-based setting.

**NHS Staff Survey.** Our second experiment is based on the 2014 UK National Health Service (NHS) Survey (Picker Institute Europe, 2015). The goal of the survey was to "*gather information that will help to improve the working lives of staff in the NHS*". Answers to the survey questions ranged from 'strongly disagree' (1) to 'strongly agree' (5). We averaged survey answers for related questions to create a dataset of continuous indices for: *job satisfaction* ($J$), *manager satisfaction* ($M$), *organization satisfaction* ($O$), and *overall health* ($Y$). The goal is to predict health $Y$ based on the remaining information. Additionally, we collected the race ($Z$) of the survey respondents. Using this data, we formulate a ground-truth causal graph shown in Figure 5.4 (equivalent to Model B in Figure 5.1). This causal graph includes correlations between all noise terms $\epsilon_J, \epsilon_M, \epsilon_O$ and we assume the following structural equations

$$O = \boldsymbol{\phi}_O(Z)^\top \boldsymbol{w}_O + \epsilon_O$$
$$M = \boldsymbol{\phi}_M(Z, O)^\top \boldsymbol{w}_M + \epsilon_M$$
$$J = \boldsymbol{\phi}_J(Z, O, M)^\top \boldsymbol{w}_J + \epsilon_J.$$

Just as in the law school example, we measure the impact of unobserved confounding by comparing this model to the unconfounded model (i.e., all noise terms are jointly independent). As there is no general efficient way to grid-search for positive definite matrices that maximize CFU for a given $p_{max}$, we make use of our optimization-based procedure for calculating maximum CFU, as described in Algorithm 1. Figure 5.5 shows the results of our method on the NHS dataset. Note that we only show positive $p_{max}$ because our optimization problem eq. (5.6) only constrains the absolute value of the off-diagonal correlations. This allows the procedure to learn whether positive or negative correlations result in greater CFU. As in the law school dataset we see an initial increase in CFU for small $p_{max}$, followed by a plateau, ending with another small increase. Just like before, all settings result in lower CFU than the two baseline techniques.

## 5.6   Conclusion

We presented two techniques to assess the impact of unmeasured confounding in causal additive noise models. We formulated unmeasured confounding as covariance between noise terms. We then introduced a grid-based approach for confounding between two variables, and an optimization-based approach for confounding in the general case. We demonstrated our approach on two real-world fairness datasets. Our techniques can also be applied for sensitivity analysis of other causal queries, which is an interesting direction for future research. Currently, our approach is limited to a specific type of unobserved confounding in ANMs assuming a certain functional form of the structural equations. However, we believe the developed tools are an important step towards making causal models suitable to address discrimination in real-world prediction problems.

<div style="text-align: right">

# 6

</div>

# Fairness and privacy

In this chapter, we go beyond the assumption that sensitive attributes are available in the clear—i.e., as clear text without encryption—to train or evaluate fair models. To avoid disparate treatment, sensitive attributes should not be considered. On the other hand, in order to avoid disparate impact, sensitive attributes must be examined—e.g., in order to learn a fair model, or to check if a given model is fair. We introduce methods from secure multi-party computation which allow us to avoid both. By encrypting sensitive attributes, we show how an outcome-based fair model may be learned, checked, or have its outputs verified and held to account, *without users revealing their sensitive attributes and without modelers having to disclose their models*.

The main content of this chapter has been published in the following paper:

> BLIND JUSTICE: FAIRNESS WITH ENCRYPTED SENSITIVE ATTRIBUTES
> Niki Kilbertus, Adrià Gascón, Matt Kusner, Michael Veale, Krishna P. Gummadi, Adrian Weller.
> [https://github.com/nikikilbertus/blind-justice]
> *International Conference on Machine Learning (ICML), 2018*

## 6.1 Introduction

To motivate the main challenge in this chapter, we recall two notions of discrimination briefly mentioned in Chapter 3. The first type, *disparate treatment* (or *direct discrimination*), occurs if individuals are treated differently according to their sensitive attributes (with all others equal). Intuitively, to avoid disparate treatment, one should not inquire about individuals' sensitive attributes, i.e., apply fairness through unawareness. *Disparate impact* (or *indirect discrimination*) occurs when the *outcomes* of decisions disproportionately benefit or harm individuals from subgroups with particular sensitive attribute settings without appropriate justification. For example, firms deploying car insurance telematics devices (Handel et al., 2014) build up high dimensional pictures of driving behavior which might easily proxy for

sensitive attributes even when they are omitted. As we have discussed in Chapter 3, most work on observational group matching criteria has thus focused on avoiding notions of disparate impact.

In order to check and enforce such requirements, the modeler is usually assumed to have access to the sensitive attributes for individuals in the training data—however, this may be undesirable for several reasons (Žliobaitė & Custers, 2016). First, individuals are unlikely to want to entrust sensitive attributes to modelers in all application domains. Where applications have clear discriminatory potential, it is understandable that individuals may be wary of providing sensitive attributes to modelers who might exploit them to negative effect, especially with no guarantee that a fair model will indeed be learned and deployed. Even if certain modelers themselves were trusted, the wide provision of sensitive data creates heightened privacy risks in the event of a data breach. Further, legal barriers may limit collection and processing of sensitive personal data. A timely example is the EU's General Data Protection Regulation (GDPR), which contains heightened prerequisites for the collection and processing of some sensitive attributes or the California Consumer Privacy Act. Unlike other data, modelers cannot justify using sensitive characteristics in fair learning with their "legitimate interests"—and instead will often need explicit, freely given consent (Veale & Edwards, 2018).

One way to address these concerns was recently proposed by Veale & Binns (2017). The idea is to involve a highly trusted third party, and may work well in some cases. However, there are significant potential difficulties: individuals must disclose their sensitive attributes to the third party (even if an individual trusts the party, she may have concerns that the data may somehow be obtained or hacked by others, e.g., Graham, 2017); and the modeler must disclose their model to the third party, which may be incompatible with their intellectual property or other business concerns.

To overcome these seemingly conflicting interests, we propose an approach to detect and mitigate disparate impact without disclosing readable access to sensitive attributes. This reflects the notion that decisions should be blind to an individual's status—depicted in courtrooms by a blindfolded Lady Justice holding balanced scales (Bennett Capers, 2012). We assume the existence of a regulator with fairness aims (such as a data protection authority or anti-discrimination agency). With recent methods from *secure multi-party computation* (MPC), we enable auditable fair learning while ensuring that both individuals' sensitive attributes and the modeler's model remain private from all other parties—including the regulator. Desirable fairness and accountability applications we enable include:

1. **Fairness certification.** Given a model and a dataset of individuals, check that the model satisfies a given observational group fairness constraint; if yes, generate a certificate.

2. **Fair model training.** Given a dataset of individuals, learn a model guaranteed and certified to be fair.

3. **Decision verification.** A malicious modeler might go through fair model training, but then use a different model in practice. To address such accountability concerns (Kroll et al., 2016), we efficiently provide for an individual to challenge a received outcome, verifying that it matches the outcome from the previously certified model.

We rely on recent theoretical developments in MPC, specifically the protocols for training linear and logistic models in secure MPC developed by Mohassel & Zhang (2017). In this work, we further extend these protocols to also admit linear constraints in order to enforce fairness requirements. These extensions may be of independent interest. Finally, we demonstrate the efficacy of our methods by evaluating them on synthetic and real-world datasets.

We note that the privacy or secrecy constraints considered in our approach are separate from other theorized, setup-dependent attacks, e.g., model extraction (Tramèr et al., 2016) or inversion (Fredrikson et al., 2015). If relevant, modelers may need to consider these separately. In particular, our approach aims at fixing the issue of providing, transmitting, and storing sensitive data "in the clear". The notion of privacy that is achieved by encrypting these data is distinct from, e.g., the guarantees provided by differential privacy. In principle, our protocols do not prevent a malicious modeler from statistically inferring information about the protected attribute given non-sensitive data. The challenge of training fair models in a differentially private fashion and potential tensions between fairness and differential privacy have been explored (Jagielski et al., 2019; Cummings et al., 2019; Bagdasaryan et al., 2019; Ding et al., 2020; Xu et al., 2020).

Besides work on differentially private fair model training, our original work has also been followed up by a line of research on fair machine learning without or only limited demographic information (Hashimoto et al., 2018; Lahoti et al., 2020; Coston et al., 2019; Chen et al., 2019; Rastegarpanah et al., 2020).

## 6.2 Fairness and privacy requirements

**Assumptions and incentives.** We assume three categories of participants: a *modeler* M, a *regulator* REG, and *users* $U_1, \ldots, U_n$. For each user $U_i$, we consider a vector of sensitive features $z_i \in \mathcal{Z} = \{0, 1\}^p$ (e.g., ethnicity or gender) which might be a source of discrimination, and a vector of non-sensitive features $x_i \in \mathcal{X} = \mathbb{R}^d$, again assumed to be discrete or real. Here, we deviate from the previous restriction of only dealing with a single binary sensitive attributes and instead allow for $p \in \mathbb{N}$ binary sensitive features. Additionally, for each user there is a non-sensitive label or target $y_i \in \mathcal{Y} = \{0, 1\}$ which the modeler M would like to

predict. Again, we assume binary for simplicity labels though our MPC approach could be extended to multi-label settings.

Modeler M wishes to train a parametric model $f_\theta : \mathcal{X} \to \mathcal{Y}$, which accurately maps features $x_i$ to labels $y_i$, in a supervised fashion. We assume M needs to keep the model private for intellectual property or other business reasons. The model $f_\theta$ does not use sensitive information $z_i$ as input to prevent disparate treatment (direct discrimination).

For each user $U_i$, the modeler M observes or is provided $x_i, y_i$. The sensitive information in $z_i$ is required to ensure $f_\theta$ meets a given observational group fairness condition $F$. While each user $U_i$ wants $f_\theta$ to meet $F$, they also wish to keep $z_i$ private from all other parties. The regulator REG aims to ensure that M deploys only models that meet fairness condition $F$. It has no incentive to collude with M. If collusion were a concern, more sophisticated cryptographic protocols would be required, which we briefly touch upon in Section 6.2. Further, the modeler M might be legally obliged to demonstrate to the regulator REG that their model meets fairness condition $F$ before it can be publicly deployed. As part of this, REG also has a positive duty to enable the training of fair models.

Later in this section, we define and address three fundamental problems in our setup: certification, training, and verification. For each problem, we present its functional goal and its privacy requirements. We refer to $D = \{(x_i, y_i)\}_{i=1}^n$ and $Z = (z_i)_{i=1}^n \in \mathbb{R}^{n \times p}$ as the non-sensitive and sensitive data, respectively. We will reference the overview Figure 6.1 throughout the description of various components of the protocols. A concise description of all three protocols is then given right before Section 6.4.

**Fairness considerations.** In this chapter we focus on the p%-rule from eq. (3.3) as the fairness criterion $F$, which we restate here as

$$\frac{P(\hat{Y} = 1 \mid Z = z)}{P(\hat{Y} = 1 \mid Z = z')} \geq \frac{p}{100} \quad \text{for } z, z' \in \mathcal{Z} . \tag{6.1}$$

A similar MPC approach could also be used for other observational group matching criteria mentioned in Table 3.3, which have been addressed with efficient standard (non-private) methods.

**Fairness certification.** Given a notion of fairness $F$, the modeler M would like to work with the regulator REG to obtain a certificate that model $f_\theta$ is fair. To do so, we propose that users send their non-sensitive data $D$ to REG; and send *encrypted* versions of their sensitive data $Z$ to both M and REG. Neither M nor REG can read the sensitive data. However, we can design a secure protocol between M and REG (described in Section 6.3) to certify if the model is fair. This setup is shown in Figure 6.1 (*Left*).

Figure 6.1 Our setup for *Fairness certification* (*Left*), *Fair model training* (*Center*), and *Decision verification* (*Right*).

While both REG and M learn the outcome of the certification, we require the following *privacy constraints*: (C1) *privacy of sensitive user data*: no one other than $U_i$ ever learns $z_i$ in the clear, (C2) *model secrecy*: only M learns $f_{\theta}$ in the clear, and (C3) *minimal disclosure of D to* REG: only REG learns $D$ in the clear.

**Fair model training.** How can a modeler M learn a fair model without access to users' sensitive data $Z$? We propose to solve this by having users send their non-sensitive data $D$ to M and to distribute encryptions of their sensitive data to M and REG as in certification. We shall describe a secure MPC protocol between M and REG to train a fair model $f_{\theta}$ privately. This setup is shown in Figure 6.1 (*Center*).

*Privacy constraints*: (C1) privacy of sensitive user data, (C2) model secrecy, and (C3) minimal disclosure of $D$ to M.

**Decision verification.** Assume that a malicious M has had model $f_{\theta}$ successfully certified by REG as above. It then swaps $f_{\theta}$ for another potentially unfair model $f_{\theta'}$ in the real world. When a user receives a decision $\hat{y}$, e.g., her mortgage is denied, she can then challenge that decision by asking REG for a verification $V$. The verification involves M and REG, and consists of verifying that $f_{\theta'}(x) = f_{\theta}(x)$, where $x$ is the user's non-sensitive data. This ensures that the user would have been subject to the same result with the certified model $f_{\theta}$, even if $f_{\theta'} \neq f_{\theta}$ and $f_{\theta'}$ is not fair. Hence, while there is no simple technical way to prevent a malicious M from deploying an unfair model, it will get caught if a user challenges a decision that would differ under $f_{\theta}$. This setup is shown in Figure 6.1 (*Right*).

*Privacy constraint*: While REG and the user learn the outcome of the verification, we require (C1) privacy of sensitive user data, and (C2) model secrecy.

**Design choices.** We include a regulator party in our setup for several reasons. Given fair learning is of most benefit to vulnerable individuals, we do not wish to deter adoption with high individual burdens. While MPC could be carried out without the involvement of a regulator, using all users as parties, this comes at a greater computational cost. With current methods, taking that approach would be unrealistic given the size of the user-base in many domains of concern, and would furthermore require all users to be online simultaneously. Introducing a regulator removes these barriers and leaves users' computational burden at a minimum level, with envisaged applications practical with only their web browsers.

In cases where users are uncomfortable sharing $D$ with either REG or M, it is trivial to extend all three tasks such that all of $x_i, y_i, z_i$ remain private throughout, with the computation cost increasing only by a factor of 2. This extension would sometimes be desirable as it restricts the view of M to the final model, prohibiting inferences about $Z$ when $D$ is known. However, this setup hinders exploratory data analysis by the modeler which might promote robust model-building, and, in the case of verification, validation by the regulator that user-provided data is correct.

## 6.3 Our solution

Our proposed solution to these three problems is to use secure Multi-Party Computation (MPC). Before we describe how it can be applied to fair learning, we first present the basic principles of MPC, as well as its limitations particularly in the context of machine learning applications.

**MPC for machine learning.** Multi-party computation protocols allow two parties $P_1$ and $P_2$ holding secret values $x_1$ and $x_2$ to evaluate an agreed-upon function $f$, via $y = f(x_1, x_2)$ in a way in which the parties (either both or one of them) learn *only y*. For example, if $f(x_1, x_2) = \mathbf{1}[x_1 < x_2]$, then the parties would learn which of their values is bigger, but nothing else. Here, the indicator $\mathbf{1}[\cdot]$ is 1 if its argument is true and 0 otherwise. This corresponds to the well-known *Yao's millionaires problem*: two millionaires want to conclude who is richer without disclosing their wealth to each other. The problem was introduced by Andrew Yao in 1982, and kicked off the area of multi-party computation in cryptography.

In our setting—instead of a simple comparison as in the millionaires problem—$f$ will be either (i) a procedure to check the fairness of a model and certify it, (ii) a machine learning training procedure with fairness constraints, or (iii) a model evaluation to verify a decision.

The two parties involved in our computation are the modeler M and the regulator REG. The inputs depend on the case (see Figure 6.1).

As generic solutions do not yet scale to real-world data analysis tasks, one typically has to tailor custom protocols to the desired functionality. This approach has been followed successfully for a variety of machine learning tasks such as logistic and linear regression (Nikolaenko et al., 2013b; Gascón et al., 2017; Mohassel & Zhang, 2017), neural network training (Mohassel & Zhang, 2017) and evaluation (Juvekar et al., 2018; Liu et al., 2017), matrix factorization (Nikolaenko et al., 2013a), and principal component analysis (Al-Rubaie et al., 2017). In the next section we review challenges beyond scalability issues that arise when implementing machine learning algorithms in MPC.

**Challenges in multi-party machine learning.** MPC protocols can be classified into two groups depending on whether the target function is represented as either a Boolean or arithmetic circuit. All protocols proceed by having the parties jointly evaluate the circuit, processing it gate by gate while keeping intermediate values hidden from both parties by means of a secret sharing scheme. While representing functions as circuits can be done without losing expressiveness, it means certain operations are impractical. In particular, algorithms that execute different branches depending on the input data will explode in size when implemented as circuits, and in some cases lose their run time guarantees (e.g., consider binary search).

Crucially, this applies to *floating-point arithmetic*. While this is work in progress, state-of-the-art MPC floating-point arithmetic implementations take more than 15 milliseconds to multiply two 64 bit numbers (Demmler et al., 2015a, Table 4), which is prohibitive for our applications. Hence, machine learning MPC protocols are limited to *fixed-point* arithmetic. Overcoming this limitation is a key challenge for the field. Another necessity for the feasibility of MPC is to approximate non-linear functions such as the sigmoid, ideally by (piecewise) linear functions.

**Input sharing.** To implement the functionality from Figure 6.1, we first need a secure procedure for the users to *secret share* a sensitive value, for example her race, with the modeler M and the regulator REG. We use *additive secret sharing*. A value $z$ is represented in a finite domain $\mathbb{Z}_q$—we use $q = 2^{64}$. To share $z$, the user samples a value $r$ from $\mathbb{Z}_q$ uniformly at random, and sends $z - r$ to M and $r$ to REG. While $z$ can be reconstructed (and subsequently operated on) inside the MPC computation by means of a simple addition, each share on its own does not reveal anything about $z$ (other than that it is in $\mathbb{Z}_q$). One can think of arithmetic sharing as a "distributed one-time pad".

In Figure 6.1, we now reinterpret the key held by REG and the encrypted $z$ by M, as their corresponding shares of the sensitive attributes and denote them by $\langle z \rangle_1$ and $\langle z \rangle_2$ respectively.

The idea of privately outsourcing computation to two non-colluding parties in this way is recurrent in MPC, and often referred to as the two-server model (Mohassel & Zhang, 2017; Gascón et al., 2017; Nikolaenko et al., 2013b; Al-Rubaie et al., 2017).

**Signing and checking a model.** We will see that *certification* and *verification* partly correspond to sub-procedures of the *fair training* task: during training we check the fairness constraint *F*, and repeatedly evaluate partial models on the training dataset (using gradient descent). Hence, *certification* and *verification* do not add technical difficulties over training, which is described in detail in Section 6.4. However, for verification, we still need to "sign" the model, i.e., REG should obtain a signature $s(\theta)$ as a result of model certification, see Figure 6.1 (*Left*). This signature is used to check in the verification phase, whether a given model $\theta'$ from M satisfies $s(\theta') = s(\theta)$ for a certified fair model $\theta$ (in which case $\theta = \theta'$ with high probability). Moreover, we need to preserve the secrecy of the model, i.e., REG should not be able to recover $\theta$ from $s(\theta)$. These properties, given that the space of models is large, calls for a cryptographic hash function, such as SHA-256.

Additionally, in our functionality, the hash of $\theta$ should be computed inside MPC, to hide $\theta$ from REG. Fortunately, cryptographic hashes such as SHA-256 are a common benchmark functionality in MPC, and their execution is highly optimized. More concretely, the overhead of computing $s(\theta)$, which needs to be done both for certification and verification is of the order of fractions of a second (Keller et al., 2013, Figure 14). While cryptographic hash functions have various applications in MPC, we believe the application to machine learning model certification is novel.

Ultimately, certification is implemented in MPC as a check that $\theta$ satisfies the criterion *F*, followed by the computation of $s(\theta)$. The regulator REG keeps the signature $s(\theta)$ of the fair model parameters for later verification. On the other hand, for verification, the MPC protocol first computes the signature of the model provided by M, and then proceeds with a prediction as long as the computed signature matches the one obtained by REG in the verification phase. An alternative solution is possible based on symmetric encryption under a shared key, as highly efficient MPC implementations of block ciphers such as AES are available (Keller et al., 2017).

**Fair training.** To realize the *fair training* functionality, we closely follow the techniques recently introduced by Mohassel & Zhang (2017). Specifically, we extend their custom MPC protocol for logistic regression to additionally handle linear constraints. This extension may be of independent interest and has applications for privacy-preserving machine learning beyond fairness. The concrete technical difficulties in achieving this goal, and how to overcome them, are presented in the next section. The formal privacy guarantees of our fair training protocol are stated in the following proposition.

**Proposition 4.** *For non-colluding* M *and* REG, *our protocol implements the fair model training functionality satisfying constraints (C1)-(C3) in Section 6.2 in the presence of a semi-honest adversary.*

The proof holds in the random oracle model, as a standard simulation argument combining several MPC primitives (Mohassel & Zhang, 2017; Gascón et al., 2017). It leverages security of arithmetic sharing, garbled circuits, and oblivious transfer protocols in the semi-honest model (Goldreich et al., 1987). Before going into details about the specific technical challenges of fair model training, we now provide a short, high-level introduction to MPC, including a description of the relevant techniques from Mohassel & Zhang (2017) used in our protocol.

**Secret sharing.** A secret sharing scheme allows one to split a value $x$ (the secret) among two parties, so that no party has unilateral access to $x$. In our setting, a user Alice will secret share a sensitive value, for example her race, among a modeler M and a regulator REG. Among prominent secret sharing schemes are *Shamir secret sharing*, *xor sharing*, *Yao sharing*, or *arithmetic multiplicative/additive sharing*. In our protocols we alternate between Yao sharing and additive sharing for efficiency. We have already describe the latter. To recap: the value $x$ is represented in a finite domain $\mathbb{Z}_q$. To share her race, Alice samples a value $r$ from $\mathbb{Z}_q$ uniformly at random, and sends $x - r$ to M and $r$ to REG. We call each of $x - r$ and $r$ a *share*, and denote them as $\langle x \rangle_1$ and $\langle x \rangle_2$. Now M and REG can recover $x$ by adding their shares, but each share on its own does not reveal anything about the value of $x$ (other than that it is smaller than $q$). The case where $q = 2$ corresponds to xor sharing.

**Function evaluation.** MPC can be classified in two groups depending on how $f$ is represented: either as a Boolean or arithmetic circuit. All protocols proceed by having the parties jointly evaluate the circuit, processing it gate by gate. For each gate $g$ for which the value for the input wires $x, y$ is shared among the parties, the parties run a subprotocol to produce the value $z = g(x, y)$ of the output wire, again shared, without revealing any information in the process. In the setting where we use arithmetic additive sharing, the two parties M and REG hold shares, $\langle x \rangle_1, \langle y \rangle_1$ and $\langle x \rangle_2, \langle y \rangle_2$, respectively. In this case, $f$ is represented as an arithmetic circuit, and hence each gate $g$ in the circuit is either an addition or a multiplication. Note that if $g$ is an addition gate, then a sharing of $z = g(x, y)$ can be obtained by having each party simply compute locally, i.e., without any interaction, $\langle z \rangle_i = \langle x \rangle_i + \langle y \rangle_i$, for $i \in \{1, 2\}$. If $g$ is a multiplication, the subprotocol to compute shares of $z$ is much more costly. Fortunately, it can be divided into an offline and an online phase.

**The preprocessing model in MPC.** In this model, two parties $P_1, P_2$ engage in an offline phase, which is data independent, and compute (and store) *shared multiplication triples* of the form $(a, b, c)$, with $c = ab$. Here, $a, b \in \mathbb{Z}_q$ are drawn uniformly at random, and

each value $a, b, c$ is shared among the parties as explained above. In the online phase, a multiplication gate $z = \mathtt{mul}(x, y)$ on shared values $x, y$ can be evaluated as follows: (1) each $P_i$ sets $\langle e \rangle_i = \langle x \rangle_i - \langle a \rangle_i$ and $\langle f \rangle_i = \langle y \rangle_i - \langle b \rangle_i$, (2) the parties exchange their shares of $e$ and $f$ and reconstruct these values locally by simply adding the shares, and (3) each $P_i$ computes $\langle z \rangle_i = ef + f\langle a \rangle_i + e\langle b \rangle_i + \langle c \rangle_i$. The correctness of this protocol can be checked by multiplying out all terms. Privacy relies on the uniform randomness of $a, b$, and hence $\langle e \rangle_i$ and $\langle f \rangle_i$ completely mask the values of $\langle x \rangle_i$ and $\langle y \rangle_i$, respectively. For a formal proof see Demmler et al. (2015b).

Hence, for each multiplication in the function to be evaluated, the parties need to jointly generate a multiplication triple in advance. For computations with many multiplications (like in our case) this can be a costly process. However, this constraint is easy to accommodate in our architecture for private fair model training, as M and REG can run the offline phase once "overnight". Arithmetic multiplication via precomputed triples is a common technique exploited by popular MPC frameworks (Demmler et al., 2015b; Damgård et al., 2012). In this setting, a number of protocols for triple generation (which we did not describe) are available (Keller et al., 2018) and under continuous improvement. These protocols are often based on either Oblivious Transfer or Homomorphic Encryption.

**The two-server model for multi-party learning.** Due to a sequence of theoretical breakthroughs mentioned above as well as the general speedup enabled by faster hardware, in the last three decades MPC has gone from being a mathematical curiosity to a technology of practical interest with commercial applications. There exist a number of openly available implementations (Demmler et al., 2015b; Zahur & Evans, 2015) for generic MPC protocols such as the ones based on arithmetic sharing (Damgård et al., 2012), garbled circuits (Yao, 1986), or GMW (Goldreich et al., 1987). These protocols have different trade-offs in terms of the number of parties they support, network requirements, and scalability for different kinds of computations. In our work, we focus on the 2-party case, as the MPC computation is done by M and REG.

Since generic protocols do not yet scale to input sizes typically encountered in machine learning applications like ours, we extend the SGD protocol from Mohassel & Zhang (2017), in which the following useful accelerating techniques are presented.

- Efficient rescaling: As our arithmetic shares represent fixed-point numbers, we need to rescale by the precision $p$ after every multiplication. This involves dividing by $2^p$, an expensive operation to do in MPC, and in particular in arithmetic sharing. Mohassel & Zhang (2017) show an elegant solution to this problem: the parties can rescale locally by dropping $p$ bits of their shares. It is not hard to see that this might produce the wrong result. However, the parameters of the arithmetic secret sharing scheme can be

set such that with a tunable arbitrarily large probability the error is at most $\pm 1$. This trick can be used for any division by a power of two.

- Alternating sharing types: As already pointed out in previous work (Demmler et al., 2015b), alternating between secret sharing schemes can provide significant acceleration for some applications. Intuitively, arithmetic operations are fast in arithmetic shares, while comparisons are fast in schemes that represent functions as Boolean circuits. Examples of the latter are the GMW protocol and Yao's garbled circuits. In our implementation, we follow this recipe and implement matrix-vector multiplication using arithmetic sharing, while for evaluating our variant of sigmoid, we rely on the protocol from Mohassel & Zhang (2017) implemented with garbled circuits using the Obliv-C framework (Zahur & Evans, 2015).

- Matrix multiplication triples: Another observation is that the idea described above for preprocessing multiplications over arithmetic shares can be reinterpreted at the level of matrices. This results in a faster online and offline phases, see Mohassel & Zhang (2017) for details.

**How to prove that a protocol is secure.** We did not provide a formal definition of security so far and instead referred the reader to Mohassel & Zhang (2017). In MPC, privacy in the case of semi-honest adversaries is argued in the simulation paradigm, see Goldreich (2004) or Lindell (2016) for formal definitions and detailed proofs. Intuitively, in this paradigm one proves that every inference that a party—in our case either REG or M—could draw from observing the execution trace of the protocol could also be drawn from the output of the execution and the party's input. This is done by proving the existence of a *simulator* that can produce an execution trace that is indistinguishable from the actual execution trace of the protocol. A crucial point is that the simulator only has access to the input and output of the party being simulated.

**Step-by-step description of protocols.** Before moving on to specific technical implementation challenges of the building blocks above, we now put the all together into a concise description of the protocols for fairness certification, fair model training, and decision verification.

**Fairness certification**

1. Users send their non-sensitive data to REG in the clear and send one share of their sensitive data to M and REG respectively.

2. M and REG engage in a secure two-party computation to which M provides its shares of sensitive attributes $\langle Z \rangle_1$ as well as model parameters $\theta$ and REG provides its shares of sensitive attributes $\langle Z \rangle_2$ as well as all non-sensitive data $X, y$ as inputs.

3. Within the secure computation, the agreed upon protocol reconstructs $Z$ and evaluates the (violation of) the fairness criterion $F(\theta, X, y, Z)$.

4. If the model $\theta$ satisfies the fairness criterion $F$ on the provided dataset, still within the joint computation a hash of the model parameters $s(\theta)$ is computed.

5. As a result of the joint computation, REG receives and indicator whether the fairness criterion is satisfied and if so, the hash of the provided model parameters. M receives no output from the computation.

**Fair model training**

1. Users send their non-sensitive data to M in the clear and send one share of their sensitive data to M and REG respectively.

2. M and REG engage in a secure two-party computation to which M provides its shares of sensitive attributes $\langle Z \rangle_1$ as well as all non-sensitive data $X, y$ and REG provides its shares of sensitive attributes $\langle Z \rangle_2$ as inputs.

3. Within the secure computation, the agreed upon protocol reconstructs $Z$ and runs through a fixed number of stochastic gradient descent epochs to minimize the following constrained optimization problem

$$\min_{\theta} \sum_{i=1}^{n} \ell_{\theta}(x_i, y_i) \quad \text{s.t.} \quad F(\theta, x_i, y_i, z_i) \leq 0 \,.$$

4. M receives the resulting model parameters $\theta$ of this optimization as output of the joint computation. REG receives no output from the computation.

**Decision verification**

1. A specific user who has used the service of M with inputs $x_i$ has received $y$ as an outcome from M. M may have arrived at decision $y$ by some unknown model denoted by ?, i.e., $y = ?(x_i)$.

2. The user sends $x_i, y$ to REG.

3. M and REG engage in a secure two-party computation to which M provides model parameters $\theta$ and REG provides $x_i, y$ as well as the signature $s$ of model parameters that have previously been certified as fair.

4. Within the secure computation, the agreed upon protocol computes the hash of the provided model parameters $\theta$ and checks whether it matches the signature $s$. If so, it verifies that applying the provided model to the provided inputs $x_i$ results in the provided outcome $y$.

5. REG receives indicators for whether both checks passed the verification as output of the joint computation. M receives no output from the computation.

## 6.4 Technical challenges of fair training

We now present our custom tailored approaches for learning and evaluating fair models with encrypted sensitive attributes. We do this via the following contributions:

- We argue that current optimization techniques for fair learning algorithms are unstable for fixed-point data, which is required by our MPC techniques.

- We describe optimization schemes that are well-suited for learning over fixed-point number representations.

- We combine tricks to approximate non-linear functions with specialized operations to make fixed-point arithmetic feasible and avoid over- and under-flows.

The optimization problem at hand is to learn a classifier $\boldsymbol{\theta}$ subject to a (convex relaxation of a) fairness constraint $F(\boldsymbol{\theta})$:

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^{n} \ell_{\boldsymbol{\theta}}(x_i, y_i) \qquad \text{subject to } F(\boldsymbol{\theta}) \leq \mathbf{0}, \tag{6.2}$$

where $\ell_{\boldsymbol{\theta}}$ is a loss term (the logistic loss in our applications). We collect user data from $U_1, \ldots, U_n$ into matrices $\boldsymbol{X} \in \mathbb{R}^{n \times d}, \boldsymbol{Z} \in \{0,1\}^{n \times p}$, and a label vector $\boldsymbol{y} \in \{0,1\}^n$.

We will focus on a convex approximation of the $p\%$-rule, see eq. (6.1), for linear classifiers following the derivation of Zafar et al. (2017b). To this end, we measure unfairness as the covariance between $\boldsymbol{Z}$ and the signed distance of the feature vectors $\boldsymbol{X}$ to the decision boundary implied by a logistic model with parameters $\boldsymbol{\theta}$, i.e., these distances are given by $\boldsymbol{X\theta}$. We then compute the covariance as

$$\mathrm{cov}(\boldsymbol{z}, \boldsymbol{x}^\top \boldsymbol{\theta}) = \mathbb{E}[(\boldsymbol{z} - \bar{\boldsymbol{z}})(\boldsymbol{x}^\top \boldsymbol{\theta} - \bar{\boldsymbol{x}}^\top \boldsymbol{\theta})] = \mathbb{E}[(\boldsymbol{z} - \bar{\boldsymbol{z}})\boldsymbol{x}^\top \boldsymbol{\theta}] - \underbrace{\mathbb{E}[(\boldsymbol{z} - \bar{\boldsymbol{z}})]}_{=0} \bar{\boldsymbol{x}}^\top \boldsymbol{\theta} = \mathbb{E}[(\boldsymbol{z} - \bar{\boldsymbol{z}})\boldsymbol{x}^\top \boldsymbol{\theta}],$$

where $\bar{\boldsymbol{z}}$ and $\bar{\boldsymbol{x}}$ are the means of all inputs $z_i$ and $x_i$ respectively. By defining $\hat{\boldsymbol{Z}}$ to be the matrix of all $\hat{z}_i := z_i - \bar{\boldsymbol{z}}$, we can now write the approximate fairness constraint as

$$F(\boldsymbol{\theta}) = \frac{1}{n} |\hat{\boldsymbol{Z}}^\top \boldsymbol{X\theta}| - \boldsymbol{c}, \tag{6.3}$$

where $\boldsymbol{c} \in \mathbb{R}^d$ is a constant vector with non-negative entries corresponding to the tightness of the fairness constraint. Note that the entries of $\boldsymbol{c}$ take the role of slack variables—the smaller the entries of $\boldsymbol{c}$, the tighter the constraint. Unlike the original $p\%$-rule, eq. (6.3) is

convex in $\boldsymbol{\theta}$. The correspondence to the $p\%$-rule can be understood as follows: Because of $\mathrm{P}(y = 1 \,|\, z) = \mathrm{P}(\mathrm{sign}(\boldsymbol{x}^\top \boldsymbol{\theta}) = 1 \,|\, z)$, whenever the $p\%$-rule is satisfied for $p = 100$, we have that $\mathrm{P}(\boldsymbol{x}^\top \boldsymbol{\theta} > 0 \,|\, z) = \mathrm{P}(\boldsymbol{x}^\top \boldsymbol{\theta} > 0 \,|\, z')$ for all $z, z'$. In that situation the covariance will also be approximately zero at least asymptotically in the large sample limit. Hence, a sensible relaxation of the 100%-rule, which corresponds to the covariance being exactly zero, we can allow the covariance to deviate from zero in a controlled fashion. In this sense, the slack or tightness variable $c$ roughly corresponds to the $p$ value in the $p\%$-rule. With $A := \nicefrac{1}{n}\, \hat{\boldsymbol{Z}}^\top \boldsymbol{X}$, our final $p\%$ *constraint* thus reads $\boldsymbol{F}(\boldsymbol{\theta}) = |A\boldsymbol{\theta}| - \boldsymbol{c}$, where the absolute value is taken element-wise.

**Current techniques.** To solve the optimization problem in eq. (6.2), with the fairness function $\boldsymbol{F}$ in eq. (6.3), Zafar et al. (2017b) use Sequential Least Squares Programming (SLSQP). This technique works by reformulating eq. (6.2) as a sequence of Quadratic Programs (QPs). After solving each QP, their algorithm uses the Han-Powell method, a quasi-Newton method that iteratively approximates the Hessian $\boldsymbol{H}$ of the objective function via the update

$$H_{t+1} = H_t + \frac{l_\Delta l_\Delta^\top}{\boldsymbol{\theta}_\Delta^\top l_\Delta} - \frac{H_t \boldsymbol{\theta}_\Delta \boldsymbol{\theta}_\Delta^\top H_t}{\boldsymbol{\theta}_\Delta^\top H_t \boldsymbol{\theta}_\Delta} \,,$$

where

$$l_\Delta = l(\boldsymbol{\theta}_{t+1}, \boldsymbol{\lambda}_{t+1}) - l(\boldsymbol{\theta}_t, \boldsymbol{\lambda}_t) \quad \text{and} \quad \boldsymbol{\theta}_\Delta = \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t \,.$$

The Lagrangian is given by

$$l(\boldsymbol{\theta}_t, \boldsymbol{\lambda}_t) = \sum_{i=1}^n \ell_{\boldsymbol{\theta}_i}(x_i, y_i) + \boldsymbol{\lambda}_t^\top \boldsymbol{F}(\boldsymbol{\theta}_t) \,.$$

There are two issues with this approach from an MPC perspective. First, solving a sequence of QPs is prohibitively time-consuming in MPC. Second, while the above Han-Powell update performs well on floating-point data, the two divisions by non-constant, non-integer numbers easily underflow or overflow with fixed-point numbers.

**Fixed-point-friendly optimization techniques.** Instead, to solve the optimization problem in eq. (6.2), we perform stochastic gradient descent and experiment with the following techniques to incorporate the constraints.

- *Lagrangian multipliers.* Here we minimize

$$\mathcal{L} := \frac{1}{n} \sum_{i=1}^n \ell_{\boldsymbol{\theta}}^{\mathrm{BCE}}(x_i, y_i) + \boldsymbol{\lambda}^\top \max\{\boldsymbol{F}(\boldsymbol{\theta}), \boldsymbol{0}\} \,,$$

using stochastic gradient descent, i.e., alternating updates

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} \mathcal{L} \quad \text{and} \quad \lambda \leftarrow \max\{\lambda + \eta_{\lambda} \nabla_{\lambda} \mathcal{L}, 0\},$$

where $\eta_{\boldsymbol{\theta}}, \eta_{\lambda}$ are the learning rates.

- *Projected gradient descent.* For this method, consider specifically the $p\%$-rule based notion $\boldsymbol{F}(\boldsymbol{\theta}) = |\boldsymbol{A}\boldsymbol{\theta}| - \boldsymbol{c}$. We first define $\hat{\boldsymbol{A}}$ as the matrix consisting of the rows of $\boldsymbol{A}$ for which $\boldsymbol{F}(\boldsymbol{\theta}) > \boldsymbol{0}$, i.e., where the constraint is active. In each step, we project the computed gradient of the binary-cross-entropy loss $\mathcal{L}^{\text{BCE}}$—either of a single example or averaged over a minibatch—back into the constraint set, i.e.,

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta_{\boldsymbol{\theta}} (\text{Id}_d - \hat{\boldsymbol{A}}^{\top} (\hat{\boldsymbol{A}} \hat{\boldsymbol{A}}^{\top})^{-1} \hat{\boldsymbol{A}}) \nabla_{\boldsymbol{\theta}} \ell_{\boldsymbol{\theta}}^{\text{BCE}} . \tag{6.4}$$

- *Interior point log barrier (Boyd & Vandenberghe, 2004).* We can approximate eq. (6.2) for the $p\%$-rule constraint $\boldsymbol{F}(\boldsymbol{\theta}) = |\boldsymbol{A}\boldsymbol{\theta}| - \boldsymbol{c}$ by:

$$\text{minimize} \ \sum_{i=1}^{n} \ell_{\boldsymbol{\theta}}^{\text{BCE}}(\boldsymbol{x}_i, y_i) - \frac{1}{t} \sum_{j=1}^{p} \left( \log(\boldsymbol{a}_j^{\top} \boldsymbol{\theta} + c_j) + \log(-\boldsymbol{a}_j^{\top} \boldsymbol{\theta} + c_j) \right),$$

where $\boldsymbol{a}_j$ is the $j$th row of $\boldsymbol{A}$. The parameter $t$ trades off the approximation of the true objective and the smoothness of the objective function. Throughout training, $t$ is increased, allowing the solution to move closer to the boundary. As the gradient of the objective has a simple closed form representation, we can perform regular (stochastic) gradient descent.

After extensive experiments, described in detail in Section 6.5, we found the Lagrangian multipliers technique to work best. It yields high accuracies, reliably stays within the constraints and is robust to hyperparameter changes such as learning rates or the batch size. For a proof of concept, in Section 6.5 we focus on the $p\%$-rule, i.e., eq. (6.3). We note that the gradients for equalized odds or equal opportunity criteria take a similarly simple form, i.e., balancing the true positive or true negative rates is simple to implement for the Lagrangian multiplier technique, but harder for projected gradient descent. However, these fairness notions are more expensive as we have to compute $\boldsymbol{Z}^{\top} \boldsymbol{X}$ for each update step, instead of pre-computing it once at the beginning of training, see Algorithm 2. We could speed up the computation again by evaluating the constraint only on the current minibatch for each update, in which case we risk violating the fairness constraint.

**MPC-friendliness.** For eq. (6.3), we can compute the gradient updates in all three methods with elementary linear algebra operations (matrix multiplications) and a single evaluation of the logistic function. While MPC is well suited for linear operations, most nonlinear

Figure 6.2 Piecewise linear approximations for the non-linear sigmoid function (in black) from Mohassel & Zhang (2017) in blue and from Faiedh et al. (2001) in orange.

functions are prohibitively expensive to evaluate in an MPC framework. Hence we tried two piecewise linear approximations for $\sigma(x)$. The first was recently suggested for machine learning in an MPC context (Mohassel & Zhang, 2017) and is simply constant 0 and 1 for $x < -0.5$ and $x > 0.5$ respectively, and linear in between. The second uses the optimal first order Chebychev polynomial on each interval $[x, x + 1]$ for $x \in \{-5, -4, \ldots, 4\}$, and is constant 0 or 1 outside of $[-5, 5]$ (Faiedh et al., 2001). While it is more accurate, we only report results for the simpler first approximation, as it yielded equal or better results in all our experiments. Both approximations are shown in Figure 6.2.

As the largest number that can be represented in fixed-point format with $m$ integer and $m$ fractional bits is roughly $2^m + 1$, overflow becomes a common problem. Since we whiten the features $X$ column-wise, we need to be careful whenever we add roughly $2^m$ numbers or more, because we cannot even represent numbers greater than $2^m$. In particular, the minibatch size has to be smaller than this limit. For large $n$, the multiplication $Z^\top X$ in the fairness function $F$ for the $p\%$-rule is particularly problematic.

Hence, we split both factors into blocks of size $b \times b$ with $b < 2^m$ and normalize the result of each blocked matrix multiplication by $b$ before adding up the blocks. We then multiply the sum by $b/n > 2^{-m}$. As long as $b$, $b/n$ (and thus also $n/b$) can be represented with sufficient precision, which is the case in all our experiments, this procedure avoids under- and overflow. Note that we require the sample size $n$ to be a multiple of $b$. In practice, we have to either discard or duplicate part of the data. Since the latter may introduce bias, we recommend subsampling. Once we have (an approximation of) $A \in \mathbb{R}^{p \times d}$, we resort to normal matrix multiplication, as typically $p, d \lesssim 100$, see Table 6.1.

Division is prohibitively expensive in MPC. Because there are no simple protocols such as the ones described for addition and multiplication, divisions have to approximated by a sequence of more basic computations that are faster to perform in MPC. Hence, we set the minibatch size to a power of two, which allows us to use fast bit shifts for divisions when averaging over minibatches. Unlike general division, bit shifts can be implemented efficiently

in MPC. To exploit the same trick when averaging over/across blocks in the blocked matrix multiplication, we choose $n$ as the largest possible power of two, see Table 6.1.

---

**Algorithm 2** Fair model training with private sensitive values using Lagrangian multipliers for $\boldsymbol{F}(\boldsymbol{\theta}) = {}^1/n|\boldsymbol{Z}^\top X| - \boldsymbol{c}$ with two parties: the modeler M and the regulator REG.

---

**Input from M:** $\langle \boldsymbol{Z} \rangle_{\mathbf{1}} \in \mathbb{Z}_q^{n \times p}$

**Input from REG:** $X \in \mathbb{Z}_q^{n \times d}$, $\boldsymbol{y} \in \mathbb{Z}_q^n$, $\langle \boldsymbol{Z} \rangle_{\mathbf{2}} \in \mathbb{Z}_q^{n \times p}$

**Input publicly known:** learning rates $\eta_{\boldsymbol{\theta}}, \eta_{\lambda}$, number of training examples $n$, minibatch size $2^s$, constraints $\boldsymbol{c} \in \mathbb{Z}_q^p$, and number of epochs $N_e$

1: $\boldsymbol{\theta} \leftarrow \mathbf{0}, \lambda \leftarrow \mathbf{0}$

2: $A \leftarrow \text{BLOCKEDMULTSHIFTAVG}(\boldsymbol{Z}^\top, X)$

3: **for** $j$ from 1 to $N_e$ **do**

4:     **for** $i$ from 1 to ${}^n/2^s$ **do**

5:         $(X_i, \boldsymbol{y}_i) \leftarrow \text{SAMPLEMINIBATCH}(X, \boldsymbol{y})$

6:         $\boldsymbol{F} \leftarrow |A\boldsymbol{\theta}| - \boldsymbol{c}$

7:         $\nabla_\lambda \leftarrow \max\{\boldsymbol{F}, \mathbf{0}\}$

8:         $\sigma \leftarrow \text{SIGMOIDAPPROX}(X_i\boldsymbol{\theta})$

9:         $\nabla_{\boldsymbol{\theta}}^{\text{BCE}} \leftarrow \text{SHIFTDIVIDE}(X_i^\top (\sigma - \boldsymbol{y}_i), 2^s)$

10:         $\nabla_{\boldsymbol{\theta}}^{\text{CON}} \leftarrow \begin{cases} A^\top \lambda, & \text{if } A > \mathbf{0} \wedge \boldsymbol{F} > \mathbf{0} \\ -A^\top \lambda, & \text{if } A < \mathbf{0} \wedge \boldsymbol{F} > \mathbf{0} \\ \mathbf{0}, & \text{if } \boldsymbol{F} \leq \mathbf{0} \end{cases}$

11:         $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta_{\boldsymbol{\theta}}(\xi_j^{\text{BCE}}\nabla_{\boldsymbol{\theta}}^{\text{BCE}} + \xi_j^{\text{CON}}\nabla_{\boldsymbol{\theta}}^{\text{CON}})$

12:         $\lambda \leftarrow \max\{\lambda + \eta_\lambda \nabla_\lambda, \mathbf{0}\}$

13: **return** Parameters $\boldsymbol{\theta}$

---

**The fair training algorithm.** Algorithm 2 describes the computations M and REG have to perform for fair model training using the Lagrangian multiplier technique and the $p\%$-rule from eq. (6.3). The hyperparameters introduced in the algorithm, such as $\xi^{\text{BCE}}, \xi^{\text{CON}}$, will be explained in the next section. We implicitly assume all computations are performed jointly on additively shared secrets by M and REG as described in Section 6.3. This means that M and REG each receive a secret share of the protected attributes $\boldsymbol{Z}$. Following the protocols outlined in Section 6.3, they can then jointly evaluate the steps in Algorithm 2. This allows them to operate on the sensitive values within the MPC computation, while preventing unilateral access to them by M and REG. The result of these computations is the same as evaluating the algorithm as described with data in the clear.

BLOCKEDMULTSHIFTAVG stands for the blocked matrix multiplication to avoid overflow for fixed-point numbers described towards the end of Section 6.4. Note that it already contains

Table 6.1 Dataset sizes and online timing results of MPC certification and training over 10 epochs with batch size 64.

|                     | Adult     | Bank      | COMPAS    | German    | SQF       |
|---------------------|-----------|-----------|-----------|-----------|-----------|
| $n$ training examples | $2^{14}$ | $2^{15}$ | $2^{12}$ | $2^9$    | $2^{16}$ |
| $d$ features        | 51        | 62        | 7         | 24        | 23        |
| $p$ sensitive attr. | 1         | 1         | 7         | 1         | 1         |
| certification       | 802 ms    | 827 ms    | 288 ms    | 250 ms    | 765 ms    |
| training            | 43 min    | 51 min    | 7 min     | 1 min     | 111 min   |

the division by $n$. The averaging within the blocked matrix multiplications as well as over the results thereof are done by fast bit shifts instead of slow MPC division circuits. This is possible, because we chose all parameters such that divisions are always by powers of two.

We found the piecewise linear approximation of the sigmoid function introduced by Mohassel & Zhang (2017)

$$\text{SIGMOIDAPPROX}(x) := \begin{cases} 0 & \text{if } x \leq -\frac{1}{2}, \\ x + \frac{1}{2} & \text{if } -\frac{1}{2} < x < \frac{1}{2}, \\ 1 & \text{if } x \geq \frac{1}{2}. \end{cases}$$

to work well in practice.

## 6.5   Experiments

The root cause for most technical difficulties pointed out in the previous section is the necessity to work with fixed-point numbers and the high computational cost of MPC. Hence, major concerns are loss of precision and infeasible running times. In this section, we show how to overcome both doubts and that fair training, certification and verification are feasible for realistic datasets.

**Experimental setup and datasets.**   We work with two separate code bases. Our Python code does not implement MPC, but allows to flexibly switch between floating and fixed-point numbers as well as exact non-linear functions and their approximations. We used it mostly for validation and empirical guidance in our design choices. The full MPC protocol is implemented in C++ on top of the Obliv-C garbled circuits framework (Zahur & Evans, 2015) and the Absentminded Crypto Kit (Doerner, 2018). This is done as described in Section 6.3 for the Lagrangian multiplier technique. It accurately mirrors the computations performed by the first implementation on encrypted data. Except for the timing results in Table 6.1, all

comparisons with floating-point numbers or non-linearities were done with the versatile Python implementation.

All our experiments use a batch size of 64, a fixed number of epochs scaling inversely with dataset size $n$ (such that we always perform roughly 15,000 gradient updates), fixed learning rates of $\eta_\theta = 10^{-4}, \eta_\lambda = 0.05$, and an annealing schedule for $1/t$ in the interior point logarithmic barrier method as described by Boyd & Vandenberghe (2004). The weights for the gradients of the regular binary cross entropy loss (BCE) and the loss from the constraint terms (CON) follow the schedules

$$\zeta_j^{\mathrm{BCE}} = \frac{N_e}{N_e + j} , \qquad \zeta_j^{\mathrm{CON}} = \frac{N_e + 10j}{N_e} .$$

Weight decay, adaptive learning rate schedules, and momentum neither consistently improved nor impaired training. Therefore, all reported numbers were achieved with vanilla SGD, for fixed learning rates, and without any regularization. After extensive testing on all datasets, we converged to a fixed-point representation with 16 bits for the integer and fractional part respectively. The smaller the number of bits, the faster the MPC implementation and the higher the risk of loss of precision or over- and underflows. We found 16 bits to be the minimally needed precision for all our experiments to work.

We consider 5 real world datasets, namely the adult (*Adult*), German credit (*German*), and bank market (*Bank*) datasets from the UCI machine learning repository (Lichman, 2013), the stop, question and frisk 2012 dataset (*SQF*),[1] and the COMPAS dataset (Angwin et al., 2016) (*COMPAS*). For practical purposes (see Section 6.4), we subsample $2^i$ examples from each dataset with the largest possible $i$, see Table 6.1. Moreover, we also run on synthetic data, generated as described by Zafar et al. (2017b, Section 4.1), as it allows us to control the correlation between the sensitive attributes and the class labels. It is thus well suited to observe how different optimization techniques handle the fairness accuracy trade-off. For comparison we use the SLSQP approach described in Section 6.4 as a baseline. We run all methods for a range of constraint values in $[10^{-4}, 10^0]$ and a corresponding range for SLSQP.

In the plots in this section, discontinuations of lines indicate failed experiments. The most common reasons are overflow and underflow for fixed-point numbers, and instability due to exploding gradients. Plots and analyses for the remaining datasets can be found in Appendix A.

**Comparing optimization techniques.**    First we evaluate which of the three optimization techniques works best in practice. Figure 6.3 shows the test set accuracy over the constraint
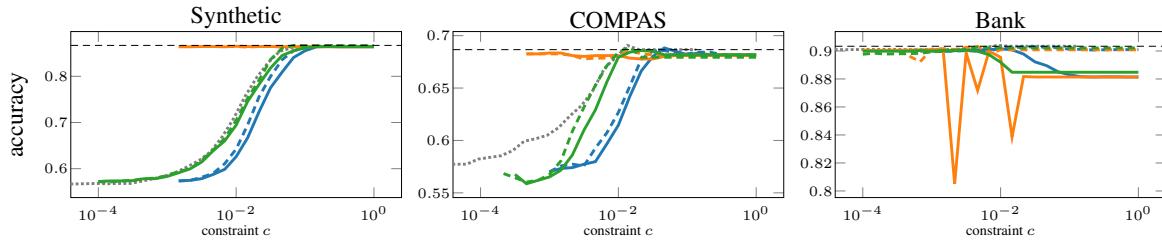
---

[1]https://perma.cc/6CSM-N7AQ

Figure 6.3 Test set accuracy over the $p\%$ value for different optimization methods (blue: iplb, orange: projected, green: Lagrange) and either no approximation (*continuous*) or a piecewise linear approximation (*dashed*) of the sigmoid using floating-point numbers. The gray dotted line is the baseline (see Section 6.4) and the black dashed line is unconstrained logistic regression (from scikit-learn).

value. By design, the synthetic dataset exhibits a clear trade-off between accuracy and fairness. The Lagrange technique closely follows the (dotted) baseline from Zafar et al. (2017b), whereas iplb performs slightly worse (and fails for small $c$). Even though the projected gradient method formally satisfies the proxy constraint for the $p\%$ rule, it does so by merely shrinking the parameter vector $\theta$, which is why it also fails for small $c$. We analyze this behavior in more detail in Appendix A.

The COMPAS dataset is the most challenging as it contains 7 sensitive attributes, one of which has only 10 positive instances in the training set. Since we enforce the fairness constraint individually for each sensitive attribute (we randomly picked one for visualization), the classifier tends to collapse to negative predictions. All three methods maintain close to optimal accuracy in the unconstrained region, but collapse more quickly than SLSQP. This example shows that the $p\%$-rule proxy itself needs careful interpretation when applied to multiple sensitive attributes simultaneously and that our SGD based approach seems particularly prone to collapse in such a scenario. On the Bank dataset, accuracy increases for iplb and Lagrange when the constraint becomes active as $c$ decreases until they match the baseline. Determining the cause of this—perhaps unintuitive—behavior requires further investigation. We currently suspect the constraint to act as a regularizer. The projected gradient method is unreliable on the Bank dataset.

Empirically, the Lagrangian multiplier technique is most robust with maximal deviations of accuracy from SLSQP of $< 4\%$ across the 6 datasets and all constraint values. We substantiate this claim in Appendix A. For the rest of this section we only report results for Lagrangian multipliers. Figure 6.3 also shows that using a piecewise linear approximation as described in Section 6.4 for the logistic function does not spoil performance.

**Fair training, certification and verification.** Figure 6.4 shows how the fractions of users with positive outcomes in the two groups ($z = 0$ is continuous and $z = 1$ is dashed) are
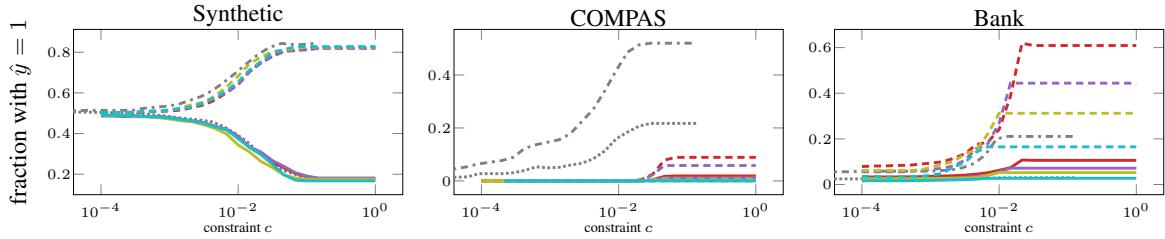
Figure 6.4 The fraction of people with $z = 0$ (*continuous/dotted*) and $z = 1$ (*dashed/dash-dotted*) who get assigned positive outcomes (red: no approx. + float, purple: no approx. + fixed, yellow: pw linear + float, turquoise: pw linear + fixed, gray: baseline). As the slack $c$ decreases, the fairness constraint is tightened and we observe dashed and continuous lines of the same color moving closer together. This indicates that as the fairness constraint is enforced, equal fractions of both populations receive positive outcomes.

gradually balanced as we decrease the fairness constraint $c$. These plots can be interpreted as the degree to which disparate impact is mitigated as the constraint is tightened. The effect is most pronounced for the synthetic dataset by construction. As discussed above, the collapse for the COMPAS dataset occurs faster than for SLSQP due to the constraints from multiple sensitive attributes. In the Bank dataset, for large $c$—before the constraint becomes active—the fractions of positive outcomes for $z = 1$ differ, which is related to the slightly suboptimal accuracy at large $c$ that needs further investigation. However, as the constraint becomes active, the fractions are balanced at a similar rate as the baseline. Overall, our Lagrangian multiplier technique with fixed point numbers and piecewise linear approximations of non-linearities robustly manages to satisfy the p%-rule proxy at similar rates as the baseline with only minor losses in accuracy on all but the challenging COMPAS dataset.

In Table 6.1 we show the online running times of 10 training epochs on a laptop computer. While training takes orders of magnitudes longer than a non-MPC implementation, our approach still remains feasible and realistic. We use the one time offline precomputation of multiplication triples described and timed in Mohassel & Zhang (2017, Table 2). As pointed out in Section 6.3, certification of a trained model requires checking whether $F(\theta) > 0$. We already perform this check at least once for each gradient update during training. It only takes a negligible fraction of the computation time, see Table 6.1. Similarly, the operations required for certification stay well below one second.

## 6.6   Conclusion

Real world fair learning has suffered from a dilemma: in order to enforce fairness, sensitive attributes must be examined; yet in many situations, users may feel uncomfortable in revealing these attributes, or modelers may be legally restricted in collecting and utilizing them.

By introducing recent methods from MPC, and extending them to handle linear constraints as required for various notions of fairness, we have demonstrated that it is practical on real-world datasets to: a) certify and sign a model as fair; b) learn a fair model; and c) verify that a fair-certified model has indeed been used; all while maintaining cryptographic privacy of all users' sensitive attributes. Connecting concerns in privacy, algorithmic fairness and accountability, our proposal empowers regulators to provide better oversight, modelers to develop fair and private models, and users to retain control over data they consider highly sensitive.

We have demonstrated the practicability of private and fair model training, certification and verification using MPC as described in Figure 6.1. Using the methods and tricks introduced in Section 6.4, we can overcome accuracy as well as over- and underflow concerns due to fixed-point numbers. Offline precomputation combined with a fast C++ implementation yield viable running times for reasonably large datasets on a laptop computer even though we acknowledge that these are still considerably slower than times obtained without protection via secure MPC. Such trade-offs and as well as potential tensions in the interplay of different notions of fairness, privacy and security pose important challenges for real-world applications, opening up new fields of research that go beyond the traditional assumptions of fair machine learning.

# 7

# Decisions versus predictions

In this chapter, we go beyond the assumption that labels are always available in the training data and show that naive predictive modeling using only the data at hand is therefore not the optimal way to arrive at decisions. To consistently learn accurate predictive models, one needs access to ground truth labels. Unfortunately, in practice, labels may only exist conditional on certain decisions—if a loan is denied, there is not even an option for the individual to pay back the loan. Hence, the observed data distribution depends on how decisions are being made. In this chapter, we show that in this *selective labels* setting, learning a predictor directly only from available labeled data is suboptimal in terms of both fairness and utility. To avoid this undesirable behavior, we propose to directly learn decision policies that maximize utility under fairness constraints and thereby take into account how decisions affect which data is observed in the future. Our results suggest the need for a paradigm shift in the context of fair machine learning from the currently prevalent idea of simply building predictive models from a single static dataset via risk minimization, to a more interactive notion of "*learning to decide*". In particular, such policies should not entirely neglect part of the input space, drawing connections to explore/exploit tradeoffs in reinforcement learning, data missingness, and potential outcomes in causal inference. Experiments on synthetic and real-world data illustrate the favorable properties of learning to decide in terms of utility and fairness.

The main content of this chapter has been published in the following paper:

## 7.1   Introduction

We start by revisiting some common settings of consequential decisions that may be (partially) automated. In pretrial release decisions, a judge may consult a learned model of the probability of recidivism to decide whether to grant bail or not. In loan decisions, a bank may decide whether or not to offer a loan based on learned estimates of the credit default probability. In fraud detection, an insurance company may flag suspicious claims based on a machine learning model's predicted probability that the claim is fraudulent. In all these scenarios, the goal of the decision maker (bank, law court, or insurance company) may be to take decisions that maximize a given utility function. Since such a utility is a function of the entire policy, rather than an individual loss term for each prediction in isolation, it may encode preferences that go beyond merely "correct or incorrect" for a predictive classification task. For example, the utility may accommodate various fairness and diversity considerations, or enforce actions that could improve the wellbeing of certain individuals and groups in the long run—even when such actions are not warranted by the directly observed outcomes. In contrast, the goal of a supervised predictive machine learning model is solely to provide accurate predictions given the available training set, typically under the implicit assumption that the training data is an i.i.d. sample from the distribution encountered during test time. Such a simplistic approach, which we refer to as *learning to predict*, ignores that once decisions are based on these predictions, they may interact with the data collection or have direct impact on the underlying distribution relevant during deployment (Perdomo et al., 2020).

Nevertheless, most work on fair machine learning does not distinguish between decisions and label predictions, which also leads to a perceived trade-off between fairness and accuracy or performance. This stems from the fact that viewing "incorrect" predictions as positively desirable, e.g., because they may promote fairness, is incompatible with the standard goal of minimizing a predictive loss in supervised learning. From that viewpoint, an accurate prediction is equivalent—or at least directly translates—to a good decision. Only recently has the distinction been made explicit, typically emphasizing that *not predicting historically recorded labels correctly* can often be part of the goal when it comes to fairness (Corbett-Davies et al., 2017; Kleinberg et al., 2017a; Mitchell et al., 2018; Valera et al., 2018). We also remind the reader of the discussion of some fundamental challenges, in particular long- versus short-term goals, in Section 2.2. This recent line of work has shown that if a predictive model achieves perfect prediction accuracy, *deterministic threshold rules*, which derive decisions deterministically from the predictive model by thresholding, indeed achieve maximum utility under various fairness constraints. At first, this lends support to focusing on deterministic threshold rules and seemingly justifies using predictions and decisions interchangeably.

However, in many practical scenarios including the ones described in the beginning of this section, the decision determines whether a label is realized or not—if bail (a loan) is denied, there is not even an option for the individual to reoffend (pay back the loan). This problem has been referred to by Lakkaraju et al. (2017) as *selective labels*. As a consequence, the labeled data used to train predictive models often depend on the decisions taken, which likely leads to suboptimal performance. For example, a racist initial policy may categorically reject applicants from a certain demographic group. Therefore, no data about these individuals is collected and there are no guarantees for how a predictive model may extrapolate into this unseen region of the input space. Indeed, deterministic threshold rules using even slightly imperfect predictive models can be far from optimal (Woodworth et al., 2017). This negative result raises the following question: *Can we do better if we learn directly to decide rather than to predict?*[1] Here, by "learning to decide" we mean learning to maximize a utility of a decision policy (rather than predictions) when deployed under test time conditions (rather than merely trained on observed data).

In the present chapter, we first articulate how the "learning to predict" approach fails in a utility maximization setting (with fairness constraints) that accommodates a variety of real-world applications, including those mentioned previously. We show that label data gathered under deterministic rules (e.g., prediction based threshold rules) are neither sufficient to improve the accuracy of the underlying predictive model, nor the utility of the decision making process. We then demonstrate how to overcome this undesirable behavior using a particular family of stochastic decision rules and introduce a simple gradient-based algorithm to learn them from data. Experiments on synthetic and real-world data illustrate our theoretical results and show that, under imperfect predictions, *learning to predict* is inferior to *learning to decide*.

**Related work.** The work most closely related to ours analyzes the long-term effects of consequential decisions informed by data-driven predictive models on underrepresented groups (Hu & Chen, 2018; Liu et al., 2018; Mouzannar et al., 2019). However, this line of work focuses mainly on the evolution of several measures of well-being under a perfect predictive model and neglecting the data collection phase. In contrast, we focus on analyzing how to improve a suboptimal decision process when labels exist only for positive decisions. Potential issues arising from neglecting the data collection process have also been highlighted in a survey of what machine learning practitioners in industry actually need to enforce fairness (Holstein et al., 2019). Dimitrakakis et al. (2019) similarly point out how attempts of

---

[1]We remark that our notion of learning to decide is not immediately related to (Bayesian) decision theory. Moreover, for this thesis, we restrict the notion of learning to predict to simple point estimates (categorical predictions) from risk minimization based on available data. Crafting "predictions" more carefully, for example by regarding a data-missingness model or proper uncertainty estimates, may not be prone to the same issues. However, most works on fairness in machine learning have been staged within our simplistic "learning to predict" framework.

fair machine learning may fail when there is uncertainty about the underlying probabilistic model of the world. They achieve fairness in such settings by deploying a Bayesian approach that aims at enforcing fairness in all possible models weighted by their probability given the current information. In contrast, we focus on analyzing how to improve a suboptimal decision process when labels exist only for positive decisions. We also build on previous work on counterfactual inference and policy learning (Athey & Wager, 2017; Ensign et al., 2018a; Gillen et al., 2018; Heidari & Krause, 2018; Joseph et al., 2016; Jung et al., 2018; Kallus, 2018; Kallus & Zhou, 2018; Lakkaraju & Rudin, 2017). In these settings, the decision typically determines which of the potential outcomes is observed and the focus is on confounders that affect both the decision and the outcome (Rubin, 2005). In contrast, in our approach the decision determines whether there will be an outcome at all, but there is no unobserved confounding. Two notable exceptions are by Kallus & Zhou (2018) and Ensign et al. (2018a), which also consider limited feedback. However, Kallus & Zhou (2018) focus on designing unbiased estimates for fairness measures, rather than learning how to decide. Ensign et al. (2018a) assume a deterministic mapping between features and labels, which allows them to reduce the problem to the apple tasting problem (Helmbold et al., 2000). Remarkably, in their deterministic setting, they also conclude that the optimal decisions should be stochastic.

Unlike in the fairness literature, where deterministic policies dominate (Corbett-Davies et al., 2017; Valera et al., 2018; Meyer et al., 2019), stochastic policies are often necessary to ensure adequate exploration (Silver et al., 2014) in contextual bandits (Dudík et al., 2011; Langford et al., 2008; Agarwal et al., 2014) and reinforcement learning (Jabbari et al., 2017; Sutton & Barto, 1998). However, the typical problem setting there differs fundamentally from ours and typically neither fairness constraints nor selective labels are taken into account. A recent notable exception is Joseph et al. (2016), initiating the study of fairness in multi-armed bandits, however, using a fairness notion orthogonal to the observational group matching criteria we consider in our work, and ignoring the selective labels problem.

## 7.2   Decisions from imperfect predictive models

First, we briefly recap notation. Let $\mathcal{X} \subseteq \mathbb{R}^d$ be the feature domain, $\mathcal{Z} = \{0, 1\}$ the range of sensitive attributes, and $\mathcal{Y} = \{0, 1\}$ the set of ground truth labels. A *decision rule* or *policy*[2] is a mapping $\pi : \mathcal{X} \times \mathcal{Z} \to \mathcal{P}(\{0, 1\})$ that maps an individual's feature vector and sensitive attribute to a probability distribution over *decisions* $d \in \{0, 1\}$. Note that instead of a deterministic predictor $\hat{Y}$ we are now considering a stochastic policy for decisions instead of predictions. We sample $x, z$ and $y$ from a ground truth distribution $P(X, Z, Y) = P(Y \mid X, Z) P(X, Z)$. Decisions $d$ are sampled from a policy $d \sim \pi(D \mid x, z)$, where we often

---

[2]We use the terms *decision rule*, *decision making process* and *policy* interchangeably in this chapter.

write $\pi(x,z)$ for $\pi(D \mid x,z)$ and $\pi(D = 1 \mid x,z)$ for the probability of a positive decision given features $x, z$. The decision determines whether the label $y \sim P(Y \mid X, Z)$ comes into existence. In loan decisions, the feature vector $x$ may include salary, education, or credit history; the sensitive attribute $z$ may indicate sex; a loan can be granted ($d = 1$) or denied ($d = 0$); and the label $y$ indicates repayment ($y = 1$) or default ($y = 0$) *upon receiving a loan*.

Inspired by Corbett-Davies et al. (2017), we measure the *utility* as the expected overall profit provided by the policy with respect to the distribution P, i.e.,

$$u_P(\pi) := \mathbb{E}_{x,z,y\sim P, d\sim\pi(x,z)} [y\,d - c\,d] = \mathbb{E}_{x,z\sim P} [\pi(D = 1 \mid x,z)(P(Y = 1 \mid x,z) - c)] , \quad (7.1)$$

where $c \in (0,1)$ reflects economic considerations of the decision maker. For example, in a loan scenario, the utility gain is $(1 - c)$ if a loan is granted and repaid, $-c$ if a loan is granted but the individual defaults, and zero if the loan is not granted. One could think of adding a term for negative decisions of the form $g(y)(1 - d)$ for some given definition of $g$, however, we would not be able to compute such a term due to the selective labels, except for constant $g$. Therefore, without loss of generality, we assume that $g(y) = 0$ for all $y$, because any non-zero constant $g$ can easily be absorbed in our framework.

For fairness considerations, we define the *f-benefit for group $z \in \{0,1\}$* with respect to the distribution P by

$$b_P^z(\pi) := \mathbb{E}_{x,y\sim P(X,Y \mid z), d\sim\pi(x,z)}[f(d,y)] ,$$

with $f : \{0,1\} \times \{0,1\} \to \mathbb{R}$. Note that common observational group matching fairness criteria can be expressed as $b_P^0(\pi) = b_P^1(\pi)$ for different choices of $f$ (and perhaps conditioning on specific values of $y$ or $d$ corresponding to criteria of separability or sufficiency respectively). For simplicity, we will focus on demographic parity (or no disparate impact), which simply amounts to $f(d,y) = d$.

Under perfect knowledge of $P(Y \mid x,z)$, the policy maximizing the above utility subject to the group benefit fairness constraint $b_P^0(\pi) = b_P^1(\pi)$ is a deterministic threshold rule (Corbett-Davies et al., 2017)[3]

$$\pi^*(D = 1 \mid x,z) = \mathbf{1}[P(Y = 1 \mid x,z) \geq c_z] , \quad (7.2)$$

where we allow for group specific cost factors $c_0, c_1$ such that $b_P^0(\pi) = b_P^1(\pi)$. Without fairness constraints, we simply have $c_0 = c_1 = c$. However, as discussed by Woodworth et al. (2017), in practice we typically do not have access to the true conditional distribution $P(Y \mid x,z)$, but instead to an imperfect predictive model $Q(Y \mid x,z)$ trained on a finite training set. Such a predictive model can similarly be used to implement a deterministic threshold

---

[3]Here, $\mathbf{1}[\bullet]$ is 1 if the predicate $\bullet$ is true and 0 otherwise.

rule as

$$\pi_Q(D = 1 \mid x, z) = \mathbf{1}[Q(Y = 1 \mid x, z) \geq c].  \tag{7.3}$$

Here, the predictor $Q(Y = 1 \mid x, z) \approx P(Y = 1 \mid x, z) - \delta_z$, with $\delta_z = c_z - c$, directly incorporates the fairness constraint, i.e., it is trained to maximize predictive power subject to the fairness constraint. In this context, Woodworth et al. (2017) have shown that this approach often leads to better performance than post-processing a potentially unfair predictor as proposed by Hardt et al. (2016b). Unfortunately, they have also shown that, because of the mismatch between $Q(Y = 1 \mid x, z)$ and $P(Y = 1 \mid x, z) - \delta_z$, the resulting policy $\pi_Q$ will usually still be suboptimal in terms of both utility and fairness. To make things worse, due to the selective labeling, the data points $x, z, y$ observed under a given policy $\pi_0$ are not i.i.d. samples from the ground truth distribution $P(X, Z, Y)$, but instead from the weighted distribution

$$P_{\pi_0}(X, Z, Y) \propto P(Y \mid X, Z) \, \pi_0(D = 1 \mid X, Z) \, P(X, Z).  \tag{7.4}$$

Consequently, if $\pi_0$ is not optimal, i.e., $\pi_0 \neq \pi^*$, the necessary i.i.d. assumption for consistency results of empirical risk minimization is violated, which may also be one reason for a common observation in fairness, namely that predictive errors are often systematically larger for minority groups (Angwin et al., 2016). In the remainder, we will say that the distributions $P_{\pi_0}(X, Z, Y)$ and $P_{\pi_0}(X, Z)$ are *induced* by the policy $\pi_0$. In the next section, we study how to learn the optimal policy, potentially subject to fairness constraints, if the data is collected from an initial faulty policy $\pi_0$.

## 7.3   From deterministic to stochastic policies

Consider a class of policies $\Pi$, within which we want to maximize utility, as defined in eq. (7.1) subject to the group benefit fairness constraint $b_P^0(\pi) = b_P^1(\pi)$. We formulate this as an unconstrained optimization with an additional penalty term, namely to maximize

$$v_P(\pi) := u_P(\pi) - \frac{\lambda}{2}(b_P^0(\pi) - b_P^1(\pi))^2  \tag{7.5}$$

over $\pi \in \Pi$ under the assumption that we do not have access to samples from the ground truth distribution $P(X, Z, Y)$, which $u_P(\pi)$ and $b_P^z(\pi)$ depend on. Instead, we only have access to samples from a distribution $P_{\pi_0}(X, Z, Y)$ induced by a given initial policy $\pi_0$ as in eq. (7.4). We first analyze this problem for deterministic threshold rules, before considering general deterministic policies, and finally also general stochastic policies.

### 7.3.1   Deterministic policies

Assume the initial policy $\pi_0$ is a given deterministic threshold rule and $\Pi$ is the set of all deterministic threshold rules, which means that each $\pi \in \Pi$ (and $\pi_0$) is of the form eq. (7.3)

for some predictive model $Q(Y \mid x, z)$. Given a hypothesis class of predictive models $\mathcal{Q}$, we reformulate eq. (7.5) to maximize

$$v_{\mathrm{P}}(\pi_Q) := u_{\mathrm{P}}(\pi_Q) - \frac{\lambda}{2}(b_{\mathrm{P}}^0(\pi_Q) - b_{\mathrm{P}}^1(\pi_Q))^2 \tag{7.6}$$

over $Q \in \mathcal{Q}$, where the utility and the benefits for $z \in \{0, 1\}$ are simply $u_{\mathrm{P}}(\pi_Q) = \mathbb{E}_{x,z,y \sim \mathrm{P}}[\mathbf{1}[Q(Y = 1 \mid x, z) \geq c](y - c)]$ and $b_{\mathrm{P}}^z(\pi_Q) = \mathbb{E}_{x,z,y \sim \mathrm{P}}[f(\mathbf{1}[Q(Y = 1 \mid x, z) \geq c], y)]$. Note that eq. (7.5) has a unique optimum $\pi^*$ (up to differences on sets of measure zero). Therefore, if $\pi^* \in \Pi$ (the set of all deterministic threshold rules), eq. (7.6) will also reach this optimum if $\mathcal{Q}$ is rich enough. However, the optimal predictor $Q^*$ may not be unique, because the utility and the benefits are not sensitive to the precise values of $Q(Y = 1 \mid x, z)$ above or below the threshold $c$.

If we only have access to samples from the distribution $\mathrm{P}_{\pi_0}$ induced by some $\pi_0 \neq \pi^*$, we may choose to simply learn a predictive model $Q_0^* \in \mathcal{Q}$ that empirically maximizes the objective $v_{\mathrm{P}_{\pi_0}}(\pi_Q)$, where the utility and the benefits are computed with respect to the induced distribution $\mathrm{P}_{\pi_0}$. However, the following negative result shows that, under mild conditions, $Q_0^*$ leads to a suboptimal deterministic threshold rule.

**Proposition 5.** *If there exists a subset $\mathcal{V} \subset \mathcal{X} \times \mathcal{Z}$ of positive measure under $\mathrm{P}$ such that $\mathrm{P}(Y = 1 \mid \mathcal{V}) \geq c$ and $\mathrm{P}_{\pi_0}(Y = 1 \mid \mathcal{V}) < c$, then there exists a maximum $Q_0^* \in \mathcal{Q}$ of $v_{\mathrm{P}_{\pi_0}}$ such that $v_{\mathrm{P}}(\pi_{Q_0^*}) < v_{\mathrm{P}}(\pi_{Q^*})$.*

*Proof.* First, note that any deterministic policy $\pi$ is fully characterized (up to differences of measure zero) by the sets $W_d(\pi) = \{(x, z) \mid \pi(D = 1 \mid x, z) = d\}$ for $d \in \{0, 1\}$. For a deterministic threshold rule $\pi_Q$, we write $W_d(Q) = \{(x, z) \mid \mathbf{1}[Q(Y = 1 \mid x, z) > c] = d\} = W_d(\pi_Q)$. By definition, we have that $v(\pi_Q) \leq v(\pi_{Q^*})$. We note that whenever the symmetric difference between the sets $W_d(Q)$ and $W_d(Q^*)$, $W_d(Q) \Delta W_d(Q^*)$, has positive inner measure (induced by $\mathrm{P}$) for $d \in \{0, 1\}$ and a $Q \in \mathcal{Q}$, we have $v(\pi_Q) \neq v(\pi_{Q^*})$ and thus $v(\pi_Q) < v(\pi_{Q^*})$. Thus it only remains to show that $W_d(Q^*) \Delta W_d(Q_0^*)$ has positive inner measure for $d \in \{0, 1\}$. Since $\mathrm{P}(Y = 1 \mid \mathcal{V}) \geq c$ by assumption, we have $\mathcal{V} \subset W_1(Q^*)$. At the same time, because of $\mathrm{P}_{\pi_0}(Y = 1 \mid \mathcal{V}) < c$ by assumption, we have $\mathcal{V} \cap W_1(\pi_0) = \varnothing$. Finally, we note that for any $Q \in \mathcal{Q}$, we have that $v_{\mathrm{P}_{\pi_0}}(Q) = v_{\mathrm{P}_{\pi_0}}(Q \cdot \chi_{W_1(\pi_0)})$, where $\chi_\bullet$ is the indicator function on the set $\bullet$. Therefore, we can choose a $Q_0^*$ maximizing $v_{\mathrm{P}_{\pi_0}}$ such that $W_1(Q_0^*) \subset W_1(\pi_0)$ and thus $\mathcal{V} \cap W_1(Q_0^*) = \varnothing$. Therefore $\mathcal{V} \subset W_1(Q_0^*) \Delta W_1(Q^*)$ and $\mathcal{V}$ has positive measure under $\mathrm{P}$ by assumption. Thus $W_d(Q_0^*) \Delta W_d(Q^*)$ has positive inner measure and we conclude $v_{\mathrm{P}}(\pi_{Q_0^*}) < v_{\mathrm{P}}(\pi_{Q^*})$.                                            □

**Lending example.** We briefly illustrate this result in a lending example based on FICO credit score data as described in Hardt et al. (2016b). Such single feature scenarios are highly
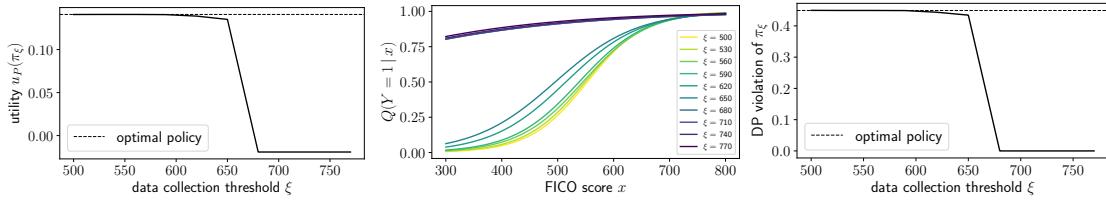
Figure 7.1 We show the utility (left) and the predictive models $Q_\xi$ learned from data collected with an initial threshold of $\xi$ (middle). Finally, we present the violation of demographic parity (right) of threshold decision rules $\pi_\xi$ learned from data collected with an initial threshold of $\xi$. Harsh data collection policies (i.e., large $\xi$)—while achieving demographic parity—render the learned policies useless in terms of utility.

relevant for score-based decision support systems where full training data and the functional form of the score are often not available (e.g., also for pretrial risk assessment). For any score that is strictly monotonic in the true success rate, the optimal policy is simply to threshold the score. This lends additional support to score-based systems.

Here, we can generate new scores for a given group via inverse transform sampling from the known cumulative distribution functions. We consider 80% white and 20% black applicants. A hypothetical new bank that has access to FICO scores $x \in \mathcal{X} := \{300, \dots, 820\}$, but not to the corresponding repayment probabilities may expect to be profitable if at least 70% of granted loans are repaid, i.e., $c = 0.7$. A risk-averse lender may initially choose a high score threshold $\xi \in \mathcal{X}$ and employ the decision rule $\mathbf{1}[x > \xi]$. After collecting repayment data $\mathcal{D}^{(\xi)} := \{(x_i, y_i)\}_{i=1}^n$ with this initial threshold, they learn a model $Q_\xi(Y = 1 \mid x)$ and then decide based on $\pi_\xi(D = 1 \mid x) = \mathbf{1}[Q_\xi(Y = 1 \mid x) > c]$.

For a range of initial data collection score thresholds $\xi \in [500, 800]$, we sample 10,000 scores from the specified population (80% white, 20% black) via inverse transform sampling given the cumulative distributions functions over scores of the two groups. The relatively large number of examples is chosen to illustrate that the negative result is not a consequence of insufficient data. We then fit an L2 regularized logistic regression model to each of these datasets using 5-fold cross validation to select the regularization parameter. This results in a predictive model $Q_\xi$ for each initial data collection threshold $\xi$. For each of these models we construct the decision rule $\pi_\xi(D = 1 \mid x) = \mathbf{1}[Q_\xi(Y = 1 \mid x) > c]$, with $c = 0.7$. We then estimate utility and fairness violation of demographic parity on a large sample from the entire population (one million examples).

In Figure 7.1 we show how the initial data collection threshold $\xi$ affects utility and fairness of the resulting predictive model-based decision rule. Conservatively high initial thresholds of $\xi \geq 650$ lead to essentially useless decisions $\pi_\xi$, because of imperfect prediction models regardless of how much data was collected. More lenient initial policies can result in near

optimal decisions with improved fairness compared to the maximum utility policy for the given cost $c$ (dashed).

A simple fix seems to present itself: Do not start with high thresholds. However, Proposition 5 tells us that it does not matter how low we set the initial threshold, if we use it to derive deterministic decisions. By deterministically thresholding the score, we inevitable reject a subset of the population (with positive measure) and thus can never learn whether some of them may actually repay. This will be highlighted in the next paragraphs, showing practical impossibility results for recovering from a bad initial policy even in a sequential training setting when using deterministic decision rules. The only way to overcome this issue in our example is not to draw a hard threshold, bat to accept applicants with some non-zero probability for *every possible score*. We will formalize this idea and its advantages in Section 7.3.2.

**Impossibility results.** Supplementing the result in Proposition 5, we will now prove that—in certain situations—a sequence of deterministic threshold rules, fails to recover the optimal policy despite it being in the hypothesis class. We assume that each threshold rule is of the form of eq. (7.3) and its associated predictive model is trained using the data gathered through the deployment of previous threshold rules. To this end, we consider a *sequential policy learning task*, which is given by a tuple $(\pi_0, \Pi', \mathcal{A})$, where:

a) $\Pi' \subset \Pi$ is the hypothesis class of policies,

b) $\pi_0 \in \Pi'$ is the initial policy, and

c) $\mathcal{A} : \Pi' \times \bigcup_{i=1}^{\infty}(\mathcal{X} \times \mathcal{Z} \times \mathcal{Y})^i \to \Pi'$ is an update rule.

The update rule $\mathcal{A}$ takes an existing policy $\pi_t$ and a dataset $\mathcal{D} \in (\mathcal{X} \times \mathcal{Z} \times \mathcal{Y})^n$ and produces an updated policy $\pi_{t+1}$, which typically aims to improve the policy in terms of the objective function $v_P(\pi)$ in eq. (7.5). In our setting, the dataset $\mathcal{D}$ is collected by deploying previous policies, i.e., from a mixture of the distributions $P_{\pi_\tau}(X, Z, Y)$ with $\tau \leq t$.

Recall that for deterministic threshold policies we can partition the space $\mathcal{X} \times \mathcal{Z} = W_0(\pi) \cup W_1(\pi)$ into regions of negative and positive decisions. Then, we say an update rule is *non-exploring on $\mathcal{D}$* if and only if $W_0(\mathcal{A}(\pi, \mathcal{D})) \subset W_0(\pi)$. Intuitively, this means that no individual who has received a negative decision under the old policy $\pi$ would receive a positive decision under the new policy $\mathcal{A}(\pi, \mathcal{D})$. Common learning algorithms for classification, such as gradient boosted trees are *error based*, i.e., they only change the decision function when they make errors on the training set. As a result, they lead to non-exploring update rules on $\mathcal{D}$ whenever they achieve zero error.

**Proposition 6.** *Let $(\pi_0, \Pi', \mathcal{A})$ be a sequential policy learning task, where $\Pi' \subset \Pi$ are deterministic threshold policies based on a class of predictive models, and let the initial policy be more strict than the*
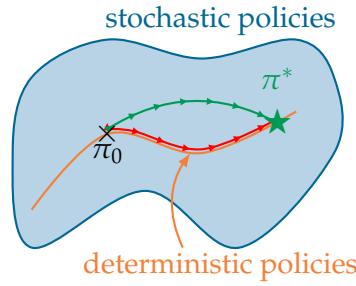
Figure 7.2 This figure illustrates how it can be impossible to find the optimal policy when the allowed set of policies is restricted to deterministic decision rules.

*optimal one, i.e., $W_0(\pi_0) \supsetneq W_0(\pi^*)$. If $\mathcal{A}$ is non-exploring on any i.i.d. sample $\mathcal{D} \sim P_{\pi_t}(X, Z, Y)$ with probability at least $1 - \delta_t$ for all $t \in \mathbb{N}$, then*

$$\Pr[\pi_T \neq \pi^*] > 1 - \sum_{t=0}^{T} \delta_t \quad \text{for any } T \in \mathbb{N}.$$

*Proof.* At each step we have

$$\Pr[\pi_t = \pi^*] = \Pr[W_0(\pi_t) = W_0(\pi^*)] \leq \Pr[W_0(\pi_t) \supset W_0(\pi^*)] \leq \delta_t + \Pr[\pi_{t-1} = \pi^*].$$

By the assumption that $\pi_0 \neq \pi^*$, we recursively get $P[\pi_t = \pi^*] \leq \sum_{i=0}^{t} \delta_i$ which concludes the proof.                                                                                               □

We can thus conclude that, for error based learning algorithms under no fairness constraints, learning within deterministic threshold policies is guaranteed to fail. Even though the optimal policy lies within the set of deterministic threshold policies, it cannot easily be approximated within this set starting from a suboptimal predictive model.

Figure 7.2 illustrates that, even though the optimal policy $\pi^*$ is deterministic, when starting from a deterministic initial policy $\pi_0$ (black cross), we cannot iteratively reach $\pi^*$ (green star) when updating solely within deterministic policies (red line following the orange line of deterministic policies). It is necessary to deploy stochastic policies (blue area) along the way to then be able to converge to the optimal policy (along the green line). We will introduce such "exploring policies" after our final impossibility result for error based learning algorithms.

**Corollary 3.** *A deterministic threshold policy $\pi \neq \pi^*$ with $\Pr[\pi(x,z) \neq y] = 0$ under P will fail to converge to $\pi^*$ under an error based learning algorithm for the underlying predictive model with probability* 1.

*Proof.* Since error based learning algorithms lead to non-exploring policies whenever

$$\sum_{(x,z,y)\in\mathcal{D}} \mathbf{1}[\pi(x,z)\neq y] = 0\,,$$

using the assumption $\Pr[\pi(x,z)\neq y] = 0$, we can use Proposition 6 with $\delta_t = 0$ for all $t\in\mathbb{N}$. □

While we have focused on deterministic threshold rules, our results readily generalize to *all* deterministic policies. An arbitrary deterministic policy $\pi$ can always be written as a threshold rule $\pi_Q$ as in eq. (7.3) with $Q(Y = 1\,|\,x,z) = \mathbf{1}[\pi(D = 1\,|\,x,z) = 1]$. To conclude, if we can only observe the outcomes of previous decisions taken by a deterministic initial policy $\pi_0$, these outcomes may be insufficient to find the (fair) deterministic decision rule that maximizes utility.

### 7.3.2 Stochastic policies

A naive but instructive way to overcome the undesirable behavior exhibited by deterministic policies discussed in the previous section, is to fully randomize initial decisions, i.e., $\pi_0(D = 1\,|\,x,z) = 1/2$ for all $x,z$. It readily follows from eq. (7.4) that then $P_{\pi_0} = P$. Hence, if the hypothesis class of predictive models $\mathcal{Q}$ is rich enough, we could learn the optimal policy $\pi^*$ from data gathered under $\pi_0$. In practice, fully randomized initial policies are unacceptable in terms of utility or unethical—it would entail releasing defendants by a coin flip. Fortunately, we will show next that full randomization is not required to learn the optimal policy. We only need to choose an initial policy $\pi_0$ such that $\pi_0(D = 1\,|\,x,z) > 0$ on any measurable subset of $\mathcal{X}\times\mathcal{Z}$ with positive probability under P, a requirement that is more acceptable for the decision maker in terms of initial utility. We refer to any policy with this property as an *exploring* policy. A policy $\pi$ is exploring, if and only if the true distribution P is absolutely continuous with respect to the induced distribution $P_\pi$. This means the data collection distribution must not ignore regions where the true distribution puts mass. We note that this condition does not strictly require randomness, but could be achieved by a pre-determined process, e.g., "$d = 1$ for every $n$-th decision". For an exploring policy $\pi_0$, we can compute the utility in eq. (7.1) and the group benefits for $z\in\{0,1\}$ via inverse propensity score weighting

$$
\begin{aligned}
u_{P_{\pi_0}}(\pi,\pi_0) &:= \mathbb{E}_{\substack{x,z,y\sim P_{\pi_0}\\ d\sim\pi(x,z)}}\left[\frac{d(y-c)}{\pi_0(D=1\,|\,x,z)}\right], \\
b^z_{P_{\pi_0}}(\pi,\pi_0) &:= \mathbb{E}_{\substack{x,z,y\sim P_{\pi_0}\\ d\sim\pi(x,z)}}\left[\frac{f(d,y)}{\pi_0(D=1\,|\,x,z)}\right].
\end{aligned}
\tag{7.7}
$$

Crucially, even though $u_{\mathrm{P}}(\pi) = u_{\mathrm{P}_{\pi_0}}(\pi, \pi_0)$ and $b_{\mathrm{P}}^z(\pi) = b_{\mathrm{P}_{\pi_0}}^z(\pi, \pi_0)$, the expectations are with respect to the induced distribution $\mathrm{P}_{\pi_0}(X, Z, Y)$, yielding the following positive result.

**Proposition 7.** *Let $\Pi$ be the set of exploring policies and let $\pi_0 \in \Pi \setminus \{\pi^*\}$. Then, the optimal objective value is*

$$v(\pi^*) = \sup_{\pi \in \Pi \setminus \{\pi^*\}} \left\{ u_{\mathrm{P}_{\pi_0}}(\pi, \pi_0) - \frac{\lambda}{2}(b_{\mathrm{P}_{\pi_0}}^0(\pi, \pi_0) - b_{\mathrm{P}_{\pi_0}}^1(\pi, \pi_0))^2 \right\}.$$

*Proof.* We already know that the supremum is upper bounded by $v(\pi^*)$, i.e., it suffices to construct a sequence of policies $\{\pi_n\}_{n \in \mathbb{N}_{>0}} \subset \Pi \setminus \{\pi^*\}$ such that $v(\pi_n) \to v(\pi^*)$ for $n \to \infty$. Using notation from the proof of Proposition 5, we define

$$\pi_n(D = 1 \mid x, z) := \begin{cases} 1 & \text{if } (x, z) \in W_1(\pi^*), \\ \frac{1}{n} & \text{otherwise}. \end{cases}$$

It is clear that $\pi_n$ is exploring, i.e., $\pi_n \in \Pi$, for all $n \in \mathbb{N}_{>0}$ as well as that $\pi_n \neq \pi^*$. To compute

$$\lim_{n \to \infty} v_{\mathrm{P}_{\pi_0}}(\pi_n, \pi_0) = \lim_{n \to \infty} \left( u_{\mathrm{P}_{\pi_0}}(\pi_n, \pi_0) - \frac{\lambda}{2}\left(b_{\mathrm{P}_{\pi_0}}^0(\pi_n, \pi_0) - b_{\mathrm{P}_{\pi_0}}^1(\pi_n, \pi_0)\right)^2 \right)$$

we look at the individual limits. For the utility we have

$$\lim_{n \to \infty} u_{\mathrm{P}_{\pi_0}}(\pi_n, \pi_0) = \lim_{n \to \infty} \mathbb{E}_{x,z,y \sim \mathrm{P}_{\pi_0}(X,Z,Y)} \left[ \frac{\pi_n(D = 1 \mid x, z)}{\pi_0(D = 1 \mid x, z)}(y - c) \right]$$

$$= \int_{W_1(\pi^*)} \frac{\mathrm{P}(Y = 1 \mid x, z) - c}{\pi_0(D = 1 \mid x, z)} \, d\mathrm{P}_{\pi_0}(x, z) +$$

$$\underbrace{\lim_{n \to \infty} \frac{1}{n} \int_{W_1(\pi^*)^{\complement}} \frac{\mathrm{P}(Y = 1 \mid x, z) - c}{\pi_0(D = 1 \mid x, z)} \, d\mathrm{P}_{\pi_0}(x, z)}_{=:C_1 \text{ with } |C_1| < \infty \text{ for any given exploring } \pi_0 \in \Pi}$$

$$= \int_{W_1(\pi^*)} (y - c) \, d\mathrm{P}(x, z, y) + \lim_{n \to \infty} \frac{C_1}{n}$$

$$= u_{\mathrm{P}}(\pi^*).$$

Similarly, for the benefit terms that are linear in both arguments, such as $f(d,y) = d$, we have for $z \in \{0,1\}$

$$
\begin{aligned}
\lim_{n\to\infty} b^z_{P_{\pi_0}}(\pi_n, \pi_0) &= \mathbb{E}_{x,y\sim P_{\pi_0}(X,Y\,|\,z)}\left[\frac{f(\pi_n(D=1\,|\,x,z),y)}{\pi_0(D=1\,|\,x,z)}\right] \\
&= \int_{W_1(\pi^*)} \frac{f(1, P(Y=1\,|\,x,z))}{\pi_0(D=1\,|\,x,z)}\, d\,P_{\pi_0}(x\,|\,z) + \\
&\quad \underbrace{\lim_{n\to\infty}\frac{1}{n}\int_{W_1(\pi^*)^{\complement}} \frac{f(1, P(Y=1\,|\,x,z))}{\pi_0(D=1\,|\,x,z)}\, d\,P_{\pi_0}(x\,|\,z)}_{=:C^z_2 \text{ with } |C^z_2|<\infty \text{ for any given exploring } \pi_0\in\Pi} \\
&= \int_{W_1(\pi^*)} f(1,y)\, d\,P(x,y\,|\,z) + \lim_{n\to\infty}\frac{C^z_2}{n} \\
&= b^z_P(\pi^*)\,.
\end{aligned}
$$

Because all the limits are finite, via the rules for sums and products of limits we get

$$
\begin{aligned}
\lim_{n\to\infty} v_{P_{\pi_0}}(\pi_n, \pi_0) &= \lim_{n\to\infty} u_{P_{\pi_0}}(\pi_n, \pi_0) - \frac{\lambda}{2}\big(\lim_{n\to\infty} b^0_{P_{\pi_0}}(\pi_n, \pi_0) - \lim_{n\to\infty} b^1_{P_{\pi_0}}(\pi_n, \pi_0)\big)^2 \\
&= u_P(\pi^*) - \frac{\lambda}{2}(b^0_P(\pi^*) - b^1_P(\pi^*))^2 \\
&= v_P(\pi^*)\,.
\end{aligned}
$$

$\square$

This shows that—unlike within deterministic threshold models—within exploring policies we can learn the optimal policy using only data from an induced distribution. Finally, we would like to highlight that not all exploring policies may be (equally) acceptable to society. For example, in lending scenarios without fairness constraints (i.e., $\lambda = 0$), it may appear wasteful to deny a loan with probability greater than zero to individuals who are believed to repay by the current model. In those cases, one may like to consider exploring policies that, given sufficient evidence, decide $d = 1$ deterministically, i.e., $\pi_0(D=1\,|\,x,z) = 1$ for some values of $x, z$. We will operationalize this notion in Section 7.4 as what we call the *semi-logistic policy*. Other settings, like the criminal justice system, call for a more general discussion about the ethics of non-deterministic decision making.

## 7.4   How to learn exploring policies

In this section, we exemplify Proposition 7 via a simple, yet practical, gradient-based algorithm to find the solution to eq. (7.5) within a (differentiable) parameterized class of exploring policies $\Pi(\Theta)$ using data gathered by a given, already deployed, exploring policy

---

**Algorithm 3** CONSEQUENTIALLEARNING: train a sequence of policies $\pi_{\boldsymbol{\theta}_t}$ of increasing $v_{\mathrm{P}}(\pi_{\boldsymbol{\theta}_t})$.

---

**Input:** cost $c$, time steps $T$, decisions $N$, iterations $M$, minibatch size $B$, penalty $\lambda$, learning rate $\alpha$

  1:  $\boldsymbol{\theta}_0 \leftarrow$ INITIALIZEPOLICY$()$

  2:  **for** $t = 0, \ldots, T-1$ **do**                                                 $\triangleright$ time steps

  3:       $\mathcal{D}^t \leftarrow$ COLLECTDATA$(\boldsymbol{\theta}_t, N)$

  4:       $\boldsymbol{\theta}_{t+1} \leftarrow$ UPDATEPOLICY$(\boldsymbol{\theta}_t, \mathcal{D}^t, M, B, \alpha)$

  5:  **return** $\{\pi_{\boldsymbol{\theta}_t}\}_{t=0}^{T}$

 

  6:  **function** COLLECTDATA$(\boldsymbol{\theta}, N)$

  7:       $\mathcal{D} \leftarrow \varnothing$

  8:       **for** $i = 1, \ldots, N$ **do**                                        $\triangleright$ $N$ decisions

  9:           $(x_i, z_i) \sim \mathrm{P}(X, Z)$ and $d_i \sim \pi_{\boldsymbol{\theta}}(x_i, z_i)$

10:           **if** $d_i = 1$ **then**                               $\triangleright$ positive decision

11:               $\mathcal{D} \leftarrow \mathcal{D} \cup \{(x_i, z_i, y_i)\}$ with $y_i \sim \mathrm{P}(Y \mid x_i, z_i)$

12:       **return** $\mathcal{D}$                                  $\triangleright$ data observed under $\pi_{\boldsymbol{\theta}}$

 

13:  **function** UPDATEPOLICY$(\boldsymbol{\theta}', \mathcal{D}, M, B, \alpha)$

14:       $\boldsymbol{\theta}^{(0)} \leftarrow \boldsymbol{\theta}'$

15:       **for** $j = 1, \ldots, M$ **do**                                    $\triangleright$ iterations

16:           $\mathcal{D}^{(j)} \leftarrow$ MINIBATCH$(\mathcal{D}, B)$                  $\triangleright$ sample minibatch

17:           $\nabla \leftarrow 0, \, n_j \leftarrow 0$

18:           **for** $(x, z, y) \in \mathcal{D}^{(j)}$ **do**                  $\triangleright$ accumulate gradients

19:               $d \sim \pi_{\boldsymbol{\theta}^{(j)}}(x, z)$

20:               **if** $d = 1$ **then**

21:                   $n_j \leftarrow n_j + 1$

22:                   $\nabla \leftarrow \nabla + \nabla_{\boldsymbol{\theta}} v(\pi_{\boldsymbol{\theta}}, \pi_{\boldsymbol{\theta}'})|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(j)}}$

23:           $\boldsymbol{\theta}^{(j+1)} \leftarrow \boldsymbol{\theta}^{(j)} + \alpha \frac{\nabla}{n_j}$

24:       **return** $\theta^M$

---

$\pi_0$. To this end, we consider a class of parameterized exploring policies $\Pi(\Theta)$ and we aim to find the policy $\pi_{\boldsymbol{\theta}^*} \in \Pi(\Theta)$ that solves the optimization problem in eq. (7.5).

We use stochastic gradient ascent (SGA) (Kiefer et al., 1952) to learn the parameters of the new policy, i.e.,

$$\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i + \alpha_i \nabla_{\boldsymbol{\theta}} v_{\mathrm{P}}(\pi_{\boldsymbol{\theta}})|_{\boldsymbol{\theta}=\boldsymbol{\theta}_i} \,,$$

where

$$\nabla_{\boldsymbol{\theta}} v_{\mathrm{P}}(\pi_{\boldsymbol{\theta}}) = \nabla_{\boldsymbol{\theta}} u_{\mathrm{P}}(\pi_{\boldsymbol{\theta}}) - \lambda(b_{\mathrm{P}}^0(\pi_{\boldsymbol{\theta}}) - b_{\mathrm{P}}^1(\pi_{\boldsymbol{\theta}}))(\nabla_{\boldsymbol{\theta}} b_{\mathrm{P}}^0(\pi_{\boldsymbol{\theta}}) - \nabla_{\boldsymbol{\theta}} b_{\mathrm{P}}^1(\pi_{\boldsymbol{\theta}})) \,,$$

and $\alpha_i > 0$ is the learning rate at step $i \in \mathbb{N}$. With the reweighting from eq. (7.7) and the log-derivative trick (Williams, 1992), we can compute the gradient of the utility and the benefits as

$$
\begin{aligned}
\nabla_{\boldsymbol{\theta}} u_{\mathrm{P}}(\pi_{\boldsymbol{\theta}}) &= \mathbb{E}_{\substack{x,z,y \sim \mathrm{P}_{\pi_0} \\ d \sim \pi_{\boldsymbol{\theta}}(x,z)}} \left[ \frac{d\,(y-c)\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}}{\pi_0(D=1 \mid x,z)} \right], \\
\nabla_{\boldsymbol{\theta}} b_{\mathrm{P}}^z(\pi_{\boldsymbol{\theta}}) &= \mathbb{E}_{\substack{x,z,y \sim \mathrm{P}_{\pi_0} \\ d \sim \pi_{\boldsymbol{\theta}}(x,z)}} \left[ \frac{f(d,y)\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}}{\pi_0(D=1 \mid x,z)} \right],
\end{aligned}
\tag{7.8}
$$

where $\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}} := \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(D \mid x,z)$ is the score function (Hyvärinen, 2005). Thus, our implementation resembles a REINFORCE algorithm with horizon one.

Note that we can obtain an expression for $\nabla_{\boldsymbol{\theta}_t} v_{\mathrm{P}}(\pi_{\boldsymbol{\theta}_t})$ by simply replacing $\pi_0$ with $\pi_{\boldsymbol{\theta}_{t-1}}$ in eq. (7.8). Thus we can estimate the gradient with samples $(x_i, z_i, y_i)$ from the distribution $\mathrm{P}_{\pi_{t-1}}$ induced by the previous policy $\pi_{t-1}$, and sample the decisions from the policy under consideration $d_i \sim \pi_{\boldsymbol{\theta}_t}$. This yields an unbiased finite sample Monte-Carlo estimator for the gradients

$$
\begin{aligned}
\nabla_{\boldsymbol{\theta}_t} u(\pi_{\boldsymbol{\theta}_t}, \pi_{\boldsymbol{\theta}_{t-1}}) &\approx \frac{1}{n_{t-1}} \sum_{i=1}^{n_{t-1}} \frac{d_i(y_i - c)}{\pi_{\boldsymbol{\theta}_{t-1}}(D=1 \mid x_i, z_i)} \nabla_{\boldsymbol{\theta}_t} \log \pi_{\boldsymbol{\theta}_t}(D = d_i \mid x_i, z_i) \,, \\
\nabla_{\boldsymbol{\theta}_t} b^z(\pi_{\boldsymbol{\theta}_t}, \pi_{\boldsymbol{\theta}_{t-1}}) &\approx \frac{1}{n_{t-1}} \sum_{i=1}^{n_{t-1}} \frac{f(d_i, y_i)}{\pi_{\boldsymbol{\theta}_{t-1}}(D=1 \mid x_i, z_i)} \nabla_{\boldsymbol{\theta}_t} \log \pi_{\boldsymbol{\theta}_t}(D = d_i \mid x_i, z_i) \,,
\end{aligned}
\tag{7.9}
$$

where $n_{t-1}$ is the number of positive decisions taken by $\pi_{\boldsymbol{\theta}_{t-1}}$. Here, it is important to notice that, while the decisions by $\pi_{\boldsymbol{\theta}_{t-1}}$ were actually taken and, as a result, (feature and label) data was gathered under $\pi_{\boldsymbol{\theta}_{t-1}}$, the decisions $d_i \sim \pi_{\boldsymbol{\theta}_t}$ are just sampled to implement SGA. The overall policy learning process is summarized in Algorithm 3, where MINIBATCH($\mathcal{D}, B$) samples a minibatch of size $B$ from the dataset $\mathcal{D}$ and INITIALIZEPOLICY() initializes the policy parameters.

Unfortunately, the above procedure has two main drawbacks. First, it may require an abundance of data from $\mathrm{P}_{\pi_0}$, which can be unacceptable in terms of utility if $\pi_0$ is far from optimal. Second, if $\pi_0(D = 1 \mid x, z)$ is small in a region where $\pi_{\boldsymbol{\theta}}$ often takes positive decisions, one may expect that an empirical estimate of the above gradient will have high
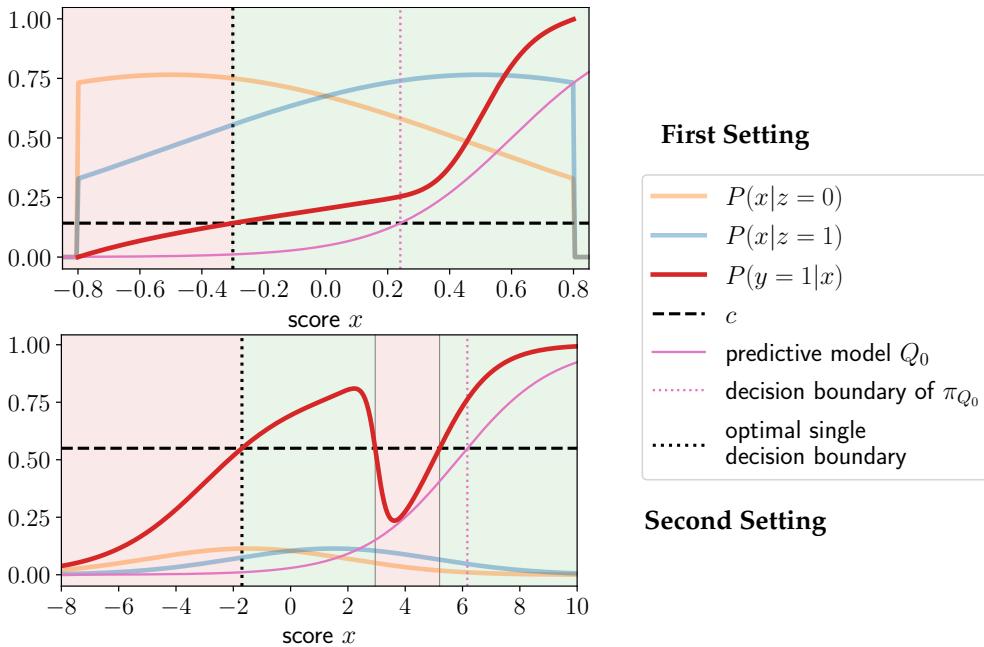
Figure 7.3 Two synthetic settings. In red, we show $P(Y = 1 \mid x)$, where the score $x$ is drawn from different distributions for the two groups (blue/orange). For given $c$ (black, dashed), the optimal policy decides $d = 1$ ($d = 0$) in the shaded green (red) regions. The vertical black, dotted line shows the best policy achievable with a single threshold on $x$. In pink, we show a possible imperfect logistic predictive model and its corresponding (suboptimal) threshold in $x$.

variance, due to similar arguments as in weighted inverse propensity scoring (Sutton & Barto, 1998). On the other hand, in most practical applications updating the model after every single decision is impractical. Typically, a fixed model will be deployed for a certain period, before it is updated using the data collected within this period. This is also a natural mode of operation for predictive models in real-world applications.

To overcome these drawbacks, we build two types of sequences of policies $\{\pi_{\theta_t}\}_{t=0}^T$:

a) the *iterative sequence* $\pi_{t+1} := \mathcal{A}(\pi_t, \mathcal{D}^t)$ with $\mathcal{D}^t \sim P_{\pi_t}(X, Z, Y)$, where only the data gathered by the immediately previous policy are used to update the current policy;

b) the *aggregated sequence* $\pi_{t+1} := \mathcal{A}(\pi_t, \bigcup_{i=0}^t \mathcal{D}^i)$ with $\mathcal{D}^i \sim P_{\pi_i}(X, Z, Y)$, where the data gathered by all previous policies are used to update the current policy.

**Remarks.** Note that in Algorithm 3 we learn each policy $\pi_t$ only using data from the previous policy $\pi_{t-1}$. This may readily be generalized to a mix of multiple previous policies $\pi_{t'}$ in eq. (7.9). Averaging multiple gradient estimators for several $t' < t$ is again an unbiased gradient estimator. To reduce variance, in practice one may consider recent policies $\pi_{t'}$ most similar to $\pi_t$.

The way in which we use weighted sampling to estimate the above gradients closely relates to the concept of weighted inverse propensity scoring (wIPS), commonly used in counterfactual learning (Bottou et al., 2013; Swaminathan & Joachims, 2015a), off-policy reinforcement learning (Sutton & Barto, 1998), and contextual bandits (Langford et al., 2008). However, a key difference is that, in wIPS, the labels $y$ are always observed. As an example, in the case of counterfactual learning one may interpret $\pi_0(x, z)$ in eq. (7.4) as a treatment assignment mechanism in a randomized controlled trial. Under this interpretation, the two most prominent differences with respect to the literature become apparent. First, we do not observe outcomes in the control group. Second, in observational studies for treatment effect estimation (Rubin, 2005), one usually estimates the direct causal effect of $d$ on $y$, i.e., $P(Y \mid do(D = d'), x, z)$, in the presence of confounders $x, z$ that affect both $d$ and $y$. This could be evaluated in a (partially) randomized controlled trial, where wIPS also comes in naturally (Pearl, 2009). In contrast, in our setting, the true label $y$ is independent of the decision $d$ and we estimate the conditional $P(Y \mid x, z)$ using data from the induced distribution $P_{\pi_0}(X, Z) \propto P(X, Z)\pi_0(D = 1 \mid x, z)$. With exploring policies, we obtain indirect access to the true data distribution $P(x, z)$ (positivity), and thus to an unbiased estimator of the conditional distribution $P(Y \mid x, z)$ (consistency).

Despite this difference, we believe that recent advances to reduce the variance of the gradients in weighted inverse propensity scoring, such as clipped-wIPS (Bottou et al., 2013), self-normalized estimator (Swaminathan & Joachims, 2015b), or doubly robust estimators (Dudík et al., 2011), may also be applicable to our setting.

Finally, we opt for the simple SGA approach on (semi-)logistic policies over, e.g., contextual bandits algorithms, because it provides a direct and fairer comparison with commonly used prediction based decision policies (e.g., logistic regression), also often trained via SGA.

While our algorithm works for any differentiable class of exploring policies, here we consider two examples of exploring policy classes in particular.

**Logistic policy.** Our first concrete parameterization of $\pi_{\boldsymbol{\theta}}$, a *logistic policy* is given by

$$\pi_{\boldsymbol{\theta}}(D = 1 \mid x, z) = \sigma(\boldsymbol{\phi}(x, z)^{\top}\boldsymbol{\theta}) \in (0, 1),$$

where $\sigma(a) := \frac{1}{1+\exp(-a)}$ is the logistic function, $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^m$ are the model parameters, and $\boldsymbol{\phi} : \mathbb{R}^d \times \{0, 1\} \to \mathbb{R}^m$ is a fixed feature map. Note that any logistic policy is an exploring policy and we can analytically compute its score function $\nabla_{\boldsymbol{\theta}_t} \log \pi_{\boldsymbol{\theta}_t}(D = 1 \mid x, z)$ as

$$\nabla_{\boldsymbol{\theta}_t} \log(\sigma(\boldsymbol{\phi}_i^{\top}\boldsymbol{\theta}_t)) = \frac{\boldsymbol{\phi}_i}{1 + e^{\boldsymbol{\phi}_i^{\top}\boldsymbol{\theta}_t}} \in \mathbb{R}^m,$$

where $\boldsymbol{\phi}_i := \boldsymbol{\phi}(x_i, z_i)$. Using this expression, we can rewrite the empirical estimator for the gradient in eq. (7.9)

$$\nabla_{\boldsymbol{\theta}_t} u(\pi_{\boldsymbol{\theta}_t}, \pi_{\boldsymbol{\theta}_{t-1}}) \approx \frac{1}{n_{t-1}} \sum_{i=1}^{n_{t-1}} \frac{1 + e^{-\boldsymbol{\phi}_i^\top \boldsymbol{\theta}_{t-1}}}{1 + e^{\boldsymbol{\phi}_i^\top \boldsymbol{\theta}_t}} d_i (y_i - c) \boldsymbol{\phi}_i,$$

$$\nabla_{\boldsymbol{\theta}_t} b^z(\pi_{\boldsymbol{\theta}_t}, \pi_{\boldsymbol{\theta}_{t-1}}) \approx \frac{1}{n_{t-1}} \sum_{i=1}^{n_{t-1}} \frac{1 + e^{-\boldsymbol{\phi}_i^\top \boldsymbol{\theta}_{t-1}}}{1 + e^{\boldsymbol{\phi}_i^\top \boldsymbol{\theta}_t}} f(d_i, y_i) \boldsymbol{\phi}_i.$$

Given the above expression, we have all the necessary ingredients to implement Algorithm 3.

**Semi-logistic policy.** As discussed in the previous section, randomizing decisions may be questionable in certain practical scenarios. For example, in loan decisions, it may appear wasteful for the bank and contestable for the applicant to deny a loan with probability greater than zero to individuals who are believed to repay by the current model. In those cases, one may consider the following modification of the logistic policy, which we refer to as *semi-logistic policy*:

$$\tilde{\pi}_{\boldsymbol{\theta}}(D = 1 \mid x, z) = \begin{cases} 1 & \text{if } \boldsymbol{\phi}(x, z)^\top \boldsymbol{\theta} \geq 0, \\ \sigma(\boldsymbol{\phi}(x, z)^\top \boldsymbol{\theta}) & \text{if } \boldsymbol{\phi}(x, z)^\top \boldsymbol{\theta} < 0. \end{cases}$$

Similarly as in the logistic policy, we can compute the score function analytically as:

$$\nabla_{\boldsymbol{\theta}} \log \tilde{\pi}_{\boldsymbol{\theta}}(D = 1 \mid x, z) = \frac{\boldsymbol{\phi}(x, z)}{1 + e^{\boldsymbol{\phi}(x, z)^\top \boldsymbol{\theta}}} \mathbf{1}[\boldsymbol{\phi}(x, z)^\top \boldsymbol{\theta} < 0],$$

and use this expression to compute an unbiased estimator for the gradient in eq. (7.9) as:

$$\nabla_{\boldsymbol{\theta}_t} u(\pi_{\boldsymbol{\theta}_t}, \pi_{\boldsymbol{\theta}_{t-1}}) \approx \frac{1}{n_{t-1}} \sum_{\substack{i=1 \\ \boldsymbol{\phi}_i^\top \boldsymbol{\theta}_t < 0}}^{n_{t-1}} \frac{d_i (y_i - c) \boldsymbol{\phi}_i}{1 + e^{\boldsymbol{\phi}_i^\top \boldsymbol{\theta}_t}} \times \begin{cases} 1 & \text{if } \boldsymbol{\phi}_i^\top \boldsymbol{\theta}_{t-1} \geq 0, \\ (1 + e^{-\boldsymbol{\phi}_i^\top \boldsymbol{\theta}_{t-1}}) & \text{if } \boldsymbol{\phi}_i^\top \boldsymbol{\theta}_{t-1} < 0. \end{cases}$$

$$\nabla_{\boldsymbol{\theta}_t} b^z(\pi_{\boldsymbol{\theta}_t}, \pi_{\boldsymbol{\theta}_{t-1}}) \approx \frac{1}{n_{t-1}} \sum_{\substack{i=1 \\ \boldsymbol{\phi}_i^\top \boldsymbol{\theta}_t < 0}}^{n_{t-1}} \frac{f(d_i, y_i) \boldsymbol{\phi}_i}{1 + e^{\boldsymbol{\phi}_i^\top \boldsymbol{\theta}_t}} \times \begin{cases} 1 & \text{if } \boldsymbol{\phi}_i^\top \boldsymbol{\theta}_{t-1} \geq 0, \\ (1 + e^{-\boldsymbol{\phi}_i^\top \boldsymbol{\theta}_{t-1}}) & \text{if } \boldsymbol{\phi}_i^\top \boldsymbol{\theta}_{t-1} < 0. \end{cases}$$

Note that the semi-logistic policy is an exploring policy and thus satisfies the assumptions of Proposition 7. Finally, in all our experiments, we directly work with the available features $x$ as inputs and add a constant offset, i.e., $\phi(x, z) = (1, x)$.

## 7.5  Experiments

We learn a sequence of policies $\{\pi_{\theta_t}\}_{t=1}^{T}$ using the following strategies:

**Optimal:** decisions are taken by the optimal deterministic threshold rule $\pi^*$ given by eq. (7.2), i.e., $\pi_t = \pi^*$ for all $t$. It can only be computed when the ground truth conditional $\mathrm{P}(Y \mid x, z)$ is known.

**Deterministic:** decisions are taken by deterministic threshold policies $\pi_t = \pi_{Q_t}$, where $Q_t$ are logistic predictive models maximizing label likelihood trained either in an iterative or aggregate sequence.

**Logistic:** decisions are taken by logistic policies $\pi_t = \pi_{\theta_t}$ trained via Algorithm 3 either in an iterative or aggregate sequence.

**Semi-logistic:** decisions are taken by semi-logistic policies $\tilde{\pi}_t = \tilde{\pi}_{\theta_t}$ trained via Algorithm 3 either in an iterative or aggregate sequence.

It is crucial that while each of the above methods decides over the same set of proposed $\{(x_i, z_i)\}_{i=1}^{N}$ at each time step $t$, depending on their decisions, they may collect labels for differing subsets and thus receive different amounts of new training data. During learning, we record the following metrics:[4]

**Utility:** the utility $u_\mathrm{P}(\pi_t)$ achieved by the current policy $\pi_t$ estimated empirically on a held-out dataset, the *test set*, sampled i.i.d. from the ground truth distribution $\mathrm{P}(X, Z, Y)$. This is the utility that the decision maker would obtain if they deployed the current policy $\pi_t$ at large in the population.

**Effective utility:** the utility realized during the learning process up to time $t$, i.e.,

$$\hat{u}(t) = \frac{1}{N \cdot t} \sum_{t' \leq t} \sum_{(x_i, z_i, y_i) \in \mathcal{D}^{t'}} (y_i - c),$$

where $\mathcal{D}^{t'}$ are the data in which the policy $\pi_{t'}(x_i, z_i)$ took positive decisions $d_i = 1$ and $N$ is the number of considered examples at each time step $t$. This is the utility accumulated by the decision maker while learning better policies.

**Fairness:** the difference in group benefits between sensitive groups $\Delta b_\mathrm{P}(\pi) := b_\mathrm{P}^0(\pi) - b_\mathrm{P}^1(\pi)$ for disparate impact: $f(d, y) = d$. A policy fully satisfies the chosen criterion if and only if $\Delta b_\mathrm{P}(\pi) = 0$. Again, we estimate fairness empirically on the test set and thus measure the level of fairness $\pi_t$ would achieve in the entire population.

Detailed parameter settings for the experiments are explained in Section B.3 in Appendix B.

---

[4]For readability we only show medians over 30 runs. Figures with 25 and 75 percentiles are in Appendix B.
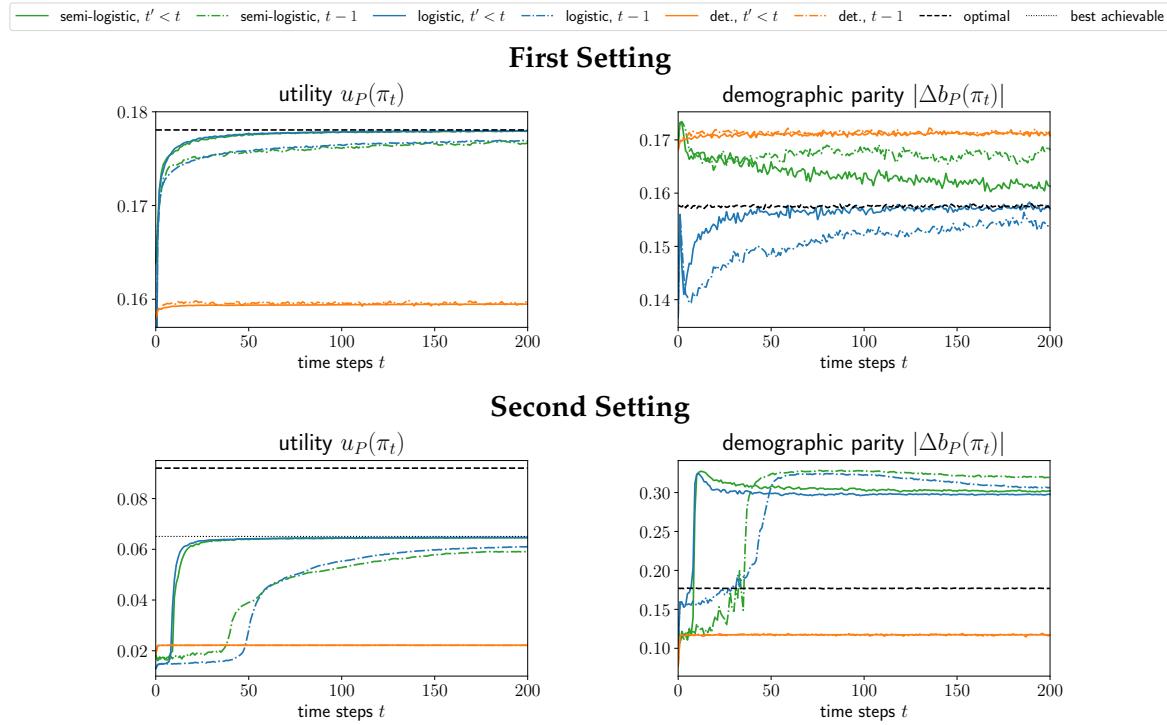
Figure 7.4 Utility and demographic parity in the synthetic settings of Figure 7.3 without enforcing fairness constraints, i.e., $\lambda = 0$.

### 7.5.1   Experiments on synthetic data

We assume that there is a single non-sensitive feature $x \in \mathbb{R}$ per individual—similar to the lending example in Section 7.3—and a sensitive attribute $z \in \{0, 1\}$. While $P(X \mid z = 0) \neq P(X \mid z = 1)$, in our experiments the policies only take $x$ as input, and *not* the sensitive attribute, which is only used for the fairness constraints. We consider two different settings, illustrated in Figure 7.3, where $z \sim \text{Ber}(0.5)$ and the distributions over $x$ differ for the two groups, see Appendix B. In the first setting, the conditional probability $P(Y = 1 \mid x)$ is strictly monotonic in the score and does not depend on $z$, but is not well calibrated, i.e., not directly proportional to $x$. In the second setting, the conditional probability $P(Y = 1 \mid x)$ crosses the cost threshold $c$ multiple times resulting in two disjoint intervals for which the optimal decision is $d = 1$ (green areas).

Figure 7.4 summarizes the results for $\lambda = 0$, i.e., without explicitly enforcing fairness constraints. Our method outperforms prediction based deterministic threshold rules in terms of utility in both settings. This can be easily understood from the evolution of policies illustrated in Figure B.2 in Appendix B. In the first setting, exploring policies locate the optimal decision boundary, whereas the deterministic threshold rules get stuck, even though $P(Y = 1 \mid x)$ is monotonic in $x$. In the second setting, our methods explore more and eventually identify the best single threshold at the black vertical dotted line in Figure 7.3. In
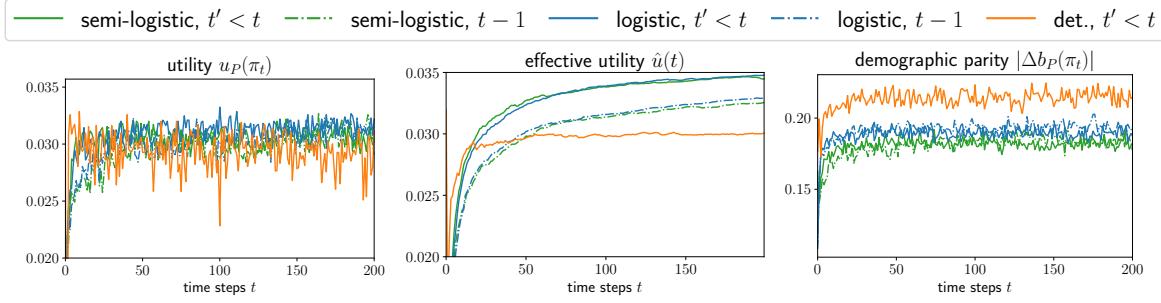
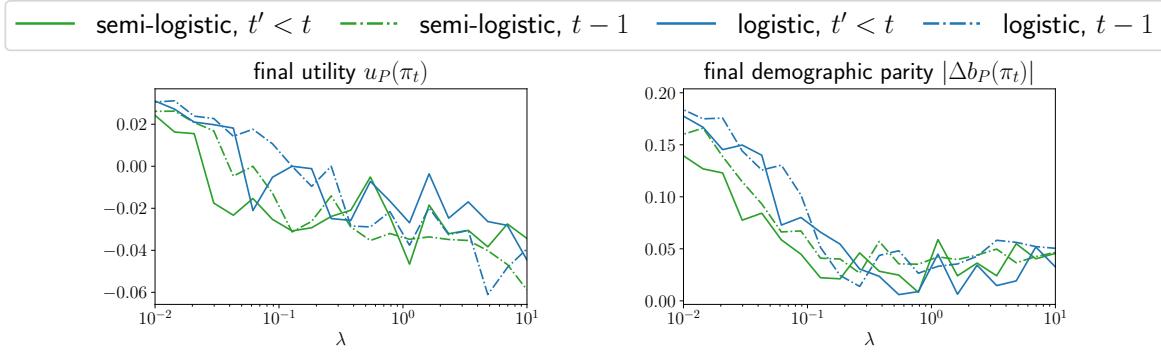Figure 7.5 Training progress on COMPAS data for $\lambda = 0$, i.e., without fairness constraints.



Figure 7.6 Fairness evaluation on COMPAS data for the final ($t = 200$) policy as a function of $\lambda$ for demographic parity. All quantities are estimated on the held-out set.

contrast, non-exploring deterministic threshold rules converge to a suboptimal threshold at $x \approx 5$, ignoring the left green region.

In the first setting, we also observe that the suboptimal predictive models amplify unfairness beyond the levels exhibited by the optimal policy. For our approach, levels of unfairness are comparable to or even below those of the optimal policy. The second setting shows that depending on the ground truth distribution, higher utility can be directly linked to larger fairness violations. In such cases, our approach allows to explicitly control for fairness. Results on utility and demographic parity under fairness constraints with different $\lambda$ are shown in Figure B.3 in Appendix B. In essence, $\lambda$ trades off utility and fairness violations to the point of perfect fairness in the ground truth distribution.

### 7.5.2 Experiments on real data

Here, we use the COMPAS recidivism dataset compiled by ProPublica (Angwin et al., 2016), which comprises of information about criminal offenders screened through the COMPAS tool in Broward County, Florida during 2013-2014. For each offender, the dataset contains a set of demographic features, the criminal history, and the risk score assigned by COMPAS. Moreover, ProPublica collected whether or not these individuals were rearrested within two

years of the screening. In our experiments, $z \in \{0, 1\}$ indicates whether individuals were identified "white", $y$ indicates rearrest, and $d \sim \pi(x, z)$ determines whether an individual is let out on parole. Again, $z$ is not used as an input. We use 80% of the data for training, where at each step $t$, we sample (with replacement) $N$ individuals, and the remaining 20% as a held-out set to evaluate each learned policy in the population of interest.

We first summarize the results for $\lambda = 0$, i.e., without fairness constraints in Figure 7.5. A slight initial utility advantage of the deterministic threshold rule is quickly overcome by our exploring policies. This is best seen when looking at *effective utility*, the average utility accumulated by the decision maker on training data up to time $t$, for which our strategies dominate after $t = 100$. Hence, early exploration not only pays off to eventually be able to take better decisions, but also reaps higher profit during training. Moreover, all strategies based on exploring policies consistently achieve lower violations of demographic parity than the deterministic threshold rules. In summary, even without fairness constraints, i.e., in a pure utility maximization setting, exploring policies achieve higher utility and simultaneously reduce unfairness compared to deterministic threshold rules.

In Figure 7.6, we show how utility and demographic parity of the final policy $\pi_{t=200}$ changes as a function of $\lambda$ when constraining demographic parity. As expected, while we are able to achieve perfect demographic parity, this comes with a drop in utility. All remaining metrics under fairness constraints are shown in Figure B.5 in Appendix B. Finally, two remarks are in order. First, for real-world data we cannot evaluate the optimal policy and do not expect it to reside in our model class. However, even when logistic models do not perfectly capture the conditional $P(Y = 1 \mid x)$, our comparisons here are "fair" in that all strategies have equal modeling capacity. Second, we take the COMPAS dataset as our (empirical) ground truth distribution even though it likely also suffered from selective labels. To learn about the real distribution underlying the dataset, we would need to actually deploy our strategy.

## 7.6 Conclusion

In this chapter, we have analyzed consequential decision making using imperfect predictive models, which are learned from data gathered by potentially biased historical decisions. First, we have articulated how this approach fails to optimize utility when starting with a non-optimal deterministic policy. Next, we have presented how directly learning to decide with exploring policies avoids this failure mode while respecting a common fairness constraint. Finally, we have introduced and evaluated a simple, yet practical gradient-based algorithm to learn fair exploring policies.

Unlike most previous work on fairness in machine learning, which phrases decision making directly as a prediction problem, we argue for a shift from "learning to predict" to "learning to

decide". In particular, we propose to not simply equate decisions with predictions obtained directly from limited available data, but to remain cognizant of how made decisions can affect and interfere with future data collection by continued exploration. Not only does this lead to improved fairness in this context, but it also establishes connections to other areas such as counterfactual inference, reinforcement learning and contextual bandits. Within reinforcement learning, it would be interesting to move beyond a static distribution P by incorporating feedback from decisions or non-static externalities. Moreover, since we have shown how shifting focus from learning predictions to learning decisions requires exploration, we hope to stimulate future research on how to explore ethically in different domains.

The crucial difference between mere predictions and actual decisions has been further highlighted in concurrent and later work (Kleinberg et al., 2018a; Rambachan et al., 2020). There, the authors argue for a social planner taking decisions to maximize a social welfare function that may also include fairness preferences. If the social planner has access to predictions from machine learning systems, it is optimal to keep the decision and prediction stage completely separate and train the machine learning pipeline to maximize predictive performance without any additional constraints. This suggests that machine learning systems should indeed be used for prediction alone, but there is a separate optimization problem in deriving decisions from these predictions. It is this second stage in which we must account for fairness considerations. However, these approaches do not consider selective labels. In general, as soon as the observed data depends on the chosen decision rule, the prediction and decision stage may not be easily decoupled anymore. We believe the importance to clearly distinguish between predictions and decisions on a technical level reflects our arguments in Chapter 1 that discrimination is usually a systemic cultural issue. Merely adjusting or improving predictive models does not suffice to make beneficial decisions.

# 8

# Conclusion

In this thesis we developed methodology and tools to overcome certain restrictive assumptions commonly made in fair machine learning for consequential decision making. Thereby we broaden the applicability of such systems and bring them closer to application. This chapter summarizes these contributions, draws conclusions, and suggests directions for future work. At the end, we again put our contributions into the broader context of socially beneficial machine learning.

## 8.1 Summary of contributions

In Chapter 4 we applied a causal lens to see beyond observational group matching criteria and put the causal data generating mechanisms in the center of the discussion. This allowed us to overcome fundamental limitations of observational fairness criteria and revealed some subtle—yet crucial—aspects relating to the meaning we assign to protected attributes. The main constituents of our conceptual framework are resolving variables and proxy variables, which play a dual role in defining types of discrimination from a skeptic or benevolent viewpoint. We developed a practical procedure to remove proxy discrimination given the structural equation model is of a certain functional form, and analyzed a similar approach for unresolved discrimination. Whilst not always feasible, the causal approach naturally creates an incentive to scrutinize the data more closely and work out plausible assumptions to be discussed alongside any conclusions regarding fairness. In particular, our framework assumes that we have access to the true underlying causal graph.

Since this is a strong assumption in itself, in Chapter 5 we also developed tools to analyze the sensitivity of counterfactually fair classifiers to potential unmeasured confounding. A grid-based approach was introduced for confounding—modeled as unobserved correlation between error terms—between two variables, and an optimization-based approach for confounding in the general multivariate case. These methods for sensitivity analysis are a step

towards extending the applicability of causal methods to assess and mitigate discrimination in real-world settings.

In Chapter 6 we addressed a dilemma of traditional approaches to fairness: in order to enforce fairness, sensitive attributes must be examined; yet in many situations, users may refuse to provide them, or modelers may be prohibited from collecting them. We adopted and improved recent methods in secure multi-party computation and demonstrated that it is practical on real-world datasets to: a) certify and sign a model as fair; b) learn a fair model; and c) verify that a fair-certified model has indeed been used; all without users ever having to disclose their sensitive attributes in the clear to anyone and without the modeler having to disclose their model to any other party. These techniques empower users to retain control over data and modelers to preserve their intellectual property rights while still being able to train fair models. Thereby, our contributions are an important step towards jointly addressing concerns in privacy, algorithmic fairness and accountability.

Chapter 7 scrutinized the traditional framework of predictive modeling for consequential decisions. Under certain conditions, predictive models can lead to optimal decisions, even in terms of fairness. However, we argued that one of the required assumptions, namely that we have access to i.i.d. labeled data, is commonly violated in practice. Often, outcomes only exist when a certain decision is made. We have shown how deterministic threshold rules from predictive models fail in this scenario. Further, we introduced an approach of directly learning decisions with exploring policies instead and demonstrated its efficacy on synthetic and real-world data via a simple, gradient-based implementation. Our results strongly point towards the importance of distinguishing predictive modeling and decision making in consequential settings. It thereby contributes both conceptually as well as methodologically to our understanding of the applicability of algorithmic fairness techniques.

## 8.2    Conclusions and directions for future work

This thesis is a first step towards being more cognizant of the confines of unrealistic technical assumptions underlying traditional approaches to fair machine learning and revealing paths to overcome these limitations. *One of the key conclusions from our findings is that as machine learning researchers, we must expand our horizon in terms of what we factor into the modeling process.*

The causal framework introduced in Chapter 4 suggests deeper scrutiny of how the data came about, forcing us to bring the data collection phase to the center of our attention. It requires us to hypothesize or ideally even validate assumptions about why the data are what the data are. Making such assumptions is difficult and typically requires a much deeper understanding of the subject matter than machine learning researchers or engineers can

be expected to bring to the table. However, from our findings we conclude that instead of avoiding the discomfort of making and communicating assumptions all together, to build socially beneficial systems we must start to become comfortable with potentially untestable modeling assumptions. As a positive side effect, we hope that such discomfort will encourage researchers and engineers to initiate a dialog with domain experts to seek further backing or refutation of their assumptions. From a technical angle, Chapter 5 also develops methodological tools to alleviate the uncertainty around potentially wrong assumptions.

Looking forward, we believe that causal modeling of societal interactions between humans and machines will play a crucial part in building fair, accountable and explainable systems. In future work, concrete technical advances could be to extend our procedures for avoiding causal measures of discrimination to larger function classes and allowing for a more fine-grained division into fair and unfair pathways. More broadly, we will also need to improve our understanding of the ontological stability of various social constructs that are relevant for fair machine learning systems, such as protected attributes. Identifying ontologically stable concepts as well as robust causal mechanisms in socioalgorithmic systems seems to be a challenge that can only be tackled effectively in an interdisciplinary endeavor. Since causal insights can only be obtained when we are willing to make assumptions (or provide existing knowledge), an interesting and important direction for future work is also to develop more flexible tools for sensitivity analysis. Concretely, how can we scale our method for unobserved confounding efficiently to large systems, how can we formulate restrictions on observed confounding more flexibly, or how can we extend the formalism to also allow for structural misspecifications of edges in the causal graph? These tools can uncover when there is "wiggle room" for causal assumptions, which alleviates the cognitive burden of making them confidently, which may otherwise lead to inaction.

In Chapter 6, it became apparent that the modes of practical data availability and transportability must also be taken into account by machine learning researchers and data scientists. As the collection, storage and transport of data become a key driver of economic success as well as subject to data privacy regulations, we must move away from the common idea that "all data is always available everywhere". Instead, data availability must become a first-class citizen in the modeling process, especially for fairness sensitive applications. Therefore, we believe that further adaptations of existing as well as developments of new cryptographic techniques for privacy preserving machine learning will be fruitful directions for future work. Specifically, it would be interesting to combine our methods from Chapter 6, which ensure data security in transport and at rest, with privacy guarantees from differential privacy. More broadly, can we leverage ideas from (partially homomorphic encryption), secure multi-party computation, or zero-knowledge proofs to hold modelers accountable by proving publicly that their models satisfy given fairness properties without having to

disclose their intellectual property? Can we practically scale such applications to large models and datasets?

Finally, just like modeling should take into account data collection, it must also consider the future impact of decisions. Chapter 7 demonstrated the importance of including downstream consequences in the selective labels setting, where future observations depend on the current model. As pointed out in the introductory chapters, this is but one possibility of a broader phenomenon sometimes referred to as *performativity* (Perdomo et al., 2020). While in our setup we assumed the ground truth distribution to be static, an interesting direction for future work would be to also allow the ground truth distribution to change over time depending on how decisions are taken. We currently lack reliable models of how humans react to and adjust their behavior when facing automated decision systems. An important question for the future will be how to navigate the fine line between building models that are robust to individual attempts of cheating the system while broadly incentivizing desirable behavior. This hints at the intersection between mechanism design, game theory, and machine learning to be a promising area for further advances in socially beneficial automated decision making.

*Tied together, our contributions on relaxing specific assumptions at different stages of the data science loop all highlight the importance of taking into account the entire interacting socialgorithmic system from the get-go.*

As a second insight, we note that our solutions predominantly rely upon combining and advancing recent ideas and concepts from *different fields of research*. Chapter 4 raises questions in the intersection of sociology, philosophy and causal inference regarding ontological stability of socially constructed quantities. In Chapter 5, we combined advances in efficient gradient-based optimization and modern, highly tuned auto differentiation frameworks to perform sensitivity analysis without one-size-fits-all identifiability assumptions. In Chapter 6, we exploited recent theoretical and implementation advances in cryptography, particularly secure multi-party computation, to reconcile fairness with privacy and accountability concerns. Finally, in Chapter 7, we borrowed ideas from bandit and reinforcement learning to uncover inadequate assumptions about the interaction between predictions and consequential decisions. Adapting concepts around exploration versus exploitation allowed us to accommodate downstream effects of decisions in our models, and raised further interesting questions about the relation between exploration, stochasticity and ethical decision-making.

*Hence, we conclude that when it comes to fair machine learning systems, bridging the gap between theory and practice appears to require a truly interdisciplinary approach that will benefit from considering advances in adjacent as well as more distant fields of research.*

## 8.3   Final remarks

Let us end by circling back to the introduction, reflecting on the goals we set ourselves and to what extent we achieved them. There are myriads of ways in which injustice and discrimination manifest themselves in our society. We have briefly elaborated on some of the broader challenges of achieving fairness when making consequential decisions in Chapter 1. While we repeatedly emphasized that machine learning certainly cannot solve these problems single-handedly, we ended Section 2.4 on an optimistic note. Hopefully, the subsequently developed concepts and techniques support this optimistic view in showing that there are paths to extend the narrow, unrealistic formalization of "fair consequential decisions", perhaps even to the point where their application in the real world does more good than harm. That said, let us now re-emerge from these formally defined settings and acknowledge that racism, gender discrimination, xenophobia, and stereotyping are still deeply ingrained in our society. While data-driven decision systems play a role, I appeal to all of us to listen to and learn from the ones suffering the hardship so we can give them agency. Smarter algorithms are no substitute for empathy, mutual respect, and activism.

# References

Abebe, R., Barocas, S., Kleinberg, J., Levy, K., Raghavan, M., and Robinson, D. G. Roles for computing in social change. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pp. 252–260, 2020.

Agarwal, A., Hsu, D., Kale, S., Langford, J., Li, L., and Schapire, R. Taming the monster: A fast and simple algorithm for contextual bandits. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pp. 1638–1646, Bejing, China, 2014. PMLR.

Al-Rubaie, M., Wu, P. Y., Chang, J. M., and Kung, S. Privacy-preserving PCA on horizontally-partitioned data. In *DSC*, pp. 280–287. IEEE, 2017.

Angrist, J. and Krueger, A. B. Instrumental variables and the search for identification: From supply and demand to natural experiments. Technical report, National Bureau of Economic Research, 2001.

Angwin, J., Larson, J., Mattu, S., and Kirchner, L. Machine bias: There is software used across the country to predict future criminals. and it is biased against blacks. *ProPublica, May*, 23, 2016. URL https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing.

Athey, S. and Wager, S. Efficient policy learning. *arXiv preprint arXiv:1702.02896*, 2017.

Awad, E., Dsouza, S., Kim, R., Schulz, J., Henrich, J., Shariff, A., Bonnefon, J.-F., and Rahwan, I. The moral machine experiment. *Nature*, 563(7729):59–64, 2018.

Bagdasaryan, E., Poursaeed, O., and Shmatikov, V. Differential privacy has disparate impact on model accuracy. In *Advances in Neural Information Processing Systems*, pp. 15479–15488, 2019.

Barber, D. Identifying graph clusters using variational inference and link to covariance parametrization. *Philosophical Transactions of the Royal Society A*, 367:4407–4426, 2009.

Barocas, S. and Selbst, A. D. Big data's disparate impact. *California Law Review*, 104, 2016a. URL https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2477899.

Barocas, S. and Selbst, A. D. Big data's disparate impact. *California Law Review*, 104:671–732, 2016b.

Barocas, S., Hardt, M., and Narayanan, A. *Fairness and Machine Learning*. fairmlbook.org, 2019. http://www.fairmlbook.org.

Benjamin, R. *Race after technology: Abolitionist tools for the new jim code*. Polity, 2019. ISBN 1509526404.

Bennett Capers, I. Blind justice. *Yale Journal of Law & Humanities*, 24:179, 2012.

Bentham, J. *An Introduction to the Principles of Morals and Legislation*. Dover Publications, 1780.

Berk, R., Heidari, H., Jabbari, S., Kearns, M., and Roth, A. Fairness in criminal justice risk assessments: The state of the art. *Sociological Methods & Research*, 2017. doi: 10.1177/0049124118782533. URL https://doi.org/10.1177/0049124118782533.

Bertrand, M. and Mullainathan, S. Are emily and greg more employable than lakisha and jamal? a field experiment on labor market discrimination. *American economic review*, 94(4): 991–1013, 2004.

Bickel, P. J., Hammel, E. A., O'Connell, J. W., et al. Sex bias in graduate admissions: Data from berkeley. *Science*, 187(4175):398–404, 1975.

Biddle, D. *Adverse impact and test validation: A practitioner's guide to valid and defensible employment testing*. Gower Publishing, Ltd., 2006.

Binns, R. Fairness in machine learning: Lessons from political philosophy. In Friedler, S. A. and Wilson, C. (eds.), *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, volume 81 of *Proceedings of Machine Learning Research*, pp. 149–159, New York, NY, USA, 23–24 Feb 2018. PMLR. URL http://proceedings.mlr.press/v81/binns18a.html.

Blair, R. J. R. A cognitive developmental approach to morality: Investigating the psychopath. *Cognition*, 57(1):1–29, 1995.

Bolukbasi, T., Chang, K.-W., Zou, J. Y., Saligrama, V., and Kalai, A. T. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in Neural Information Processing Systems*, pp. 4349–4357, 2016.

Bonchi, F., Hajian, S., Mishra, B., and Ramazzotti, D. Exposing the probabilistic causal structure of discrimination. *International Journal of Data Science and Analytics*, 3(1):1–21, 2017.

Bongers, S., Forré, P., Peters, J., Schölkopf, B., and Mooij, J. M. Foundations of structural causal models with cycles and latent variables. *arXiv preprint arXiv:1611.06221*, 2016.

Bottou, L., Peters, J., nonero Candela, J. Q., Charles, D. X., Chickering, D. M., Portugaly, E., Ray, D., Simard, P., and Snelson, E. Counterfactual reasoning and learning systems: The example of computational advertising. *Journal of Machine Learning Research*, 14:3207–3260, 2013.

Boyd, S. and Vandenberghe, L. *Convex optimization*. Cambridge university press, 2004.

Broussard, M. *Artificial unintelligence: How computers misunderstand the world*. MIT Press, 2018.

Buolamwini, J. and Gebru, T. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on Fairness, Accountability and Transparency*, pp. 77–91, 2018.

Calders, T. and Verwer, S. Three naive bayes approaches for discrimination-free classification. *Data Mining and Knowledge Discovery*, 21(2):277–292, 2010.

Caliskan, A., Bryson, J. J., and Narayanan, A. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186, 2017.

Calmon, F., Wei, D., Vinzamuri, B., Ramamurthy, K. N., and Varshney, K. R. Optimized pre-processing for discrimination prevention. In *Advances in Neural Information Processing Systems*, pp. 3992–4001, 2017.

Chen, J., Kallus, N., Mao, X., Svacha, G., and Udell, M. Fairness under unawareness: Assessing disparity when protected class is unobserved. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, FAT* '19, pp. 339–348, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450361255. doi: 10.1145/3287560. 3287594. URL https://doi.org/10.1145/3287560.3287594.

Chiappa, S. and Gillam, T. Path-specific counterfactual fairness. *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

Chouldechova, A. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2):153–163, 2017.

Chouldechova, A. and Roth, A. A snapshot of the frontiers of fairness in machine learning. *Communications of the ACM*, 63(5):82–89, 2020.

Cinelli, C., Kumor, D., Chen, B., Pearl, J., and Bareinboim, E. Sensitivity analysis of linear structural causal models. In *Proceedings of the 36th International Conference on Machine Learning*, 2019.

Corbett-Davies, S., Pierson, E., Feller, A., and Sharad, G. A computer program used for bail and sentencing decisions was labeled biased against blacks. it's actually not that clear. *The Washington Post, October*, 2016. URL https://www.washingtonpost.com/news/monkey-cage/wp/2016/10/17/can-an-algorithm-be-racist-our-analysis-is-more-cautious-than-propublicas.

Corbett-Davies, S., Pierson, E., Feller, A., Goel, S., and Huq, A. Algorithmic decision making and the cost of fairness. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 797–806. ACM, 2017.

Cornia, N. and Mooij, J. M. Type-ii errors of independence tests can lead to arbitrarily large errors in estimated causal effects: An illustrative example. In *Proceedings of the Workshop on Causal Inference (UAI)*, pp. 35–42, 2014.

Coston, A., Ramamurthy, K. N., Wei, D., Varshney, K. R., Speakman, S., Mustahsan, Z., and Chakraborty, S. Fair transfer learning with missing protected attributes. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '19, pp. 91–98, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450363242. doi: 10.1145/3306618.3314236. URL https://doi.org/10.1145/3306618.3314236.

Crawford, K. The hidden biases in big data. *Harvard Business Review*, 1, 2013.

Cummings, R., Gupta, V., Kimpara, D., and Morgenstern, J. On the compatibility of privacy and fairness. In *Adjunct Publication of the 27th Conference on User Modeling, Adaptation and Personalization*, pp. 309–315, 2019.

Damgård, I., Pastro, V., Smart, N. P., and Zakarias, S. Multiparty computation from somewhat homomorphic encryption. In *CRYPTO*, volume 7417 of *Lecture Notes in Computer Science*, pp. 643–662. Springer, 2012.

Dembroff, R., Kohler-Hausmann, I., and Sugarman, E. What taylor swift and beyoncé teach us about sex and causes. *University of Pennsylvania Law Review, Forthcoming*, 2020.

Demmler, D., Dessouky, G., Koushanfar, F., Sadeghi, A., Schneider, T., and Zeitouni, S. Automated synthesis of optimized circuits for secure computation. In *ACM Conference on Computer and Communications Security*, pp. 1504–1517. ACM, 2015a.

Demmler, D., Schneider, T., and Zohner, M. ABY – a framework for efficient mixed-protocol secure two-party computation. In *NDSS*. The Internet Society, 2015b.

Despart, Z. Proposed bail lawsuit settlement includes child care, phones, rides for poor defendants, 2019. URL https://www.houstonchronicle.com/news/houston-texas/houston/article/Proposed-bail-lawsuit-settlement-includes-child-13764225.php.

Dieterich, W., Mendoza, C., and Brennan, T. Compas risk scales: Demonstrating accuracy equity and predictive parity. Technical report, Technical report, Northpointe, July 2016. http://www. northpointeinc. com/northpointe-analysis, 2016.

Dimitrakakis, C., Liu, Y., Parkes, D., and Radanovic, G. Bayesian fairness. In *AAAI*, 2019.

Ding, J., Zhang, X., Li, X., Wang, J., Yu, R., and Pan, M. Differentially private and fair classification via calibrated functional mechanism. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 622–629, 2020.

Dobbie, W., Goldin, J., and Yang, C. S. The effects of pretrial detention on conviction, future crime, and employment: Evidence from randomly assigned judges. *American Economic Review*, 108(2):201–40, 2018.

Doerner, J. Absentminded crypto kit. https://bitbucket.org/jackdoerner/absentminded-crypto-kit, 2018.

Dong, J., Roth, A., Schutzman, Z., Waggoner, B., and Wu, Z. S. Strategic classification from revealed preferences. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pp. 55–70, 2018.

Dorie, V., Harada, M., Carnegie, N. B., and Hill, J. A flexible, interpretable framework for assessing sensitivity to unmeasured confounding. *Statistics in Medicine*, 35(20):3453–3470, 2016.

Dougherty, C. Google photos mistakenly labels black people 'gorillas'. *The New York Times, Bits, July*, 2015. URL https://bits.blogs.nytimes.com/2015/07/01/google-photos-mistakenly-labels-black-people-gorillas.

Drton, M. and Richardson, T. Iterative conditional fitting for Gaussian ancestral graph models. *Proceedings of the Twentieth Conference on Uncertainty in Artificial Intelligence*, pp. 130–137, 2004.

Dudík, M., Langford, J., and Li, L. Doubly Robust Policy Evaluation and Learning. In *Proceedings of the 28th International Conference on Machine Learning*, pp. 1097–1104. Omnipress, 2011.

Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, ITCS '12, pp. 214–226. ACM, 2012.

Edwards, H. and Storkey, A. Censoring representations with an adversary. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2016.

Ensign, D., Friedler, S. A., Neville, S., Scheidegger, C., and Venkatasubramanian, S. Decision making with limited feedback: Error bounds for recidivism prediction and predictive policing. *JMLR*, 2018a.

Ensign, D., Friedler, S. A., Neville, S., Scheidegger, C., and Venkatasubramanian, S. Runaway feedback loops in predictive policing. *FAT*, 2018b.

Eubanks, V. *Automating inequality: How high-tech tools profile, police, and punish the poor*. St. Martin's Press, 2018.

Faiedh, H., Gafsi, Z., and Besbes, K. Digital hardware implementation of sigmoid function and its derivative for artificial neural networks. *Proceeding of the 13th International Conference on Microelectronics, 2001.*, pp. 189 – 192, 11 2001.

Feldman, M., Friedler, S. A., Moeller, J., Scheidegger, C., and Venkatasubramanian, S. Certifying and removing disparate impact. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 259–268, 2015.

Flores, A. W., Bechtel, K., and Lowenkamp, C. T. False positives, false negatives, and false analyses: A rejoinder to machine bias: There's software used across the country to predict future criminals. and it's biased against blacks. *Fed. Probation*, 80:38, 2016.

Fredrikson, M., Jha, S., and Ristenpart, T. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, pp. 1322–1333, 2015.

Friedler, S. A., Scheidegger, C., and Venkatasubramanian, S. On the (im) possibility of fairness. *arXiv preprint arXiv:1609.07236*, 2016.

Friedler, S. A., Scheidegger, C., Venkatasubramanian, S., Choudhary, S., Hamilton, E. P., and Roth, D. A comparative study of fairness-enhancing interventions in machine learning. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, FAT* '19, pp. 329–338, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450361255. doi: 10.1145/3287560.3287589. URL https://doi.org/10.1145/3287560.3287589.

Friedman, B. and Nissenbaum, H. Bias in computer systems. *ACM Transactions on Information Systems (TOIS)*, 14(3):330–347, 1996.

Gascón, A., Schoppmann, P., Balle, B., Raykova, M., Doerner, J., Zahur, S., and Evans, D. Privacy-Preserving Distributed Linear Regression on High-Dimensional Data. *Proceedings on Privacy Enhancing Technologies*, 2017(4):345–364, October 2017.

Gebhard, T., Kilbertus, N., Parascandolo, G., Harry, I., and Schölkopf, B. Convwave: Searching for gravitational waves with fully convolutional neural nets. In *Workshop on Deep Learning for Physical Sciences (DLPS) at the 31st Conference on Neural Information Processing Systems*, December 2017. URL https://dl4physicalsciences.github.io/files/nips_dlps_2017_13.pdf.

Gebhard, T. D., Kilbertus, N., Harry, I., and Schölkopf, B. Convolutional neural networks: A magic bullet for gravitational-wave detection? *Phys. Rev. D*, 100:063015, Sep 2019. doi: 10.1103/PhysRevD.100.063015. URL https://link.aps.org/doi/10.1103/PhysRevD.100.063015.

Gillen, S., Jung, C., Kearns, M., and Roth, A. Online learning with an unknown fairness metric. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 31*, pp. 2600–2609. Curran Associates, Inc., 2018. URL http://papers.nips.cc/paper/7526-online-learning-with-an-unknown-fairness-metric.pdf.

Goldreich, O. *The Foundations of Cryptography – Volume 2, Basic Applications*. Cambridge University Press, 2004.

Goldreich, O., Micali, S., and Wigderson, A. How to play any mental game or A completeness theorem for protocols with honest majority. In *STOC*, pp. 218–229. ACM, 1987.

Graham, C. NHS cyber attack: Everything you need to know about 'biggest ransomware' offensive in history. *Telegraph, May 20, 2017.* URL http://www.telegraph.co.uk/news/2017/05/13/nhs-cyber-attack-everything-need-know-biggest-ransomware-offensive/.

Green, B. and Hu, L. The myth in the methodology: Towards a recontextualization of fairness in machine learning. *Machine Learning: The Debates workshop at the 35th International Conference on Machine Learning (ICML)*, 2018.

Grgić-Hlača, N., Zafar, M. B., Gummadi, K. P., and Weller, A. The case for process fairness in learning: Feature selection for fair decision making. In *NIPS Symposium on Machine Learning and the Law, Barcelona, Spain*, volume 8, 2016.

Grgić-Hlača, N., Redmiles, E. M., Gummadi, K. P., and Weller, A. Human perceptions of fairness in algorithmic decision making: A case study of criminal risk prediction. In *WWW*, 2018a.

Grgić-Hlača, N., Zafar, M. B., Gummadi, K. P., and Weller, A. Beyond distributive fairness in algorithmic decision making: Feature selection for procedurally fair learning. In *AAAI*, 2018b.

Grgić-Hlača, N., Weller, A., and Redmiles, E. M. Dimensions of diversity in human perceptions of algorithmic fairness. *arXiv preprint arXiv:2005.00808*, 2020.

Handel, P., Skog, I., Wahlstrom, J., Bonawiede, F., Welch, R., Ohlsson, J., and Ohlsson, M. Insurance telematics: Opportunities and challenges with the smartphone solution. *IEEE Intelligent Transportation Systems Magazine*, 6(4):57–70, 2014.

Hanna, R. N. and Linden, L. L. Discrimination in grading. *American Economic Journal: Economic Policy*, 4(4):146–68, 2012.

Hardt, M., Megiddo, N., Papadimitriou, C., and Wootters, M. Strategic classification. In *Proceedings of the 2016 ACM Conference on Innovations in Theoretical Computer Science*, ITCS '16, pp. 111–122, New York, NY, USA, 2016a. Association for Computing Machinery. ISBN 9781450340571. doi: 10.1145/2840728.2840730. URL https://doi.org/10.1145/2840728.2840730.

Hardt, M., Price, E., Srebro, N., et al. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems*, pp. 3315–3323, 2016b.

Hashimoto, T., Srivastava, M., Namkoong, H., and Liang, P. Fairness without demographics in repeated loss minimization. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1929–1938, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL http://proceedings.mlr.press/v80/hashimoto18a.html.

Hébert-Johnson, U., Kim, M., Reingold, O., and Rothblum, G. Multicalibration: Calibration for the (Computationally-identifiable) masses. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1939–1948, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL http://proceedings.mlr.press/v80/hebert-johnson18a.html.

Heidari, H. and Krause, A. Preventing disparate treatment in sequential decision making. In *IJCAI*, pp. 2248–2254, 2018.

Helmbold, D. P., Littlestone, N., and Long, P. M. Apple tasting. *Information and Computation*, 161(2):85–139, 2000.

Holstein, K., Wortman Vaughan, J., Daumé, H., Dudik, M., and Wallach, H. Improving fairness in machine learning systems: What do industry practitioners need? In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, CHI '19, pp. 1–16, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450359702. doi: 10.1145/3290605.3300830. URL https://doi.org/10.1145/3290605.3300830.

Hoyer, P. O., Janzing, D., Mooij, J. M., Peters, J., and Schölkopf, B. Nonlinear causal discovery with additive noise models. *Advances in Neural Information Processing Systems*, 21:689–696, 2008.

Hron, J., Krauth, K., Jordan, M. I., and Kilbertus, N. Exploration in two-stage recommender systems. *arXiv preprint arXiv:2009.08956*, 2020.

Hu, L. and Chen, Y. A short-term intervention for long-term fairness in the labor market. In *World Wide Web Conference*, WWW '18, pp. 1389–1398, Republic and Canton of Geneva, Switzerland, 2018. International World Wide Web Conferences Steering Committee. ISBN 978-1-4503-5639-8.

Hu, L. and Kohler-Hausmann, I. What's sex got to do with machine learning? In *ACM Conference on Fairness, Accountability, and Transparency*, 2020.

Hyvärinen, A. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6:695–709, 2005.

Ingold, D. and Soper, S. Amazon doesn't consider the race of its customers. should it? *Bloomberg, April*, 2016. URL https://www.bloomberg.com/graphics/2016-amazon-same-day.

Jabbari, S., Joseph, M., Kearns, M., Morgenstern, J., and Roth, A. Fairness in reinforcement learning. In Precup, D. and Teh, Y. W. (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1617–1626, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. URL http://proceedings.mlr.press/v70/jabbari17a.html.

Jagielski, M., Kearns, M., Mao, J., Oprea, A., Roth, A., Malvajerdi, S. S., and Ullman, J. Differentially private fair learning. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 3000–3008, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL http://proceedings.mlr.press/v97/jagielski19a.html.

Joseph, M., Kearns, M., Morgenstern, J. H., and Roth, A. Fairness in learning: Classic and contextual bandits. In *Advances in Neural Information Processing Systems*, pp. 325–333, 2016.

Jung, J., Shroff, R., Feller, A., and Goel, S. Algorithmic decision making in the presence of unmeasured confounding. *arXiv preprint arXiv:1805.01868*, 2018.

Juvekar, C., Vaikuntanathan, V., and Chandrakasan, A. Gazelle: A Low Latency Framework for Secure Neural Network Inference. *IACR Cryptology ePrint Archive*, 2018:73, 2018.

Kahneman, D., Knetsch, J. L., and Thaler, R. Fairness as a constraint on profit seeking: Entitlements in the market. *The American economic review*, pp. 728–741, 1986.

Kallus, N. Balanced policy evaluation and learning. In *Advances in Neural Information Processing Systems*, pp. 8909–8920, 2018.

Kallus, N. and Zhou, A. Residual unfairness in fair machine learning from prejudiced data. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2439–2448, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.

Kamiran, F., Žliobaitė, I., and Calders, T. Quantifying explainable discrimination and removing illegal discrimination in automated decision making. *Knowledge and information systems*, 35(3):613–644, 2013.

Kamishima, T., Akaho, S., Asoh, H., and Sakuma, J. Fairness-aware classifier with prejudice remover regularizer. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 35–50. Springer, 2012.

Kay, M., Matuszek, C., and Munson, S. A. Unequal representation and gender stereotypes in image search results for occupations. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, pp. 3819–3828. ACM, 2015.

Kearns, M. and Roth, A. *The ethical algorithm: The science of socially aware algorithm design*. Oxford University Press, 2019.

Kearns, M., Neel, S., Roth, A., and Wu, Z. S. Preventing fairness gerrymandering: Auditing and learning for subgroup fairness. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2564–2572, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL http://proceedings.mlr.press/v80/kearns18a.html.

Keller, M., Scholl, P., and Smart, N. P. An architecture for practical actively secure MPC with dishonest majority. In *ACM Conference on Computer and Communications Security*, pp. 549–560. ACM, 2013.

Keller, M., Orsini, E., Rotaru, D., Scholl, P., Soria-Vazquez, E., and Vivek, S. Faster secure multi-party computation of AES and DES using lookup tables. In *ACNS*, volume 10355 of *Lecture Notes in Computer Science*, pp. 229–249. Springer, 2017.

Keller, M., Pastro, V., and Rotaru, D. Overdrive: Making SPDZ great again. In *EUROCRYPT (3)*, volume 10822 of *Lecture Notes in Computer Science*, pp. 158–189. Springer, 2018.

Kiefer, J., Wolfowitz, J., et al. Stochastic estimation of the maximum of a regression function. *The Annals of Mathematical Statistics*, 23(3):462–466, 1952.

Kilbertus, N., Carulla, M. R., Parascandolo, G., Hardt, M., Janzing, D., and Schölkopf, B. Avoiding discrimination through causal reasoning. In *Advances in Neural Information Processing Systems*, pp. 656–666, 2017.

Kilbertus, N., Gascon, A., Kusner, M., Veale, M., Gummadi, K., and Weller, A. Blind justice: Fairness with encrypted sensitive attributes. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2630–2639, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018a. PMLR. URL http://proceedings.mlr.press/v80/kilbertus18a.html.

Kilbertus, N., Parascandolo, G., and Schölkopf, B. Generalization in anti-causal learning. *arXiv preprint arXiv:1812.00524*, 2018b.

Kilbertus, N., Ball, P. J., Kusner, M. J., Weller, A., and Silva, R. The sensitivity of counterfactual fairness to unmeasured confounding. In *Proceedings of the 35th Conference on Uncertainty*

*in Artificial Intelligence (UAI)*, pp. 213. AUAI Press, July 2019. URL http://auai.org/uai2019/proceedings/papers/213.pdf.

Kilbertus, N., Kusner, M. J., and Silva, R. A class of algorithms for general instrumental variable models. In *Advances in Neural Information Processing Systems 33*, 2020a.

Kilbertus, N., Rodriguez, M. G., Schölkopf, B., Muandet, K., and Valera, I. Fair decisions despite imperfect predictions. In Chiappa, S. and Calandra, R. (eds.), *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pp. 277–287, Online, 26–28 Aug 2020b. PMLR. URL http://proceedings.mlr.press/v108/kilbertus20a.html.

Kim, J. S., Chen, J., and Talwalkar, A. Model-agnostic characterization of fairness trade-offs. In *Proceedings of the 37th International Conference on Machine Learning*, Proceedings of Machine Learning Research. PMLR, 2020.

Kim, M., Reingold, O., and Rothblum, G. Fairness through computationally-bounded awareness. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 31*, pp. 4842–4852. Curran Associates, Inc., 2018. URL http://papers.nips.cc/paper/7733-fairness-through-computationally-bounded-awareness.pdf.

Kingma, P. D. and Welling, M. Auto-encoding variational bayes. In *Proceedings of the International Conference on Learning Representations (ICLR)*, volume 1, 2014.

Kleinberg, J., Lakkaraju, H., Leskovec, J., Ludwig, J., and Mullainathan, S. Human decisions and machine predictions. *The quarterly journal of economics*, 133(1):237–293, 2017a.

Kleinberg, J., Mullainathan, S., and Raghavan, M. Inherent Trade-Offs in the Fair Determination of Risk Scores. In Papadimitriou, C. H. (ed.), *8th Innovations in Theoretical Computer Science Conference (ITCS 2017)*, volume 67 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pp. 43:1–43:23, Dagstuhl, Germany, 2017b. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. ISBN 978-3-95977-029-3. doi: 10.4230/LIPIcs.ITCS.2017.43. URL http://drops.dagstuhl.de/opus/volltexte/2017/8156.

Kleinberg, J., Ludwig, J., Mullainathan, S., and Rambachan, A. Algorithmic fairness. In *Aea papers and proceedings*, volume 108, pp. 22–27, 2018a.

Kleinberg, J., Ludwig, J., Mullainathan, S., and Sunstein, C. R. Discrimination in the age of algorithms. *Journal of Legal Analysis*, 10, 2018b.

Kohler-Hausmann, I. Eddie murphy and the dangers of counterfactual causal thinking about detecting racial discrimination. *Nw. UL Rev.*, 113:1163, 2018.

Kroll, J. A., Huey, J., Barocas, S., Felten, E. W., Reidenberg, J. R., Robinson, D. G., and Yu, H. Accountable algorithms. *University of Pennsylvania Law Review*, 165, 2016.

Kusner, M. J., Loftus, J., Russell, C., and Silva, R. Counterfactual fairness. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 30*, pp. 4069–4079. Curran Associates, Inc., 2017. URL http://papers.nips.cc/paper/6995-counterfactual-fairness.pdf.

Lahoti, P., Gummadi, K. P., and Weikum, G. Operationalizing individual fairness with pairwise fair representations. *Proc. VLDB Endow.*, 13(4):506–518, December 2019. ISSN 2150-8097. doi: 10.14778/3372716.3372723. URL https://doi.org/10.14778/3372716.3372723.

Lahoti, P., Beutel, A., Chen, J., Lee, K., Prost, F., Thain, N., Wang, X., and Chi, E. H. Fairness without demographics through adversarially reweighted learning. *arXiv preprint arXiv:2006.13114*, 2020.

Lakkaraju, H. and Rudin, C. Learning Cost-Effective and Interpretable Treatment Regimes. In Singh, A. and Zhu, J. (eds.), *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pp. 166–175, Fort Lauderdale, FL, USA, 20–22 Apr 2017. PMLR.

Lakkaraju, H., Kleinberg, J., Leskovec, J., Ludwig, J., and Mullainathan, S. The selective labels problem: Evaluating algorithmic predictions in the presence of unobservables. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 275–284. ACM, 2017.

Langford, J., Strehl, A., and Wortman, J. Exploration scavenging. In *Proceedings of the 25th International Conference on Machine Learning*, ICML '08, pp. 528–535, New York, NY, USA, 2008. ACM.

Lee, P. Learning from tay's introduction. *Official Microsoft Blog, March*, 2016. URL https://blogs.microsoft.com/blog/2016/03/25/learning-tays-introduction.

Lichman, M. UCI machine learning repository, 2013. URL http://archive.ics.uci.edu/ml.

Lindell, Y. How To Simulate It – A Tutorial on the Simulation Proof Technique. *IACR Cryptology ePrint Archive*, 2016:46, 2016.

Liu, J., Juuti, M., Lu, Y., and Asokan, N. Oblivious neural network predictions via minionn transformations. In *CCS*, pp. 619–631. ACM, 2017.

Liu, L. T., Dean, S., Rolf, E., Simchowitz, M., and Hardt, M. Delayed impact of fair machine learning. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 3150–3158, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL http://proceedings.mlr.press/v80/liu18c.html.

Loftus, J. R., Russell, C., Kusner, M. J., and Silva, R. Causal reasoning for algorithmic fairness. *arXiv preprint arXiv:1805.05859*, 2018.

Louizos, C., Swersky, K., Li, Y., Welling, M., and Zemel, R. The variational fair autoencoder. In *International Conference on Learning Representations*, 2016.

Lum, K. and Isaac, W. To predict and serve? *Significance*, 13(5):14–19, 2016. doi: 10.1111/j.1740-9713.2016.00960.x. URL https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/j.1740-9713.2016.00960.x.

Madras, D., Creager, E., Pitassi, T., and Zemel, R. Learning adversarially fair and transferable representations. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 3384–3393, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL http://proceedings.mlr.press/v80/madras18a.html.

Meyer, M. N., Heck, P. R., Holtzman, G. S., Anderson, S. M., Cai, W., Watts, D. J., and Chabris, C. F. Objecting to experiments that compare two unobjectionable policies or treatments. *Proceedings of the National Academy of Sciences*, 2019. ISSN 0027-8424. doi: 10.1073/pnas.1820701116.

Miller, J., Milli, S., and Hardt, M. Strategic classification is causal modeling in disguise. In *Proceedings of the 37th International Conference on Machine Learning*, Proceedings of Machine Learning Research. PMLR, 2020.

Milli, S., Miller, J., Dragan, A. D., and Hardt, M. The social cost of strategic classification. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pp. 230–239, 2019.

Mitchell, S., Potash, E., and Barocas, S. Prediction-based decisions and fairness: A catalogue of choices, assumptions, and definitions. *arXiv preprint arXiv:1811.07867*, 2018.

Mohamed, S., Png, M.-T., and Isaac, W. Decolonial ai: Decolonial theory as sociotechnical foresight in artificial intelligence. *Philosophy & Technology*, pp. 1–26, 2020.

Mohassel, P. and Zhang, Y. SecureML: A system for scalable privacy-preserving machine learning. In *IEEE Symposium on Security and Privacy (SP)*, pp. 19–38, 2017.

Moulin, H. *Fair division and collective welfare*. MIT press, 2004.

Mouzannar, H., Ohannessian, M. I., and Srebro, N. From fair decision making to social equality. In *FAT*, 2019.

Muñoz, C., Smith, M., and Patil, D. Big data: A report on algorithmic systems, opportunity, and civil rights. *Washington, DC: Executive Office of the President, White House*, 2016. URL https://obamawhitehouse.archives.gov/sites/default/files/microsites/ostp/2016_0504_data_discrimination.pdf.

Nabi, R. and Shpitser, I. Fair inference on outcomes. In *AAAI Conference on Artificial Intelligence*, volume 2018, pp. 1931. NIH Public Access, 2018.

Nikolaenko, V., Ioannidis, S., Weinsberg, U., Joye, M., Taft, N., and Boneh, D. Privacy-preserving matrix factorization. In *ACM Conference on Computer and Communications Security*, pp. 801–812. ACM, 2013a.

Nikolaenko, V., Weinsberg, U., Ioannidis, S., Joye, M., Boneh, D., and Taft, N. Privacy-preserving ridge regression on hundreds of millions of records. In *IEEE Symposium on Security and Privacy*, pp. 334–348. IEEE Computer Society, 2013b.

Noble, S. U. *Algorithms of Oppression: How Search Engines Reinforce Racism*. NYU Press, 2018. ISBN 9781479849949.

O'Neil, C. *Weapons of math destruction: How big data increases inequality and threatens democracy*. Broadway Books, 2016.

Parascandolo, G., Kilbertus, N., Rojas-Carulla, M., and Schölkopf, B. Learning independent causal mechanisms. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 4036–4044, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL http://proceedings.mlr.press/v80/parascandolo18a.html.

Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.

Pearl, J. *Causality*. Cambridge University Press, 2009.

Perdomo, J. C., Zrnic, T., Mendler-Dünner, C., and Hardt, M. Performative prediction. In *Proceedings of the 37th International Conference on Machine Learning*, Proceedings of Machine Learning Research. PMLR, 2020.

Peters, J., Janzing, D., and Schölkopf, B. *Elements of causal inference: foundations and learning algorithms*. MIT press, 2017.

Picker Institute Europe. National Health Service national staff survey, 2014, 2015. URL http://doi.org/10.5255/UKDA-SN-7776-1.

Pleiss, G., Raghavan, M., Wu, F., Kleinberg, J., and Weinberger, K. Q. On fairness and calibration. In *Advances in Neural Information Processing Systems*, pp. 5680–5689, 2017.

Qureshi, B., Kamiran, F., Karim, A., Ruggieri, S., and Pedreschi, D. Causal inference for social discrimination reasoning. *Journal of Intelligent Information Systems*, 54(2):425–437, November 2019. doi: 10.1007/s10844-019-00580-x. URL https://doi.org/10.1007/s10844-019-00580-x.

Rambachan, A., Kleinberg, J., Mullainathan, S., and Ludwig, J. An economic approach to regulating algorithms. Technical report, National Bureau of Economic Research, 2020.

Rastegarpanah, B., Crovella, M., and Gummadi, K. P. Fair inputs and fair outputs: The incompatibility of fairness in privacy and accuracy. *arXiv preprint arXiv:2005.09209*, 2020.

Rawls, J. Justice as fairness. *The philosophical review*, 67(2):164–194, 1958.

Rawls, J. *A theory of justice*. Harvard university press, 2009.

Richardson, T. Markov properties for acyclic directed mixed graphs. *Scandinavian Journal of Statistics*, 30:145–157, 2003.

Robins, J. M., Rotnitzky, A., and Scharfstein, D. O. Sensitivity analysis for selection bias and unmeasured confounding in missing data and causal inference models. In *Statistical Models in Epidemiology, the Environment, and Clinical Trials*, 2000.

Roemer, J. E. *Theories of distributive justice*. Harvard University Press, 1998.

Roemer, J. E. *Equality of opportunity*. Harvard University Press, 2009.

Rosenbaum, P. *Observational Studies*. Springer, 2002.

Rosenbaum, P. R. and Rubin, D. B. The central role of the propensity score in observational studies for causal effects. *Biometrika*, pp. 41–55, 1983.

Rubin, D. Causal inference using potential outcomes. *Journal of the American Statistical Association*, 100(469):322–331, 2005.

Russell, C., Kusner, M. J., Loftus, J., and Silva, R. When worlds collide: integrating different counterfactual assumptions in fairness. In *Advances in Neural Information Processing Systems*, pp. 6414–6423, 2017.

Sandel, M. J. *Justice: What's the right thing to do?* Macmillan, 2010.

Shpitser, I. Counterfactual graphical models for longitudinal mediation analysis with unobserved confounding. *Cognitive Science*, 37:1011–1035, 2013.

Silva, R., Chu, W., and Ghahramani, Z. Hidden common cause relations in relational learning. *Advances in Neural Information Processing Systems*, 20:1345–1352, 2007.

Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., and Riedmiller, M. Deterministic policy gradient algorithms. In *Proceedings of the 31st International Conference on Machine Learning*, pp. 387–395. JMLR.org, 2014.

Snow, J. Amazon's face recognition falsely matched 28 members of congress with mugshots. *American Civil Liberties Union of Northern California, July*, 2018. URL https://www.aclunc.org/blog/amazon-s-face-recognition-falsely-matched-28-members-congress-mugshots.

Sokolic, J., Rodrigues, M. R., Qiu, Q., and Sapiro, G. Learning to identify while failing to discriminate. In *ICCV Workshops*, pp. 2537–2544, 2017.

Speicher, T., Heidari, H., Grgić-Hlača, N., Gummadi, K. P., Singla, A., Weller, A., and Zafar, M. B. A unified approach to quantifying algorithmic unfairness: Measuring individual & group unfairness via inequality indices. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 2239–2248. ACM, 2018.

Sutton, R. S. and Barto, A. G. *Introduction to Reinforcement Learning*. MIT Press, Cambridge, MA, USA, 1st edition, 1998. ISBN 0262193981.

Swaminathan, A. and Joachims, T. Counterfactual risk minimization: Learning from logged bandit feedback. In *ICML*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 814–823. JMLR.org, 2015a.

Swaminathan, A. and Joachims, T. The self-normalized estimator for counterfactual learning. In *Proceedings of the 28th International Conference on Neural Information Processing Systems*, pp. 3231–3239, Cambridge, MA, USA, 2015b. MIT Press.

Sweeney, L. Discrimination in online ad delivery. *Queue*, 11(3):10, 2013.

Tramèr, F., Zhang, F., Juels, A., Reiter, M. K., and Ristenpart, T. Stealing machine learning models via prediction apis. In *USENIX Security Symposium*, pp. 601–618, 2016.

Träuble, F., Creager, E., Kilbertus, N., Goyal, A., Locatello, F., Schölkopf, B., and Bauer, S. Is independence all you need? on the generalization of representations learned from correlated data. *arXiv preprint arXiv:2006.07886*, 2020.

Turiel, E. *The culture of morality: Social development, context, and conflict*. Cambridge University Press, 2002.

Valera, I., Singla, A., and Gomez-Rodriguez, M. Enhancing the accuracy and fairness of human decision making. In *Neural Information Processing Systems*, 2018.

VanderWeele, T. J. and Robinson, W. R. On causal interpretation of race in regressions adjusting for confounding and mediating variables. *Epidemiology*, 25(4):473, 2014.

Veale, M. and Binns, R. Fairer machine learning in the real world: Mitigating discrimination without collecting sensitive data. *Big Data & Society*, 4(2), 2017.

Veale, M. and Edwards, L. Clarity, Surprises, and Further Questions in the Article 29 Working Party Draft Guidance on Automated Decision-Making and Profiling. *Computer Law & Security Review*, 2018. doi: 10.1016/j.clsr.2017.12.002.

Wang, H., Grgić-Hlača, N., Lahoti, P., Gummadi, K. P., and Weller, A. An empirical study on learning fairness metrics for compas data with human supervision. *arXiv preprint arXiv:1910.10255*, 2019.

Wick, M., panda, s., and Tristan, J.-B. Unlocking fairness: a trade-off revisited. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32*, pp. 8783–8792. Curran Associates, Inc., 2019. URL http://papers.nips.cc/paper/9082-unlocking-fairness-a-trade-off-revisited.pdf.

Wightman, L. F. LSAC national longitudinal bar passage study. *LSAC Research Report Series.*, 1998.

Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.

Woodworth, B., Gunasekar, S., Ohannessian, M. I., and Srebro, N. Learning non-discriminatory predictors. In Kale, S. and Shamir, O. (eds.), *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*, pp. 1920–1953, Amsterdam, Netherlands, 07–10 Jul 2017. PMLR.

Xu, D., Du, W., and Wu, X. Removing disparate impact of differentially private stochastic gradient descent on model accuracy. *arXiv preprint arXiv:2003.03699*, 2020.

Yaari, M. E. and Bar-Hillel, M. On dividing justly. *Social choice and welfare*, 1(1):1–24, 1984.

Yao, A. C.-C. How to Generate and Exchange Secrets (Extended Abstract). In *FOCS*, pp. 162–167. IEEE Computer Society, 1986.

Yona, G. and Rothblum, G. Probably approximately metric-fair learning. In *International Conference on Machine Learning*, pp. 5666–5674, 2018.

Young, H. P. *Equity: in theory and practice*. Princeton University Press, 1995.

Zafar, M. B., Valera, I., Gómez Rodriguez, M., and Gummadi, K. P. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *Proceedings of the 26th International Conference on World Wide Web*, pp. 1171–1180, 2017a.

Zafar, M. B., Valera, I., Rogriguez, M. G., and Gummadi, K. P. Fairness Constraints: Mechanisms for Fair Classification. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, pp. 962–970, 2017b.

Zahur, S. and Evans, D. Obliv-c: A language for extensible data-oblivious computation. *IACR Cryptology ePrint Archive*, 2015:1153, 2015.

Zemel, R. S., Wu, Y., Swersky, K., Pitassi, T., and Dwork, C. Learning fair representations. *Proceedings of the International Conference of Machine Learning*, 28:325–333, 2013. URL http://jmlr.org/proceedings/papers/v28/zemel13.pdf.

Zhang, J. and Bareinboim, E. Fairness in decision-making–the causal explanation formula. In *32nd AAAI Conference on Artificial Intelligence*, 2018.

Zhang, L. and Wu, X. Anti-discrimination learning: a causal modeling-based framework. *International Journal of Data Science and Analytics*, pp. 1–16, 2017. ISSN 2364-4168.

Žliobaitė, I. and Custers, B. Using sensitive personal data may be necessary for avoiding discrimination in data-driven decision models. *Artificial Intelligence and Law*, 24(2):183–201, 2016.

# A

# Additional experiments for Chapter 6

## A.1 Results on remaining datasets

Analogously to Figures 6.3 and Figure 6.4 we report the results on test accuracy as well as the mitigation of disparate impact for the Lagrangian multiplier method in Figure A.1. In the Adult dataset we are able to mitigate disparate impact with slightly worse accuracy as compared to the baseline. Note that the German dataset contains only 512 training and 200 test examples, which explains the discrete jumps in accuracy in minimal steps of $1/200 = 0.005$. Hence, even though the Lagrangian multiplier technique here consistently removes disparate impact to a similar extent as the baseline, interpretations of results on such small datasets require great care. For the much larger stop, question and frisk dataset we again observe the curious initial increase in accuracy similar to our observations for the Bank dataset. In this dataset about 93% of all samples have positive labels, which explains the near optimal accuracy when collapsing to always predict 1, which happens for the baseline as well for our method at a similar rate as $c$ decreases.

## A.2 Disadvantages of other optimization methods

In Section 6.5 we suggest the Lagrangian multiplier technique for fair model training using fixed-point numbers. Here we substantiate this suggestion with further empirical evidence. Figure A.2 shows analogous results to Figure 6.4 and the second row of Figure A.1. These plots reveal the shortcomings of the interior point logarithmic barrier and the projected gradient methods.

**Interior Point Logarithmic Barrier method.** While the interior point logarithmic barrier method does balance the fractions of people being assigned positive outcomes between the two different demographic groups when the constraint is tightened, it soon breaks down entirely due to overflow and underflow errors. The number of failed runs was substantially
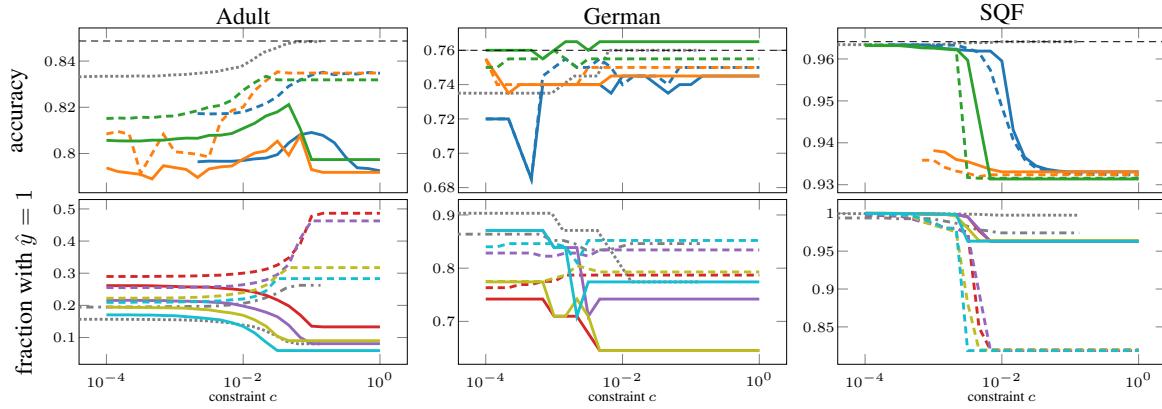
Figure A.1 **First row:** The color code is blue: iplb, orange: projected, green: Lagrange with *continuous* lines for no approximation and *dashed* lines for piecewise linear approximation. The gray dotted line is the baseline and the dashed black line marks unconstrained logistic regression. **Second row:** *Continuous/dotted* lines correspond to $z = 0$ and *dashed/dash-dotted* lines to $z = 1$. The color code is (red: no approx. + float, purple: no approx. + fixed, yellow: pw linear + float, turquoise: pw linear + fixed, gray: baseline).



Figure A.2 We plot the fraction of people with $z = 0$ (*continuous/dotted*) and with $z = 1$ (*dashed/dash-dotted*) who get assigned positive outcomes over the constraint $c$ for 6 different datasets. The different colors correspond to (red: no approximation + floats, purple: no approximation + fixed-point, yellow: piecewise linear + floats, turquoise: piecewise linear + fixed-point, gray: baseline).

higher than for the Lagrangian multiplier technique. As explained by Boyd & Vandenberghe (2004), when we increase the parameter $t$ of the interior point logarithmic barrier method during training, the barrier becomes steeper, approaching the function

$$I_-(x) = \begin{cases} 0 & \text{for } x \leq 0, \\ \infty & \text{for } x > 0. \end{cases}$$

From this it becomes obvious that when facing tight constraints, the gradients might change from almost zero to large values within a single update of the parameters $\boldsymbol{\theta}$. Moreover, iplb requires careful tuning and scheduling of $t$. Hence, the interior point logarithmic barrier method, while achieving good results over some domains, is not well suited for MPC.

**Projected gradient method.** In Figure A.2, we observe that the projected gradient method seems to fail in most cases, since it does not actually balance the fractions of positive outcomes across the sensitive groups. There is a simple explanation why it can satisfy the constraint $F(\boldsymbol{\theta}) \leq 0$ for the $p$%-rule even with small $c$ and still retain near optimal accuracy. Note that the accuracy only depends on the direction of $\boldsymbol{\theta}$, i.e., it is invariant to arbitrary rescaling of $\boldsymbol{\theta}$. Since the constraint $F(\boldsymbol{\theta}) = |A\boldsymbol{\theta}| - c \leq 0$ is always satisfied for $\boldsymbol{\theta} = 0$, dividing any $\boldsymbol{\theta}$ by a large enough factor will result in a classifier that achieves equal accuracy and satisfies the constraint (by continuity). However, minimizing the loss in the original logistic regression optimization problem (or equivalently maximizing the likelihood), which is not invariant under rescaling of $\boldsymbol{\theta}$, counteracts shrinking $\boldsymbol{\theta}$ as it enforces high confidence of decisions, i.e., large $\boldsymbol{\theta}$. The projection method produces high accuracy classifiers with small weights that formally fulfill the fairness constraint, but do not properly mitigate disparate impact as measured by the true $p$%-rule instead of the computational proxy. It also often fails for small constraint values, as the projection matrix in eq. (6.4) turns out to become near singular producing over- and underflow errors.

# Additional experiments for Chapter

## B.1    Experiments on synthetic data

**Setup.** The precise setup for the two different synthetic settings, illustrated in Figure 7.3, is as follows. The only feature $x$ is a scalar score and $z \sim \text{Ber}(0.5)$. In the first setting, $x$ is sampled from a normal distribution $\mathcal{N}(\mu = 0.5 - z, \sigma = 1)$ truncated to $x \in [-0.8, 0.8]$, and the conditional probability $\text{P}(Y \mid x)$ is strictly monotonic in the score and does not explicitly depend on $s$. As a result, for any $c$, there exists a single decision boundary for the score that results in the optimal policy, which is contained in the class of logistic policies. Note, however, that the score is not well calibrated, i.e., $\text{P}(Y \mid x)$ is not directly proportional to $x$.

In the second setting, $x \sim \mathcal{N}(\mu = 3(0.5 - z), \sigma = 3.5)$. Here, the conditional probability $\text{P}(Y \mid x)$ crosses the cost threshold $c$ multiple times, resulting in two disjoint intervals of scores for which the optimal decision is $d = 1$ (green areas). Consequently, the optimal policy cannot be implemented by a deterministic threshold rule based on a logistic predictive model. We show the best achievable single decision threshold in Figure 7.3.

**Repeated figure.** First, in Figure B.1 we again show the contents of Figure 7.4 in the main text, but added effective utility and also show shaded regions for the 25th and 75th percentile over 30 runs.

**Evolution of policies.** In Figure B.2 we show for a representative run at $\lambda = 0$ how the different policies evolve in the two synthetic settings over time. The two columns correspond to the two different synthetic settings. For all policies, we show snapshots at a fixed number of logarithmically spaced time steps between $t = 0$ and $t = 200$. For deterministic threshold rules, we show the logistic function of the underlying predictive model. The vertical dashed line corresponds to the decision boundary in $x$. For the logistic and semi-logistic policies, the lines correspond to $\pi_t(D = 1 \mid x)$, i.e., to the probability of giving a positive decision for a given input $x$. Note that the semi-logistic policies have a discontinuity, because we do
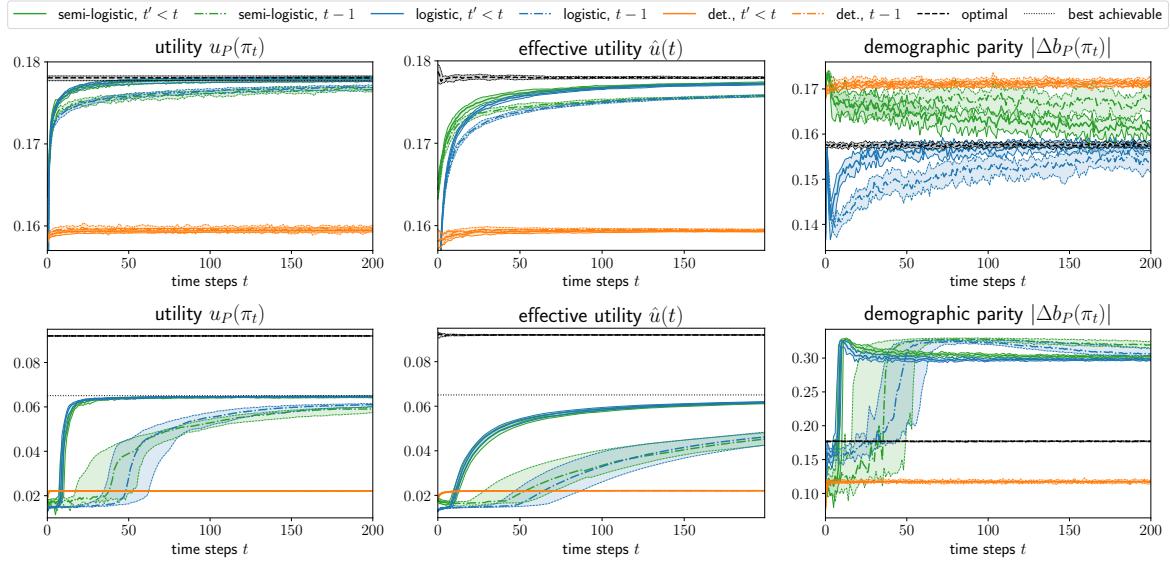
Figure B.1 Utility, effective utility, and demographic parity in the synthetic settings of Figure 7.3.

not randomize when the model believes $d = 1$ is a favorable decision with more than 50% certainty. For reference, we also show the true conditional distribution, the cost parameter, as well as the best achievable single decision boundary.

In the first setting, the exploring policies locate the optimal decision boundary, whereas the deterministic threshold rules, which are based on learned predictive models, do not, even though $P(Y = 1 \mid x)$ is monotonic in $x$ and has a sigmoidal shape. The predictive models focus on fitting the rightmost part of the conditional well, but ignore the left region, from which they never receive data.

In the second setting, our methods explore more and eventually take positive decisions for $x$ right of the vertical dotted line in Figure 7.3, which is indeed the best achievable single threshold policy. In contrast, non-exploring deterministic threshold rules again suffer from the same issue as in the first setting and converge to a suboptimal threshold at $x \approx 5$. They ignore the left green region in Figure 7.3 and do not overcome the dip of $P(Y = 1 \mid x)$ below $c$, because they never receive data for $x \leq 4$.

**Adding fairness constraints.** Figure B.3 shows how all metrics at the final time step $t = 200$ evolve as $\lambda$ is increased over the range $[10^{-0.5}, 10^4]$. We use the benefit function for demographic parity in the fairness constraint, i.e., $f(d, y) = d$. The first row corresponds to the first setting and the second row corresponds to the second setting. In both cases, our approach achieves perfect fairness for sufficiently large $\lambda$ at the expected cost of a drop in (effective) utility.
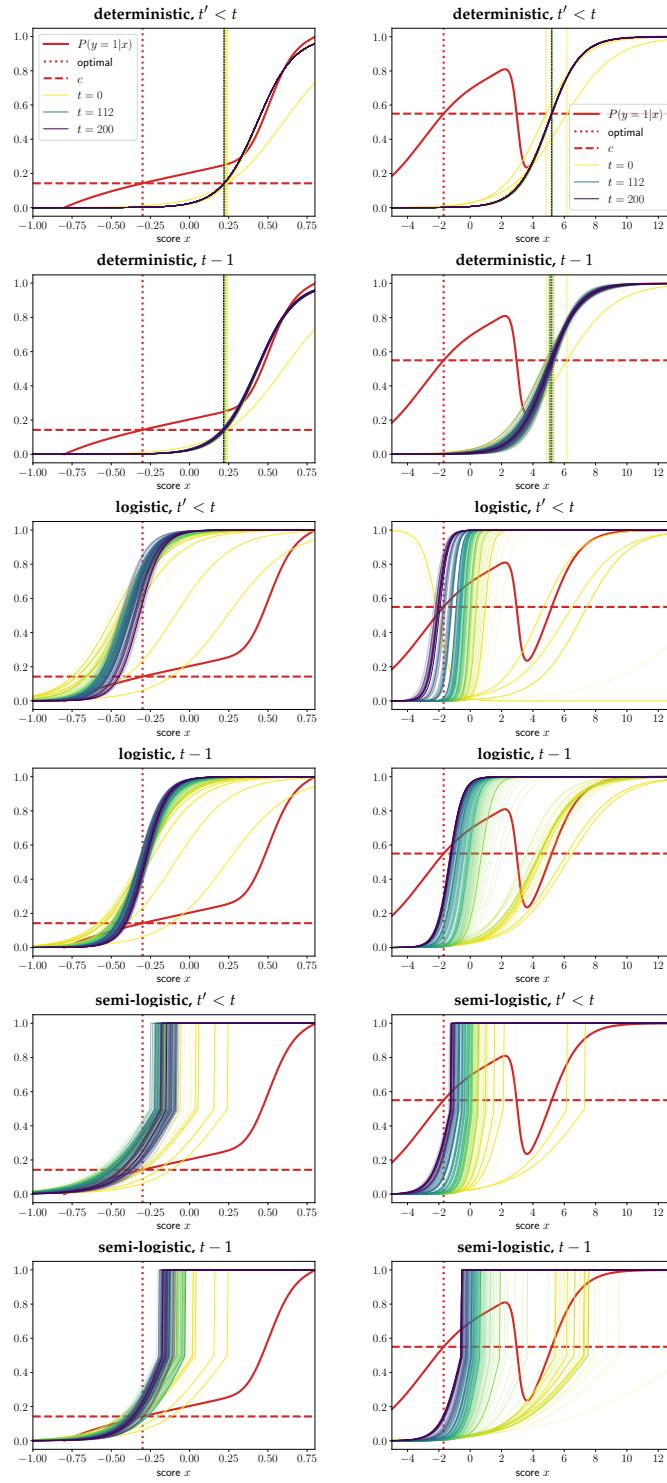
Figure B.2 Learned predictive models for deterministic threshold rules and learned decision rules for the (semi-)logistic policies. The columns correspond to the two synthetic settings. We overlay the ground truth distribution $P(Y = 1 \mid x)$ (red line), cost parameter $c$ (dashed, red), and optimal single decision boundary in $x$ within our model class (dotted, red). We describe the plots in detail in the text.
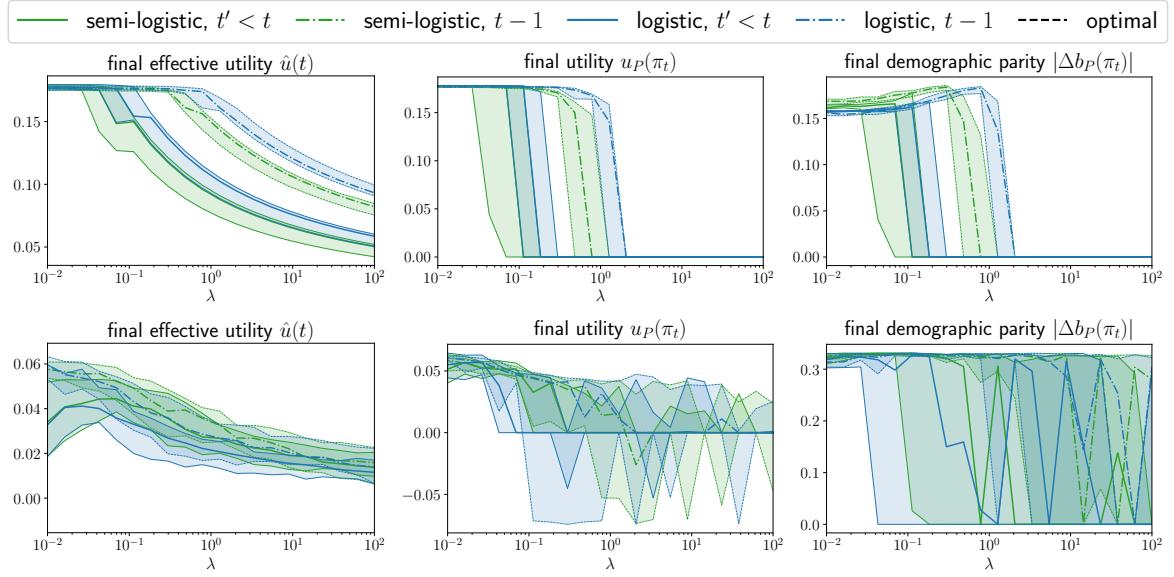
**Figure B.3** We show utility, effective utility, and demographic parity (columns) at the final time step $t = 200$ as a function of $\lambda$ where we constrain demographic parity ($f(d, y) = d$). The first row corresponds to the first setting and the second row corresponds to the second setting.

## B.2    Experiments on real data

First, in Figure B.4 we again show the contents of Figure 7.5 in the main text with shaded regions for the 25th and 75th percentile over 30 runs. Analogously to Figure B.3, we show the effect of enforcing fairness constraints in the COMPAS dataset in Figure B.5. The overall trends are similar to the results we have observed in the synthetic settings, reinforcing the applicability of our approach on real-world data.
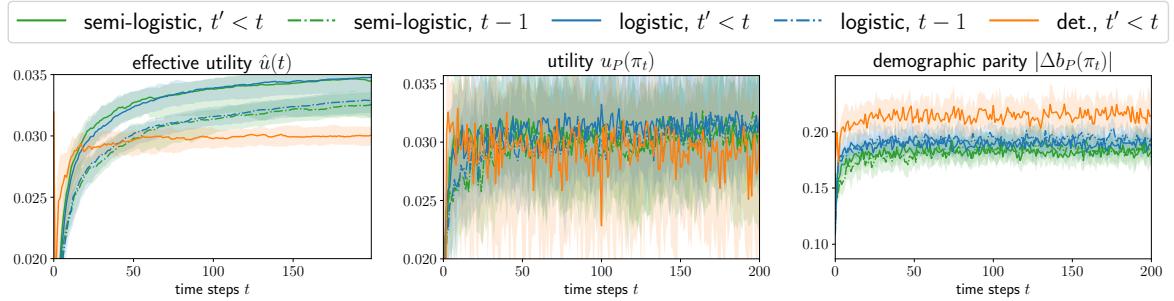


**Figure B.4** Performance on COMPAS data. We show the training progress for $\lambda = 0$, where all metrics are estimated on the held-out dataset.
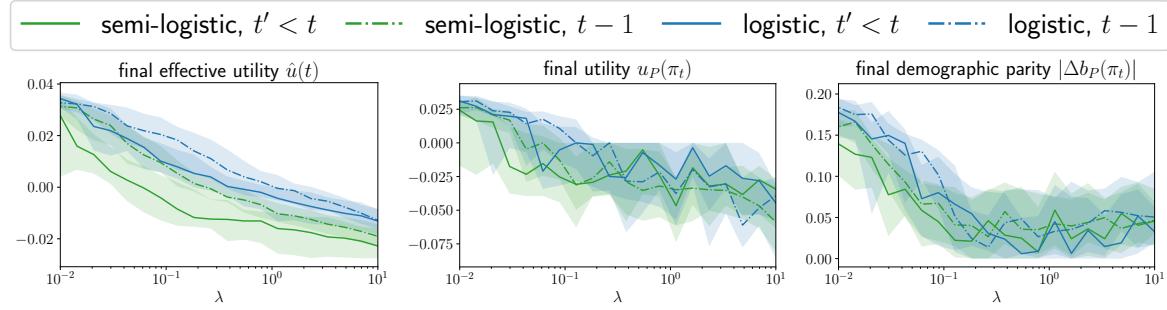
Figure B.5 We show (effective) utility and demographic parity (columns) for the COMPAS dataset at the final time step $t = 200$ estimated on the held-out dataset as a function of $\lambda$.

## B.3  Parameter settings

The parameters used for the different experiments have been found by few iterations of manual trial. The number of time steps is $T = 200$ for all datasets. For the first synthetic setting we used $\alpha = 1$, $B = 256$, $M = 128$, $N = B \cdot M$, and $c \approx 0.142$ (chosen such that the optimal decision boundary is at $x = -0.3$). For the second synthetic setting we used $\alpha = 0.5$, $B = 128$, $M = 32$, $N = B \cdot M$, and $c = 0.55$. Here we also decay the learning rate by a factor of 0.8 every 30 time steps. For the COMPAS dataset we used $\alpha = 0.1$, $B = 64$, $M = 40 \cdot B$, $N = B^2$, and $c = 0.6$. While the initialization for the synthetic settings can be seen in Figure B.2, for COMPAS we trained a logistic predictive model on 500 i.i.d. examples for initializing policies and predictive models.