

FADE: FAir Double Ensemble Learning for Observable and Counterfactual Outcomes

Alan Mishler^{1,2*} and Edward Kennedy¹

^{1*}Department of Statistics & Data Science, Carnegie Mellon University, Pittsburgh, PA, 15213, USA.

²J. P. Morgan AI Research, New York, NY, USA.

*Corresponding author(s). E-mail(s): contact@alanmishler.com;
Contributing authors: edward@stat.cmu.edu;

Abstract

Methods for building fair predictors often involve tradeoffs between fairness and accuracy and between different fairness criteria, but the nature of these tradeoffs varies. Recent work seeks to characterize these tradeoffs in specific problem settings, but these methods often do not accommodate users who wish to improve the fairness of an existing benchmark model without sacrificing accuracy, or vice versa. These results are also typically restricted to observable accuracy and fairness criteria. We develop a flexible framework for fair ensemble learning that allows users to efficiently explore the fairness-accuracy space or to improve the fairness or accuracy of a benchmark model. Our framework can simultaneously target multiple observable or counterfactual fairness criteria, and it enables users to combine a large number of previously trained and newly trained predictors. We provide theoretical guarantees that our estimators converge at fast rates. We apply our method on both simulated and real data, with respect to both observable and counterfactual accuracy and fairness criteria. We show that, surprisingly, multiple unfairness measures can sometimes be minimized simultaneously with little impact on accuracy, relative to unconstrained predictors or existing benchmark models.

Keywords: fairness, counterfactual fairness, causal inference, in-processing, post-processing

1 Introduction

Classification and regression models are increasingly widely used to inform or render decisions in domains such as healthcare, criminal justice, education, hiring, and consumer finance. Given the high-stakes nature of such decisions, it is important to ensure that these models are both accurate, to maximize their overall benefits and minimize their overall harms; and fair, so that the benefits and harms do not accrue disproportionately to already (under)privileged groups. In recent years, there have been many well-publicized cases of algorithmic systems whose performance varies over sensitive features such as race and gender in ways that appear to harm marginalized populations (Angwin and Larson, 2016; Buolamwini and Gebru, 2018; Obermeyer et al, 2019).

In response to concerns such as these, the algorithmic fairness community has developed a wide array of methods for removing or minimizing unfairness in models. In some cases, the most accurate models under consideration do not satisfy a chosen fairness criterion, so there is a fairness-accuracy tradeoff (Friedler et al, 2019; Menon and Williamson, 2018; Zhao and Gordon, 2019). Many methods therefore aim to maximize predictive accuracy subject to a bound on some quantitative unfairness criterion (Zafar et al, 2017; Donini et al, 2018; Woodworth et al, 2017). Some methods adopt a complementary perspective, seeking to minimize unfairness subject to an accuracy constraint (Zafar et al, 2017; Coston et al, 2021). Many strict versions of fairness criteria are pairwise unsatisfiable in real-world settings, so there may also be fairness-fairness tradeoffs (Chouldechova, 2017; Kleinberg et al, 2017; Kim et al, 2020).

In some cases, however, tradeoffs are small to nonexistent: model fairness can be increased with minimal loss of accuracy, or vice versa (Dutta et al, 2020; Coston et al, 2021; Rodolfa et al, 2021). Recently there has been growing interest in characterizing these tradeoffs both theoretically and empirically for specific problems and specific classes of models (Berk et al, 2017; Kim et al, 2020; Liu and Vicente, 2021). Current methods for illuminating these tradeoffs do not readily accommodate users who wish to improve the fairness and/or accuracy of an existing benchmark model rather than exploring the fairness-accuracy space. Additionally, these methods are designed to handle *observable* accuracy and fairness criteria, i.e. criteria that depend on observable outcomes. They do not address *counterfactual* accuracy and fairness criteria, which depend on counterfactual outcomes and which are relevant to many settings in which algorithms are used to support decision making. In general, there are very few methods designed to build predictors that satisfy counterfactual versions of common fairness criteria like equalized odds (Mishler et al, 2021).

To address these limitations, we propose *FAir Double Ensemble learning (FADE)*, a simple and flexible framework that builds predictors as weighted combinations of basis functions that are chosen by the user. Within this framework, we develop three methods: (1) minimizing risk subject to fairness constraints, (2) minimizing unfairness subject to a risk constraint, and (3) efficiently generating a large class of unfairness-penalized predictors. The

weights in method (3) have a closed-form expression that varies smoothly over a vector of unfairness penalty parameters, allowing users to trace out paths in fairness-accuracy spaces. It is computationally extremely fast to compute and evaluate thousands or even tens of thousands of models of this form. These methods accommodate users who wish to improve the fairness of an existing model without sacrificing accuracy, or vice versa, or who wish to understand fairness-accuracy and fairness-fairness tradeoffs in their problem.

In sum, FADE has the following properties:

- It allows users to target specific fairness and accuracy constraints as well as to efficiently explore fairness-accuracy and fairness-fairness tradeoffs.
- It can target a range of both observable and counterfactual accuracy and fairness criteria, separately or simultaneously.
- It enables users to combine previously trained and newly trained predictors, thereby collapsing the distinction between in-processing and post-processing approaches to building fair predictors.
- In the context of counterfactual accuracy and fairness, it utilizes doubly robust estimators, which yield fast convergence rates under relatively weak nonparametric assumptions. The excess risk and excess unfairness of our estimators, suitably defined, shrink to 0 at up to \sqrt{n} rates even when relevant nuisance parameters are estimated at slower rates that are typical in nonparametric machine learning.

The remainder of the paper is organized as follows. In Section 2, we discuss background and related work. In Section 3, we formalize our problem and define three estimands, all of which are formulated as optimization problems. The first estimand involves minimizing risk subject to a fairness constraint, the second involves minimizing unfairness subject to a risk constraint, and the third utilizes penalty terms rather than constraints to trace out paths in fairness-accuracy space. Section 4 provides a set of assumptions necessary to relate counterfactual quantities to observable data. In Section 5, we define our estimators and give theoretical guarantees for their performance. We illustrate FADE on simulated data in Section 6 and on real data in Sections 7 and 8.¹ In a counterfactual setting, FADE substantially improves both the fairness and accuracy of the COMPAS recidivism predictor (Section 7). In an observable setting, FADE yields many predictors that perform comparably to or better than other fairness methods on an income prediction task, while allowing users much more flexibility in the final model form (Section 8). We conclude in Section 9.

2 Background and related work

We use the terms “predictor” and “model” interchangeably to refer to any mapping from covariates to outputs that is intended to estimate an unknown

¹Code to reproduce the results will be made available in the Github repository [amishler/FAir-Double-Ensemble-Learning](https://github.com/amishler/FAir-Double-Ensemble-Learning).

quantity, whether that quantity is an unobserved label or an as-yet unrealized outcome. We use “accuracy” or “risk” to refer to any measure that tracks how well a predictor estimates the target quantity, such as mean-squared error or 0-1 error, and “performance” to refer to a model’s joint accuracy and fairness characteristics.

2.1 Ways of achieving fairness

The fairness literature generally distinguishes three approaches for developing fair predictors. *Pre-processing* approaches transform the input data to remove bias (Calmon et al, 2017; Feldman et al, 2015; Kamiran and Calders, 2012; Zemel et al, 2013). *In-processing* or *in-training* approaches enforce fairness via constraints or regularization terms during the learning process (Donini et al, 2018; Kamishima et al, 2012; Woodworth et al, 2017; Zafar et al, 2017). *Post-processing* approaches learn functions to map the outputs of existing predictors to new outputs (Hardt et al, 2016; Pleiss et al, 2017; Kim et al, 2019). Our approach enables users to combine previously existing predictors with newly trained predictors or other basis functions, essentially collapsing the distinction between in-processing and post-processing.

2.2 Observational and counterfactual fairness

Many popular fairness criteria place restrictions on the joint distribution of predictions, outcomes, and a sensitive feature. For example, the criterion of *independence*, also known as *statistical parity* or *demographic parity*, requires that the predictions be marginally independent of the sensitive feature (Calders et al, 2009; Barocas et al, 2018), while *separation* or *equalized odds* requires that they be independent conditional on the outcome (Hardt et al, 2016). Criteria such as equalized odds that depend on the outcome may be defined with respect to observable or potential (aka “counterfactual”) outcomes. Counterfactual versions of these criteria are appropriate in risk assessment settings, i.e. settings in which the model is meant to estimate the risk of an adverse outcome absent a specific intervention (Coston et al, 2020). In these settings, the potential outcomes of interest are the outcomes that would occur if, possibly counterfactually, a decision variable were set to some baseline level (Neyman, 1923; Holland, 1986). Examples arise in areas such as healthcare, where doctors must predict who would develop complications without further treatment; criminal justice, where judges must predict who would recidivate if released pretrial; and consumer finance, where banks must predict who would default if issued a loan.

A distinct set of causal fairness criteria consider counterfactuals with respect to the sensitive feature rather than with respect to a decision variable (Kilbertus et al, 2017; Kusner et al, 2017; Nabi and Shpitser, 2018; Zhang and Bareinboim, 2018; Nabi et al, 2019; Wang et al, 2019). These criteria consider questions like “what would the risk prediction be if the defendant had been of a different race their whole life?” rather than “what would the outcome be if

this person were released pretrial?” We do not consider these criteria here; see (Mishler et al, 2021), Section 3.2 for further discussion.

Most of the existing fairness literature is concerned with observable fairness criteria and accuracy measures. To our knowledge, only two papers have developed methods to satisfy the type of counterfactual criteria described above. Mishler et al (2021) developed a post-processing method that maximizes accuracy while satisfying (approximate) counterfactual equalized odds or related fairness criteria. Their approach takes as input a binary classifier and outputs a randomized binary classifier. In contrast, our method applies to both classification and regression, and it combines in-processing and post-processing.

Coston et al (2021) developed a method to minimize various unfairness measures subject to an accuracy constraint. They considered both the observable setting and a *selective labels* setting, when the outcome of interest is observed only for a non-representative subset of the population. Although the terminology differs, this is essentially equivalent to the counterfactual setting that we consider. Their method involves iterative optimization and outputs a randomized predictor that is constructed as a distribution over a set of models.

In contrast to both of the above methods, our framework can handle any of the following: (1) minimizing risk subject to fairness constraints, (2) minimizing unfairness subject to an accuracy constraint, and (3) efficiently producing a large set of models that vary in their risk and fairness properties. Method (3) utilizes a set of closed-form solutions that are extremely fast to compute. Our methods also output deterministic rather than randomized predictors, though in a classification setting these can be turned into randomized classifiers by treating the output in $[0, 1]$ as a probability and sampling from the corresponding Bernoulli distribution, as described in Section 8.1. Our methods apply to a large class of observable and counterfactual accuracy and fairness criteria.

2.3 Fairness-accuracy and fairness-fairness tradeoffs

Within a candidate set of models, the most accurate model and the most fair model may not be the same model, in which case there is a fairness-accuracy tradeoff. The shape of this tradeoff depends on the model set, the accuracy and fairness criteria, and the distribution of the data (Dutta et al, 2020). While some papers emphasize the unavoidable existence of such tradeoffs (Calders et al, 2009; Corbett-Davies et al, 2017; Menon and Williamson, 2018; Woodworth et al, 2017; Zhao and Gordon, 2019), other papers have found that in practical settings they are sometimes so small as to be irrelevant; that is, relative to a baseline model, it may be possible to substantially improve a given fairness criterion with little to no decrement in accuracy, or vice versa (Coston et al, 2021; Rodolfa et al, 2021).

Different fairness criteria may also trade off with one another. In their strictest form, many fairness criteria are mutually unsatisfiable in real-world conditions (Chouldechova, 2017; Kleinberg et al, 2017). In practice, many methods make use of continuous-valued relaxations of these criteria, which

may be more or less simultaneously satisfiable, to a degree that again depends on the modeling choices and data distribution.

Recent work aims to characterize fairness-accuracy and fairness-fairness tradeoffs both theoretically (Dutta et al, 2020; Kim et al, 2020) and empirically (Berk et al, 2017; Liu and Vicente, 2021). Like Berk et al (2017), our penalized predictor method uses fairness regularization terms to trace out different paths in fairness-accuracy space; however, their results consider observable accuracy and fairness measures, whereas ours encompass both observable and counterfactual measures. We also consider a (broad) class of fairness criteria that yield closed-form solutions, and we provide theoretical guarantees for our methods.

In some cases, users have clear accuracy or fairness constraints that they wish their models to satisfy. These constraints might derive from moral, legal, or business considerations. For example, a business might wish to ensure that a hiring algorithm generates positive recommendations for roughly equal percentages of male and female applicants, in order to avoid potential disparate impact. Conversely, a business might wish to improve the fairness of an existing model without sacrificing accuracy (profit). We provide an explicit correspondence between our constrained and penalized predictors and show how the set of penalized models can be “seeded” with models that target specific fairness or accuracy constraints.

Our method also makes it easy for users and auditors to understand whether a model in use could be made more fair without a substantial loss of accuracy, or vice versa. This is useful both for improving model performance and for understanding whether a particular level of unfairness can be justified as a type of “business necessity,” or whether fairness can be improved without compromising accuracy (Coston et al, 2021).

3 Setup and estimands

Our data is of the form $Z = (A, X, S, D, Y) \sim \mathbb{P}$, for sensitive feature $A \in \{0, 1\}$, additional covariates $X \in \mathcal{X}$, previously trained predictor(s) $S \in \mathcal{S}$, decision or treatment $D \in \mathcal{D}$, and outcome or label $Y \in [\ell_y, u_y]$ with bounds ℓ_y, u_y . If no previously trained predictors are available, then we have $S = \emptyset$. We denote by Y_i^0 the potential outcome $Y_i^{D=0}$, that is, the outcome or label that would be observed for individual i if, possibly contrary to fact, the decision were set to $D_i = 0$. For example, Y^0 could indicate whether an individual would recidivate if released pretrial. We assume that Y^0 also lies in $[\ell_y, u_y]$. In settings where we are interested in observational rather than counterfactual fairness, we may have $D = \emptyset$. We assume that $Z \subseteq \mathcal{Z} \subset \mathbb{R}^p$ for compact \mathcal{Z} .

Let $W = (A, X, S) \in \mathcal{W}$ represent the collected covariates. We let \tilde{Y} denote either Y and Y^0 as appropriate, since we are interested in both observational and counterfactual fairness and accuracy measures. We refer to $\tilde{Y} = Y$ as the *observable* setting and $\tilde{Y} = Y^0$ as the *counterfactual setting*. Broadly speaking, we seek functions of the form $f : \mathcal{W} \mapsto [\ell_y, u_y]$ that are both accurate and fair in predicting \tilde{Y} . Our goals are (1) to enable users to target specific fairness or

accuracy constraints, and (2) to trace out the fairness and accuracy properties of a large set of models, both in order to understand setting-specific fairness-accuracy and fairness-fairness tradeoffs and in order to maximize the user’s ability to choose a desirable model.

Remark 1 (Additional notation) We let $\|\cdot\|$ denote an appropriate L_2 norm. That is, for any random variable $f(Z)$ taking values in \mathbb{R} , $\|f(Z)\| = (\int (f(Z))^2 d\mathbb{P}(Z))^{1/2}$ denotes the $L_2(\mathbb{P})$ norm, while for a non-random vector $v \in \mathbb{R}^k$, $\|v\| = (\sum_{j=1}^k v_j^2)^{1/2}$ denotes the Euclidean L_2 norm. For a random vector $f(Z)$ taking values in \mathbb{R}^k , $\|f(Z)\| = (\sum_{j=1}^k \|f_j(Z)\|^2)^{1/2}$.

3.1 FADE summary

The “Ensemble learning” part of “FAir Double Ensemble learning” has its usual sense, referring to an ensemble of predictors. The “Double” part captures several features of our approach: (1) it combines in-processing and post-processing; (2) it accommodates both observable and counterfactual outcomes; (3) it (optionally) has two stages, first learning predictors and then learning their ensemble weights; and (4) in the counterfactual setting, it utilizes doubly robust estimators, which also appear in the literature under the heading “double machine learning” (Tsiatis, 2006; Chernozhukov et al, 2018). These features are illustrated in the remainder of this section and in Section 5.

3.2 Accuracy and fairness measures

The risk (accuracy) measure we consider is the MSE:

$$\text{Risk}(f) = \mathbb{E}[(f(W) - \tilde{Y})^2]$$

We consider (un)fairness measures $\text{UF}(f)$ that can be expressed in the form

$$\text{UF}(f) = |\mathbb{E}[g(W, \tilde{Y})f(W)]| \tag{1}$$

where $g(W, \tilde{Y})$ is a bounded *fairness function* that depends only on W and \tilde{Y} , not on D . This accommodates a broad range of measures, including measures described by the following proposition. All proofs are given in the Appendix.

Proposition 1 *Let $\alpha_0, \alpha_1 \in \mathbb{R}$ and let h_0, h_1 be mappings from $\{0, 1\} \times \tilde{Y}$ to $\{0, 1\}$. Consider an unfairness measure*

$$\text{UF}(f) = |\alpha_0 \mathbb{E}[f(W) \mid h_0(A, \tilde{Y}) = 1] - \alpha_1 \mathbb{E}[f(W) \mid h_1(A, \tilde{Y}) = 1]|$$

and assume that $\mathbb{P}(h_0(A, \tilde{Y}) = 1) > 0, \mathbb{P}(h_1(A, \tilde{Y}) = 1) > 0$. Then there exists a fairness function $g(W, \tilde{Y})$ such that $\text{UF}(f) = |\mathbb{E}[g(W, \tilde{Y})f(W)]|$, namely

$$g(W, \tilde{Y}) = \alpha_0 \frac{h_0(A, \tilde{Y})}{\mathbb{E}[h_0(A, \tilde{Y})]} - \alpha_1 \frac{h_1(A, \tilde{Y})}{\mathbb{E}[h_1(A, \tilde{Y})]}$$

That is, (1) is compatible with any fairness measure that can be expressed as a (weighted) difference of average predictions conditioned on events that are a function of the sensitive feature and the outcome. Functions of this form must in general be estimated, since they depend on unknown expected values. We focus in this paper on the following measures, which we refer to equivalently as disparities. We first express each in a canonical form, and then we identify the corresponding fairness function $g(W, \tilde{Y})$, dropping the arguments (W, \tilde{Y}) for convenience.

Definition 1 The *rate disparity (rate-diff)* is

$$|\mathbb{E}[f(W) | A = 0] - \mathbb{E}[f(W) | A = 1]|$$

with fairness function

$$g^{\text{rate}} = \frac{1 - A}{\mathbb{E}[1 - A]} - \frac{A}{\mathbb{E}[A]}$$

Definition 2 For $\tilde{Y} \in \{0, 1\}$, the *generalized False Positive Rate disparity (FPR-diff)* is

$$|\mathbb{E}[f(W) | A = 0, \tilde{Y} = 0] - \mathbb{E}[f(W) | A = 1, \tilde{Y} = 0]|$$

with fairness function

$$g^{\text{FPR}} = \frac{(1 - \tilde{Y})(1 - A)}{\mathbb{E}[(1 - \tilde{Y})(1 - A)]} - \frac{(1 - \tilde{Y})A}{\mathbb{E}[(1 - \tilde{Y})A]}$$

Definition 3 For $\tilde{Y} \in \{0, 1\}$, the *generalized False Negative Rate disparity (FNR-diff)* is

$$|\mathbb{E}[1 - f(W) | A = 0, \tilde{Y} = 1] - \mathbb{E}[1 - f(W) | A = 1, \tilde{Y} = 1]|$$

with fairness function

$$g^{\text{FNR}} = \frac{\tilde{Y}A}{\mathbb{E}[\tilde{Y}A]} - \frac{(1 - \tilde{Y})(1 - A)}{\mathbb{E}[\tilde{Y}(1 - A)]}$$

These definitions are closely related to common fairness criteria described in Section 2.2. The criterion of *independence* requires predictions $f(W)$ to be independent of the sensitive feature A . Rate-diff measures violations of this criterion (Calders and Verwer, 2010). Equal opportunity requires the false negative rates to be equal across the two groups, while equalized odds requires both the false positive and the false negative rates to be equal (Hardt et al, 2016). FNR-diff therefore measures violations of equal opportunity, while FPR-diff and FNR-diff together measure violations of equalized odds.

For continuous-valued predictors, it may be challenging to attain full (conditional) independence. Hence it is common to focus only on average conditional predictions (e.g. Corbett-Davies et al, 2017).

3.3 Predictor classes

We consider predictors that lie in the linear span of a set of basis functions $b = b(W) = (b_1(W), \dots, b_k(W))$, where each function $b_j(W)$ maps from \mathcal{W} to \mathbb{R} . That is, for a given b we seek predictors in the set \mathcal{F}_b , where

$$\mathcal{F}_b = \{b^T \beta : \beta \in \mathbb{R}^k\}$$

Predictors of this form are commonly referred to as (a linear) ensemble, stacked predictors, or aggregated predictors (Breiman, 1996; Juditsky and Nemirovski, 2000; Tsybakov, 2003; Polley and Van Der Laan, 2010). In the context of our method, we refer to these as FADE predictors. The vector b is determined by the user. It can include for example previously trained predictors S , newly trained predictors, or arbitrary orthogonal basis functions such as trigonometric functions or polynomials. Our approach makes it easy for users to search across a range of different bases b .

In this paper, we consider a regime in which $k < n$, where n is the sample size, since this simplifies estimation. In practice, users might wish to use bases of dimension $k \geq n$, such as spline bases or kernel basis functions. We briefly consider these and other possibilities in Appendix F. In our asymptotic analyses, we also generally assume that the basis is eventually fixed, meaning $k \not\rightarrow \infty$. We intend to analyze settings in which $k \geq n$ and/or $k \rightarrow \infty$ in future work.

Depending on b , the set \mathcal{F}_b may be relatively rich. For example, b could be a truncated orthonormal basis of the space $L_2(\mathcal{W})$ of square-integrable functions, in which case \mathcal{F}_b could approximate L_2 to a degree chosen by the user.

Our theoretical analyses utilize the following two assumptions about the basis.

Assumption 1 (PSD outer product). *Uniformly in n , the eigenvalues of $\mathbb{E}[bb^T]$ are bounded above and away from 0.*

This assumption asserts that the basis functions $b_1(W), \dots, b_k(W)$ are not too collinear. It means that $\mathbb{E}[bb^T]$ is always positive semi-definite. In a regime in which the basis is eventually fixed, this assumption simply requires that the basis functions are never perfectly collinear.

Assumption 2 (Bounded basis norm). *Uniformly in n , $\sup_{w \in \mathcal{W}} \|b(w)\| < \infty$, where $\|b(w)\|$ is the Euclidean L_2 norm.*

When $k \not\rightarrow \infty$, this assumption simply requires the norm $\|b(w)\|$ to be finite over the covariate space. In a setting in which k is allowed to grow to infinity, this assumption can be relaxed to one which controls the growth rate of $\sup_{w \in \mathcal{W}} \|b(w)\|$. See Remark 7 in Section 5.

3.4 Estimands

We first define two estimands that are solutions to constrained least squares problems. These estimands represent users who have clear target fairness or accuracy constraints. We then show that these estimands can be equivalently expressed via penalized least squares problems that admit closed-form solutions. These solutions are indexed by an unfairness penalty parameter; by varying this parameter, we may trace out curves in accuracy-fairness space over \mathcal{F}_b .

Suppose that there are t fairness measures that can be expressed via fairness functions $g_j, j = 1, \dots, t$. For a given k -dimensional basis b , define the risk-minimization (*risk-min*) parameter β_r^* and the unfairness-minimization (*unfair-min*) parameter β_u^* as the solutions to the following optimization problems:

$$\begin{aligned} \beta_r^* &= \arg \min_{\beta \in \mathbb{R}^k} \mathbb{E}[(b^T \beta - \tilde{Y})^2] \\ &\text{subject to } (\mathbb{E}[g_j b^T \beta])^2 \leq \epsilon_j^2, \quad j = 1, \dots, t \\ \beta_u^* &= \arg \min_{\beta \in \mathbb{R}^k} \sum_{j=1}^t \alpha_j (\mathbb{E}[g_j b^T \beta])^2 \\ &\text{subject to } \mathbb{E}[(b^T \beta - \tilde{Y})^2] \leq \epsilon \end{aligned}$$

for user-chosen constraints $\epsilon_j \geq 0, \epsilon > 0$, and weights α_j . That is, β_r^* indexes the most accurate predictor in \mathcal{F}_b among those that satisfy t specified fairness constraints, and β_u^* indexes the most fair predictor among those that satisfy a specified risk constraint.

We constrain $\epsilon > 0$ because otherwise we'd be insisting on a perfectly accurate predictor, which is generally impossible in practice. The risk-min problem is always feasible with $\epsilon_j \geq 0$, since the predictor defined by $\beta = 0$ always satisfies the fairness constraints. Under Assumption 1, β_r^* is unique, since the objective is strictly convex. The unfair-min problem may be infeasible, if there is no predictor in \mathcal{F}_b whose risk is less than or equal to ϵ . This may not be an issue in practice, if ϵ represents (an estimate of) the risk of an existing benchmark model. With slight modifications, all our subsequent results would carry through if this constraint were explicitly expressed with respect to a benchmark model; for the sake of simplicity, however, we leave it in this form. If the unfairness-minimization problem is feasible and $\sum_{j=1}^t \alpha_j \mathbb{E}[g_j b] \mathbb{E}[g_j b]^T$ is positive definite, then β_u^* is unique, since the objective is strictly convex.

Note that the fairness constraints in the risk-min problem can be equivalently written in affine form, as $|\mathbb{E}[g_j f(W)]| \leq \epsilon_j$, meaning that β_r^* is the solution to a quadratic program. We express the constraints in squared form for notational consistency with the penalized estimand, which is defined as follows.

For any $\lambda = (\lambda_1, \dots, \lambda_t)$, with all $\lambda_j \geq 0, j = 1, \dots, t$, define the penalized-minimization (*penalized-min*) estimand β_λ^* as:

$$\beta_\lambda^* = \arg \min_{\beta \in \mathbb{R}^k} \mathbb{E}[(b^T \beta - \tilde{Y})^2] + \sum_{j=1}^t \lambda_j (\mathbb{E}[g_j b^T \beta])^2$$

This can be written in an equivalent closed form:

$$\beta_\lambda^* = \left(\mathbb{E}[bb^T] + \sum_{j=1}^t \lambda_j \mathbb{E}[g_j b] \mathbb{E}[g_j b]^T \right)^{-1} \mathbb{E}[\tilde{Y} b]$$

Under Assumption 1, the matrix inverse always exists, since each matrix $\mathbb{E}[g_j b] \mathbb{E}[g_j b]^T$ is positive semi-definite and $\mathbb{E}[bb^T]$ is positive definite, and hence by Weyl's inequality the entire matrix is positive definite.

We now establish a correspondence between the constrained and penalized forms. Let \mathcal{I} denote the set of active fairness constraints at β_r^* , that is, $\mathcal{I} = \{j \in \{1, \dots, t\} : (\mathbb{E}[g_j b^T \beta_r^*])^2 = \epsilon_j^2\}$.

Assumption 3 (LICQ). *The set of vectors $\{\mathbb{E}[g_j b] \mathbb{E}[g_j b]^T \beta_r^* : j \in \mathcal{I}\}$ is linearly independent.*

Assumption 3, which is expected to hold in practical settings, is the Linear Independence Constraint Qualification (LICQ), which yields an injective mapping from the constrained to the penalized form.

Proposition 2 *Under Assumption 1, for any β_r^* there exists a $\lambda \in \mathbb{R}_{0+}^t$ such that $\beta_\lambda^* = \beta_r^*$. If Assumption 3 holds, then this λ is unique.*

We will utilize the penalized form to efficiently construct a large set of predictors that vary in their accuracy and fairness properties. We will see that we can exploit an empirical analogue of Proposition 2 to “seed” this set with models that target specified fairness constraints.

Proposition 3 *Fix $\lambda \in \mathbb{R}_{0+}^t, \lambda \neq 0$. Under Assumption 1, $\beta_\lambda^* = \beta_r^*$ with fairness constraints $\epsilon_j^2 = (\mathbb{E}[g_j b^T \beta_\lambda^*])^2$.*

Proposition 3 expresses the converse direction of the relationship between β_r^* and β_λ^* . The constraints ϵ_j that define β_r^* are arguably easier to reason about than the penalties λ_j that define β_λ^* . This proposition therefore facilitates interpretation of penalized estimands in terms of their corresponding constrained forms.

An analogous penalized form can be written that corresponds to β_u^* , with results that match Propositions 2 and 3. For estimation purposes, however, we focus in this paper on β_λ^* , so we do not develop that form here.

Remark 2 (Predictor truncation) Any $\beta \in \mathbb{R}^k$ indexes a predictor $b^T \beta \in \mathcal{F}_b$. Since Y^0 and Y are bounded in $[\ell_u, u_y]$, however, the resulting predictor may be truncated to lie in $[\ell_u, u_y]$, if the bounds are known.

Our final estimands consist of the risk and (un)fairness properties of any fixed predictor f_β :

$$\begin{aligned} \text{Risk}(f_\beta) &= \mathbb{E}[(f_\beta - \tilde{Y})^2] \\ \text{UF}_j(f_\beta) &= \mathbb{E}[g_j f_\beta], \quad j = 1, \dots, t \end{aligned}$$

In particular, once we have computed some estimate $\hat{\beta}$ of β_r^* or β_u^* , or β_λ^* , it is of interest to estimate the risk and fairness of the resulting predictor $f_{\hat{\beta}}$.

4 Identification

When $\tilde{Y} = Y^0$, i.e. when the risk and fairness functions are defined with respect to counterfactual rather than observable outcomes, we require assumptions in order to identify these quantities in terms of the observed data. For ease of notation, we first define three nuisance parameters that appear in the estimands and associated estimators.

$$\begin{aligned} \pi &= \pi(W) = \mathbb{P}(D = 1 \mid W) \\ \mu_0 &= \mu_0(W) = \mathbb{E}[Y \mid W, D = 0] \\ \nu_0 &= \nu_0(W) = \mathbb{E}[Y^2 \mid W, D = 0] \end{aligned}$$

$\pi(W)$ is the propensity score, while μ_0 and ν_0 are regressions with respect to the observed outcome and the squared observed outcome. In a classification setting with $Y \in \{0, 1\}$, we have $Y^2 = Y$, so $\nu_0 = \mu_0$. We make the following common causal inference assumptions:

Assumption 4 (Consistency). $Y = DY^1 + (1 - D)Y^0$.

Assumption 5 (Positivity). $\exists \delta \in (0, 1)$ s.t. $\mathbb{P}(\pi(W) \leq 1 - \delta) = 1$.

Assumption 6 (Ignorability). $Y^0 \perp\!\!\!\perp D \mid W$.

Consistency signifies that for each individual, the treatment received matches the outcome that is observed, meaning for example that one person's treatment status does not affect other people's outcomes. Positivity requires

that within covariate strata W , individuals have some chance of not receiving treatment, meaning that there is no stratum of measure > 0 in which all individuals are guaranteed to receive treatment. Finally, ignorability precludes unmeasured confounders that affect both treatment status and the potential outcome. Positivity and ignorability may be satisfied in randomized experiments in which treatment is assigned (conditionally) at random, or in observational studies given an appropriate set of covariates W .

Note that these assumptions may not hold exactly in practice. For example, in an observational study, the measured covariates may not be sufficient to fully deconfound D and Y^0 . The enterprise of *sensitivity analysis* in causal inference parameterizes violations of these assumptions and models their effect on downstream estimation (Rosenbaum, 1987; Liu et al, 2013; Richardson et al, 2014; Bonvini and Kennedy, 2021). We leave sensitivity analysis in our setting for future work.

For convenience, we also define the following:

$$\begin{aligned}\phi &= \phi(Z) = \frac{1-D}{1-\pi}(Y - \mu_0) + \mu_0 \\ \underline{\phi} &= \underline{\phi}(Z) = \frac{1-D}{1-\pi}(Y^2 - \nu_0) + \nu_0\end{aligned}$$

Under the identifying assumptions, these are the uncentered influence functions for $\mathbb{E}[Y^0]$ and $\mathbb{E}[(Y^0)^2]$, respectively (Bickel et al, 1993; van der Laan and Robins, 2003; Tsiatis, 2006; Kennedy, 2016).

Proposition 4 *Under Assumptions 5–6, the counterfactual risk, FPR-diff, and FNR-diff for any function $f : \mathcal{W} \mapsto \mathbb{R}$ are identified as follows.*

$$\mathbb{E}[(f - Y^0)^2] = \mathbb{E}[(f - \mu_0)^2] + \text{var}(Y^0) \quad (2)$$

$$= \mathbb{E}[f^2 - 2f\mu_0 + \nu_0] \quad (3)$$

$$\mathbb{E}[g^{cFPR} f(W)] = \mathbb{E} \left[\left\{ \frac{(1-\mu_0)(1-A)}{\mathbb{E}[(1-\mu_0)(1-A)]} - \frac{(1-\mu_0)A}{\mathbb{E}[(1-\mu_0)A]} \right\} f(W) \right]$$

$$\mathbb{E}[g^{cFNR} f(W)] = \mathbb{E} \left[\left\{ \frac{\mu_0 A}{\mathbb{E}[\mu_0 A]} - \frac{(1-\mu_0)(1-A)}{\mathbb{E}[\mu_0(1-A)]} \right\} f(W) \right]$$

These expressions also hold if μ_0 is replaced with ϕ and ν_0 is replaced with $\underline{\phi}$.

We do not include the rate-diff, since this involves only the sensitive feature and the decision variable, not outcomes, and is therefore trivially identified.

Remark 3 (Multiple risk expressions) Expressions (2) and (3) show that when estimating the risk-min parameter β_r^* , we can either minimize an estimate of $\mathbb{E}[(f - \mu_0)^2]$ or an estimate of $\mathbb{E}[f^2 - 2f\mu_0]$; the terms $\text{var}(Y^0)$ and ν_0 are constant with respect to f and so drop out of the minimization. The nuisance parameter ν_0 will only be required when we wish to estimate the actual risk of a given predictor, as well as when estimating the unfair-min parameter $\hat{\beta}_u$, since that involves a constraint on the

actual risk. Note that ν_0 would not be required to solve unfair-min if the accuracy constraint were defined with respect to an existing benchmark model, since the two ν_0 terms in the constraint would cancel out.

5 Estimation

We require a training set $\mathcal{D}_{\text{train}}$, which is used to construct estimates $\hat{\beta}$ of the optimal weights β_r^* or β_u^* or β_λ^* , and a test set $\mathcal{D}_{\text{test}}$, which is used to estimate the risk and fairness values of the resulting predictor(s) $f_{\hat{\beta}}$. If the user wishes to train new basis predictors, then an additional dataset $\mathcal{D}_{\text{learn}}$ is also required. This is not needed if the user is only aggregating arbitrary basis functions, like trigonometric functions, or previously existing predictors S .

In order to obtain fast rates for our estimators, in the counterfactual setting we split $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$ into separate folds for estimating the nuisance parameters and the target parameters. The sample splitting scheme is shown in Figure 1. For simplicity, we illustrate a single split, but in practice cross-fitting can be used within each dataset.

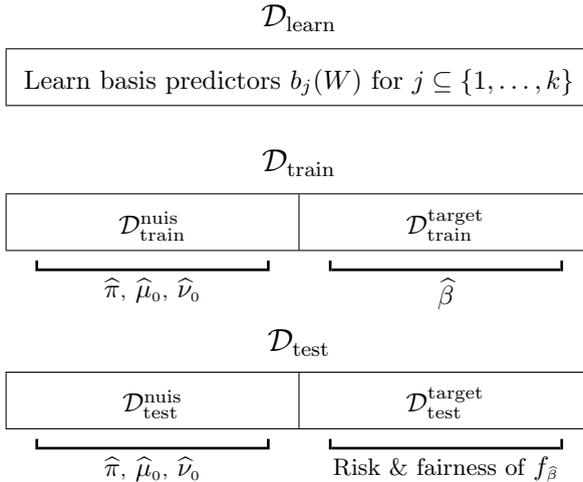


Fig. 1: Sample splitting scheme. $\mathcal{D}_{\text{learn}}$ is not needed if all the basis functions already exist. $\mathcal{D}_{\text{train}}^{\text{nuis}}$ and $\mathcal{D}_{\text{test}}^{\text{nuis}}$ are used to estimate nuisance parameters, while $\mathcal{D}_{\text{train}}^{\text{target}}$ and $\mathcal{D}_{\text{test}}^{\text{target}}$ are used to estimate target parameters. $\hat{\beta}$ represents a weight vector that is an estimate of β_r^* or β_u^* or β_λ^* , while $f_{\hat{\beta}}$ represents the predictor indexed by $\hat{\beta}$. Splitting $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$ is only required in the counterfactual setting, since there are no nuisance parameters in the observable setting. In practice, cross-fitting may be used within both $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$.

We solve empirical versions of the identified minimization problems that define the estimands. Let $\widehat{\phi}, \widehat{\underline{\phi}}$ denote estimates of ϕ and $\underline{\phi}$ constructed from estimates $\widehat{\pi}, \widehat{\mu}_0, \widehat{\nu}_0$.

For any fixed function $f : \mathcal{Z} \mapsto \mathbb{R}$, let $\mathbb{P}_n(f(Z)) = n^{-1} \sum_{i=1}^n f(Z)$ and $\mathbb{P}(f) = \int f d\mathbb{P}(Z)$ denote the sample and population expectations of f , so that for example $\mathbb{P}(\phi) = \mathbb{E}[\phi]$ while $\mathbb{P}(\widehat{\phi}) = \mathbb{E}[\widehat{\phi} \mid \mathcal{D}_{\text{train}}]$ or $\mathbb{E}[\widehat{\phi} \mid \mathcal{D}_{\text{test}}]$ is the expected value of $\widehat{\phi}(Z)$ once the relevant nuisance function estimate $\widehat{\phi}$ has been constructed. That is, $\mathbb{P}(\widehat{\phi})$ is a random variable that depends on the nuisance data, while $\mathbb{P}_n(\widehat{\phi})$ is a random variable that depends on both the nuisance and target data. In order to avoid excess notation, for quantities like $\mathbb{P}(\widehat{\phi})$ and $\mathbb{P}_n(\widehat{\phi})$, we will rely on context to make it clear whether $\widehat{\phi}$ depends on $\mathcal{D}_{\text{train}}$ or $\mathcal{D}_{\text{test}}$.

For notational convenience, let \widehat{g}_j with no arguments denote $g_j(W, Y)$ in the observable setting and $g_j(W, \widehat{\phi})$ in the counterfactual setting. That is $\widehat{g}_j = g_j$ in the observable setting, since there is no nuisance quantity to estimate, but $\widehat{g}_j \neq g_j$ in the counterfactual setting. The occasional use of \widehat{g}_j for both settings allows us to concisely state certain conditions and results.

Assumption 7 (Bounded propensity estimator). $\exists \gamma \in (0, 1)$ s.t. $\mathbb{P}(\widehat{\pi}(A, X, S) \leq 1 - \gamma) = 1$.

Assumption 7 is the empirical analogue of the positivity assumption (5). It can be trivially satisfied by truncating $\widehat{\pi}$ at $1 - \delta$, the positivity threshold in Assumption 5.

Assumption 8 (Consistent nuisance estimators). $\|\widehat{\pi} - \pi\| = o_{\mathbb{P}}(1)$ and $\|\widehat{\mu}_0 - \mu_0\| = o_{\mathbb{P}}(1)$ and $\|\widehat{\nu}_0 - \nu_0\| = o_{\mathbb{P}}(1)$.

This assumption is reasonable if nonparametric methods are used to construct the nuisance parameter estimates. With slight procedural modifications, this assumption can be relaxed to require consistency in the influence function estimators $\widehat{\phi}$ and $\widehat{\underline{\phi}}$ without necessarily requiring consistency in each of the nuisance parameter estimators. For simplicity, we do not address this.

5.1 Constrained FADE estimators

The risk-min and unfair-min estimators $\widehat{\beta}_r$ and $\widehat{\beta}_u$ are defined for the observable and counterfactual settings in Tables 1 and 2.

As with the corresponding estimands β_r^* and β_u^* , the optimization problem that defines $\widehat{\beta}_r$ is always feasible, while the problem that defines $\widehat{\beta}_u$ may not be, if there is no predictor in \mathcal{F}_b with estimated risk less than or equal to ϵ . If $\mathbb{P}_n[bb^T]$ is positive definite, then $\widehat{\beta}_r$ is unique, since the objective is strictly convex. Under Assumption 1, this will hold with probability approaching 1 in n , or

Observable ($\tilde{Y} = Y$)	Counterfactual ($\tilde{Y} = Y^0$)
$\hat{\beta}_r = \arg \min_{\beta \in \mathbb{R}^k} \mathbb{P}_n[(b^T \beta - Y)^2]$ $\text{s.t. } \left(\mathbb{P}_n[g_j(W, Y)b^T \beta] \right)^2 \leq \epsilon_j^2, j = 1, \dots, t$	$\hat{\beta}_r = \arg \min_{\beta \in \mathbb{R}^k} \mathbb{P}_n[(b^T \beta - \hat{\phi})^2]$ $\text{s.t. } \left(\mathbb{P}_n[g_j(W, \hat{\phi})b^T \beta] \right)^2 \leq \epsilon_j^2, j = 1, \dots, t$

Table 1: Definition of the unfair-min estimator $\hat{\beta}_r$ in the observable and counterfactual settings.

Observable ($\tilde{Y} = Y$)	Counterfactual ($\tilde{Y} = Y^0$)
$\hat{\beta}_u = \arg \min_{\beta \in \mathbb{R}^k} \sum_{j=1}^t \alpha_j \left(\mathbb{P}_n[g_j(W, Y)b^T \beta] \right)^2$ $\text{s.t. } \mathbb{P}_n \left[(b^T \beta - Y)^2 \right] \leq \epsilon^2$	$\hat{\beta}_u = \arg \min_{\beta \in \mathbb{R}^k} \sum_{j=1}^t \alpha_j \left(\mathbb{P}_n[g_j(W, \hat{\phi})b^T \beta] \right)^2$ $\text{s.t. } \mathbb{P}_n \left[(b^T \beta)^2 - 2(b^T \beta)\hat{\phi} + \hat{\phi} \right] \leq \epsilon^2$

Table 2: Definition of the risk-min estimator $\hat{\beta}_u$ in the observable and counterfactual settings.

with probability 1 if, say, at least one of the covariates in W is continuously distributed. If the problem that defines $\hat{\beta}_u$ is feasible and $\sum_{j=1}^t \alpha_j \mathbb{P}_n[\hat{g}_j b] \mathbb{P}_n[\hat{g}_j b]^T$ is positive definite, then $\hat{\beta}_u$ is unique, since the objective is strictly convex.

We next consider the excess risk and the excess unfairness for the constrained predictors. In the counterfactual setting, we require assumptions on the rate at which the nuisance parameters are estimated.

Assumption 9 (Nuisance parameter rates).

$$\begin{aligned} \|\hat{\pi} - \pi\| \|\hat{\mu}_0 - \mu_0\| &= o_{\mathbb{P}}(1/\sqrt{n}) \\ \|\hat{\pi} - \pi\| \|\hat{\nu}_0 - \nu_0\| &= o_{\mathbb{P}}(1/\sqrt{n}) \end{aligned}$$

Assumption 9 says that the product of errors in the nuisance parameter estimators goes to 0 faster than \sqrt{n} . That can be satisfied for example if the nuisance parameters are estimated at faster than $n^{1/4}$ rates, which can be achieved in nonparametric settings under appropriate smoothness or sparsity conditions (Györfi et al, 2002; Raskutti et al, 2011).

Definition 4 The *excess risk* is defined for $\widehat{\beta}_r$ and $\widehat{\beta}_u$ as:

$$\mathbb{P}[(b^T \widehat{\beta}_r - \widetilde{Y})^2] - \mathbb{P}[(b^T \beta_r^* - \widetilde{Y})^2] \quad (\text{risk-min})$$

$$\mathbb{P}[(b^T \widehat{\beta}_u - \widetilde{Y})^2] - \epsilon^2 \quad (\text{unfair-min})$$

Theorem 1 (Excess risk in the constrained setting) *Under Assumptions 1–2 for the observable setting, and Assumptions 1–2 and 4–9 for the counterfactual setting:*

$$\mathbb{P}[(b^T \widehat{\beta}_r - \widetilde{Y})^2] - \mathbb{P}[(b^T \beta_r^* - \widetilde{Y})^2] = O_{\mathbb{P}}(1/\sqrt{n}) \quad (\text{risk-min})$$

$$\mathbb{P}[(b^T \widehat{\beta}_u - \widetilde{Y})^2] - \epsilon^2 = O_{\mathbb{P}}(1/\sqrt{n}) \quad (\text{unfair-min})$$

Definition 5 The *excess unfairness* for $\widehat{\beta}_r$ and $\widehat{\beta}_u$ is defined as:

$$\max_{j=1, \dots, t} \left\{ \left((\mathbb{P}[g_j b^T \widehat{\beta}_r])^2 - \epsilon_j^2 \right)_+ \right\} \quad (\text{risk-min})$$

$$\sum_{j=1}^t \alpha_j \left\{ (\mathbb{P}[g_j b^T \widehat{\beta}_u])^2 - (\mathbb{P}[g_j b^T \beta_u^*])^2 \right\} \quad (\text{unfair-min})$$

where $(\cdot)_+ = \max\{\cdot, 0\}$ denotes the positive part function.

Theorem 2 (Excess unfairness in the constrained setting) *Under Assumptions 1–2 for the observable setting, and Assumptions 1–2 and 4–9 for the counterfactual setting:*

$$\max_{j=1, \dots, t} \left\{ \left((\mathbb{P}[g_j b^T \widehat{\beta}_r])^2 - \epsilon_j^2 \right)_+ \right\} = O_{\mathbb{P}}(1/\sqrt{n}) \quad (\text{risk-min})$$

$$\sum_{j=1}^t \alpha_j \left\{ (\mathbb{P}[g_j b^T \widehat{\beta}_u])^2 - (\mathbb{P}[g_j b^T \beta_u^*])^2 \right\} = O_{\mathbb{P}}(1/\sqrt{n}) \quad (\text{unfair-min})$$

These results show that if a user has specific fairness or risk constraints in mind, in the observable setting, they can generate a predictor in a rich linear space that is asymptotically guaranteed to meet these constraints, while minimizing the corresponding risk or unfairness. In the counterfactual setting, they can do the same thing with the same fast rate guarantees, as long as the nuisance parameters are estimated at fast enough rates.

Of course, any particular estimates $\widehat{\beta}_r, \widehat{\beta}_u$ may violate their target risk and fairness constraints by arbitrary amounts, since the constraints used to compute them are themselves estimated. Suppose that $\widehat{\beta}_r$ was evaluated on the test set, and one of its estimated unfairness values was found to exceed the constraint ϵ_j by an unacceptable amount. To remedy this, the user could lower the value of ϵ_j and compute a new $\widehat{\beta}_r$ under this more stringent constraint. They could repeat this process until they found a $\widehat{\beta}_r$ with acceptable estimated fairness. Since $\widehat{\beta}_r$ is the solution to a quadratic program, however, this is computationally costly, and there is no guarantee that additional searching will yield improvements. A predictor that is more fair with respect to one

fairness constraint may be *less* fair with respect to other constraints, or may incur unacceptable additional risk.

Ideally, the user might wish to treat the fairness constraints as tuning parameters, selecting a large set of constraint vectors $(\epsilon_1, \dots, \epsilon_t) \in \mathbb{R}_{0+}^t$, computing $\widehat{\beta}_r$ for each vector, and comparing the risk and fairness properties of all the resulting predictors. In the next section, we use the closed-form penalized estimators to accomplish something equivalent to this, with trivial additional computational cost.

5.2 Penalized FADE estimators

In the observable and counterfactual settings, the estimator $\widehat{\beta}_\lambda$ takes the following equivalent forms, which mirror the two expressions given for β_λ^* :

$$\begin{aligned} \widehat{\beta}_\lambda &= \arg \min_{\beta \in \mathbb{R}^k} \mathbb{P}_n[(b^T \beta - Y)^2] + \sum_{j=1}^t \lambda_j (\mathbb{P}_n[g_j b^T \beta])^2 && \text{(Observable)} \\ &= \left(\mathbb{P}_n(bb^T) + \sum_{j=1}^t \lambda_j \mathbb{P}_n(g_j b) \mathbb{P}_n(g_j b)^T \right)^{-1} \mathbb{P}_n(bY) \end{aligned}$$

$$\begin{aligned} \widehat{\beta}_\lambda &= \arg \min_{\beta \in \mathbb{R}^k} \mathbb{P}_n[(b^T \beta - \widehat{\phi})^2] + \sum_{j=1}^t \lambda_j (\mathbb{P}_n[\widehat{g}_j b^T \beta])^2 && \text{(Counterfactual)} \\ &= \left(\mathbb{P}_n(bb^T) + \sum_{j=1}^t \lambda_j \mathbb{P}_n(\widehat{g}_j b) \mathbb{P}_n(\widehat{g}_j b)^T \right)^{-1} \mathbb{P}_n(b\widehat{\phi}) \end{aligned}$$

assuming that the relevant matrix inverse exists. A sufficient condition for it to exist is that $\mathbb{P}_n[bb^T]$ is positive definite, which, as discussed above, will happen with probability 1 or approaching 1 under Assumption 1.

The procedure we propose is given in Figure 2. The user first chooses a large set of vectors $\Lambda_n \subset \mathbb{R}_{0+}^t$, which we assume may depend on sample size. They then compute the solution set $\widehat{\mathcal{B}}_n = \{\widehat{\beta}_\lambda : \lambda \in \Lambda_n\}$, estimate the risk and fairness properties of each $f_\beta : \beta \in \widehat{\mathcal{B}}_n$, and select a predictor with a favorable performance profile.

1. Pick a large set of vectors $\Lambda_n \subset \mathbb{R}_{0+}^t$.
2. Compute the solution set $\widehat{\mathcal{B}}_n = \{\widehat{\beta}_\lambda : \lambda \in \Lambda_n\}$.
3. Compute the estimated risk and fairness properties of each $f_\beta : \beta \in \widehat{\mathcal{B}}_n$.
4. Select a predictor f_β with favorable risk and fairness properties.

Fig. 2: Penalized FADE estimation procedure.

Propositions 2 and 3 established a correspondence between the constrained and penalized estimands, so each $\widehat{\beta}_\lambda$ may be regarded as an estimate either

of the penalized-min estimand β_λ^* or of some risk-min estimand β_r^* . The value of the penalized perspective is that Step 2 in this procedure can be carried out extremely efficiently. Since each matrix $\mathbb{P}_n(\hat{g}_j b)\mathbb{P}_n(\hat{g}_j b)^T$ has rank 1, the overall matrix inverse can be computed by computing $\mathbb{P}_n(bb^T)^{-1}$ and then applying a series of simple algebraic operations, per the Sherman-Morrison update formula. This is expressed in the following proposition.

Proposition 5 *Let*

$$\bar{\lambda}_j = (\lambda_1, \dots, \lambda_j), \text{ so that } \bar{\lambda}_t = \lambda$$

$$m_j = \mathbb{P}_n(\hat{g}_j b)$$

$$\hat{\mathbf{Q}}_0 = \mathbb{P}_n(bb^T)^{-1}$$

$$\hat{\mathbf{Q}}_1(\lambda_1) = \hat{\mathbf{Q}}_0 - \frac{\lambda_1 \hat{\mathbf{Q}}_0 m_1 m_1^T \hat{\mathbf{Q}}_0}{1 + \lambda_1 m_1^T \hat{\mathbf{Q}}_0 m_1}$$

$$\hat{\mathbf{Q}}_j(\bar{\lambda}_j) = \hat{\mathbf{Q}}_{j-1}(\bar{\lambda}_{j-1}) - \frac{\lambda_j \hat{\mathbf{Q}}_{j-1}(\bar{\lambda}_{j-1}) m_j m_j^T \hat{\mathbf{Q}}_{j-1}(\bar{\lambda}_{j-1})}{1 + \lambda_j m_j^T \hat{\mathbf{Q}}_{j-1}(\bar{\lambda}_{j-1}) m_j}, \text{ for } j = 2, \dots, t$$

Then

$$\hat{\beta}_\lambda = \begin{cases} \hat{\mathbf{Q}}_t(\lambda_t)^{-1} \mathbb{P}_n(b\hat{\phi}) & \text{(Counterfactual)} \\ \hat{\mathbf{Q}}_t(\lambda_t)^{-1} \mathbb{P}_n(bY) & \text{(Observable)} \end{cases}$$

Proposition 5 says that to compute the set $\hat{\mathcal{B}}_n$ requires only a single matrix inversion, to compute $\hat{\mathbf{Q}}_0$. Each vector m_j also only needs to be computed once. The remaining operations are algebraic. Since $\hat{\mathbf{Q}}_0$ is a $k \times k$ matrix and each m_j is a vector of length k , if b is a relatively small basis, then $\hat{\mathbf{Q}}_0$ will be fast to compute, and all the remaining algebraic operations will be fast. In our simulations and real data analyses, we show that we can get good results with a very small number of basis functions (e.g. 4 to 6), which yield extremely fast computations.

How should Λ_n be chosen in Step 1? Since $\Lambda_n \subset \mathbb{R}_{0+}^t$, one simple possibility is to take a one-dimensional grid of points between 0 and some arbitrary large number and then construct the t -dimensional Cartesian product. Since $\hat{\beta}_\lambda$ is smooth in λ , and since the risk and fairness measures are smooth in $\hat{\beta}$, we can expect that such a grid will enable us to move smoothly around the fairness-accuracy space, and that we won't be missing desirable predictors that lie in between the grid points².

Another possibility is to “seed” Λ_n with values that correspond to a particular $\hat{\beta}_r$. That is, fix some constraints ϵ_j and solve the dual of the risk-min program that defines $\hat{\beta}_r$. The dual solution λ^* indexes a penalized-min problem whose solution is β_r , so Λ_n can then be constructed as a grid around this λ^* . See the proof of Proposition 2 for a more detailed explanation.

²The movement won't be entirely smooth if predictions are truncated to lie in $[\ell_y, u_y]$.

This “seeding” approach provides a way to ensure that the set $\{\widehat{\beta}_\lambda : \lambda \in \Lambda_n\}$ includes estimators that in some sense target reasonable constraints, particularly for users with specific constraints in mind. This approach requires solving just a single constrained optimization problem, to establish a point of reference in fairness-accuracy space.

The procedure we have described allows users to efficiently construct and evaluate a very large set of models that fall in different points in fairness-accuracy space. In sections 6, 7, and 8, we show that this procedure enables us to find high-performing models in both observable and counterfactual settings, with simulated and real data. With minimal searching over possible bases, we are able to find models that substantially outperform existing models and methods with respect to both fairness and accuracy.

Remark 4 (Penalized version of unfair-min) Since β_λ^* is constructed as a penalized equivalent of β_r^* , the seeding approach to constructing Λ_n that we have described allows users to target particular fairness constraints but not particular risk constraints. It is straightforward to develop an analogous procedure around a penalized version of β_u^* that allows users to seed Λ_n with estimators that target particular risk constraints. In practice, it is not likely to matter much, since the construction of $\widehat{\beta}_\lambda$ should allow users to flexibly explore the fairness-accuracy space and find an estimator that accommodates their desired constraints, if one exists in the span of the chosen basis.

Remark 5 (Arbitrary FADE weights) An even simpler and plausibly just as effective alternative to computing the collection $\{\widehat{\beta}_\lambda : \lambda \in \Lambda_n\}$ is to simply define an arbitrary set $\mathcal{B} \subset \mathbb{R}^k$, perhaps constrained to lie in the simplex or in an L_1 box around the origin. That is, the user could simply evaluate arbitrary sets of basis weights to see if any of them yields a reasonable predictor. This set could be similarly constructed as a grid around a particular $\widehat{\beta}_r$ or $\widehat{\beta}_u$, if users have specific fairness or accuracy constraints they wish to target.

Remark 6 (Resemblance to ridge regression) $\widehat{\beta}_\lambda$ resembles a ridge regression estimator. In ridge regression and other regularized estimators, however, the penalty tuning parameter λ is expected to go to 0 as $n \rightarrow \infty$. In our setting, λ serves to enforce fairness rather than to modulate the variance-bias tradeoff, so there is no reason for it to shrink with n . Without unfairness penalties, the predictor won’t automatically get more fair as the data gets larger.

We now develop theoretical guarantees for the penalized FADE estimators. Let $h(n)$ denote the rate at which the product of nuisance parameter errors $\|\widehat{\pi} - \pi\| \|\widehat{\mu}_0 - \mu_0\|$ grows or converges, and let $\underline{h}(n)$ denote the rate for $\|\widehat{\pi} - \pi\| \|\widehat{\nu}_0 - \nu_0\|$. That is,

$$\begin{aligned} \|\widehat{\pi} - \pi\| \|\widehat{\mu}_0 - \mu_0\| &= O_{\mathbb{P}}(h(n)) \\ \|\widehat{\pi} - \pi\| \|\widehat{\nu}_0 - \nu_0\| &= O_{\mathbb{P}}(\underline{h}(n)) \end{aligned}$$

Under ideal conditions, Assumption 9 will hold, so that the product of nuisance parameter errors decay faster than $1/\sqrt{n}$, but the subsequent results do not require this.

Assumption 10 (Compact superset Λ). *For all n , $\Lambda_n \subseteq \Lambda \subset \mathbb{R}^t$ for some compact set Λ .*

Definition 6 For any $\lambda \in \mathbb{R}_{0+}^t$, the *excess risk* for $\hat{\beta}_\lambda$ is

$$\mathbb{P}[(b^T \hat{\beta}_\lambda - \tilde{Y})^2] - \mathbb{P}[(b^T \beta_\lambda^* - \tilde{Y})^2]$$

Theorem 3 (Uniform rate for excess risk in the penalized setting) *Under Assumptions 1–2 for the observable setting; and Assumptions 1–2, 4–9, and 10 for the counterfactual setting:*

$$\sup_{\lambda \in \Lambda} \left\{ \mathbb{P} \left[\left(b^T \hat{\beta}_\lambda - \tilde{Y} \right)^2 \right] - \mathbb{P} \left[\left(b^T \beta_\lambda^* - \tilde{Y} \right)^2 \right] \right\} = O_{\mathbb{P}}(\sqrt{1/n}) + O_{\mathbb{P}}(h(n))$$

In other words, the excess risk goes to 0 uniformly at $\sqrt{1/n}$ or the nuisance rate $h(n)$, whichever is slower. We have a similar result for the excess unfairness, which is defined as follows.

Definition 7 For any $\lambda \in \mathbb{R}_{0+}^t$, the *excess unfairness* for $\hat{\beta}_\lambda$ is

$$\left\{ \max_{j \in 1, \dots, t} \left(\mathbb{P} \left[g_j b^T \hat{\beta}_\lambda \right] - \mathbb{P} \left[g_j b^T \beta_\lambda^* \right] \right) \right\}$$

We have defined excess unfairness as the max over j , but it makes little difference if we define it instead as the sum over j . Note that here we haven't used the squared unfairness.

Theorem 4 (Uniform rate for excess unfairness in the penalized setting) *Under Assumptions 1–2 for the observable setting; and Assumptions 1–2, 4–9, and 10 for the counterfactual setting:*

$$\sup_{\lambda \in \Lambda} \left\{ \max_{j \in 1, \dots, t} \left(\mathbb{P} \left[g_j b^T \hat{\beta}_\lambda \right] - \mathbb{P} \left[g_j b^T \beta_\lambda^* \right] \right) \right\} = O_{\mathbb{P}}(\sqrt{1/n}) + O_{\mathbb{P}}(h(n))$$

Remark 7 (Allowing $k \rightarrow \infty$) We can obtain similar theoretical results in a regime in which the basis dimension k is allowed to grow to ∞ , if we require that $\sup_{w \in \mathcal{W}} \|b(w)\| = O(\sqrt{k})$ and that $k \log(k)/n \rightarrow 0$. The first requirement is a stronger version of Assumption 2, while the second insists that k not grow too fast in n . Under these additional requirements, we attain a rate of $O_{\mathbb{P}}(\sqrt{k/n}) + O_{\mathbb{P}}(\sqrt{k} \cdot h(n))$ in Theorems 3 and 4. These results extend the results of Belloni et al (2015) to a setting with nuisance parameters and penalty terms. As illustrated in that paper, these requirements are weak enough to allow the basis to asymptotically span rich function spaces such as the space of square integrable functions.

5.3 Risk and unfairness of a fixed predictor

Once a user has constructed a predictor or a set of candidate predictors, they will naturally wish to estimate the risk and fairness properties of those predictors, for example before choosing one to deploy in a decision-making context. The risk and unfairness of a fixed predictor f_β are estimated as

$$\begin{aligned}\widehat{\text{Risk}}(f_\beta) &= \begin{cases} \mathbb{P}_n[(f_\beta - Y)^2] & \text{(Observable)} \\ \mathbb{P}_n[f_\beta^2 - 2f_\beta\widehat{\phi} + \widehat{\phi}] & \text{(Counterfactual)} \end{cases} \\ \widehat{\text{UF}}_j(f_\beta) &= \begin{cases} \mathbb{P}_n[g_j(W, Y)f_\beta] & \text{(Observable)} \\ \mathbb{P}_n[g_j(W, \widehat{\phi})f_\beta] & \text{(Counterfactual)} \end{cases}\end{aligned}$$

for $j = 1, \dots, t$.

Theorem 5 (Asymptotic normality of risk and unfairness estimators) *Consider fairness functions $g_j \in \{g^{\text{rate}}, g^{\text{FPR}}, g^{\text{FNR}}\}$. Under Assumptions 1–2 for the observable (Obs.) setting; and Assumptions 1–2, 4–9, and 10 for the counterfactual (Count.) setting:*

$$\begin{aligned}\sqrt{n} \left(\widehat{\text{Risk}}(f_\beta) - \text{Risk}(f_\beta) \right) &\xrightarrow{d} \begin{cases} N \left(0, \text{var} \left((f_\beta - Y)^2 \right) \right) & \text{(Obs.)} \\ N \left(0, \text{var} \left(f_\beta^2 - 2f_\beta\phi + \phi \right) \right) & \text{(Count.)} \end{cases} \\ \sqrt{n} \left(\widehat{\text{UF}}_j(f_\beta) - \text{UF}_j(f_\beta) \right) &\xrightarrow{d} \begin{cases} N \left(0, \text{var} \left(g_j(W, Y)f_\beta \right) \right) & \text{(Obs.)} \\ N \left(0, \text{var} \left(\mathbb{P}(\gamma_0)^{-1}\eta_0 - \mathbb{P}(\gamma_1)^{-1}\eta_1 \right) \right) & \text{(Count.)} \end{cases}\end{aligned}$$

where, for $a \in \{0, 1\}$,

$$\begin{aligned}\gamma_a &= \begin{cases} (1 - \phi)\mathbb{1}\{A = a\} & \text{(for } g^{\text{FPR}}) \\ \phi\mathbb{1}\{A = a\} & \text{(for } g^{\text{FNR}}) \end{cases} \\ \eta_a &= \gamma_a \left(f_\beta - \frac{\mathbb{P}[\gamma_a f_\beta]}{\mathbb{P}[\gamma_a]} \right)\end{aligned}$$

Under Theorem 5, asymptotically valid confidence intervals can be constructed, and asymptotically valid hypothesis tests conducted, for $\text{Risk}(f_\beta)$ and $\text{UF}_j(f_\beta)$.

In the next three sections, we illustrate FADE on simulated and real data, in both observable and counterfactual settings. Given the correspondence between the constrained and penalized FADE forms established in Propositions 2 and 3, we do not conduct separate simulations for these two settings. We emphasize the penalized form because of its substantial computational advantages over the constrained form.

6 Simulations

We illustrate the penalized FADE procedure in the counterfactual setting, i.e. when $\widetilde{Y} = Y^0$. As in a real data setting, each estimator $\widehat{\beta}_\lambda$ is constructed

using only observable data, but unlike in a real data setting, we use the known values of Y^0 to evaluate the resulting predictors.

All computations in this and subsequent sections were carried out on a 2013 MacBook Pro with a 2.4 GHz dual-core processor and 8GB of RAM.

6.1 Data-generating process

The data generating process is as follows, for data $Z = (A, X, D, Y^0, Y^1, Y)$.

$$\begin{aligned} \mathbb{P}(A = 1) &= 0.3 \\ X \mid A &\sim N(A * (1, -0.8, 4, 2)^T, I_4) \\ \mathbb{P}(D = 1 \mid A, X) &= \min\{0.975, \text{expit}((A, X)^T(0.2, -1, 1, -1, 1))\} \\ \mathbb{P}(Y^0 = 1 \mid A, X, D) &= \text{expit}((A, X)^T(-5, 2, -3, 4, -5)) \\ \mathbb{P}(Y^1 = 1 \mid A, X, D) &= \text{expit}((A, X)^T(1, -2, 3, -4, 5)) \\ Y &= (1 - D)Y^0 + DY^1 \end{aligned}$$

where I_4 denotes the 4×4 identity matrix. $A = 1$ represents the minority group. There are no previously trained predictors; i.e. $S = \emptyset$, so the collected covariates consist of $W = (A, X)$. This data generating process satisfies the identifying assumptions, Assumptions 4–6: the last line expresses the consistency assumption; the propensity score $\pi(A, X) = \mathbb{P}(D = 1 \mid A, X)$ is upper bounded at 0.975 to satisfy positivity; and $Y^a \perp\!\!\!\perp D \mid W$ for $a \in \{0, 1\}$, satisfying ignorability.

The two groups $A = 0$ and $A = 1$ differ in the distribution of covariates (Figure 3) and decisions and outcomes (Table 3). The minority group experiences a positive decision ($D = 1$) 18% of the time, while the majority group experiences it 50% of the time. Outcomes Y^0 and Y are higher for the minority group, with a larger disparity for potential outcomes than for observable outcomes. All measures are computed on a dataset of size 50,000.

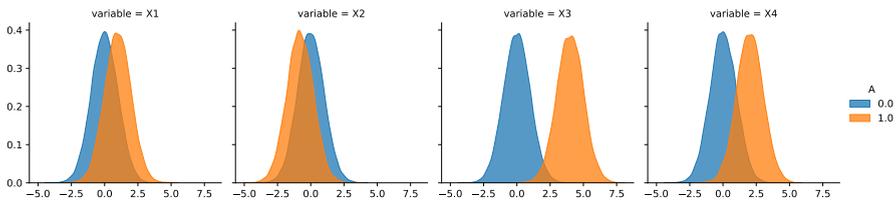


Fig. 3: Conditional covariate distributions for the two groups $A = 0$ and $A = 1$ in the simulated data. Curves are kernel density estimates.

As a reference point for our method, in Table 4 we compute the performance of the counterfactual Bayes-optimal predictor $f(A, X) = \mathbb{E}[Y^0 \mid A, X]$, which is defined in the data-generating process. We include both MSE, with is the

A	$\mathbb{E}[D A]$	$\mathbb{E}[Y^0 A]$	$\mathbb{E}[Y A]$
0	0.50	0.50	0.67
1	0.18	0.76	0.71

Table 3: Distribution of decisions and outcomes for groups $A = 0$ and $A = 1$ in the simulated data.

measure directly targeted by our method, as well as area under the curve (AUC). The Bayes-optimal predictor is highly accurate, with an MSE of 0.05 and an AUC of 0.98. The MSE of 0.05 is a lower bound on the risk achievable by any predictor. Unsurprisingly, given the difference in the distribution of outcomes across the two groups, the Bayes-optimal predictor has a large rate-diff ($|\mathbb{E}[f | A = 0] - \mathbb{E}[f | A = 1]|$). The differences in generalized false positive and false negative rates, however, are relatively small.

Predictor	MSE	AUC	rate-diff	FPR-diff	FNR-diff
Bayes-optimal	0.05	0.98	0.26	0.07	0.05

Table 4: Risk and fairness measures with respect to Y^0 for the Bayes-optimal predictor $\mathbb{E}[Y^0 | A, X]$ in the simulated data. The predictor is highly accurate, with low MSE and high AUC. It has a relatively large rate disparity but small disparities in generalized false positive and false negative rates.

6.2 Base predictors and nuisance models

We now investigate the performance of the penalized FADE procedure. We randomly sample three iid datasets of size $n = 1000$, representing $\mathcal{D}_{\text{learn}}$, $\mathcal{D}_{\text{train}}^{\text{nuis}}$, and $\mathcal{D}_{\text{train}}^{\text{target}}$. We train four base predictors on $\mathcal{D}_{\text{train}}$, with A, X as covariates and Y as the outcome. We train only on data in which $D = 0$: under the ignorability assumption, $\mathbb{E}[Y | A, X, D = 0] = \mathbb{E}[Y^0 | A, X]$, so this results in predictors which are designed to estimate Y^0 . The predictors consist of a random forest, a gradient boosted (GB) classifier, a Gaussian Naive Bayes model, and a ridge regression, all chosen for convenience and ease of computation. (In practice, a logistic regression would be a natural choice for a base predictor. Since the actual regression function $\mathbb{E}[Y^0 | A, X]$ is logistic, however, we do not use this model in order to avoid making the problem too easy, and to simulate a real data setting in which it is unlikely that the true regression function is known up to a finite dimensional parameter.) In addition to these predictors, we include a mean predictor, which always predicts the conditional sample mean $\mathbb{P}_n(Y | D = 0)$. This plays essentially the same role as an intercept in ordinary linear regression.

We use random forest classifiers to estimate the propensity and outcome models $\hat{\pi}$ and $\hat{\mu}_0$ on $\mathcal{D}_{\text{train}}^{\text{nuis}}$. All models were trained with their default tuning parameters using the scikit-learn library in Python. Predictors $\hat{\beta}_\lambda$ are computed using $\mathcal{D}_{\text{train}}^{\text{target}}$.

After a set of penalty vectors $\Lambda_n \subset \mathbb{R}_{0+}^t$ is chosen and the corresponding model coefficients $\hat{\mathcal{B}}_n = \{\hat{\beta}_\lambda : \lambda \in \Lambda_n\}$ are computed, we estimate the risk of fairness properties of every $f_\beta : \beta \in \hat{\mathcal{B}}_n$. In order to understand the true range of risk and fairness values that our method produces, we use a large test set $\mathcal{D}_{\text{test}}$ of size 10,000, and in place of $\hat{\phi}$, the nuisance quantity that would be required in a real data setting, we use the known values of Y^0 to compute the risk and fairness estimates. (For comparison purposes, estimates were also computed using μ_0 instead of Y^0 ; the results were virtually identical.) Since the true Y^0 is used, there is no need to split $\mathcal{D}_{\text{test}}$ into $\mathcal{D}_{\text{test}}^{\text{nuis}}$ and $\mathcal{D}_{\text{test}}^{\text{target}}$.

Table 5 shows the performance of the base predictors as well as the ordinary unpenalized least squares (OLS) solution, i.e. the predictor f_{β_λ} with $\lambda = 0$. The OLS predictor is the (estimated) MSE-minimal aggregation of the five base predictors, computed without regard for fairness. The OLS weights are $[-0.27, 0.09, 0.40, -0.11, 0.94]$; each base predictor appears to make a nontrivial contribution. The MSE of the base predictors ranges from 0.08 to 0.27. The four non-constant predictors improve substantially on the mean predictor with respect to both MSE and AUC. The mean predictor necessarily has a value of 0 for all three disparities, while the disparities of the other base predictors vary between 0.10 and 0.59. As expected, the OLS predictor has lower MSE than any of the base predictors. The performance of the OLS predictor is similar to the performance of the Bayes-optimal predictor in Table 5: both have a small MSE, a relatively large rate disparity, and relatively small error rate disparities.

Predictor	MSE	AUC	rate-diff	FPR-diff	FNR-diff
Mean	0.27	0.50	0.00	0.00	0.00
Random Forest	0.09	0.95	0.28	0.21	0.11
GB Classifier	0.08	0.96	0.31	0.22	0.10
Naive Bayes	0.17	0.84	0.52	0.59	0.40
Ridge	0.09	0.98	0.22	0.09	0.10
OLS	0.07	0.98	0.26	0.10	0.08

Table 5: Performance of the five base predictors and the ordinary least squares (OLS) predictor in the simulated data. The OLS weights are $[-0.27, 0.09, 0.40, -0.11, 0.94]$. The OLS predictor substantially improves on the MSE of the base predictors. The performance profile of the OLS predictor is close to the profile of the Bayes-optimal predictor in Table 4.

6.3 Results: one unfairness penalty

We now compute a set of unfairness-penalized predictors, applying a single unfairness penalty at a time. Let

$$\Lambda_{n,1} = \{0, 0.001, 0.01, 1, 10, 20, 50, 100, 500, 1000, 2000\}.$$

For each $\lambda \in \Lambda_{n,1}$, we compute $\widehat{\beta}_\lambda$ for each fairness function $g \in \{g^{\text{rate}}, g^{\text{FPR}}, g^{\text{FNR}}\}$. The value $\lambda = 0$ corresponds in each case to the OLS solution, so this yields a total of $(|\Lambda_{n,1}| - 1) * 3 + 1 = 31$ predictors.

The risk and unfairness values for each predictor are plotted in Figure 4. The disparity corresponding to the targeted constraint is represented by a solid line, while the other two disparities and the MSE are represented by dashed lines. We emphasize that the values in this figure are computed on $\mathcal{D}_{\text{test}}$, after the predictors $\widehat{\beta}_\lambda$ are computed on $\mathcal{D}_{\text{train}}$.

As expected, as λ increases, the targeted disparity of the resulting predictor generally decreases. The decrease is monotonic, except at one point: $\lambda = 1$ for FPR-diff, which may be a result of sampling noise. The rate difference decreases from 0.26 to 0.04. The FPR disparity decreases from 0.10 to 0, then remains at 0.01. The FNR disparity decreases from 0.08 to 0.04. When the rate disparity is penalized, the decrease in the target disparity is accompanied by a slight increase in MSE, from 0.07 to 0.09, as well as small increases in FPR-diff and FNR-diff. When FPR-diff is penalized, all three disparities fall together, while the increase in MSE is miniscule, from 0.067 to 0.071. The same is true when FNR-diff is targeted: the MSE increases from 0.067 to 0.069.

These results illustrate (1) that the penalty term successfully controls the target disparity, (2) that an increase in fairness need not come at the cost of a substantial decrease in accuracy, and (3) that a decrease in one disparity need not produce an increase in other disparities. In the second two panels, the FADE predictors are uniformly more fair than the OLS predictor, with essentially no change in accuracy. Additionally, in these two panels, even though only one disparity was penalized at a time, all three disparities decrease as λ increases.

6.4 Results: multiple unfairness penalties

We now apply all three unfairness penalties simultaneously. Define $\Lambda = \Lambda_{n,1} \times \Lambda_{n,1} \times \Lambda_{n,1} \subset \mathbb{R}_{0+}^3$. The collection $\widehat{\mathcal{B}}_n = \{\widehat{\beta}_\lambda : \lambda \in \Lambda\}$ now indexes $|\Lambda_{n,1}|^3 = 1331$ predictors. We use the same base predictors and nuisance predictors as in the previous section. The process of training the predictors, computing $\widehat{\mathcal{B}}_n$, and estimating the risk and fairness for each predictor, took less than 10 seconds.

Figure 5 plots each of the three disparities against MSE, for each of the 1331 predictors, as well as the base predictors and the OLS predictor. As expected, the OLS predictor has the smallest MSE. Fewer than 1331 dots are visible in each panel, due to the fact that many of the predictors substantially overlap in fairness-accuracy space. Nevertheless, the predictors span a wide range of

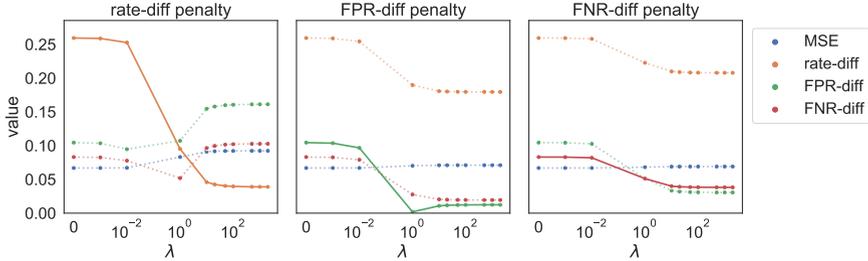


Fig. 4: Risk and fairness for predictors subject to one of three penalties, in the simulated data. The x-axis represents the unfairness penalty coefficient λ , with the leftmost point ($\lambda = 0$) in each panel corresponding to the OLS solution. The y-axis represents the MSE and the disparity values of the resulting predictor $\hat{\beta}_\lambda$, computed on an independent test set of size 10,000 using the known values of Y^0 . Solid lines indicate the metric that is penalized in training.

performance profiles. For all three disparities, predictors exist that take the disparity to 0, with relatively small increase in MSE relative to the OLS predictor. For rate-diff and FNR-diff, these predictors notably do not appear in Figure 4, where the lowest value for these two disparities are 0.04. We only discover these predictors by applying multiple penalties simultaneously. All the FADE predictors are substantially more accurate than the mean predictor.

Figure 6 plots the same 1331 predictors with respect to each pair of disparities, with color indicating MSE. This figure illustrates the interplay of three metrics at once, and is of interest to users who wish to control two disparities simultaneously. For example, users who wish to target (counterfactual) equalized odds would be interested in the bottom panel that plots FNR-diff and FPR-diff.

These views once again reveal a wide range of predictor behavior. Unsurprisingly, many of the highest MSE predictors are close to the origin, but the relationship between MSE and distance to the origin is far from monotonic. In all three panels, there is a line of predictors stretching from the OLS predictor that represent improvements in both disparities with minimal increase in MSE. In the bottom panel, for example, there are predictors with FPR-diff close to 0, FNR-diff under 0.05, and MSE under 0.10. These predictors approximately satisfy equalized odds, and they represent an increase in MSE of less than 0.03 relative to the OLS predictor.

In order to examine this more precisely, Table 6 shows the performance of the predictors with the minimum L_2 distance from the origin in the fairness-accuracy subspace defined by MSE as well as zero to three disparities. For example, the “MSE + rate” row represents the predictor with the smallest L_2 norm in the (MSE, rate-diff) vector, while the “MSE + rate + FPR + FNR” row represents the predictor with the smallest L_2 norm in the (MSE, rate-diff, FPR-diff, FNR-diff) vector. FPR-diff and FNR-diff can be minimized, singly

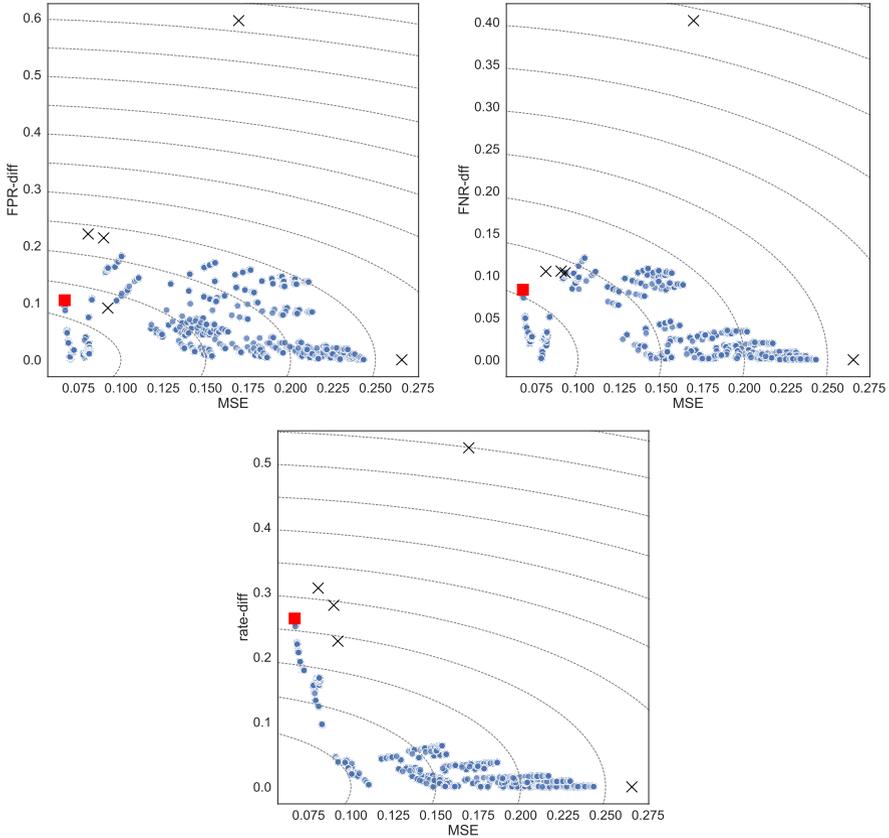


Fig. 5: Disparity and MSE values for each of 1331 predictors in the simulated data. Black “X”s represent the base predictors, with the mean predictor at the bottom right of each panel. The red square is the OLS predictor, and the blue dots are the FADE predictors. Radius lines indicate distance from the origin. Despite substantial overlap, the predictors span a wide range of fairness and accuracy values. For each disparity, many predictors exist which take that disparity to 0, at a small cost in MSE relative to the OLS predictor.

or jointly, with no increase in MSE relative to the OLS predictor. Rate-diff can be substantially reduced with a relatively small increase in MSE. Perhaps surprisingly, all three disparities can be jointly minimized, to 0.06 (rate-diff), 0.03 (FPR-diff), and 0.02 (FNR-diff), with only a 0.06 increase in MSE and a 0.02 decrease in AUC relative to the unpunished OLS predictor.

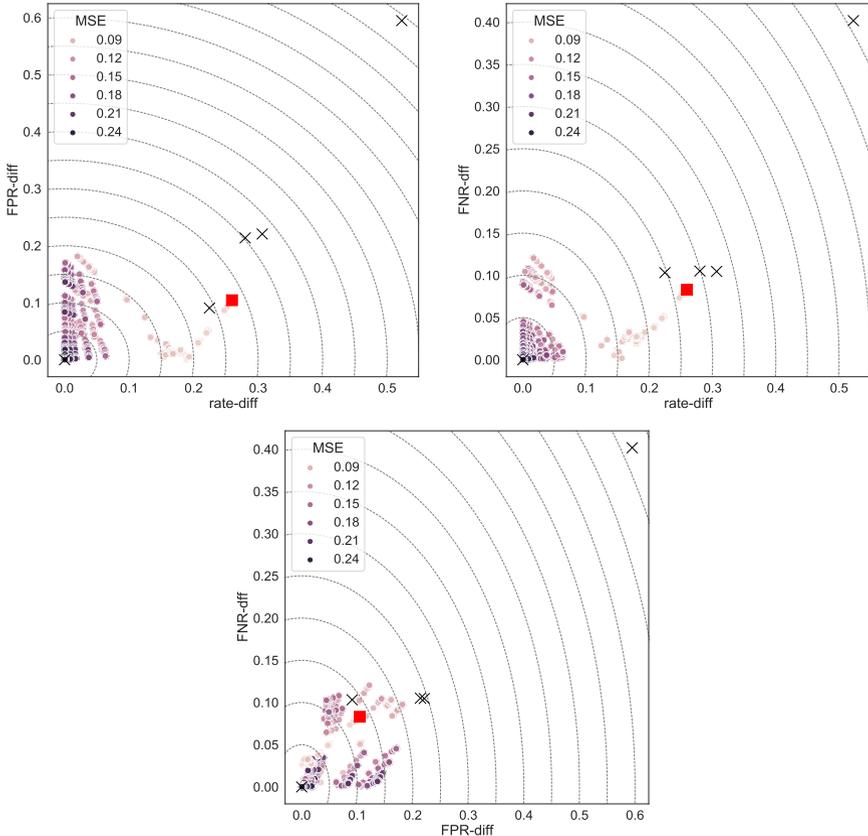


Fig. 6: Pairs of disparity values and MSE values for each of 1331 predictors in the simulated data. Black “X”s represent the base predictors, with the mean predictor at the origin in each panel. The red square is the OLS predictor, and the dots are the FADE predictors. Radius lines indicate distance from the origin. Each pair of disparities can be jointly decreased with minimal increase in MSE relative to the OLS predictor.

7 Results: recidivism risk prediction

We next illustrate FADE on the COMPAS dataset gathered by ProPublica (Angwin and Larson, 2016; Angwin et al, 2016). The dataset comprises public arrest records, criminal records, and COMPAS scores from a single county in Florida, spanning 2013–2016. COMPAS is a collection of tools developed by the company Equivant (formerly Northpointe) designed to assess the risk of recidivism. We utilize the COMPAS scores for general, as opposed to violent, recidivism. The scores consist of risk deciles, coded 1–10, which we normalize to the range $[0.1, 1]$. COMPAS takes as input up to 137 features (Northpointe,

Predictor	MSE	AUC	rate-diff	FPR-diff	FNR-diff
MSE (OLS)	0.07	0.98	0.26	0.10	0.08
MSE + rate	0.09	0.95	0.04	0.16	0.10
MSE + FPR	0.07	0.98	0.19	0.00	0.03
MSE + FNR	0.07	0.98	0.18	0.01	0.02
MSE + rate + FPR	0.12	0.90	0.04	0.05	0.09
MSE + rate + FNR	0.10	0.95	0.04	0.16	0.08
MSE + FPR + FNR	0.07	0.98	0.18	0.01	0.02
MSE + rate + FPR + FNR	0.13	0.96	0.06	0.03	0.02

Table 6: Performance of the FADE predictors that minimize the Euclidean norm of MSE and zero or more disparities, in the simulated data. The top row represents the OLS predictor, included again for reference. All three disparities can be minimized, singly or jointly, with no impact or a small impact on MSE.

2015; Rudin et al, 2020), which are unavailable in this data set. We utilize just three features as covariates: an indicator for defendant age greater than 45, an indicator for defendant age less than 25, and the number of prior arrests, ranging from 0 to 29. Previous work has found that predictors trained using just these covariates perform similarly to COMPAS (Angelino et al, 2018).

The sensitive feature is race, restricted to defendants who are coded African-American ($n = 3175$) or Caucasian ($n = 2013$). The decision variable D represent pretrial release, with $D = 0$ if defendants are released and $D = 1$ if they are detained. The outcome of interest Y^0 is rearrest within two years, should a defendant be released pretrial. Since it difficult to assess the plausibility of the positivity and ignorability assumptions without consulting with domain experts, we conducted analyses in both the counterfactual and observable setting. The results and conclusions were largely the same, so we only include the counterfactual results.

We split the data into five datasets, each with approximately 1040 rows: $\mathcal{D}_{\text{learn}}$, $\mathcal{D}_{\text{train}}^{\text{nuis}}$, $\mathcal{D}_{\text{train}}^{\text{target}}$, $\mathcal{D}_{\text{test}}^{\text{nuis}}$, and $\mathcal{D}_{\text{test}}^{\text{target}}$. As base predictors, we used the four model types from the previous section as well as a logistic regression. We used random forest classifiers for the nuisance predictors in both the training and test data. Table 7 gives the estimated performance of the five base predictors, COMPAS, and the OLS predictor, which spans COMPAS and the base predictors. Previous work found differences in a binarized version of COMPAS for both observable (Angwin and Larson, 2016) and counterfactual (Mishler, 2019) false positive vs false negative rates for African-American vs. Caucasian defendants. Those differences appear here in the generalized error rates. COMPAS also has a large rate disparity. Perhaps surprisingly, the base predictors all yield smaller disparities than COMPAS, even though they generally also have smaller MSE.

We compute FADE predictors using the same sets of penalty vectors Λ as in the previous section. Figure 7 shows disparities and MSE values for all 1331 predictors. Most predictors fall within a narrow range of MSE values that also includes COMPAS, so the primary value of aggregation here is in reducing

disparities. Nearly all the FADE predictors improve on COMPAS in terms of both risk and fairness. The top row of Figure 7 shows that all three disparities can be individually reduced to 0 with minimal cost in MSE relative to the OLS predictor, and the bottom row shows that these improvements also extend over pairs of disparities.

	MSE	rate-diff	FPR-diff	FNR-diff
Mean	0.26	0.00	0.00	0.00
Random Forest	0.28	0.05	0.03	0.02
Logistic	0.22	0.06	0.06	0.00
GB Classifier	0.23	0.06	0.01	0.05
Ridge	0.22	0.05	0.05	0.00
COMPAS	0.24	0.15	0.15	0.08
OLS	0.22	0.09	0.09	0.05

Table 7: Estimated performance of the five base predictors, COMPAS, and the OLS predictor in the COMPAS dataset. The OLS weights are $[0.39, 0.12, 0.79, 0.10, -1.08, 0.93]$. The OLS predictor does not perform substantially better than the base predictors. COMPAS has different false positive and false negative rates for African-American vs Caucasian defendants, as well as a rate disparity. The base predictors all have smaller disparities than COMPAS, and—perhaps surprisingly—generally smaller MSE.

8 Results: income prediction

Finally, we apply our method in the observable setting, using the Adult dataset (Dua and Graff, 2017). This dataset comprises demographic variables derived from the 1994 U.S. Census. We consider sex as a sensitive feature, coded 0 or 1, and we utilize as covariates a set of indicator variables representing age by decade, and a set of indicator variables representing the number of years of education. The classification task is to predict whether an individual’s income is over \$50K/year, for example for the purpose of deciding whether to issue a loan.

We randomly split the data into four datasets: $\mathcal{D}_{\text{learn}}$ and $\mathcal{D}_{\text{train}}$, consisting of 14,653 and 14,652 rows; and $\mathcal{D}_{\text{test}}$ and $\mathcal{D}_{\text{validate}}$, consisting of 9,768 and 9,769 rows. $\mathcal{D}_{\text{validate}}$ is used to compare the performance of the selected best predictors to predictors that are developed using existing fairness methods.

This dataset does not contain any previously trained predictors. We use the same five base predictor types as in the COMPAS analysis. Additionally, we use $\mathcal{D}_{\text{train}}$ to train three “fair” predictors with other fairness methods: *adversarial debiasing* (Zhang et al, 2018), *reductions* (Agarwal et al, 2018), and a *meta-algorithm* (Celis et al, 2020). All predictors were trained using the Python library aif360, a set of tools that provide access to a range of fairness methods via a consistent interface (Bellamy et al, 2018). The three chosen

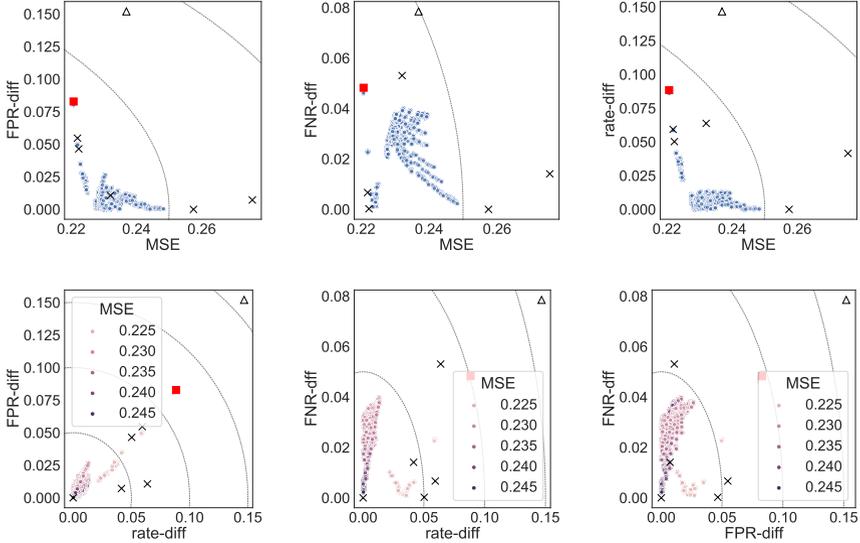


Fig. 7: Disparity against MSE (top row) or pairs of disparities colored by MSE (bottom row), for each of 1331 predictors in the COMPAS data. The black triangle represents COMPAS, the black “X”s represent the other base predictors, the red square is the OLS predictor, and the blue dots are the FADE predictors. Radius lines indicate distance from the origin. Most of the predictors improve on COMPAS in terms of both MSE and the relevant disparity. For each disparity, many predictors exist which take that disparity to 0, at a small cost in MSE relative to the OLS predictor.

methods yield binary classifiers, and they are all designed to minimize *rate-diff* in an observable setting with a binary outcome. Since, for a binary classifier \hat{Y} , MSE is equal to classification error $\mathbb{P}(\hat{Y} \neq Y)$, we may regard these three predictors as the result of methods that seek to minimize MSE among specific classes of binary predictors.

We construct FADE predictors using the five base predictors (“base5”) or the five base predictors and the three fair predictors (“base8”). Table 8 gives the performance of the base predictors, the fair predictors, and the two OLS predictors. Compared to the base predictors, two of the three fairness methods result in a substantially lower *rate-diff*, which is the disparity they aim to minimize. The base predictors and the OLS predictors have lower MSE, higher AUC, and higher disparities compared to these two fair predictors.

We compute FADE predictors using the same sets of penalty vectors Λ as in the previous two sections. Figure 7 shows disparities and MSE values for all 1331 predictors. Fairness-accuracy tradeoffs are most evident for *rate-diff* and *FPR-diff* in the top row of this figure, where only the five base predictors are

Predictor	MSE	AUC	rate-diff	FPR-diff	FNR-diff
Mean	0.18	0.50	0.00	0.00	0.00
Random Forest	0.14	0.81	0.19	0.14	0.24
Logistic	0.14	0.81	0.20	0.14	0.25
GB Classifier	0.14	0.82	0.19	0.13	0.24
Ridge	0.14	0.81	0.17	0.13	0.16
Adversarial	0.21	0.67	0.07	0.00	0.03
Reductions	0.22	0.62	0.01	0.03	0.07
Meta	0.30	0.68	0.18	0.26	0.26
OLS - base5	0.14	0.82	0.19	0.13	0.24
OLS - base8	0.14	0.81	0.20	0.14	0.25

Table 8: Estimated performance in the Adult dataset of five base predictors, three predictors trained using existing fairness methods, and the two OLS predictors, which aggregate only the five base predictors or all eight predictors. The OLS weights are $[0.03, 0.29, 0.65, 0.03, 0.01]$, for base5, and $[-0.01, 0.26, 0.56, 0.18, 0.04, -0.02, -0.03, 0]$, for base8. Only two of the three fairness methods successfully control their targeted disparity, rate-diff. The Meta predictor has a rate-diff which is comparable to the base predictors which are trained without regard to fairness. The OLS predictors perform comparably to the base predictors.

used. In the bottom row, where the fair predictors are included as basis functions, the tradeoffs essentially disappear: all three disparities can be reduced to 0 with virtually no cost in MSE.

8.1 Model validation

Using the performance estimates from the test data, we select the FADE predictors that minimize the distance from the origin in each of seven fairness-accuracy subspaces, for both the base5 and base8 predictors. We then compute risk and fairness estimates for each of these predictors, as well as the three fair predictors, on the validation data (Table 9). The estimates on the test and validation data differed by no more than approximately 0.005.

Both the base5 and base8 FADE predictors are substantially more fair than the OLS predictors, while incurring very small increases in MSE. The high AUC values confirm that these are accurate predictors. All the FADE predictors have small disparities compared to the OLS predictors. Explicitly minimizing multiple disparities simultaneously is not necessarily more costly in terms of performance than minimizing a single disparity.

The FADE predictors have substantially lower MSE and higher AUC than the fair predictors, and they are in many cases more fair. The fair predictors achieve values of 0.08, 0.01, and 0.17 for rate-diff, the disparity they aim to minimize. The base5 FADE predictors that include rate-diff in their criteria

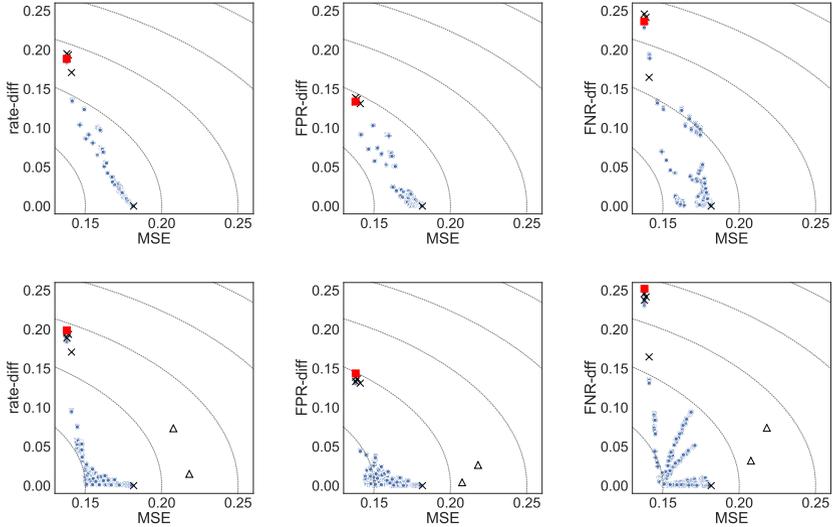


Fig. 8: Disparity against MSE for FADE predictors based on the five base predictors (top row) or the five base predictors plus the three fair predictors, for each of 1331 FADE predictors in the Adult data. Black “X”s represent the base predictors, the red square is the OLS predictor, and the blue dots are the FADE predictors. Radius lines indicate distance from the origin. For the sake of legibility, the Meta predictor, which has an MSE of 0.30, is excluded. The top row exhibits small but clear fairness-accuracy tradeoffs for rate-diff and FPR-diff. The bottom row shows that with the inclusion of the fair predictors, each disparity can be taken to 0 with almost no cost in MSE relative to the OLS model.

achieve rate-diffs of 0.04, 0.04, 0.06, and 0.02. The corresponding base8 FADE predictors achieve rate-diffs of 0.01, 0.03, 0.01, and 0.03.

The base5 results show that FADE yields predictors that perform comparably to or better than existing fairness methods. The base8 results highlight the flexibility of our approach: multiple predictors can be aggregated, regardless of whether or not they are trained with fairness properties in mind, with different weights to target different disparities. In this case, including the fair predictors in the aggregation improves both accuracy and fairness.

Each of the fair prediction methods contains tuning parameters that can be adjusted to return different predictors, as well as settings that allow them to target different fairness constraints, such as equalized odds. However, each method can only target a single fairness constraint at once. Additionally, these methods can take substantial time to run. A single run of the Meta method took roughly 5 seconds, the Reductions method ran in approximately 15 seconds, and the Adversarial method, which relies on neural nets, took roughly

a minute. By contrast, we were able to train the base predictors and compute and evaluate 1331 FADE predictors in 25 seconds.

The three fair predictors are binary by construction, whereas our method returns continuous predictors. Of course, these continuous predictors can be treated as binary, either by thresholding the output (yielding a deterministic classifier) or by treating the output as a probability and sampling from a corresponding Bernoulli distribution (yielding a randomized classifier). It is fast to compute estimates of the accuracy and fairness values from either of these two binarized classifiers and choose one which minimizes the criteria of interest. For example, the base5 FADE predictors include a predictor which, when thresholded at 0.5 to yield a deterministic binary classifier, achieves a classification error of 0.24 and disparities of 0 (to two digits). This has very slightly higher classification error than the Adversarial and Reductions predictors, but it exactly achieves equalized odds and demographic parity.

	Predictor	MSE	AUC	rate-diff	FPR-diff	FNR-diff
base5	MSE (OLS-base5)	0.14	0.82	0.19	0.13	0.24
	MSE + rate	0.16	0.73	0.04	0.02	0.10
	MSE + FPR	0.15	0.80	0.09	0.06	0.13
	MSE + FNR	0.16	0.75	0.10	0.09	0.01
	MSE + rate + FPR	0.16	0.73	0.04	0.02	0.10
	MSE + rate + FNR	0.16	0.75	0.06	0.05	0.01
	MSE + FPR + FNR	0.16	0.75	0.06	0.05	0.01
	MSE + rate + FPR + FNR	0.17	0.73	0.02	0.02	0.00
base8	MSE (OLS-base8)	0.14	0.81	0.20	0.14	0.25
	MSE + rate	0.15	0.79	0.01	0.03	0.02
	MSE + FPR	0.14	0.79	0.06	0.01	0.10
	MSE + FNR	0.15	0.79	0.05	0.01	0.01
	MSE + rate + FPR	0.15	0.79	0.03	0.00	0.01
	MSE + rate + FNR	0.15	0.79	0.01	0.03	0.01
	MSE + FPR + FNR	0.15	0.79	0.04	0.00	0.01
	MSE + rate + FPR + FNR	0.15	0.79	0.03	0.01	0.01
fair	Adversarial	0.21	0.67	0.08	0.01	0.04
	Reductions	0.22	0.62	0.01	0.03	0.06
	Meta	0.30	0.68	0.17	0.26	0.24

Table 9: Performance of the predictors that minimize the Euclidean norm of MSE and zero to three disparities, in the Adult data. The OLS predictors and the three fair predictors are included again for reference. Predictors are selected on the test data and evaluated on the validation data. The base5 FADE predictors aggregate the five base predictors, and the base8 FADE predictors aggregate all eight predictors. All three disparities can be minimized, singly or jointly, with only a small impact on MSE, and a small impact on AUC. Aggregated predictors are more accurate than the fair predictors, and have comparable or smaller values of rate-diff, the disparity that the fair predictors aim to minimize.

9 Conclusion

We developed a framework, FAir Double Ensemble learning (FADE), for constructing fair predictors and for exploring fairness-accuracy and fairness-fairness tradeoffs. This framework is extremely flexible, allowing users to combine arbitrary sets of predictors, including previously trained predictors and newly trained ones, regardless of whether they are designed to satisfy fairness constraints or not. FADE thereby collapses the distinction between in-processing and post-processing approaches to building fair predictors. FADE can accommodate a wide range of disparities and allows users to minimize multiple disparities simultaneously. FADE also accommodates both observable and counterfactual outcomes, joining a very small set of existing methods for targeting counterfactual accuracy and counterfactual versions of fairness criteria like equalized odds.

Within the FADE framework, we developed three methods. The first two “constrained” methods allow users to minimize mean squared error subject to explicit fairness constraints, or minimize unfairness subject to an explicit constraint on the mean squared error. The third “penalized” method allows users to efficiently construct large sets of predictors and evaluate their risk and fairness properties. The penalized method enables users to efficiently explore fairness-accuracy and fairness-fairness tradeoffs in their problem setting, and it enables them to find a predictor with a favorable risk and fairness profile. Our theoretical results show that FADE predictors converge to optimal behavior at fast rates, and our empirical results show that in many cases, disparities can be substantially reduced with no tangible loss of accuracy—or even an increase in accuracy—relative to the unpenalized least squares solution or an existing benchmark predictor.

Although our penalized approach is designed to minimize mean squared error and to penalize certain classes of disparities, the resulting predictors can naturally be evaluated with respect to any accuracy or fairness metric. For example, users might wish to consider only binary classifiers, so they may wish to evaluate classification error on thresholded versions of the FADE predictors. The penalized approach provides a principled way to explore various fairness-accuracy spaces, even if the fairness and/or accuracy metrics of interest aren’t explicitly represented in the penalized expression.

Finally, the efficiency of our penalized method relies on the particular closed form of the parameterized predictors, which arises as a result of the mean squared error and the squared fairness terms. However, any quadratic function that involves a positive definite matrix has a closed form solution. This form could be preserved under different accuracy metrics and fairness terms by, for example, adding a regularization term $\beta^T M \beta$ for some positive-definite matrix M . This suggests that our approach could be adapted to explicitly target other accuracy and/or fairness metrics.

Acknowledgments. We are grateful to Alexandra Chouldechova, Aaditya Ramdas, Cosma Shalizi, Ilya Shpitser, and Larry Wasserman for comments on

earlier versions of this work. This work was completed while Alan Mishler was a PhD student at Carnegie Mellon University.

Appendix A Proof preliminaries

For convenience, we collect all the assumptions that appear in the paper, with their short descriptors:

1. For all n , $\mathbb{E}[bb^T]$ is positive definite. (PSD outer product)
2. Uniformly in n , $\sup_{w \in \mathcal{W}} \|b(w)\| < \infty$. (Bounded basis norm)
3. The set $\{\mathbb{E}[g_j b] \mathbb{E}[g_j b]^T \beta_r^* : j \in \mathcal{I}\}$ is linearly independent. (LICQ - population)
4. $Y = DY^1 + (1 - D)Y^0$. (Consistency)
5. $\exists \delta \in (0, 1)$ s.t. $\mathbb{P}(\pi(W) \leq 1 - \delta) = 1$. (Positivity)
6. $Y^0 \perp\!\!\!\perp D \mid W$. (Ignorability)
7. $\exists \gamma \in (0, 1)$ s.t. $\mathbb{P}(\hat{\pi}(A, X, S) \leq 1 - \gamma) = 1$. (Bounded propensity estimator)
8. $\|\hat{\pi} - \pi\| = o_{\mathbb{P}}(1)$ and $\|\hat{\mu}_0 - \mu_0\| = o_{\mathbb{P}}(1)$ and $\|\hat{\nu}_0 - \nu_0\| = o_{\mathbb{P}}(1)$. (Consistent nuisance estimators)
9. $\|\hat{\pi} - \pi\| \|\hat{\mu}_0 - \mu_0\| = o_{\mathbb{P}}(1/\sqrt{n})$. (Nuisance parameter rates)
 $\|\hat{\pi} - \pi\| \|\hat{\nu}_0 - \nu_0\| = o_{\mathbb{P}}(1/\sqrt{n})$.
10. $\Lambda_n \subseteq \Lambda \subset \mathbb{R}^t$ for some compact Λ . (Compact Λ)

Recall that for any function $f : \mathcal{Z} \mapsto \mathbb{R}$, we defined $\mathbb{P}_n(f) = n^{-1} \sum_{j=1}^n f = \int f d\mathbb{P}_n(Z)$ and $\mathbb{P}(f) = \int f d\mathbb{P}(Z)$ as the sample and population expectations of f , so that for example $\mathbb{P}(\hat{\phi}) = \mathbb{E}[\hat{\phi} \mid \mathcal{D}_{\text{train}}]$ or $\mathbb{E}[\hat{\phi} \mid \mathcal{D}_{\text{test}}]$ is the expected value of $\hat{\phi}(Z)$ once the relevant nuisance function estimate $\hat{\phi}$ has been constructed.

We state several lemmas that are used in the proofs for the constrained and penalized settings. The first is a restatement of Lemma 2 in [Kennedy et al \(2020\)](#).

Lemma 1 (Kennedy, 2020). *Let $\hat{f} : \mathcal{Z} \mapsto \mathbb{R}$ be a function estimated on a nuisance dataset $\mathcal{D}^{\text{nuis}}$ independent of \mathbb{P}_n , and let $f : \mathcal{Z} \mapsto \mathbb{R}$ be another function. Assume $\text{var}(\hat{f} - f \mid \mathcal{D}^{\text{nuis}}) < \infty$. Then*

$$(\mathbb{P}_n - \mathbb{P})(\hat{f} - f) = O_{\mathbb{P}} \left(\frac{\|\hat{f} - f\|}{\sqrt{n}} \right)$$

Lemma 2 (Double robustness). *Let $f : \mathcal{W} \mapsto \mathbb{R}^p$ for any p be a function with $\|f(W)\| \leq M < \infty$ for some M . Under Assumption 5 (positivity),*

$$\begin{aligned}\|f(W)(\widehat{\phi} - \phi)\| &= O_{\mathbb{P}}(\|\mu_0 - \widehat{\mu}_0\| \|\widehat{\pi} - \pi\|) \\ \|f(W)(\widehat{\underline{\phi}} - \underline{\phi})\| &= O_{\mathbb{P}}(\|\nu_0 - \widehat{\nu}_0\| \|\widehat{\pi} - \pi\|)\end{aligned}$$

It follows immediately that

$$\begin{aligned}\mathbb{P}\left(f(W)(\widehat{\phi} - \phi)\right) &= O_{\mathbb{P}}(\|\mu_0 - \widehat{\mu}_0\| \|\widehat{\pi} - \pi\|) \\ \mathbb{P}\left(f(W)(\widehat{\underline{\phi}} - \underline{\phi})\right) &= O_{\mathbb{P}}(\|\nu_0 - \widehat{\nu}_0\| \|\widehat{\pi} - \pi\|)\end{aligned}$$

Proof

$$\begin{aligned}\mathbb{P}\left(f(W)(\widehat{\phi} - \phi)\right) &= \mathbb{P}\left(f(W)\left(\frac{1-D}{1-\widehat{\pi}}(Y - \widehat{\mu}_0) + \widehat{\mu}_0 - \frac{1-D}{1-\pi}(Y - \mu_0) - \mu_0\right)\right) \\ &= \mathbb{P}\left(f(W)\left(\frac{1-D}{1-\widehat{\pi}}(\mu_0 - \widehat{\mu}_0) + \widehat{\mu}_0 - \frac{1-D}{1-\pi}(\mu_0 - \mu_0) - \mu_0\right)\right) \\ &= \mathbb{P}\left(f(W)\left(\frac{1-\pi}{1-\widehat{\pi}}(\mu_0 - \widehat{\mu}_0) + \widehat{\mu}_0 - \mu_0\right)\right) \\ &= \mathbb{P}\left(f(W)\left(\frac{(\mu_0 - \widehat{\mu}_0)(\widehat{\pi} - \pi)}{1-\pi}\right)\right) \\ &\leq \frac{1}{\delta} \mathbb{P}(f(W)(\mu_0 - \widehat{\mu}_0)(\widehat{\pi} - \pi)) \\ &\leq \frac{1}{\delta} \|f(W)\| \|\mu_0 - \widehat{\mu}_0\| \|\widehat{\pi} - \pi\| \\ &= O_{\mathbb{P}}(\|\mu_0 - \widehat{\mu}_0\| \|\widehat{\pi} - \pi\|)\end{aligned}$$

where the second and third lines use iterated expectation, conditioning on W ; the fifth line uses Assumption 5 (positivity); and the sixth line uses the Cauchy-Schwarz inequality. \square

Appendix B Proofs of propositions

B.1 Proof of Proposition 1

Proof Let $\alpha_0, \alpha_1 \in \mathbb{R}$, let h_0, h_1 be mappings from $\{0, 1\} \times \widetilde{Y}$ to $\{0, 1\}$, and let $g(W, \widetilde{Y}) = \alpha_0 \frac{h_0(A, \widetilde{Y})}{\mathbb{E}[h_0(A, \widetilde{Y})]} - \alpha_1 \frac{h_1(A, \widetilde{Y})}{\mathbb{E}[h_1(A, \widetilde{Y})]}$. Then

$$\begin{aligned}\mathbb{E}[f(W)h_0(A, \widetilde{Y})] &= \mathbb{E}[\mathbb{E}[f(W)h_0(A, \widetilde{Y}) \mid h_0(A, \widetilde{Y})]] \\ &= \mathbb{E}[f(W) \mid h_0(A, \widetilde{Y}) = 1] \mathbb{P}(h_0(A, \widetilde{Y}) = 1) \\ &= \mathbb{E}[f(W) \mid h_0(A, \widetilde{Y}) = 1] \mathbb{E}[h_0(A, \widetilde{Y})] \\ \implies \mathbb{E}[f(W) \mid h_0(A, \widetilde{Y}) = 1] &= \frac{\mathbb{E}[f(W)h_0(A, \widetilde{Y})]}{\mathbb{E}[h_0(A, \widetilde{Y})]}\end{aligned}$$

where $\mathbb{E}[h_0(A, \tilde{Y})] > 0$ by assumption. By similar reasoning for h_1 , it follows that

$$|\alpha_0 \mathbb{E}[f(W) \mid h_0(A, \tilde{Y}) = 1] - \alpha_1 \mathbb{E}[f(W) \mid h_1(A, \tilde{Y}) = 1]| = |\mathbb{E}[g(W, \tilde{Y})f(W)]|$$

as desired. \square

B.2 Proof of Proposition 2

Proof Define the Lagrangian $L(\beta, v)$ of the risk-min program that defines β_r^* :

$$L(\beta, v) = \mathbb{E}[(b^T \beta - \tilde{Y})^2] + \sum_{j=1}^t v_j \left\{ \left(\mathbb{E}[g_j b^T \beta] \right)^2 - \epsilon_j^2 \right\}$$

Under Assumption 1, the risk-min program is a strictly convex quadratic program, and it is feasible by construction. Therefore, the dual program has at least one solution $\lambda \in \mathbb{R}_{0+}^t$, and strong duality holds, so the primal solution β_r^* and λ jointly satisfy the KKT conditions. In particular, $\nabla_{\beta} L(\beta_r^*, \lambda) = 0$ (stationarity). By computing $\nabla_{\beta} L(\beta, u)$, we see that this equality holds iff

$$\begin{aligned} \beta_r^* &= \left(\mathbb{E}[bb^T] + \sum_{j=1}^t \lambda_j \mathbb{E}[g_j b] \mathbb{E}[g_j b]^T \right)^{-1} \mathbb{E}[\tilde{Y} b] \\ &= \beta_{\lambda}^* \end{aligned}$$

If Assumption 3 holds, then λ is unique by Theorem 2 in Wachsmuth (2013). \square

B.3 Proof of Proposition 3

Proof Assumption 1 again ensures that β_{λ}^* exists and is unique. The KKT conditions for the risk-min program are satisfied by setting $\epsilon_j^2 = (\mathbb{P}_n[\hat{g}_j b^T \hat{\beta}_{\lambda}])^2$, which then implies that

$$\begin{aligned} \beta_{\lambda}^* &= \arg \min_{\beta \in \mathbb{R}^k} \mathbb{E}[(b^T \beta - \tilde{Y})^2] \\ &\text{subject to } (\mathbb{E}[g_j b^T \beta])^2 \leq (\mathbb{P}_n[\hat{g}_j b^T \hat{\beta}_{\lambda}])^2, \quad j = 1, \dots, t \\ &= \beta_r^* \end{aligned}$$

\square

B.4 Proof of Proposition 4

Proof Beginning with the risk, note that

$$\begin{aligned} \text{var}(Y^0) &= \mathbb{E}[(Y^0)^2] - (\mathbb{E}[Y^0])^2 \\ &= \mathbb{E}\{\mathbb{E}[(Y^0)^2 \mid W]\} - (\mathbb{E}\{\mathbb{E}[Y^0 \mid W]\})^2 \\ &= \mathbb{E}\{\mathbb{E}[Y^2 \mid W, D = 0]\} - (\mathbb{E}\{\mathbb{E}[Y \mid W, D = 0]\})^2 \\ &= \mathbb{E}[\nu_0] - (\mathbb{E}[\mu_0])^2 \end{aligned}$$

where the second line uses iterated expectation and the third line follows from the consistency and ignorability assumptions. We then have

$$\begin{aligned} \mathbb{E}[(f - Y^0)^2] &= \mathbb{E}\{\mathbb{E}[f^2 - 2Y^0 + (Y^0)^2 \mid W]\} \\ &= \mathbb{E}[f^2 - 2\mu_0 + \nu_0] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}[(f - \mu_0)^2] + (\mathbb{E}[\nu_0] - \mathbb{E}[\mu_0^2]) \\
&= \mathbb{E}[(f - \mu_0)^2] + \text{var}(Y^0)
\end{aligned}$$

The last and third-to-last lines give the equalities in the proposition. Turning to the FPR-diff, by Definition 2 and Proposition 1, we have

$$\begin{aligned}
\mathbb{E}[g^{\text{cFPR}} f(W)] &= \mathbb{E} \left[\left\{ \frac{(1 - Y^0)(1 - A)}{\mathbb{E}[(1 - Y^0)(1 - A)]} - \frac{(1 - Y^0)A}{\mathbb{E}[(1 - Y^0)A]} \right\} f(W) \right] \\
&= \mathbb{E} \left[\mathbb{E} \left\{ \frac{(1 - Y^0)(1 - A)}{\mathbb{E}[(1 - Y^0)(1 - A)]} - \frac{(1 - Y^0)A}{\mathbb{E}[(1 - Y^0)A]} \mid W \right\} f(W) \right] \\
&= \mathbb{E} \left[\left\{ \frac{(1 - \mu_0)(1 - A)}{\mathbb{E}[(1 - \mu_0)(1 - A)]} - \frac{(1 - \mu_0)A}{\mathbb{E}[(1 - \mu_0)A]} \right\} f(W) \right]
\end{aligned}$$

where the last line again follows from the consistency and ignorability assumptions. The result for g^{cFNR} follows by identical reasoning. \square

Appendix C Proof of Theorem 5

We prove this theorem first, since the result is used in the proofs of the other theorems. In the observable setting, the theorem follows immediately from the central limit theorem, so the subsequent derivations are for the counterfactual setting.

C.1 Asymptotic normality of the risk estimator

For any fixed predictor f_β , we have

$$\begin{aligned}
\widehat{\text{Risk}}(f_\beta) - \text{Risk}(f_\beta) &= \mathbb{P}_n[f_\beta^2 - (2b^T \beta)\widehat{\phi} + \widehat{\phi}] - \mathbb{P}[f_\beta^2 - (2b^T \beta)\phi + \phi] \\
&= (\mathbb{P}_n - \mathbb{P}) \{f_\beta^2 - 2f_\beta\phi + \phi\} + \\
&\quad (\mathbb{P}_n - \mathbb{P}) \{2f_\beta(\phi - \widehat{\phi}) + (\widehat{\phi} - \phi)\} + \\
&\quad \mathbb{P} \{2f_\beta(\phi - \widehat{\phi}) + (\widehat{\phi} - \phi)\}
\end{aligned}$$

The second term of the last equality is $O_{\mathbb{P}}(\|\widehat{\mu}_0 - \mu_0\| \|\widehat{\pi} - \pi\|/\sqrt{n}) = o_{\mathbb{P}}(1/\sqrt{n})$ by Lemma 1, Lemma 2, and Assumption 9. The third term is $o_{\mathbb{P}}(1/\sqrt{n})$ by Lemma 2 and Assumption 9. We therefore have

$$\widehat{\text{Risk}}(f_\beta) - \text{Risk}(f_\beta) = (\mathbb{P}_n - \mathbb{P}) \{f_\beta^2 - 2f_\beta\phi + \phi\} + o_{\mathbb{P}}(1/\sqrt{n}) \quad (\text{C1})$$

and the result follows by the central limit theorem.

C.2 Asymptotic normality of the unfairness estimators

Since g^{rate} does not depend on the outcome, we have $\widehat{g}^{\text{ind}} = g^{\text{rate}}$, and the result follows immediately from the central limit theorem. We now prove the

result for g^{FPR} in the counterfactual setting. We have

$$\mathbb{P}_n(\widehat{g}_j f_\beta) - \mathbb{P}(g_j f_\beta) = \left\{ \frac{\mathbb{P}_n[\widehat{\gamma}_0 f_\beta]}{\mathbb{P}_n[\widehat{\gamma}_0]} - \frac{\mathbb{P}_n[\widehat{\gamma}_1 f_\beta]}{\mathbb{P}_n[\widehat{\gamma}_1]} \right\} - \left\{ \frac{\mathbb{P}[\gamma_0 f_\beta]}{\mathbb{P}[\gamma_0]} - \frac{\mathbb{P}[\gamma_1 f_\beta]}{\mathbb{P}[\gamma_1]} \right\} \quad (\text{C2})$$

Considering just the $\widehat{\gamma}_0$ and γ_0 terms, we have

$$\begin{aligned} \frac{\mathbb{P}_n[\widehat{\gamma}_0 f_\beta]}{\mathbb{P}_n[\widehat{\gamma}_0]} - \frac{\mathbb{P}[\gamma_0 f_\beta]}{\mathbb{P}[\gamma_0]} &= \frac{\mathbb{P}_n[\widehat{\gamma}_0 f_\beta] \mathbb{P}[\gamma_0] - \mathbb{P}[\gamma_0] \mathbb{P}_n[\widehat{\gamma}_0]}{\mathbb{P}_n[\widehat{\gamma}_0] \mathbb{P}[\gamma_0]} \\ &= \frac{\mathbb{P}[\gamma_0] (\mathbb{P}_n[\widehat{\gamma}_0 f_\beta] - \mathbb{P}[\gamma_0 f_\beta]) - \mathbb{P}[\gamma_0 f_\beta] (\mathbb{P}_n[\widehat{\gamma}_0] - \mathbb{P}[\gamma_0])}{\mathbb{P}_n[\widehat{\gamma}_0] \mathbb{P}[\gamma_0]} \\ &= \mathbb{P}_n[\widehat{\gamma}_0]^{-1} \left\{ \underbrace{(\mathbb{P}_n[\widehat{\gamma}_0 f_\beta] - \mathbb{P}[\gamma_0 f_\beta])}_{(1)} - \frac{\mathbb{P}[\gamma_0 f_\beta]}{\mathbb{P}[\gamma_0]} \underbrace{(\mathbb{P}_n[\widehat{\gamma}_0] - \mathbb{P}[\gamma_0])}_{(2)} \right\} \quad (\text{C3}) \end{aligned}$$

Terms (1) and (2) in (C3) can be expanded as follows:

$$\begin{aligned} (1) &= (\mathbb{P}_n - \mathbb{P})\gamma_0 f_\beta + (\mathbb{P}_n - \mathbb{P})((\widehat{\gamma}_0 - \gamma_0) f_\beta) + \mathbb{P}((\widehat{\gamma}_0 - \gamma_0) f_\beta) \\ (2) &= (\mathbb{P}_n - \mathbb{P})\gamma_0 + (\mathbb{P}_n - \mathbb{P})(\widehat{\gamma}_0 - \gamma_0) + \mathbb{P}(\widehat{\gamma}_0 - \gamma_0) \end{aligned}$$

In both these expressions, the second term is $O_{\mathbb{P}}(\|\widehat{\phi} - \phi\|/\sqrt{n}) = o_{\mathbb{P}}(1/\sqrt{n})$ by Lemma 1 and Assumption 9, and the third term is $o_{\mathbb{P}}(1/\sqrt{n})$ by Lemma 2 and Assumption 9. Under Assumption 7, $\mathbb{P}_n[\widehat{\gamma}_0]^{-1}$ is bounded, while $\mathbb{P}[\gamma_0 f_\beta]/\mathbb{P}[\gamma_0]$ is bounded under Assumption 5. Therefore, we can rewrite (C2) as

$$\begin{aligned} &\mathbb{P}_n[\widehat{\gamma}_0]^{-1} (\mathbb{P}_n - \mathbb{P}) \left\{ \gamma_0 \left(f_\beta - \frac{\mathbb{P}[\gamma_0 f_\beta]}{\mathbb{P}[\gamma_0]} \right) \right\} + o_{\mathbb{P}}(1/\sqrt{n}) \\ &= \mathbb{P}_n[\widehat{\gamma}_0]^{-1} (\mathbb{P}_n - \mathbb{P}) \eta_0 + o_{\mathbb{P}}(1/\sqrt{n}) \end{aligned}$$

We can therefore rewrite $\mathbb{P}_n(\widehat{g}_j f_\beta) - \mathbb{P}(g_j f_\beta)$ as

$$\mathbb{P}_n[\widehat{\gamma}_0]^{-1} (\mathbb{P}_n - \mathbb{P}) \eta_0 + \mathbb{P}_n[\widehat{\gamma}_1]^{-1} (\mathbb{P}_n - \mathbb{P}) \eta_1 + o_{\mathbb{P}}(1/\sqrt{n})$$

Note that the analysis of term (2) in (C3) yields that $\mathbb{P}_n[\widehat{\gamma}_0] - \mathbb{P}[\gamma_0] = o_{\mathbb{P}}(1)$. Applying the central limit theorem to the vector (η_0, η_1) , followed the continuous mapping theorem, Slutsky's theorem, and the delta method, we have

$$\sqrt{n} (\mathbb{P}_n(\widehat{g}_j f_\beta) - \mathbb{P}(g_j f_\beta)) \xrightarrow{d} N(0, \text{var} (\mathbb{P}(\gamma_0)^{-1} \eta_0 - \mathbb{P}(\gamma_1)^{-1} \eta_1))$$

as desired. The result for g^{FNR} follows by identical reasoning.

Appendix D Proofs for constrained FADE

We state two additional lemmas that are only used in the constrained setting. The first lemma gives sufficient conditions under which the optimal value of an estimated convex problem converges at a particular rate to the optimal value of the target convex program. It is an adaptation of Theorem 3.5 in Shapiro (1991) that follows immediately from Theorems 2.1 and 3.4 in that same paper.

Lemma 3. (Shapiro, 1991) *Let Θ be a compact subset of \mathbb{R}^k . Let $C(\Theta)$ denote the set of continuous real-valued functions on Θ , with $\mathcal{L} = C(\Theta) \times \dots \times C(\Theta)$ the r -dimensional Cartesian product. Let $\psi(\theta) = (\psi_0, \dots, \psi_r) \in \mathcal{L}$ be a vector of convex functions. Consider the quantity α^* defined as the solution to the following convex optimization program:*

$$\begin{aligned} \alpha^* &= \min_{\theta \in \Theta} \psi_0(\theta) \\ &\text{subject to } \psi_j(\theta) \leq 0, \quad j = 1, \dots, r \end{aligned}$$

Assume that Slater's condition holds, so that there is some $\theta \in \Theta$ for which the inequalities are satisfied and non-affine inequalities are strictly satisfied, i.e. $\psi_j(\theta) < 0$ if ψ_j is non-affine. Now consider a sequence of approximating programs, for $n = 1, 2, \dots$:

$$\begin{aligned} \hat{\alpha}_n &= \min_{\theta \in \Theta} \hat{\psi}_{0n}(\theta) \\ &\text{subject to } \hat{\psi}_{jn}(\theta) \leq 0, \quad j = 1, \dots, r \end{aligned}$$

with $\hat{\psi}_n(\theta) := (\hat{\psi}_{0n}, \dots, \hat{\psi}_{rn}) \in \mathcal{L}$. Assume that $f(n)(\hat{\psi}_n - \psi)$ converges in distribution to a random element $W \in \mathcal{L}$ for some real-valued function $f(n)$. Then:

$$f(n)(\hat{\alpha}_n - \alpha_0) \rightsquigarrow L$$

for a particular random variable L . It follows that $\hat{\alpha}_n - \alpha_0 = O_{\mathbb{P}}(1/f(n))$.

D.1 Intermediate result

The next lemma applies Lemma 3 to the risk-min and unfair-min settings. For analytical purposes, we suppose that for each k , the quantities $\beta_r^*, \hat{\beta}_r, \beta_u^*, \hat{\beta}_u$ are constrained to lie in some (arbitrarily large) compact set $\Theta_k \subseteq \mathbb{R}^k$. Since $k \not\rightarrow \infty$, ultimately Θ_k is fixed to some set Θ . For example, Θ could be given by box constraints defined by the largest and smallest numbers the machine can represent. Since this is a device for asymptotic analysis, we do not express it in the actual optimization. Under Assumption 2, it follows that $b^T \beta$ is uniformly bounded in Θ . (Recall that in practice, the output of any predictor will be truncated to lie in $[\ell_y, u_y]$.)

Per Proposition 4 and Remark 3, we can write the objective function for the risk-min parameter β_r^* equivalently as $\mathbb{P}[(b^T \beta)^2 - 2(b^T \beta)\phi]$ in the counterfactual setting (or $\mathbb{P}[(b^T \beta)^2 - 2(b^T \beta)Y]$ in the observable setting), since the term $\mathbb{P}[\phi^2]$ (or $\mathbb{P}[Y^2]$) drops out of the minimization. We utilize this form for analysis.

Denote by $\psi_0, \dots, \psi_{t+1}$ and $\widehat{\psi}_0, \dots, \widehat{\psi}_{t+1}$ the population and empirical risk and unfairness functions, each of which is a mapping from Θ to \mathbb{R} . For the counterfactual setting, these are given by

$$\begin{aligned} \psi_0(\beta) &= \mathbb{P}[(b^T \beta)^2 - 2(b^T \beta)\phi] & \widehat{\psi}_0 &= \mathbb{P}_n[(b^T \beta)^2 - 2(b^T \beta)\widehat{\phi}] \\ \psi_j(\beta) &= (\mathbb{P}[g_j b^T \beta])^2 & \widehat{\psi}_j(\beta) &= (\mathbb{P}_n[\widehat{g}_j b^T \beta])^2, \quad j = 1, \dots, t \\ \psi_{t+1}(\beta) &= \mathbb{P}[(b^T \beta)^2 - (2b^T \beta)\phi + \underline{\phi}] & \psi_{t+1}(\beta) &= \mathbb{P}_n[(b^T \beta)^2 - (2b^T \beta)\widehat{\phi} + \underline{\widehat{\phi}}] \end{aligned} \tag{D4}$$

The observable setting substitutes Y for ϕ , Y^2 for ϕ , and g_j for \widehat{g}_j . Let $\mathcal{C}(\Theta)$ denote the set of continuous real-valued functions on Θ , with $\mathcal{L}(\Theta) = \mathcal{C}(\Theta) \times \dots \times \mathcal{C}(\Theta)$ the Cartesian product (with suitable dimension). Let $\psi_{(\bullet)}, \widehat{\psi}_{(\bullet)} : \Theta \mapsto \mathcal{L}(\Theta)$ be the vectors of functions that define the population and empirical optimization problem, for $\bullet \in \{r, u\}$, representing the risk-min and unfair-min problems. That is, for risk-min, define

$$\begin{aligned} \psi_{(r)} &= (\psi_0(\beta), \psi_1(\beta), \dots, \psi_t(\beta))^T \\ \widehat{\psi}_{(r)} &= (\widehat{\psi}_0(\beta), \widehat{\psi}_1(\beta), \dots, \widehat{\psi}_t(\beta))^T \end{aligned}$$

and for unfair-min, define

$$\begin{aligned} \psi_{(u)} &= \left(\sum_{j=1}^t \alpha_j \psi_j(\beta), \psi_{t+1}(\beta) \right)^T \\ \widehat{\psi}_{(u)} &= \left(\sum_{j=1}^t \alpha_j \widehat{\psi}_j(\beta), \widehat{\psi}_{t+1}(\beta) \right)^T \end{aligned}$$

The first element in each of $\psi_{(r)}$ and $\psi_{(u)}$ is the objective function, and the remaining elements are the constraint functions.

Lemma 4 (Convergence rates of estimated functions). *Under Assumptions 1 and 2 for the observable setting, and Assumptions 1 and 4–6 for the counterfactual setting, there exist random elements C_r, C_u taking values in the appropriate space $\mathcal{L}(\Theta)$ such that*

$$\begin{aligned} \sqrt{n}(\widehat{\psi}_{(r)} - \psi_{(r)}) &\xrightarrow{d} C_r \\ \sqrt{n}(\widehat{\psi}_{(u)} - \psi_{(u)}) &\xrightarrow{d} C_u \end{aligned}$$

where the convergence is in L_2 norm.

Proof We will utilize the fact that the class $\{b^T \beta : \beta \in \Theta\}$ is \mathbb{P} -Donsker, since $b^T \beta$ is parametric and Lipschitz in β under Assumption 2.

In the observable setting, we have

$$\widehat{\psi}_{(r)} - \psi_{(r)} = (\mathbb{P}_n - \mathbb{P}) \left((b^T \beta)^2 - 2(b^T \beta)Y, g_1 b^T \beta, \dots, g_t b^T \beta \right)^T$$

so that the result follows immediately from the central limit theorem and the Donsker condition. We now turn to the counterfactual setting. First, consider the objective function ψ_0 .

$$\begin{aligned} \widehat{\psi}_0(\beta) - \psi_0(\beta) &= \mathbb{P}_n \left\{ (b^T \beta)^2 - 2(b^T \beta)\widehat{\phi} \right\} - \mathbb{P} \left\{ (b^T \beta)^2 - 2(b^T \beta)\phi \right\} \\ &= (\mathbb{P}_n - \mathbb{P}) \left\{ (b^T \beta)^2 \right\} - \left\{ \mathbb{P}_n(b^T \beta\widehat{\phi}) - \mathbb{P}(b^T \beta\phi) \right\} \\ &= (\mathbb{P}_n - \mathbb{P}) \left\{ (b^T \beta)^2 - \phi \right\} + (\mathbb{P}_n - \mathbb{P})(2(b^T \beta)(\phi - \widehat{\phi})) + \mathbb{P}(2(b^T \beta)(\phi - \widehat{\phi})) \end{aligned}$$

The second term is $O_{\mathbb{P}}(\|\widehat{\mu}_0 - \mu_0\| \|\widehat{\pi} - \pi\|/\sqrt{n}) = o_{\mathbb{P}}(1/\sqrt{n})$ by Lemma 1, Lemma 2, and Assumption 9. The third term is $o_{\mathbb{P}}(1/\sqrt{n})$ by Lemma 2 and Assumption 9. We therefore have

$$\widehat{\psi}_0(\beta) - \psi_0(\beta) = (\mathbb{P}_n - \mathbb{P}) \left\{ (b^T \beta)^2 - \phi \right\} + o_{\mathbb{P}}(1/\sqrt{n}) \quad (\text{D5})$$

We now consider the unfairness functions $\psi_j, j = 1, \dots, t$. We have

$$\begin{aligned} \widehat{\psi}_j(\beta) - \psi_j(\beta) &= \left\{ \mathbb{P}_n(\widehat{g}_j b^T \beta) + \mathbb{P}(g_j b^T \beta) \right\} \left\{ \mathbb{P}_n(\widehat{g}_j b^T \beta) - \mathbb{P}(g_j b^T \beta) \right\} \\ &= \left\{ \mathbb{P}_n(\widehat{g}_j b^T \beta) + \mathbb{P}(g_j b^T \beta) \right\} \left(\mathbb{P}_n(\widehat{\gamma}_0)^{-1}, \mathbb{P}_n(\widehat{\gamma}_1)^{-1} \right) (\mathbb{P}_n - \mathbb{P}) \begin{pmatrix} \eta_0 \\ \eta_1 \end{pmatrix} + o_{\mathbb{P}}(1/\sqrt{n}) \end{aligned} \quad (\text{D6})$$

where the second line follows the derivation in Section C.2, coupled with the fact that $\mathbb{P}_n(\widehat{g}_j b^T \beta) + \mathbb{P}(g_j b^T \beta) = o_{\mathbb{P}}(1)$. Finally, the analysis of ψ_{t+1} is already given in Section C.1:

$$\widehat{\psi}_{t+1} - \psi_{t+1} = (\mathbb{P}_n - \mathbb{P}) \left((b^T \beta)^2 - 2(b^T \beta)\phi + \underline{\phi} \right) + o_{\mathbb{P}}(1/\sqrt{n}) \quad (\text{D7})$$

Suppose we have a single fairness function g_j . Combining (D5), (D6), and (D7), we have shown that $\widehat{\psi}_{(r)} - \psi_{(r)}$ can be written as

$$\widehat{\psi}_{(r)} - \psi_{(r)} = M(\mathbb{P}_n - \mathbb{P}) \begin{pmatrix} (b^T \beta)^2 - \phi \\ \eta_0 \\ \eta_1 \end{pmatrix} + \begin{pmatrix} o_{\mathbb{P}}(1/\sqrt{n}) \\ o_{\mathbb{P}}(1/\sqrt{n}) \end{pmatrix}, \text{ where}$$

$$M = \begin{bmatrix} 1 & & 0 \\ 0 & \left\{ \mathbb{P}_n(\widehat{g}_j b^T \beta) + \mathbb{P}(g_j b^T \beta) \right\} \mathbb{P}_n(\widehat{\gamma}_0)^{-1} - \left\{ \mathbb{P}_n(\widehat{g}_j b^T \beta) + \mathbb{P}(g_j b^T \beta) \right\} \mathbb{P}_n(\widehat{\gamma}_1)^{-1} & 0 \end{bmatrix}$$

Applying the central limit theorem, Slutsky's theorem, the continuous mapping theorem, and the delta method, we have that $\sqrt{n}(\widehat{\psi}_{(r)}(\beta) - \psi_{(r)}(\beta))$ converges to a normal distribution for any fixed β . Under the Donsker condition, this convergence is uniform over β , and $\sqrt{n}(\widehat{\psi}_{(r)} - \psi_{(r)})$ converges to a Gaussian process. Equivalent reasoning applies in the case of multiple fairness functions, and to $\sqrt{n}(\widehat{\psi}_{(u)} - \psi_{(u)})$. \square

We now prove Theorems 1 and 2. The two proofs proceed along similar lines. We will again utilize the fact that $\{b^T\beta : \beta \in \Theta\}$ is \mathbb{P} -Donsker, so that the empirical process $\{\sqrt{n}(\mathbb{P}_n - \mathbb{P})(b^T\beta) : \beta \in \Theta\}$ converges to a Gaussian process.

D.2 Proof of Theorem 1 (Excess risk in constrained FADE)

Proof We consider the risk-min problem first. We expand the excess risk by adding and subtracting the objective function at the solution $\widehat{\beta}_r$:

$$\begin{aligned} & \mathbb{P} \left[(b^T \widehat{\beta}_r)^2 - 2(b^T \widehat{\beta}_r) \widehat{\phi} \right] - \mathbb{P} \left[(b^T \beta_r^*)^2 - 2(b^T \beta_r^*) \phi \right] \\ &= \mathbb{P} \left[(b^T \widehat{\beta}_r)^2 - 2(b^T \widehat{\beta}_r) \widehat{\phi} \right] - \mathbb{P}_n \left[(b^T \widehat{\beta}_r)^2 - 2(b^T \widehat{\beta}_r) \widehat{\phi} \right] + \\ & \quad \mathbb{P}_n \left[(b^T \widehat{\beta}_r)^2 - 2(b^T \widehat{\beta}_r) \widehat{\phi} \right] - \mathbb{P} \left[(b^T \beta_r^*)^2 - 2(b^T \beta_r^*) \phi \right] \end{aligned}$$

The second term is $O_{\mathbb{P}}(1/\sqrt{n})$ by Lemma 4 and Shapiro's theorem. The first term is just $\psi_0(\widehat{\beta}_r) - \widehat{\psi}_0(\widehat{\beta}_r)$, which is $O_{\mathbb{P}}(1/\sqrt{n})$ by (D5) in the proof of Lemma 4 coupled with the Donsker condition. Hence, the excess risk is $O_{\mathbb{P}}(1/\sqrt{n})$, as claimed.

We now turn to the unfair-min problem. The excess risk is

$$\begin{aligned} & \mathbb{P}[(b^T \widehat{\beta}_u)^2 - 2(b^T \widehat{\beta}_u) \phi + \underline{\phi}] - \epsilon^2 \\ & \leq \mathbb{P}[(b^T \widehat{\beta}_u)^2 - 2(b^T \widehat{\beta}_u) \phi + \underline{\phi}] - \mathbb{P}_n[(b^T \widehat{\beta}_u)^2 - 2(b^T \widehat{\beta}_u) \widehat{\phi} + \widehat{\underline{\phi}}] \\ &= -(\mathbb{P}_n - \mathbb{P}) \left[(b^T \widehat{\beta}_u)^2 + (2b^T \widehat{\beta}_u) \phi + \underline{\phi} \right] + \\ & \quad (\mathbb{P}_n - \mathbb{P}) \left[(2b^T \widehat{\beta}_u) (\widehat{\phi} - \phi) + (\widehat{\underline{\phi}} - \underline{\phi}) \right] + \\ & \quad \mathbb{P} \left[(2b^T \widehat{\beta}_u) (\widehat{\phi} - \phi) + (\widehat{\underline{\phi}} - \underline{\phi}) \right] \end{aligned}$$

The first term is $O_{\mathbb{P}}(1/\sqrt{n})$ by the central limit theorem and the Donsker condition. The second term is $O_{\mathbb{P}}(\|\widehat{\phi} - \phi\|/\sqrt{n}) = o_{\mathbb{P}}(1/\sqrt{n})$ by Lemma 1, Lemma 2, and Assumption 8. The last term is $o_{\mathbb{P}}(1/\sqrt{n})$ by Lemma 2. The excess risk is therefore $O_{\mathbb{P}}(1/\sqrt{n})$, as claimed. \square

D.3 Proof of Theorem 2 (Excess unfairness in constrained FADE)

Proof We consider the unfair-min problem first. We expand the excess unfairness by adding and subtracting the objective function at the solution $\widehat{\beta}_u$:

$$\begin{aligned} & \sum_{j=1}^t \alpha_j (\mathbb{P}[g_j b^T \widehat{\beta}_u])^2 - \sum_{j=1}^t \alpha_j (\mathbb{P}[g_j b^T \beta_u^*])^2 \\ &= \sum_{j=1}^t \alpha_j \left\{ (\mathbb{P}[g_j b^T \widehat{\beta}_u])^2 - (\mathbb{P}_n[\widehat{g}_j b^T \widehat{\beta}_u])^2 \right\} + \sum_{j=1}^t \alpha_j \left\{ (\mathbb{P}_n[\widehat{g}_j b^T \widehat{\beta}_u])^2 - (\mathbb{P}[g_j b^T \beta_u^*])^2 \right\} \end{aligned} \tag{D8}$$

Again, the second term is $O_{\mathbb{P}}(1/\sqrt{n})$ by Lemma 4 and Shapiro's theorem. The first term is equal to $\sum_{j=1}^t \alpha_j (\psi_j(\widehat{\beta}) - \widehat{\psi}_j(\widehat{\beta}))$, which is $O_{\mathbb{P}}(1/\sqrt{n})$ by (D6) in the proof

of Lemma 4 coupled with the Donsker condition. The excess unfairness is therefore $O_{\mathbb{P}}(1/\sqrt{n})$, as claimed.

We now turn to the risk-min problem. The excess unfairness for constraint j is

$$\begin{aligned} & (\mathbb{P}[g_j b^T \widehat{\beta}_u])^2 - \epsilon^2 \\ & \leq (\mathbb{P}[g_j b^T \widehat{\beta}_u])^2 - (\mathbb{P}_n[\widehat{g}_j b^T \widehat{\beta}_u])^2 \\ & = O_{\mathbb{P}}(1/\sqrt{n}) \end{aligned}$$

where the last line simply uses the analysis for term (1) from (D8). The excess unfairness is therefore $O_{\mathbb{P}}(1/\sqrt{n})$, as claimed. \square

Appendix E Proofs for penalized FADE

Throughout this section, let

$$\begin{aligned} \widehat{\mathbf{Q}}_{\lambda} &= \mathbb{P}_n(bb^T) + \sum_{j=1}^t \lambda_j \mathbb{P}_n(\widehat{g}_j b) \mathbb{P}_n(\widehat{g}_j b)^T \\ \mathbf{Q}_{\lambda} &= \mathbb{P}(bb^T) + \sum_{j=1}^t \lambda_j \mathbb{P}(g_j b) \mathbb{P}(g_j b)^T \end{aligned}$$

so that

$$\begin{aligned} \beta_{\lambda} &= \mathbf{Q}_{\lambda}^{-1} \mathbb{P}(b\widetilde{Y}) \\ \widehat{\beta}_{\lambda} &= \begin{cases} \widehat{\mathbf{Q}}_{\lambda}^{-1} \mathbb{P}_n(bY) & \text{(Observable)} \\ \widehat{\mathbf{Q}}_{\lambda}^{-1} \mathbb{P}_n(b\widehat{\phi}) & \text{(Counterfactual)} \end{cases} \end{aligned}$$

Under the assumptions of Theorems 3 and 4, we prove several preliminary results that are used in the theorem proofs.

Lemma 5. *(Bounded norm for $\widehat{\mathbf{Q}}_{\lambda}^{-1}$).*

$$\mathbb{P}(\|\widehat{\mathbf{Q}}_{\lambda}^{-1}\|) \leq C \rightarrow 1 \text{ for some constant } C$$

Proof This follows from Assumption 1 plus the consistency of $\widehat{\mathbf{Q}}_{\lambda}$ for \mathbf{Q}_{λ} . It follows that $\|\widehat{\mathbf{Q}}_{\lambda}^{-1}\| = O_{\mathbb{P}}(1)$. \square

Lemma 6. *Fix a $\lambda \in \Lambda$. Then*

$$\|\widehat{\beta}_{\lambda} - \beta_{\lambda}^*\| = \begin{cases} O_{\mathbb{P}}(\sqrt{1/n}) & \text{(Observable)} \\ O_{\mathbb{P}}(\sqrt{1/n}) + O_{\mathbb{P}}(h(n)) & \text{(Counterfactual)} \end{cases}$$

Proof In the observable setting, we have

$$\widehat{\beta}_{\lambda} - \beta_{\lambda}^* = (\widehat{\mathbf{Q}}_{\lambda}^{-1} - \mathbf{Q}_{\lambda}^{-1}) \mathbb{P}(bY)$$

$$\begin{aligned}
&= \widehat{\mathbf{Q}}_\lambda^{-1}(\mathbf{Q}_\lambda - \widehat{\mathbf{Q}}_\lambda)\mathbf{Q}^{-1}\mathbb{P}(bY) \\
&= \widehat{\mathbf{Q}}_\lambda^{-1}(\mathbf{Q}_\lambda - \widehat{\mathbf{Q}}_\lambda)\beta_\lambda^* \\
&= \widehat{\mathbf{Q}}_\lambda^{-1}\left\{(\mathbb{P}_n - \mathbb{P})(bb^T\beta_\lambda^*) + \right. \\
&\quad \left. \sum_{j=1}^t \left[\lambda_j(\mathbb{P}_n - \mathbb{P})(bg_j)\mathbb{P}(g_j b^T\beta_\lambda^*) + \mathbb{P}_n(bg_j)(\mathbb{P}_n - \mathbb{P})(g_j b^T\beta_\lambda^*) \right] \right\}
\end{aligned}$$

The norm of each term in the braces is $O_{\mathbb{P}}(1/\sqrt{n})$ by the central limit theorem. By Lemma 5, $\widehat{\mathbf{Q}}_\lambda^{-1}$ doesn't contribute to the rate, so $\|\widehat{\beta}_\lambda - \beta_\lambda^*\| = O_{\mathbb{P}}(1/\sqrt{n})$ as claimed.

In the counterfactual setting, we have

$$\begin{aligned}
\widehat{\beta}_\lambda - \beta_\lambda^* &= \widehat{\mathbf{Q}}_\lambda^{-1}\mathbb{P}_n(b\widehat{\phi}) - \mathbf{Q}_\lambda^{-1}\mathbb{P}(b\phi) \\
&= (\widehat{\mathbf{Q}}_\lambda^{-1} - \mathbf{Q}_\lambda^{-1})\mathbb{P}(b\phi) + \widehat{\mathbf{Q}}_\lambda^{-1}(\mathbb{P}_n(b\widehat{\phi}) - \mathbb{P}(b\phi)) \\
&= \widehat{\mathbf{Q}}_\lambda^{-1}(\mathbf{Q}_\lambda - \widehat{\mathbf{Q}}_\lambda)\mathbf{Q}^{-1}\mathbb{P}(b\phi) + \widehat{\mathbf{Q}}_\lambda^{-1}(\mathbb{P}_n(b\widehat{\phi}) - \mathbb{P}(b\phi)) \\
&= \underbrace{\widehat{\mathbf{Q}}_\lambda^{-1}(\mathbf{Q}_\lambda - \widehat{\mathbf{Q}}_\lambda)\beta_\lambda^*}_{(1)} + \underbrace{\widehat{\mathbf{Q}}_\lambda^{-1}(\mathbb{P}_n(b\widehat{\phi}) - \mathbb{P}(b\phi))}_{(2)} \tag{E9}
\end{aligned}$$

The norm of term (2) in (E9) is $O_{\mathbb{P}}(1/\sqrt{n}) + O_{\mathbb{P}}(h(n))$ by Lemma 2 and Lemma 5. For term (1), ignoring the leading $\widehat{\mathbf{Q}}_\lambda^{-1}$ for now, we have

$$\begin{aligned}
(\mathbf{Q}_\lambda - \widehat{\mathbf{Q}}_\lambda)\beta_\lambda^* &= \underbrace{(\mathbb{P}_n - \mathbb{P})(bb^T\beta_\lambda^*)}_{(a)} + \\
&\quad \sum_{j=1}^t \lambda_j \left[\underbrace{(\mathbb{P}_n(b\widehat{g}_j) - \mathbb{P}(bg_j))\mathbb{P}(g_j b^T\beta_\lambda^*)}_{(b)} + \underbrace{\mathbb{P}_n(b\widehat{g}_j)(\mathbb{P}_n(\widehat{g}_j b^T\beta_\lambda^*) - \mathbb{P}(gb^T\beta_\lambda^*))}_{(c)} \right]
\end{aligned}$$

The norm of term (a) is $O_{\mathbb{P}}(\sqrt{1/n})$ by the central limit theorem. Terms (b) and (c) decompose as follows:

$$\begin{aligned}
(b) &= \mathbb{P}(g_j b^T\beta_\lambda^*) \left\{ (\mathbb{P}_n - \mathbb{P})(bg_j) + (\mathbb{P}_n - \mathbb{P})(b(\widehat{g}_j - g_j)) + \mathbb{P}(b(\widehat{g}_j - g_j)) \right\} \\
(c) &= \mathbb{P}_n(b\widehat{g}_j) \left\{ (\mathbb{P}_n - \mathbb{P})(g_j b^T\beta_\lambda^*) + (\mathbb{P}_n - \mathbb{P})(\widehat{g}_j - g_j)(b^T\beta_\lambda^*) + \mathbb{P}((\widehat{g}_j - g_j)b^T\beta_\lambda^*) \right\}
\end{aligned}$$

The norms of the first term in braces in each of these two expressions is $O_{\mathbb{P}}(1/\sqrt{n})$ by the central limit theorem. The norm of the second term is $o_{\mathbb{P}}(1/\sqrt{n})$ by Lemma 1 and Assumption 8. The norm of the third term is $O_{\mathbb{P}}(1/\sqrt{n}) + O_{\mathbb{P}}(h(n))$ by Lemma 2. Using Lemma 5, the consistency of $\mathbb{P}_n(b\widehat{g}_j)$ for $\mathbb{P}(bg_j)$, and the boundedness of $\mathbb{P}(g_j b^T\beta_\lambda^*)$, we have

$$\|\widehat{\beta}_\lambda - \beta_\lambda^*\| = O_{\mathbb{P}}(\sqrt{1/n}) + O_{\mathbb{P}}(h(n))$$

as claimed. \square

We now prove the two theorems. We will use the fact that under Assumption 1, $\widehat{\beta}_\lambda - \beta_\lambda^*$ is Lipschitz in λ , and Λ is compact, so the set $\{\widehat{\beta}_\lambda - \beta_\lambda^* : \lambda \in \Lambda\}$ is Donsker.

E.1 Proof of Theorem 3 (Excess risk in penalized FADE)

Fix a $\lambda \in \Lambda$. We have

$$\begin{aligned} \mathbb{P} \left[\left(b^T \widehat{\beta}_\lambda - \widetilde{Y} \right)^2 \right] - \mathbb{P} \left[\left(b^T \beta_\lambda^* - \widetilde{Y} \right)^2 \right] &= \|b^T \widehat{\beta}_\lambda - \widetilde{Y}\|^2 - \|b^T \beta_\lambda^* - \widetilde{Y}\|^2 \\ &= \left(\|b^T \widehat{\beta}_\lambda - \widetilde{Y}\| - \|b^T \beta_\lambda^* - \widetilde{Y}\| \right) \left(\|b^T \widehat{\beta}_\lambda - \widetilde{Y}\| + \|b^T \beta_\lambda^* - \widetilde{Y}\| \right) \end{aligned}$$

Since $\widehat{\beta}_\lambda$ is consistent for β_λ^* (by Lemma 6), the second factor is $O_{\mathbb{P}}(1)$, so we can just consider the first factor.

$$\begin{aligned} \|b^T \widehat{\beta}_\lambda - \widetilde{Y}\| - \|b^T \beta_\lambda^* - \widetilde{Y}\| &\leq \|b^T \widehat{\beta}_\lambda - b^T \beta_\lambda^*\| \\ &= O_{\mathbb{P}} \left(\|\widehat{\beta}_\lambda - \beta_\lambda^*\| \right) \\ &= O_{\mathbb{P}}(\sqrt{1/n}) + O_{\mathbb{P}}(h(n)) \end{aligned}$$

where the first line uses the reverse triangle inequality, the second line uses Assumption 2, and the third line uses Lemma 6. Under the Donsker condition, the convergence is uniform over Λ :

$$\sup_{\lambda \in \Lambda} \left\{ \mathbb{P} \left[\left(b^T \widehat{\beta}_\lambda - \widetilde{Y} \right)^2 \right] - \mathbb{P} \left[\left(b^T \beta_\lambda^* - \widetilde{Y} \right)^2 \right] \right\} = O_{\mathbb{P}}(\sqrt{1/n}) + O_{\mathbb{P}}(h(n))$$

E.2 Proof of Theorem 4 (Excess unfairness in penalized FADE)

Fix a $\lambda \in \Lambda$. The excess unfairness for g_j is

$$\begin{aligned} \mathbb{P} \left[g_j b^T \widehat{\beta}_\lambda \right] - \mathbb{P} \left[g_j b^T \beta_\lambda^* \right] &= \mathbb{P} [g_j b^T (\widehat{\beta}_\lambda - \beta_\lambda^*)] \\ &\leq \mathbb{P} [|g_j b^T (\widehat{\beta}_\lambda - \beta_\lambda^*)|] \\ &= \sqrt{(\mathbb{P} [|g_j b^T (\widehat{\beta}_\lambda - \beta_\lambda^*)|])^2} \\ &\leq \sqrt{\mathbb{P} [(g_j b^T (\widehat{\beta}_\lambda - \beta_\lambda^*))^2]} \\ &= O_{\mathbb{P}} \left(\sqrt{\mathbb{P} [(\widehat{\beta}_\lambda - \beta_\lambda^*)^2]} \right) \\ &= \|\widehat{\beta}_\lambda - \beta_\lambda^*\| \\ &= O_{\mathbb{P}}(\sqrt{1/n}) + O_{\mathbb{P}}(h(n)) \end{aligned}$$

where the last line uses Lemma 6. Under the Donsker condition, the convergence is uniform over Λ :

$$\sup_{\lambda \in \Lambda} \left\{ \max_{j \in \{1, \dots, t\}} \left(\mathbb{P} \left[g_j b^T \widehat{\beta}_\lambda \right] - \mathbb{P} \left[g_j b^T \beta_\lambda^* \right] \right) \right\} = O_{\mathbb{P}}(\sqrt{1/n}) + O_{\mathbb{P}}(h(n))$$

Appendix F Bases with dimension $k \geq n$

We can generalize our estimators slightly to accommodate case where where $k \geq n$, meaning the dimension of the basis is greater than the sample size, as is the case for example with smoothing splines or RKHSs. We simply add an appropriate penalty matrix term $\lambda_0 \beta^T \mathbf{K} \beta$ to the penalized estimator expression or to the objective function for the constrained estimators, where \mathbf{K} is a $k \times k$ smoothing matrix. In the former case, for example, the estimator $\widehat{\beta}_\lambda$ becomes

$$\arg \min_{\beta \in \mathbb{R}^k} \mathbb{P}_n[(b^T \beta - \widehat{\phi})^2] + \lambda_0 \beta^T \mathbf{K} \beta + \sum_{j=1}^t \lambda_j (\mathbb{P}_n[\widehat{g}_j b^T \beta])^2 \quad (\text{F10})$$

$$= \left(\mathbb{P}_n(bb^T) + \lambda_0 \mathbf{K} + \sum_{j=1}^t \lambda_j \mathbb{P}_n(\widehat{g}_j b) \mathbb{P}_n(\widehat{g}_j b)^T \right)^{-1} \mathbb{P}_n(b \widehat{\phi}) \quad (\text{F11})$$

For instance, in a smoothing spline setting, b represents a spline basis, and $\mathbf{K}_{ij} = \int_{\mathcal{W}} b_i''(w) b_j''(w) dw$. In an RKHS, we'd have $b_i = \sum_{j=1}^n k(\cdot, w_j)$ and $\mathbf{K}_{ij} = k(w_i, w_j)$. In a ridge setting we'd have $\mathbf{K} = I$. The penalty term ensures the invertibility of the large matrix in (F11), and it preserves the fast computability of a large set of solutions $\widehat{\beta}_\lambda$.

This penalty term may also be useful to prevent overfitting even in cases where $k < n$, if k is close to n or if the basis is very expressive.

References

- Agarwal A, Beygelzimer A, Dudik M, et al (2018) A reductions approach to fair classification. In: Dy J, Krause A (eds) Proceedings of the 35th International Conference on Machine Learning, vol 80. PMLR, pp 60–69, URL <http://proceedings.mlr.press/v80/agarwal18a.html>
- Angelino E, Larus-Stone N, Alabi D, et al (2018) Learning certifiably optimal rule lists for categorical data. Journal of Machine Learning Research 18:1–78. URL <http://jmlr.org/papers/v18/17-716.html>
- Angwin J, Larson J (2016) Bias in criminal risk scores is mathematically inevitable, researchers say. ProPublica URL <https://www.propublica.org/article/bias-in-criminal-risk-scores-is-mathematically-inevitable-researchers-say>
- Angwin J, Larson J, Mattu S, et al (2016) Machine bias. ProPublica URL <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>
- Barocas S, Hardt M, Narayanan A (2018) Fairness and Machine Learning. <http://www.fairmlbook.org>. URL <http://www.fairmlbook.org>

- Bellamy RKE, Dey K, Hind M, et al (2018) AI fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias. Arxiv preprint <https://arxiv.org/abs/1810.01943>
- Belloni A, Chernozhukov V, Chetverikov D, et al (2015) Some new asymptotic theory for least squares series: Pointwise and uniform results. Arxiv preprint <https://arxiv.org/abs/1212.0442>
- Berk R, Heidari H, Jabbari S, et al (2017) A convex framework for fair regression. Arxiv preprint <https://arxiv.org/abs/1706.02409>
- Bickel PJ, Klaassen CA, Ritov Y, et al (1993) Efficient and adaptive estimation for semiparametric models. Johns Hopkins series in the mathematical sciences, Johns Hopkins University Press, Baltimore, MD
- Bonvini M, Kennedy EH (2021) Sensitivity analysis via the proportion of unmeasured confounding. Journal of the American Statistical Association pp 1–11. <https://doi.org/10.1080/01621459.2020.1864382>
- Breiman L (1996) Stacked regressions. Machine Learning 24(1):49–64. <https://doi.org/10.1007/BF00117832>
- Buolamwini J, Gebru T (2018) Gender shades: Intersectional accuracy disparities in commercial gender classification. In: Friedler SA, Wilson C (eds) Proceedings of the 1st Conference on Fairness, Accountability and Transparency, Proceedings of Machine Learning Research, vol 81. PMLR, New York, NY, USA, pp 77–91, URL <http://proceedings.mlr.press/v81/buolamwini18a.html>
- Calders T, Verwer S (2010) Three naive bayes approaches for discrimination-free classification. Data Mining and Knowledge Discovery 21(2):277–292. <https://doi.org/10.1007/s10618-010-0190-x>
- Calders T, Kamiran F, Pechenizkiy M (2009) Building classifiers with independence constraints. In: 2009 IEEE International Conference on Data Mining Workshops. IEEE, pp 13–18, <https://doi.org/10.1109/ICDMW.2009.83>
- Calmon F, Wei D, Vinzamuri B, et al (2017) Optimized pre-processing for discrimination prevention. In: Guyon I, Luxburg UV, Bengio S, et al (eds) Advances in Neural Information Processing Systems. Curran Associates, Inc., NIPS 2017, pp 3992–4001, URL <http://papers.nips.cc/paper/6988-optimized-pre-processing-for-discrimination-prevention.pdf>
- Celis LE, Huang L, Keswani V, et al (2020) Classification with fairness constraints: A meta-algorithm with provable guarantees. Arxiv preprint <https://arxiv.org/abs/1806.06055>

- Chernozhukov V, Chetverikov D, Demirer M, et al (2018) Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal* 21(1):C1–C68. <https://doi.org/10.1111/ectj.12097>
- Chouldechova A (2017) Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data* 5(2):153–163. <https://doi.org/10.1089/big.2016.0047>
- Corbett-Davies S, Pierson E, Feller A, et al (2017) Algorithmic decision making and the cost of fairness. In: *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. Association for Computing Machinery, New York, NY, USA, KDD '17, p 797–806, <https://doi.org/10.1145/3097983.3098095>
- Coston A, Mishler A, Kennedy EH, et al (2020) Counterfactual risk assessments, evaluation, and fairness. In: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. Association for Computing Machinery, New York, NY, USA, FAT* '20, p 582–593, <https://doi.org/10.1145/3351095.3372851>
- Coston A, Rambachan A, Chouldechova A (2021) Characterizing fairness over the set of good models under selective labels. *Arxiv preprint* <https://arxiv.org/abs/2101.00352>
- Donini M, Oneto L, Ben-David S, et al (2018) Empirical risk minimization under fairness constraints. In: Bengio S, Wallach H, Larochelle H, et al (eds) *Advances in Neural Information Processing Systems 31*. Curran Associates, Inc., NeurIPS 2018, pp 2791–2801, URL <http://papers.nips.cc/paper/7544-empirical-risk-minimization-under-fairness-constraints.pdf>
- Dua D, Graff C (2017) UCI machine learning repository. URL <http://archive.ics.uci.edu/ml>
- Dutta S, Wei D, Yueksel H, et al (2020) Is there a trade-off between fairness and accuracy? A perspective using mismatched hypothesis testing. In: III HD, Singh A (eds) *Proceedings of the 37th International Conference on Machine Learning, Proceedings of Machine Learning Research*, vol 119. PMLR, pp 2803–2813, URL <https://proceedings.mlr.press/v119/dutta20a.html>
- Feldman M, Friedler SA, Moeller J, et al (2015) Certifying and removing disparate impact. In: *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. Association for Computing Machinery, New York, NY, USA, KDD '15, p 259–268, <https://doi.org/10.1145/2783258.2783311>

- Friedler SA, Scheidegger C, Venkatasubramanian S, et al (2019) A comparative study of fairness-enhancing interventions in machine learning. In: Proceedings of the Conference on Fairness, Accountability, and Transparency - FAT* '19. ACM Press, pp 329–338, <https://doi.org/10.1145/3287560.3287589>, URL <http://dl.acm.org/citation.cfm?doid=3287560.3287589>
- Györfi L, Kohler M, Krzyzak A, et al (2002) A Distribution-Free Theory of Nonparametric Regression. Springer Series in Statistics, Springer-Verlag, New York, NY
- Hardt M, Price E, Price E, et al (2016) Equality of opportunity in supervised learning. In: Lee DD, Sugiyama M, Luxburg UV, et al (eds) Advances in Neural Information Processing Systems 29. Curran Associates, Inc., NIPS 2016, pp 3315–3323, URL <http://papers.nips.cc/paper/6374-equality-of-opportunity-in-supervised-learning.pdf>
- Holland PW (1986) Statistics and causal inference. *Journal of the American Statistical Association* 81(396):968. <https://doi.org/10.2307/2289069>
- Juditsky A, Nemirovski A (2000) Functional aggregation for nonparametric regression. *The Annals of Statistics* 28(3). <https://doi.org/10.1214/aos/1015951994>
- Kamiran F, Calders T (2012) Data preprocessing techniques for classification without discrimination. *Knowledge and Information Systems* 33(1):1–33. <https://doi.org/10.1007/s10115-011-0463-8>
- Kamishima T, Akaho S, Asoh H, et al (2012) Fairness-aware classifier with prejudice remover regularizer. In: Flach PA, De Bie T, Cristianini N (eds) *Machine Learning and Knowledge Discovery in Databases*, vol 7524. Springer Berlin Heidelberg, pp 35–50, https://doi.org/10.1007/978-3-642-33486-3_3
- Kennedy EH (2016) Semiparametric theory and empirical processes in causal inference. In: He H, Wu P, Chen DGD (eds) *Statistical Causal Inferences and Their Applications in Public Health Research*. Springer International Publishing, p 141–167, https://doi.org/10.1007/978-3-319-41259-7_8
- Kennedy EH, Balakrishnan S, G'Sell M (2020) Sharp instruments for classifying compliers and generalizing causal effects. *Annals of Statistics* 48(4):2008–2030. <https://doi.org/10.1214/19-AOS1874>
- Kilbertus N, Rojas Carulla M, Parascandolo G, et al (2017) Avoiding discrimination through causal reasoning. In: Guyon I, Luxburg UV, Bengio S, et al (eds) *Advances in Neural Information Processing Systems* 30. Curran Associates, Inc., NIPS 2017, pp 656–666, URL <http://papers.nips.cc/paper/6668-avoiding-discrimination-through-causal-reasoning.pdf>

- Kim JS, Chen J, Talwalkar A (2020) FACT: A diagnostic for group fairness trade-offs. In: III HD, Singh A (eds) Proceedings of the 37th International Conference on Machine Learning, Proceedings of Machine Learning Research, vol 119. PMLR, pp 5264–5274, URL <https://proceedings.mlr.press/v119/kim20a.html>
- Kim MP, Ghorbani A, Zou J (2019) Multiaccuracy: Black-box post-processing for fairness in classification. In: Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society. Association for Computing Machinery, New York, NY, USA, AIES '19, p 247–254, <https://doi.org/10.1145/3306618.3314287>
- Kleinberg J, Mullainathan S, Raghavan M (2017) Inherent trade-offs in the fair determination of risk scores. In: Papadimitriou CH (ed) 8th Innovations in Theoretical Computer Science Conference (ITCS 2017), Leibniz International Proceedings in Informatics (LIPIcs), vol 67. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, Dagstuhl, Germany, pp 43:1–43:23, <https://doi.org/10.4230/LIPIcs.ITCS.2017.0>
- Kusner MJ, Loftus J, Russell C, et al (2017) Counterfactual fairness. In: Guyon I, Luxburg UV, Bengio S, et al (eds) Advances in Neural Information Processing Systems 30. Curran Associates, Inc., NIPS 2017, pp 4066–4076, URL <http://papers.nips.cc/paper/6995-counterfactual-fairness.pdf>
- van der Laan MJ, Robins JM (2003) Unified Methods for Censored Longitudinal Data and Causality. Springer Series in Statistics, Springer, New York, NY, <https://doi.org/10.1007/978-0-387-21700-0>
- Liu S, Vicente LN (2021) Accuracy and fairness trade-offs in machine learning: A stochastic multi-objective approach. Arxiv preprint <https://arxiv.org/abs/2008.01132>
- Liu W, Kuramoto SJ, Stuart EA (2013) An introduction to sensitivity analysis for unobserved confounding in nonexperimental prevention research. *Prevention Science* 14(6):570–580. <https://doi.org/10.1007/s11121-012-0339-5>
- Menon AK, Williamson RC (2018) The cost of fairness in binary classification. In: Friedler SA, Wilson C (eds) Proceedings of the 1st Conference on Fairness, Accountability and Transparency, Proceedings of Machine Learning Research, vol 81. PMLR, New York, NY, pp 107–118, URL <http://proceedings.mlr.press/v81/menon18a.html>
- Mishler A (2019) Modeling Risk and Achieving Algorithmic Fairness Using Potential Outcomes. In: Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society - AIES '19. ACM Press, Honolulu, HI, pp 555–556, <https://doi.org/10.1145/3306618.3314323>

- Mishler A, Kennedy EH, Chouldechova A (2021) Fairness in risk assessment instruments: Post-processing to achieve counterfactual equalized odds. In: Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency. ACM, pp 386–400, <https://doi.org/10.1145/3442188.3445902>
- Nabi R, Shpitser I (2018) Fair inference on outcomes. In: Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence. Association for the Advancement of Artificial Intelligence, pp 1931–1940, URL <https://aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16683>
- Nabi R, Malinsky D, Shpitser I (2019) Learning optimal fair policies. In: Chaudhuri K, Salakhutdinov R (eds) Proceedings of the 36th International Conference on Machine Learning, Proceedings of Machine Learning Research, vol 97. PMLR, pp 4674–4682, URL <http://proceedings.mlr.press/v97/nabi19a.html>
- Neyman J (1923) Justification of applications of the calculus of probabilities to the solutions of certain questions in agricultural experimentation. Excerpts english translation (Reprinted). *Statistical Science* 5:463–472
- Northpointe (2015) Practitioners guide to COMPAS core. URL http://www.northpointeinc.com/downloads/compas/Practitioners-Guide-COMPAS-Core-_031915.pdf
- Obermeyer Z, Powers B, Vogeli C, et al (2019) Dissecting racial bias in an algorithm used to manage the health of populations. *Science* 366(6464):447–453. <https://doi.org/10.1126/science.aax2342>
- Pleiss G, Raghavan M, Wu F, et al (2017) On fairness and calibration. In: Guyon I, Luxburg UV, Bengio S, et al (eds) *Advances in Neural Information Processing Systems*, vol 30. Curran Associates, Inc., URL <https://proceedings.neurips.cc/paper/2017/file/b8b9c74ac526fffb2d39ab038d1cd7-Paper.pdf>
- Polley EC, Van Der Laan MJ (2010) Super learner in prediction. UC Berkeley Division of Biostatistics Working Paper Series 266. URL <https://biostat.berpress.com/ucbbiostat/paper266>
- Raskutti G, Wainwright MJ, Yu B (2011) Minimax rates of estimation for high-dimensional linear regression over ℓ_q -balls. *IEEE transactions on information theory* 57(10):6976–6994. URL <https://doi.org/10.1109/TIT.2011.2165799>
- Richardson A, Hudgens MG, Gilbert PB, et al (2014) Nonparametric bounds and sensitivity analysis of treatment effects. *Statistical Science* 29(4). <https://doi.org/10.1214/14-STS499>

- Rodolfa KT, Lamba H, Ghani R (2021) Empirical observation of negligible fairness-accuracy trade-offs in machine learning for public policy. Arxiv preprint <https://arxiv.org/abs/2012.02972>
- Rosenbaum PR (1987) Sensitivity analysis for certain permutation inferences in matched observational studies. *Biometrika* 74(1):13–26. <https://doi.org/10.2307/2336017>
- Rudin C, Wang C, Coker B (2020) The age of secrecy and unfairness in recidivism prediction. *Harvard Data Science Review* 2(1). <https://doi.org/10.1162/99608f92.6ed64b30>
- Shapiro A (1991) Asymptotic analysis of stochastic programs. *Annals of Operations Research* 30(1):169–186. <https://doi.org/10.1007/BF02204815>
- Tsiatis AA (2006) *Semiparametric Theory and Missing Data*. Springer, New York, NY, <https://doi.org/10.1007/0-387-37345-4>
- Tsybakov AB (2003) Optimal rates of aggregation. In: Schölkopf B, Warmuth MK (eds) *Learning Theory and Kernel Machines, Lecture Notes in Computer Science*, vol 2777. Springer Berlin Heidelberg, p 303–313, https://doi.org/10.1007/978-3-540-45167-9_23
- Wachsmuth G (2013) On LICQ and the uniqueness of lagrange multipliers. *Operations Research Letters* 41(1):78–80. <https://doi.org/10.1016/j.orl.2012.11.009>
- Wang Y, Sridhar D, Blei DM (2019) Equal opportunity and affirmative action via counterfactual predictions. Arxiv preprint <https://arxiv.org/abs/1905.10870>
- Woodworth B, Gunasekar S, Ohannessian MI, et al (2017) Learning non-discriminatory predictors. In: Kale S, Shamir O (eds) *Proceedings of the 2017 Conference on Learning Theory, Proceedings of Machine Learning Research*, vol 65. PMLR, Amsterdam, Netherlands, pp 1920–1953, URL <http://proceedings.mlr.press/v65/woodworth17a.html>
- Zafar MB, Valera I, Gomez Rodriguez M, et al (2017) Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In: *Proceedings of the 26th International Conference on World Wide Web. International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, WWW '17*, p 1171–1180, <https://doi.org/10.1145/3038912.3052660>
- Zemel R, Wu Y, Swersky K, et al (2013) Learning fair representations. In: Dasgupta S, McAllester D (eds) *Proceedings of the 30th International Conference on Machine Learning, Proceedings of Machine Learning*

Research, vol 28. PMLR, Atlanta, Georgia, USA, pp 325–333, URL <https://proceedings.mlr.press/v28/zemel13.html>

Zhang BH, Lemoine B, Mitchell M (2018) Mitigating unwanted biases with adversarial learning. In: Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society. ACM, pp 335–340, <https://doi.org/10.1145/3278721.3278779>

Zhang J, Bareinboim E (2018) Fairness in decision-making – the causal explanation formula. In: Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence. Association for the Advancement of Artificial Intelligence, pp 2037–2045, URL <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16949/15911>

Zhao H, Gordon G (2019) Inherent tradeoffs in learning fair representations. In: Wallach H, Larochelle H, Beygelzimer A, et al (eds) Advances in Neural Information Processing Systems, vol 32. Curran Associates, Inc., URL <https://proceedings.neurips.cc/paper/2019/file/b4189d9de0fb2b9cce090bd1a15e3420-Paper.pdf>