# FDA: GENERATING FAIR SYNTHETIC DATA WITH PROVABLE TRADE-OFF BETWEEN FAIRNESS AND FAITHFULNESS

Anonymous authors

Paper under double-blind review

#### ABSTRACT

We propose a novel framework called FDA for generating Fair synthetic data through Data Augmentation, offering the first method with provable trade-off guarantee between fairness and faithfulness. Unlike other existing methods, our approach utilizes a novel joint model that consists of two sub-models: one focused on enforcing strict fairness constraints while the other dedicated to preserving fidelity to the original data, coupled with a tuning mechanism that provides explicit control over the trade-off between fairness and faithfulness. Specifically, our FDA framework enables explicit quantification of the extent to which the generated fair synthetic data preserve faithfulness to the original data, while achieving an intermediate level of fairness determined by a user specified parameter  $\alpha \in [0, 1]$ . Theoretically, we show that the resulting fair synthetic data converge to the original data in probability when  $\alpha$  tends to 1, thereby implying convergence in distribution. Our framework can be also combined with some GAN-based fair models, such as DECAF, to further improve the utility of the resulting synthetic data in downstream analysis, while carefully balancing fairness. Furthermore, we obtain an upper bound of the unfairness measurement for downstream models trained on the generated fair synthetic data, which can help users to choose appropriate  $\alpha$ . Finally, we perform numerical experiments on benchmark data to validate our theoretical contributions and to compare our FDA with other methods.

033

006

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027

028

#### 1 INTRODUCTION

Artificial intelligence (AI) and machine learning (ML) algorithms have increasingly been used to
improve decision-making in almost all aspects of our lives (Zhao et al., 2018; Bogen & Rieke,
2018; Cohen et al., 2020; Mukerjee et al., 2002; Angwin et al., 2016; Berk et al., 2021). However,
there is mounting evidence showing that the developed algorithms may inherit biases and injustices
from historical data, leading to unfair decisions that discriminate against certain populations (Dastin,
2018; Datta et al., 2018; Lu et al., 2020; de Vassimon Manela et al., 2021). If not properly addressed,
biased or unfair decision-making may lead to violations of equality and anti-discrimination laws
(Krishnamurthy, 2021; Wachter et al., 2021). The emerging field of algorithmic fairness seeks to
address this urgent issue by mitigating the bias and discrimination in the AL and ML systems.

Broadly speaking, the bias mitigation methods can be categorized into three types: pre-processing, 043 in-processing, and post-processing. For a comprehensive overview of these methods, we encour-044 age readers to refer to recent review papers, such as (Pessach & Shmueli, 2022; Hort et al., 2022; 045 Mehrabi et al., 2021; Caton & Haas, 2024), and the extensive references cited within these works. 046 Pre-processing methods modify the biased training data, with the goal that any downstream model 047 trained on debiased data would achieve desired fairness requirements. In-processing methods mod-048 ify the algorithms by enforcing fairness constraints during training, with the goal that the trained algorithms achieve the desired fairness requirements on all real-life data. Post-processing methods modify the predictions based on a trained unfair model, with the goal that the final predictions sat-051 isfy certain fairness requirements. In recent years, as a pre-processing method, fair synthetic data generation has gained significant momentum (Feldman et al., 2015a; Zhang et al., 2017b; Calmon 052 et al., 2017; Xu et al., 2018; Zemel et al., 2013; Xu et al., 2019; van Breugel et al., 2021; Rajabi & Garibay, 2022). For example, Xu et al. (2019) proposed FairGAN, a GAN-based method to create

Table 1: Overview of related fair synthetic data generation methods summarized according to the following features: (1) supports trade-off between fairness and data utility; (2) allows continuous 056 labels; (3) allows categorical labels; (4) provides explicit quantification of loss of faithfulness to the original data, when meeting user-specified fairness requirement; (5) provides theoretical guarantee on the convergence of the generated fair synthetic data to the original data in probability and distribution; (6) provides theoretical analysis on the fairness of downstream models using debiased 060 synthetic data. 061

Model	Reference	(1)	(2)	(3)	(4)	(5)	(6)
FDA	THIS PAPER	-	<ul> <li>Image: A second s</li></ul>	<ul> <li>Image: A set of the set of the</li></ul>	<ul> <li>Image: A set of the set of the</li></ul>	-	<ul> <li>Image: A start of the start of</li></ul>
DECAF	VAN BREUGEL ET AL. (2021)	×	-	-	×	×	×
FAIRGAN	XU ET AL. (2018)	-	-	<ul> <li>Image: A second s</li></ul>	×	×	×
OPPDP	CALMON ET AL. (2017)	-	X	<b>√</b>	×	X	×
TABFAIRGAN	RAJABI & GARIBAY (2022)	-	X	<b>√</b>	×	X	×

067 068

071

073

synthetic data that satisfy group fairness; van Breugel et al. (2021) proposed DECAF, which trains 069 a graphical causal model using GANs and allows desired fairness constraints imposed via the associated causal graphs. Despite the successes of these earlier works, an important question remains unaddressed: To what extent does the fair synthetic data represent the statistical properties of the original data, preserving its utility for downstream analysis and modeling?

Generally, the task of generating a fair synthetic data that preserves the properties of the original data 074 creates a tension. On one hand, achieving fairness requires modification of the unfair data, which 075 may inadvertently impact the faithfulness of the synthetic data to the original data. On the other 076 hand, the synthetic data should faithfully represent the statistical properties of the original data, in 077 order to preserve its utility for downstream analysis and modeling. Therefore, achieving fairness 078 requires sacrificing some level of faithfulness, and vice versa. Due to the inherent competing nature 079 between these two goals, a trade-off between fairness and faithfulness is necessary characterized when generating fair synthetic data. In practice, striking the right balance involves careful consider-081 ation of the goals, stakeholders' priorities, and ethical implications in a given application.

082 **Contribution.** In this paper, we propose FDA, a Fair synthetic data generation framework through 083 Data Augmentation. FDA is built upon a joint modeling framework consisting of a fair model  $\mathcal{M}_{\text{fair}}$ 084 and a faithful model  $\mathcal{M}_{\text{faithful}}$ , coupled with a tuning mechanism to achieve a provable trade-off be-085 tween fairness and faithfulness in the generated synthetic data. This allows an explicit quantification of the extent to which the generated fair synthetic data preserve faithfulness to the original data, 087 while meeting specific fairness requirement controlled by  $\alpha \in [0, 1]$ , a user specified bias reduction parameter that quantifies the amount of biases removed from the original unfair data. Theoretically, setting  $\alpha = 0$ , the resulting synthetic data satisfy absolute fairness, with maximum reduction of the faithfulness to the original data. Conversely, setting  $\alpha = 1$ , the resulting synthetic data achieves per-090 fect similarity to the original data and is proved to converge to the original data in probability, further 091 implying convergence in distribution. When setting  $\alpha \in (0, 1)$ , users can achieve an intermediate 092 level of fairness while maintaining a certain level of faithfulness to the original data, both quantified by  $\alpha$ . Our framework can be combined with GAN-based fair models such as DECAF, to further 094 improve the data utility of the resulting fair synthetic data in downstream analysis, while achieving an intermediate level of fairness. In contrast to black-box methods that require time-consuming 096 training, our FDA framework generates synthetic data directly from the predictive distributions defined by our chosen joint model, which follows simple Gaussian distributions. Furthermore, to guide 098 users to choose appropriate unfairness reduction parameter  $\alpha$ , we provide theoretical analysis on the 099 fairness of downstream models trained on the generated fair synthetic data. As far as we know, our FDA is the first method to provide all these desired features, with a comparison of our FDA with 100 other methods summarized in Table 1. 101

102 **Notations.** For any positive integer K, let  $[K] = \{1, \dots, K\}$ . For any  $p \ge 1$ , let  $\mathcal{M}_p(\mathbb{R})$  be the 103 space of all probability measures on  $\mathbb{R}$  with finite p-th moment. For any two random variables U 104 and V,  $\mu_{U|V}$  denotes the conditional distribution of U given V;  $\mu_U$  and  $\mu_V$  denote their marginal 105 distributions respectively. We write  $U \stackrel{d}{=} V$  when U and V are equal in distribution. For a random 106 sequence  $\{U_i\}_{i=1}^{\infty}$ , we write  $U_n \xrightarrow{p} U(U_n \xrightarrow{d} U)$  when  $U_n$  converges in probability (in distribution) to U as  $n \to \infty$ . We use  $\Delta^{K-1}$  as the probability simplex in  $\mathbb{R}^K$ . 107

# <sup>108</sup> 2 PRELIMINARIES

110 A sequence of triplets  $\mathcal{D} = \{Y_i, X_i, S_i\}_{i=1}^n$  is observed, where for each  $i, Y_i \in \mathbb{R}$  denotes the 111 outcome,  $S_i \in [K]$  denotes the sensitive attribute, and  $X_i \in \mathbb{R}^d$  denotes other attributes. Assuming 112 that  $\mathcal{D}$  is sampled from the distribution  $\mathcal{P}_{\mathcal{D}}$ , which violates certain fairness requirements, rendering 113  $\mathcal{D}$  unfair, our objective is to generate fair synthetic data denoted as  $\hat{\mathcal{D}}$  based on  $\mathcal{D}$ . Specifically, 114 we want to ensure that  $\hat{\mathcal{D}}$  satisfies  $\alpha$ -reduction of unfairness (given in Theorem 3.5), for any fixed 115  $\alpha \in [0, 1]$ .

As discussed in van Breugel et al. (2021), predictive fairness measures such as equalized odds are not compatible in the context of fair data, as the aim is to ensure the fairness in the synthetic data distribution, rather than achieving fair algorithmic predictions. Consequently, we follow van Breugel et al. (2021) to focus on *Demographic Parity* (DP) and formally extend it to the context of fair data.

**Definition 2.1** (Demographic parity). The distribution  $\mathcal{P}_{\mathcal{D}}$ , from which  $\mathcal{D}$  is sampled, is said to satisfy *demographic parity (DP)*, if it satisfies  $(Y|S = s_1) \stackrel{d}{=} (Y|S = s_2)$ , for any  $s_1, s_2 \in [K]$ . In other words, for any  $T \subseteq \mathbb{R}$ ,  $\mathbb{P}(Y \in T|S = s_1) = \mathbb{P}(Y \in T|S = s_2)$ .

124 Then, any discrepancy between the distribution of  $Y|S = s_1$  and that of  $Y|S = s_2$ , for any 125  $s_1, s_2 \in [K]$  indicates violation of DP. Note, our proposed framework can be applied to the condi-126 tional fairness notion (Barocas et al., 2023) (see Remark 3.2). As discussed in Chzhen & Schreuder 127 (2022), various distance measures, e.g., Wasserstein distance, total variation and Kolmogorov-128 Smirnov distance, have been used to evaluate this discrepancy empirically and thereby quantify 129 the violation of DP. In this paper, we follow Chzhen & Schreuder (2022) and others (Gouic et al., 130 2020; Chzhen et al., 2020; Jiang et al., 2020; Gaucher et al., 2022; Xian et al., 2022) to use Wasser-131 stein distance due to its effectiveness to explicitly quantify unfairness measurement, faithfulness to the original data as well as data utility in downstream models using the same unit measurements 132 comparing to other distance measures. Specifically, we define the unfairness measure of the original 133 distribution of  $\mathcal{D}$  as the sum of the weighted distances between  $\{\mu_{Y|S=s}\}_{s\in[K]}$  and their com-134 mon barycenter (Villani, 2021; Santambrogio, 2015), w.r.t. the Wasserstein-2 distance<sup>1</sup>, denoted by 135  $W_2(\mu_{Y|S=s}, \cdot)$ , as defined below. 136

**Definition 2.2** (Unfairness measure). We define the *unfairness* of the distribution  $\mathcal{P}_{\mathcal{D}}$ , from which the dataset  $\mathcal{D}$  is sampled, as follows<sup>2</sup>

139 140 141

152

153

$$\mathcal{UF}(\mathcal{P}_{\mathcal{D}}) := \min_{\nu \in \mathcal{M}_2(\mathbb{R})} \sum_{s=1}^K \omega_s W_2(\mu_{Y|S=s}, \nu), \qquad (1)$$

142 for any given weights<sup>3</sup>  $(\omega_1, \cdots, \omega_K) \in \Delta^{K-1}$ .

144 It is easy to see that  $\mathcal{UF}(\mathcal{P}_{\mathcal{D}}) = 0$  if and only if there is a minimizer  $\nu$  in equation 1 such that 145  $\mu_{Y|S=s} = \nu$  for all  $s \in [K]$ , that is, it satisfies the DP constraint:  $(Y|S = s_1) \stackrel{d}{=} (Y|S = s_2)$  for 147 all  $s_1, s_2 \in [K]$ . Conversely, a larger value of this measure<sup>4</sup> indicates a more severe violation of DP 148 constraint.

**Problem statement.** For any biased dataset  $\mathcal{D}$ , and a user-specified *bias reduction factor*  $\alpha \in [0, 1]$ , we substitute the observed  $Y_i$  values with their synthetic counterparts to produce fair synthetic data, denoted as  $\hat{\mathcal{D}} = {\hat{Y}_i, X_i, S_i}_{i=1}^n$ . Here  $\hat{\mathcal{D}}$  satisfies the following bias reduction guarantee:

•  $\mathcal{UF}(\mathcal{P}_{\hat{\mathcal{D}}}) = \alpha \mathcal{UF}(\mathcal{P}_{\mathcal{D}})$  (Theorem 3.5). where  $\mathcal{P}_{\hat{\mathcal{D}}}$  denotes the distribution of  $\hat{\mathcal{D}}$ .

In the meanwhile, we can assess the loss of faithfulness between the synthetic data  $\hat{D}$  and the original data D by calculating Wasserstein-2 distance between  $\mu_{\hat{Y}}$  and  $\mu_Y$ , which is a closed-form function of  $\alpha$  (Theorem 3.6). These findings enable users to choose an appropriate  $\alpha$  by considering the explicit trade-off between fairness and faithfulness.

<sup>&</sup>lt;sup>1</sup>See Section A.1 for a formal definition of Wasserstein-p distance.

<sup>&</sup>lt;sup>2</sup>One may use  $W_2^2(\mu_{Y|S=s},\nu)$  when defining  $\mathcal{UF}(\mathcal{P}_{\mathcal{D}})$ ; then  $\alpha$  needs to be replaced by  $\alpha^2$ .

<sup>&</sup>lt;sup>3</sup>See Appendix A.3 for a discussion on how to select the weights.

<sup>&</sup>lt;sup>4</sup>See Section A.2 for a discussion on how to evaluate this unfairness measure empirically.

Importantly, under our FDA framework, we can generate high quality synthetic data with a theoretical guarantee that as  $\alpha$  approaches 1, the synthetic  $\hat{Y}$  converges to the original Y in probability (and consequently in distribution) conditional on the features X, S.



Figure 1: A graphical representation of our FDA synthetic data generation framework.

## 178 3 The proposed FDA method

166

167

169 170 171

172

173 174 175

176 177

179

204

205

**Overview.** The main idea is to simulate synthetic data from the predictive distribution defined under a joint model, denoted as  $\mathcal{M}_{FDA}$ , using both observed  $\mathcal{D}$  and additionally augmented data. A similar approach was previously considered in Jiang et al. (2022) to generate privacy preserving synthetic data. In this work, we build upon their method to address the challenge of generating fair synthetic data, providing provable theoretical guarantees on the trade-off between fairness and utility, as well as the convergence of the synthetic data to the original data in both probability and distribution.

186 Specifically,  $\mathcal{M}_{\text{FDA}}$  consists of two sub-models: (i) a fair model (see equation 2), denoted as  $\mathcal{M}_{\text{fair}}$ , 187 which specifies certain relationship between  $Y_i$  and the feature vector  $X_i$  and the sensitive attribute 188  $S_i$ , such that it imposes exact fairness constraint; and (ii) a faithful model,  $\mathcal{M}_{\text{faithful}}$ , which generates 189  $Z_i$  (see equation 3 for details) given  $Y_i$ , for each  $i \in [n]$ , such that  $Z_i$  are noisy copies of  $Y_i$  with 190 accuracy level controlled by tuning parameters. The fair synthetic data are then generated as samples 191 from the corresponding predictive distributions that are defined by the model  $\mathcal{M}_{\text{FDA}}$ .

192 By design, these  $Z_i$  contain information about  $Y_i$  so that  $\mathcal{M}_{\text{faithful}}$  plays the role of enforcing the resulting synthetic data to be close (and thus faithful) to the original data. In contrast,  $\mathcal{M}_{\text{fair}}$  plays 193 the role of imposing the desired fairness requirement, e.g., the DP constraint. Under such a frame-194 work, both models  $\mathcal{M}_{\text{faithful}}$  and  $\mathcal{M}_{\text{fair}}$  influence the resulting synthetic data, with their respective 195 contributions determined by the values of the tuning parameters introduced in  $\mathcal{M}_{faithful}$  (can be seen 196 in equation 3 and discussed thereafter). As a result, the tuning parameters control the relative influence from  $\mathcal{M}_{\text{faithful}}$  and  $\mathcal{M}_{\text{fair}}$ , and thus balance between the two competing goals, fairness and 198 faithfulness to the original data. Given the proposed framework is general, the data considered can 199 be any type and need not to be limited to be binary or discrete as considered by many others. Next, 200 we present details of  $\mathcal{M}_{\text{faithful}}$  and  $\mathcal{M}_{\text{fair}}$  when  $Y_i$ 's are continuous; the discussion on how to use the 201 proposed framework when  $Y_i$ 's are categorical is given in Remark 3.4. 202

**Fair model.** The fair model  $\mathcal{M}_{\text{fair}}$  can be written in the following general form:

$$\mathcal{M}_{\text{fair}}: Y_i = f(X_i, S_i, \beta) + \varepsilon_i, \ \varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \eta^2),$$

(2)

where f can be any fair predictor, under which  $Y_i$  satisfies the DP constraint, and  $\beta$  is the associated model parameter. Notably, the simplest choice of  $\mathcal{M}_{\text{fair}}$  is a constant mean model (CMM) with  $f = \beta_0$  being a constant function. Alternatively, one could also choose f to be a GAN-based fair predictor, e.g., DECAF (discussed in Remark 3.3).

However, since many dependencies among the variables in the original data  $\mathcal{D}$  are intentionally omitted in  $\mathcal{M}_{\text{fair}}$  to fulfill the fairness requirement, loss of information and utility is inevitable. In other words, generating synthetic data solely from  $\mathcal{M}_{\text{fair}}$  could result in considerable loss of faithfulness to the original data. As pointed out by many others (Feldman et al., 2015a; Xu et al., 2018; Chzhen & Schreuder, 2022; Zhao & Gordon, 2022; Tran et al., 2022), a sustainable solution would be to accept some degree of compromise on fairness in order to preserve the utility of the original data in the downstream analysis. We achieve this goal by introducing a faithful model as a submodel in our joint modeling framework, which plays the role of mitigating loss of information
 about the original data.

**Faithful model.** The faithful model  $\mathcal{M}_{\text{faithful}}$  takes the following form:

 $\hat{Y}_i \sim p\left(Y_i | Z_i, X_i, S_i; \hat{\beta}, \hat{\eta}^2, \hat{\sigma}^2\right)$ 

$$\mathcal{M}_{\text{faithful}}: Z_{im} = Y_i + e_{im}, \ e_{im} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2), \tag{3}$$

for  $m \in [M]$ . The specification of this model allows us to generate noisy copies of  $Y_i$ , i.e.,  $Z_i = (Z_{i1}, \dots, Z_{iM})^{\top}$ , so that  $Z_i$  contains information about  $Y_i$  with their faithfulness to the original  $Y_i$  controlled by the tuning parameters  $\sigma^2$  and M. By increasing the number of copies, M and/or decreasing the additive noise variance  $\sigma^2$ , the generated  $Z_i$  contains more information about  $Y_i$ . In practice, one can fix the value of M and adjust the value of  $\sigma^2$ , or vice versa. As discussed next in Remark 3.1, the ratio  $\sigma^2/M$  determines both the levels of fairness in the generated synthetic data as well as its faithfulness to the original data.

With  $\mathcal{M}_{\text{fair}}$  and  $\mathcal{M}_{\text{faithful}}$  chosen, we create our augmented dataset  $\mathcal{D}_{\text{DA}} = \{\mathcal{D}, \{Z_i\}_{i=1}^n\}$  given the simulated  $Z_i$ 's, and fit our joint model  $\mathcal{M}_{\text{FDA}}$  given below,

$$\underbrace{p(Y,Z \mid X,S;\beta,\eta^2,\sigma^2)}_{\mathcal{M}_{\text{FDA}}} = \underbrace{p(Y \mid X,S;\beta,\eta^2)}_{\mathcal{M}_{\text{fair}}} \underbrace{p(Z \mid Y;\sigma^2)}_{\mathcal{M}_{\text{faithful}}}.$$
(4)

Note, under the model  $\mathcal{M}_{\text{faithful}}$ , where Z depends solely on Y, the conditional independence of Z on X and S given Y implies that  $p(Z \mid Y, X, S) = p(Z \mid Y)$ .

Followed by the joint model in equation 4, the synthetic values of  $Y_i$ , denoted as  $\hat{Y}_i$  for  $i \in [n]$ , are then drawn from the corresponding predictive distribution as follows:

236

237

231 232 233

219 220 221

240 241

242 243  $\propto \underbrace{p(Y_i, Z_i \mid X_i, S_i; \hat{\beta}, \hat{\eta}^2, \hat{\sigma}^2)}_{\mathcal{M}_{\text{FDA}}} = \underbrace{p(Y_i \mid X_i, S_i; \hat{\beta}, \hat{\eta}^2)}_{\mathcal{M}_{\text{fair}}} \underbrace{p(Z_i \mid Y_i; \hat{\sigma}^2)}_{\mathcal{M}_{\text{faithful}}},$ 

where  $\hat{\beta}, \hat{\eta}^2, \hat{\sigma}^2$  are the estimates of model parameters  $\beta, \eta^2, \sigma^2$ , respectively. Figure 2 presents a graphical representation of our data augmented joint modeling framework and fair synthetic data generation process.

247 With the intentional choice of Gaussian models for both  $\mathcal{M}_{fair}$  and  $\mathcal{M}_{faithful}$ , it can be easily shown 248 that the predictive distribution in equation 5 corresponds to the following Gaussian distribution (the 249 proof is given in Appendix B.1):

250 251

253

 $\mathcal{N}\left(\frac{\frac{\hat{\sigma}^2}{M}f(X_i, S_i, \hat{\beta}) + \frac{\sum_{j=1}^M Z_{ij}}{M}\hat{\eta}^2}{\frac{\hat{\sigma}^2}{M} + \hat{\eta}^2}, \frac{\frac{\hat{\sigma}^2}{M}\hat{\eta}^2}{\frac{\hat{\sigma}^2}{M} + \hat{\eta}^2}\right).$ (6)

(5)

In summary, given the estimates of the model parameters for both  $\mathcal{M}_{\text{fair}}$  and  $\mathcal{M}_{\text{faithful}}$ , i.e.,  $\hat{\beta}, \hat{\eta}^2, \hat{\sigma}^2$ , our synthetic dataset  $\hat{\mathcal{D}} = {\hat{Y}_i, X_i, S_i}_{i=1}^n$  can be conveniently obtained by generating  $\hat{Y}_i$ , for  $i \in [n]$ from the corresponding predictive distribution given in equation 6. Algorithm 1 summarizes the key steps using our FDA framework.

258 As shown in equation 6, it is clear that both  $\mathcal{M}_{\text{fair}}$  and  $\mathcal{M}_{\text{faithful}}$  contribute information to the synthetic values of  $Y_i$ 's. At one extreme, when M = 0 or/and  $\sigma^2 = \infty$ , the Gaussian distribution in equation 6 259 260 reduces to  $\mathcal{N}\left(f(X_i, S_i, \hat{\beta}), \hat{\eta}^2\right)$ ; that is, the synthetic data  $\hat{Y}_i$  will be generated based on the fair 261 262 model  $\mathcal{M}_{\text{fair}}$  alone, resulting in synthetic data  $\mathcal{D}$  that satisfies exact DP constraint. At the other 263 extreme, when  $M = \infty$  or/and  $\sigma^2 = 0$ , the Gaussian distribution in equation 6 degenerates to a point mass at  $Y_i$  with  $\overline{Z}_i = \frac{1}{M} \sum_{l=1}^M Z_{il} \stackrel{d}{=} Y_i$ ; that is, the information in  $\mathcal{M}_{\text{faithful}}$  will completely 264 265 override that of  $\mathcal{M}_{\text{fair}}$ , resulting in  $\hat{Y}_i = Y_i$ . Thus, when  $M, \sigma^2 \in (0, \infty)$ , an intermediate level of 266 fairness and faithfulness to the original data will be achieved. 267 Remark 3.1 (Choosing tuning parameters in practice). To enhance the practical utility of our FDA

*Remark* 3.1 (Choosing tuning parameters in practice). To enhance the practical utility of our FDA framework, we further introduce a bias reduction factor, denoted as  $\alpha = \frac{\eta^2}{\eta^2 + \lambda}$  with  $\lambda = \sigma^2/M$ . This allows us to express the unfairness measure in the synthetic dataset,  $\mathcal{UF}(\mathcal{P}_{\hat{\mathcal{D}}})$ , in closed form

Algorit	hm 1 FDA fair synthetic data generation algorithm
Input:	original dataset $\mathcal{D} = \{X_i, Y_i, S_i\}_{i=1}^n$ and user specified $\lambda$ .
1. Gene	rate noisy copies $\{Z_i\}_{i=1}^n$ of $Y_i$ under the chosen faithful model $\mathcal{M}_{\text{faithful}}$ .
2. Given	the augmented data $\mathcal{D}_{DA} = \{\mathcal{D}, \{Z_i\}_{i=1}^n\},\$
( <b>2a).</b> Fi	t the joint model defined in equation 4 and obtain the parameter estimates $\hat{\eta}^2$ , $\hat{\beta}$ and $\hat{\sigma}^2$ ;
(2b). Sa	mple random draws $\hat{Y}_i$ from (6) for each $i \in [n]$ .
3. Outp	ut fair synthetic datasets $\hat{\mathcal{D}} = \{\hat{Y}_i, X_i, S_i\}_{i=1}^n$ .

as a function of this bias reduction factor. At the same time, we can quantify the Wasserstein-2 distance between the distribution of  $Y_i$ 's in the original dataset  $\mathcal{D}$  and that of  $\hat{Y}_i$ 's in the synthetic dataset  $\hat{\mathcal{D}}$  as a function of  $\alpha$ . By balancing the reduction of unfairness with the preservation of faithfulness to the original data, users can select an appropriate value of  $\alpha$ , which in turns can determine the values of the tuning parameters M and  $\sigma^2$  in the faithful model  $\mathcal{M}_{\text{faithful}}$  given the relationship  $\lambda = \sigma^2/M = \frac{(1-\alpha)\eta^2}{\eta^2}$ . Further details are provided in Theorem 3.5, Theorem 3.6, and the subsequent discussions following these theorems.

287 *Remark* 3.2 (FDA beyond DP fairness). The FDA framework can be applied to generate fair syn-288 thetic data to satisfy the conditional statistical fairness (see Appendix A.4). In this case, one simply 289 needs to choose  $\mathcal{M}_{\text{fair}}$  such that it satisfies the conditional fairness criterion. For example, the con-290 stant mean model for  $\mathcal{M}_{\text{fair}}$  mentioned earlier (i.e., CMM) naturally satisfies this fairness notion.

291 *Remark* 3.3 (FDA when combined with GAN-based fair models). We can let  $\mathcal{M}_{\text{fair}}$  be a GAN-based 292 fair data generator, e.g., DECAF. In this case, the generated synthetic data are guaranteed to achieve 293 a higher level of fidelity to the original data compared to using DECAF alone (see Theorem 3.7).

**Remark 3.4** (FDA framework when  $Y_i$ 's are categorical). In the case of categorical labels, a sample drawn from equation 6 is continuous. Then we need to map these continuous values back to discrete categories. One option is to round the continuous value to the nearest integer, or alternatively, one could define specific ranges for each category and map the continuous value to a category based on these predefined ranges. For example, when  $Y \in \{0, 1\}$ , one can use a threshold value 0.5 as a cutoff to distinguish the two categories.

300 301

#### 3.1 THEORETICAL ANALYSIS OF THE RESULTING SYNTHETIC DATA

Quantifying the fairness and faithfulness of synthetic data relative to the original data has been a longstanding challenge in previous fairness studies. While some prior works provide analytic bounds to characterize these properties, our work is the first to offer a closed-form solution, enabling an exact quantification of the trade-off between fairness and faithfulness. In this section, we proide a detailed discussion on how the synthetic data  $\hat{D}$  generated by our FDA framework achieves an  $\alpha$ -reduction in unfairness, where  $\alpha$  is the bias reduction factor introduced in Remark 3.1, and how its faithfulness to the original data can be explicitly expressed as a function of  $\alpha$ .

Theorem 3.5 (Unfairness reduction guarantee). With the bias reduction factor  $\alpha = \frac{\eta^2}{\eta^2 + \lambda}$ , where  $\eta^2$ is determined by  $\mathcal{M}_{fair}$  and  $\lambda = \sigma^2/M$ , determined by the tuning parameters  $\sigma^2$  and M in  $\mathcal{M}_{faithful}$ , the distribution of the synthetic data  $\hat{\mathcal{D}}$  generated by our FDA framework achieves the  $\alpha$ -reduction of unfairness as follows,

314

$$\mathcal{UF}(\mathcal{P}_{\hat{\mathcal{D}}}) = \alpha \mathcal{UF}(\mathcal{P}_{\mathcal{D}}) \,. \tag{7}$$

where  $\mathcal{P}_{\hat{D}}$  and  $\mathcal{P}_{\mathcal{D}}$  represent the distributions, from which the synthetic data  $\hat{\mathcal{D}}$  and the original data  $\mathcal{D}$  are sampled, respectively.

By selecting the tuning parameters  $\sigma^2$  and M to achieve the desired value of  $\alpha$ , Theorem 3.5 guarantees the  $\alpha$ -reduction of unfairness on  $\mathcal{P}_{\hat{\mathcal{D}}}$  for the resulting synthetic data  $\hat{\mathcal{D}}$  (see proof in Appendix B.4). In practice, when  $\eta$  is unknown, an estimate of  $\eta$ , denoted as  $\hat{\eta}$ , can be obtained by fitting the assumed  $\mathcal{M}_{\text{fair}}$  first. Then, the values of  $\sigma^2$  and M can be chosen such that their ratio  $\lambda = \hat{\eta}^2 (1 - \alpha) / \alpha$ . Meanwhile, the following theorem quantifies the Wasserstein-2 distance between the conditional distribution of  $\hat{Y}_i$ 's in  $\hat{\mathcal{D}}$  and that of the original  $Y_i$ 's in  $\mathcal{D}$ , when the synthetic data meets the desired fairness requirement characterized by any user-specified  $\alpha$ .

**Theorem 3.6** (Faithfulness quantification). For the synthetic data  $\hat{D}$  such that its distribution  $\mathcal{P}_{\hat{D}}$ 325 satisfies the  $\alpha$ -reduction of unfairness, the Wasserstein-2 distance between  $\mu_{\hat{Y}|X,S}$  and  $\mu_{Y|X,S}$  is given by

$$W_2(\mu_{\hat{Y}|X,S},\mu_{Y|X,S}) = \sqrt{(1-\alpha)^2 \mathbb{E}\left[\left(Y - f(X,S,\beta)\right)^2 | X,S\right] + (1-\alpha^2) \eta^2} \,. \tag{8}$$

where  $\mu_{\hat{Y}|X,S}$  and  $\mu_{Y|X,S}$  denote the conditional distribution of  $\hat{Y}$ 's given X,S in  $\hat{\mathcal{D}}$  and that of  $Y_i$ 's in  $\mathcal{D}$ , respectively. Particularly, when  $M \to \infty$ , and/or  $\sigma^2 \to 0$ , we have 332

$$\hat{Y}|X, S \xrightarrow{p} Y|X, S, \text{ and consequently, } \hat{Y}|X, S \xrightarrow{d} Y|X, S,$$
(9)

for any choice of  $\mathcal{M}_{fair}$ .

324

326

327 328

330

331

333 334 335

336 337

338

339 340 341

342 343 344

362

364

365 366

367

Theorem 3.6 (see proof in Appendix B.2) quantifies the closeness of the generated synthetic valuse  $\hat{Y}_i$ 's to the original values  $Y_i$ 's with respect to the user-specified unfairness reduction factor  $\alpha$ . When  $\lambda = 0$  (i.e.,  $M = \infty$  or  $\sigma^2 = 0$ ),  $\alpha = 1$  and  $\hat{Y}$  is identical to Y; that is,  $\hat{\mathcal{D}} = \mathcal{D}$ and  $\mathcal{P}_{\hat{\mathcal{D}}}$  is as unfair as  $\mathcal{P}_{\mathcal{D}}$ . Conversely, when  $\lambda = \infty$  (i.e., M = 0 and/or  $\sigma^2 = \infty$ ),  $\alpha = 0$ and the maximal discrepancy between  $\hat{Y}$  and Y is achieved, i.e., equation 8 attains it maximum

$$\sqrt{\mathbb{E}\left[\left(Y - f(X, S, \beta)\right)^2 | X, S\right]} + \eta^2$$
, resulting in exact DP for  $\mathcal{P}_{\hat{\mathcal{D}}}$  at the cost of faithfulness

Theorem 3.6 and Theorem 3.5 together quantify the compromise one must make in order to meet 345 the desired fairness requirement with respect to  $\alpha$ . Given a specific choice of  $\mathcal{M}_{\text{fair}}$ , a larger value 346 of  $\alpha$  results in fairer synthetic data  $\mathcal{D}$  but a greater discrepancy between the distributions of  $Y_i$ 's 347 and  $Y_i$ 's conditional on features. In practice, one can select a suitable  $\alpha$  by balancing between the 348 goals of reducing the unfairness and enhancing faithfulness. Importantly, achieving this goal can 349 be facilitated by choosing appropriate values of the tuning parameters  $\sigma^2$  and M in  $\mathcal{M}_{\text{faithful}}$ . It 350 is worth noting that the faithfulness in  $\hat{D}$  does depend on the predictor f in the fair model  $\mathcal{M}_{\text{fair.}}$ 351 Ideally, if one can find a fair predictor f in  $\mathcal{M}_{\text{fair}}$  that satisfies the DP constraint with minimal loss 352 of faithfulness, then the generated synthetic data obtained using our FDA framework achieve the 353 highest level of faithfulness to the original data while simultaneously meeting the desired fairness 354 requirement with respect to  $\alpha$ .

355 Next, we demonstrate that the introduction of our faithful model  $\mathcal{M}_{faithful}$  within our joint FDA 356 framework ensures that the generated synthetic data are guaranteed to be more faithful to the original 357 data compared to using the fair model  $\mathcal{M}_{\text{fair}}$  alone. 358

**Theorem 3.7** (Faithfulness improvement guarantee). For the synthetic data  $\hat{\mathcal{D}} := \{\hat{Y}_i, X_i, S_i\}_{i=1}^n$ 359 generated from the fair model  $\mathcal{M}_{fair}$  alone (hence satisfying the exact DP constraint), the 360 *Wasserstein-2 distance between*  $\mu_{\tilde{Y}|X,S}$  *and*  $\mu_{Y|X,S}$  *is given by,* 361

$$W_{2}(\mu_{\tilde{Y}|X,S},\mu_{Y|X,S}) = \sqrt{\mathbb{E}\left[\left(Y - f(X,S,\beta)\right)^{2} | X,S\right] + \eta^{2}}.$$
 (10)

*Comparing equation 8 and equation 10, it is easy to see that,* 

$$W_2(\mu_{\hat{Y}|X,S}, \mu_{Y|X,S}) < W_2(\mu_{\tilde{Y}|X,S}, \mu_{Y|X,S})$$
(11)

368 for any fixed M > 0 and  $\sigma^2 > 0$ , where  $\mu_{\hat{Y}|X,S}$  denotes the conditional distribution of the synthetic 369  $\hat{Y}_i$ 's generated by our joint FAD framework using both  $\mathcal{M}_{fair}$  and  $\mathcal{M}_{faithful}$ . 370

Theorem 3.7 (see proof in Appendix B.3) proves that the synthetic values  $\hat{Y}_i$ 's obtained using FDA 371 372 are closer to the original values  $Y_i$ 's than the synthetic values  $\tilde{Y}_i$ 's obtained from using  $\mathcal{M}_{\text{fair}}$  alone. 373 This finding provides provable guarantee that our joint FDA framework can be combined with any 374 existing fair data generation model to further improve the faithfulness of the generated synthetic 375 data. As we have discussed in the introduction, the generated synthetic dataset  $\mathcal{D}$  is used as the training dataset to train downstream models. For any trained downstream model g(X, S) to predict 376 Y, the model error could affect the downstream fairness. Theoretically, the upper bound of the 377 fairness violation of the downstream model is given by the following proposition.

**Proposition 3.8.** For any downstream model g(X, S) to predict Y, if  $W_2(\mu_{g(X,S)|s}, \mu_{\hat{Y}|s}) \leq \delta$  for all  $s \in [K]$ , we have

381

382

384 385 386

387

388

389

390

 $\min_{\nu \in \mathcal{P}_p(\mathbb{R})} \sum_{s=1}^K \omega_s W_2(\mu_{g(X,S)|s}, \nu) \le \alpha \mathcal{UF}(\mathcal{P}_{\mathcal{D}}) + \delta, \qquad (12)$ 

for any given weights  $(\omega_1, \cdots, \omega_K) \in \Delta^{K-1}$ .

Proposition 3.8 (see proof in Appendix B.5) establishes a uniform upper bound on fairness violations in the downstream models trained on the generated synthetic dataset. This upper bound can help users to select an appropriate  $\alpha$  to ensure the desired level of downstream model fairness. Additionally, Proposition 3.8 suggests that the downstream model error, which is captured by  $\delta$ , can negatively affect the downstream model's fairness: a smaller error in the downstream model corresponds to improved fairness guarantees in the downstream models.

391 392 393

394

#### 4 EXPERIMENTS

We demonstrate the novel features of our FDA framework in generating fair synthetic data and compare them with the baseline methods listed in Table 1 based on real data experiments. Further, we show that how our FDA framework can be combined with DECAF, a GAN-based fair data generator to improve the utility in downstream models trained on the resulting fair synthetic data.

399 All the baseline methods in Table 1 require intensive training and are extremely sensitive to the 400 architecture and hyperparameters of the models. For example, different constructed causal graphs 401 used by DECAF lead to varying utility in the resulting synthetic data. Therefore, to ensure a fair 402 comparison, we run experiments on the UCI Adult dataset (Dua & Graff, 2020), which is the only 403 dataset used by all the baseline methods. We also use the same model specifications as in the 404 original code provided by the authors to maintain consistency in our comparison. The UCI Adult dataset contains over 65,000 samples with 11 attributes. It is known to exhibit bias between gender 405 and income (Feldman et al., 2015b; Zhang et al., 2017a). Thus, we treat gender as the sensitive 406 attribute and *income* as the binary target variable representing whether an adult's income is over 407 \$50K or not. Additional important details regarding the experimental setup, including the dataset 408 split, the architecture of the downstream model, and its training process, are provided in Appendix 409 C. Note that, we also conducted experiments on the COMPAS dataset using our FDA method. Due 410 to space limitations, these results are included in Appendix D. 411

- 412
- 413 414

# 4.1 ACHIEVING TRADE-OFF BETWEEN FAIRNESS AND FAITHFULNESS IN DEBIASED SYNTHETIC DATA BY FDA

415 In this section, we show how our FDA framework facilitates the trade-off between absolute fair-416 ness ( $\alpha = 0$ ) and perfect data faithfulness ( $\alpha = 1$ ) by varying the bias reduction factor  $\alpha$  within 417 (0, 1). Specifically, synthetic Adult datasets are generated using the baseline methods and our FDA 418 under different values of  $\alpha$ . We repeat the experiments 10 times for each method. Figure 2 shows 419 (1) the estimated Wasserstein-2 distance between the synthetic and original data distributions, de-420 noted by  $\hat{W}_2(\mu_{\hat{Y}}, \mu_Y)$ , and (2) the estimated unfairness measure in the synthetic dataset, denoted 421 by  $\mathcal{UF}(\mathcal{P}_{\mathcal{D}})$ , obtained by each method, where the solid line represents the average of the 10 exper-422 iments and shadowed areas indicates variation<sup>5</sup>. As expected, when  $\alpha \to 0$ ,  $\widehat{\mathcal{UF}}(\mathcal{P}_{\hat{\mathcal{D}}}) \to 0$ . This 423 is when synthetic data achieves perfect fairness but the worst faithfulness to the original data, with 424  $\hat{W}_2(\mu_{\hat{Y}},\mu_Y)$  far away from 0. Conversely, when  $\alpha \to 1$ ,  $\hat{W}_2(\mu_{\hat{Y}},\mu_Y) \to 0$ . This is when the 425 synthetic data distribution converges to the original data distribution, achieving perfect data faithful-426 ness; in the meanwhile, the synthetic data achieves zero reduction of biases compared to the original 427 data. None of the baseline methods offer this tuning mechanism, highlighting a unique advantage of 428 our approach. It is also worth noting that, unlike other methods, our FDA shows very small variation 429 that it is almost invisible on the plots, showing the stability of our FDA method.

<sup>&</sup>lt;sup>5</sup>We observe similar patterns with the total variation-based unfairness measure detailed in the Appendix C.2, experiments on COMPAS dataset is given in Appendix D.



Figure 2: Faithfulness and fairness of the synthetic datasets by FairGAN, OPPDP, TabFairGAN, DECAF and FDA under varying values of  $\alpha$ :  $\hat{W}_2(\mu_{\hat{Y}}, \mu_Y)$  ( $\downarrow$  more faithful),  $\widehat{\mathcal{UF}}(\mathcal{P}_{\hat{\mathcal{D}}})$  ( $\downarrow$  more fair). The shadowed areas along each line represent the variations on 10 repetitions of experiments.

4.2 Improving data faithfulness and downstream utility using FDA

In this section, we conduct experiments on the Adult dataset to demonstrate how our FDA framework
 when combined with DECAF, can enhance the faithfulness of synthetic data to the original data
 and thereby enhances the performance of downstream models trained on the synthetic data while
 achieving an intermediate level of fairness.

454 Specifically, in our FDA framework, we let  $\mathcal{M}_{\text{fair}}$  be the GAN-based model DECAF (van Breugel 455 et al., 2021) (we term the corresponding model FDA-DECAF). To emphasize that, in our original FDA framework,  $\mathcal{M}_{fair}$  is a constant mean model, namely CMM, we term the corresponding 456 model FDA-CMM. We generate debiased synthetic data using DECAF, FDA-DECAF and FDA-457 CMM under different values of  $\alpha$  and afterwards evaluate the utility and fairness of a downstream 458 model trained on the resulting debiased synthetic datasets using these methods. Again, to ensure 459 fair comparisons, we follow van Breugel et al. (2021) and focus on the same downstream MLP 460 model. When evaluating the utility and fairness of the downstream model trained on the debiased 461 synthetic data, we focus on the same metrics: (1) utility: we evaluate the predictive performance 462 of the model using accuracy, precision, recall, and AUROC; (2) fairness: we assess the fairness of 463 the downstream model trained on the synthetic data using the total variation distance based measure 464  $|\mathbb{P}(\hat{Y}|S=1) - \mathbb{P}(\hat{Y}|S=0)|.$ 465

Figure 3 shows the results<sup>6</sup> of our experiments, repeated 10 times. As expected, FDA-DECAF leads 466 to consistently better utility in the downstream prediction for any  $\alpha \in (0,1)$  when compared with 467 DECAF. This finding is supported by our theoretical result given in Theorem 3.7. By introducing 468 a joint modeling approach (including  $\mathcal{M}_{\text{fair}}$  and  $\mathcal{M}_{\text{faithful}}$ ), coupled with the tuning mechanism to 469 allow for an intermediate level of fairness, our FDA framework enables the resulting synthetic data 470 to maintain a certain level of faithfulness to the original data, thereby enhancing the downstream 471 prediction performance. When  $\alpha \to 1$ , our FDA framework allows the utility of the downstream 472 model trained on the resulting synthetic data to fully recover the utility using the original data. In 473 contrast, DECAF enforces exact fairness in the synthetic data, which can lead to significant loss of 474 utility in the downstream model as shown in Figure 3.

It is also worth noting that, despite the greater efforts required to implement the FDA-DECAF method, which involves constructing causal graphs and intensive model tuning, FDA-DECAF does not perform better than our FDA-CMM, which only involves sampling from Gaussian distributions.
When achieving similar utility-fairness trade-off, our FDA-CMM method offers significant advantages in terms of computational efficiency, ease of implementation, stability and interpretability.

480 481 482

485

445

446

447 448

449

## 5 DISCUSSION

In this paper, we introduce the FDA framework to generate debiased synthetic data, aiming to achieve intermediate levels of fairness and faithfulness as controlled by a user-specified unfairness

<sup>&</sup>lt;sup>6</sup>See Appendix C.3 for more implementation details.



Figure 3: Utility and fairness of the downstream models trained on the synthetic datasets obtained by DECAF, FDA-DECAF and FDA-CMM under varying values of  $\alpha$ : DP $\downarrow$  (top left), Precision $\uparrow$ (top right), Recall $\uparrow$  (bottom left), AUROC $\uparrow$  (bottom right). The shadowed areas along each line represent the variations across 10 repeated experiments.

512

513

514

517

518 reduction factor ( $\alpha$ ). This is the first work to offer a provable trade-off characterization between 519 these two competing objectives. Our experimental results demonstrate the effectiveness of FDA in 520 achieving this trade-off, and validate our theoretical guarantees: achieving perfect data faithfulness 521 when setting  $\alpha = 1$ , which is an important feature not present in other methods, and absolute fair-522 ness when setting  $\alpha = 0$ . Moreover, our FDA framework improves the faithfulness of synthetic data 523 when integrated with DECAF, a GAN-based fair data generator, as both proven theoretically and 524 demonstrated empirically in our experiments. This enhancement in faithfulness, in turn, boosts the 525 utility of downstream models trained on the resulting synthetic data.

526 **Social implications.** Setting  $\alpha$  at different values allows users to balance the trade-off between abso-527 lute fairness and perfect data faithfulness, adapting to the specific needs and priorities of various ap-528 plications. This capability is crucial, as it enables decision-makers to make informed choices about 529 this trade-off without having to sacrifice one aspect for the other. Our FDA framework represents 530 a significant advancement in synthetic data generation, offering a transparent and robust approach 531 that involves only sampling from Gaussian distributions shown in Equation equation 6 and relies on no assumptions. Such simplicity contrasts sharply with the labor-intensive and frequently unstable 532 training processes of complex black-box methods. These features of our FDA framework make it 533 more accessible, allowing a broader range of practitioners to adopt it without requiring specialized 534 training. This could potentially transform how data is managed to ensure equity and fairness. 535

Future work. The current FDA framework addresses the generation of fair synthetic data for onedimensional labels. Extending this framework to multi-dimensional labels, including mixed labels
of continuous and categorical types, is computationally straightforward; however, the challenge lies
in the theoretical analysis required to establish the relationship between fairness and faithfulness in
this context. Addressing this challenge will be the focus of our future work.

# 540 REFERENCES

546

552

566

- 542 Martial Agueh and Guillaume Carlier. Barycenters in the wasserstein space. SIAM Journal on Mathematical Analysis, 43(2):904–924, 2011.
- Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Machine bias. In *Ethics of data and analytics*, pp. 254–264. Auerbach Publications, 2016.
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and machine learning: Limitations* and opportunities. MIT Press, 2023.
- Richard Berk, Hoda Heidari, Shahin Jabbari, Michael Kearns, and Aaron Roth. Fairness in criminal justice risk assessments: The state of the art. *Sociological Methods & Research*, 50(1):3–44, 2021.
- Miranda Bogen and Aaron Rieke. Help wanted: An examination of hiring algorithms, equity, and bias. 2018.
- Flavio Calmon, Dennis Wei, Bhanukiran Vinzamuri, Karthikeyan Natesan Ramamurthy, and Kush R
   Varshney. Optimized pre-processing for discrimination prevention. Advances in neural information processing systems, 30, 2017.
- Simon Caton and Christian Haas. Fairness in machine learning: A survey. ACM Computing Surveys, 56(7):1–38, 2024.
- Lénaïc Chizat, Pierre Roussillon, Flavien Léger, François-Xavier Vialard, and Gabriel Peyré. Faster
   wasserstein distance estimation with the sinkhorn divergence. In *Proceedings of the 34th Inter- national Conference on Neural Information Processing Systems*, pp. 2257–2269, 2020.
- Evgenii Chzhen and Nicolas Schreuder. A minimax framework for quantifying risk-fairness trade off in regression. *The Annals of Statistics*, 50(4):2416–2442, 2022.
- Evgenii Chzhen, Christophe Denis, Mohamed Hebiri, Luca Oneto, and Massimiliano Pontil. Fair
   regression with wasserstein barycenters. *Advances in Neural Information Processing Systems*, 33:7321–7331, 2020.
- Lee Cohen, Zachary C Lipton, and Yishay Mansour. Efficient candidate screening under multiple
   tests and implications for fairness. In *1st Symposium on Foundations of Responsible Computing* (FORC 2020). Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 2020.
- Jeffrey Dastin. Amazon scraps secret AI recruiting tool that showed bias against women. In *Ethics of data and analytics*, pp. 296–299. Auerbach Publications, 2018.
- Amit Datta, Anupam Datta, Jael Makagon, Deirdre K Mulligan, and Michael Carl Tschantz. Dis crimination in online advertising: A multidisciplinary inquiry. In *Conference on Fairness, Ac- countability and Transparency*, pp. 20–34. PMLR, 2018.
- 579 Daniel de Vassimon Manela, David Errington, Thomas Fisher, Boris van Breugel, and Pasquale
   580 Minervini. Stereotype and skew: Quantifying gender bias in pre-trained and fine-tuned language
   581 models. In *Proceedings of the 16th Conference of the European Chapter of the Association for* 582 *Computational Linguistics: Main Volume*, pp. 2232–2242, 2021.
- DC Dowson and BV666017 Landau. The fréchet distance between multivariate normal distributions. *Journal of multivariate analysis*, 12(3):450–455, 1982.
- 586 Dheeru Dua and Casey Graff. Uci machine learning repository. 2020. URL https://archive.
   587 ics.uci.edu/ml.
- Pavel Dvurechensky, Alexander Gasnikov, and Alexey Kroshnin. Computational optimal transport: Complexity by accelerated gradient descent is better than by sinkhorn's algorithm. In *International conference on machine learning*, pp. 1367–1376. PMLR, 2018.
- Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubra manian. Certifying and removing disparate impact. In *proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 259–268, 2015a.

594 595 596 597 598	Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasub- ramanian. Certifying and removing disparate impact. In <i>Proceedings of the 21th ACM SIGKDD</i> <i>International Conference on Knowledge Discovery and Data Mining</i> , KDD '15, pp. 259–268, New York, NY, USA, 2015b. Association for Computing Machinery. ISBN 9781450336642. doi: 10.1145/2783258.2783311. URL https://doi.org/10.1145/2783258.2783311.
599 600 601	Solenne Gaucher, Nicolas Schreuder, and Evgenii Chzhen. Fair learning with wasserstein barycenters for non-decomposable performance measures. <i>arXiv preprint arXiv:2209.00427</i> , 2022.
602 603	Thibaut Le Gouic, Jean-Michel Loubes, and Philippe Rigollet. Projection to fairness in statistical learning. <i>arXiv preprint arXiv:2005.11720</i> , 2020.
604 605 606	Max Hort, Zhenpeng Chen, Jie M Zhang, Federica Sarro, and Mark Harman. Bia mitigation for machine learning classifiers: A comprehensive survey. <i>arXiv preprint arXiv:2207.07068</i> , 2022.
607 608 609	Bei Jiang, Adrian E Raftery, Russell J Steele, and Naisyin Wang. Balancing inferential integrity and disclosure risk via model targeted masking and multiple imputation. <i>Journal of the American Statistical Association</i> , 117(537):52–66, 2022.
610 611 612	Ray Jiang, Aldo Pacchiano, Tom Stepleton, Heinrich Jiang, and Silvia Chiappa. Wasserstein fair classification. In <i>Uncertainty in artificial intelligence</i> , pp. 862–872. PMLR, 2020.
613 614	Vivek Krishnamurthy. AI and human rights law. Artificial Intelligence and the Law in Canada (Toronto: LexisNexis Canada, 2021), 2021.
615 616 617 618	Kaiji Lu, Piotr Mardziel, Fangjing Wu, Preetam Amancharla, and Anupam Datta. Gender bias in neural natural language processing. <i>Logic, Language, and Security: Essays Dedicated to Andre Scedrov on the Occasion of His 65th Birthday</i> , pp. 189–202, 2020.
619 620	Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. <i>ACM Computing Surveys (CSUR)</i> , 54(6):1–35, 2021.
621 622 623	Amitabha Mukerjee, Rita Biswas, Kalyanmoy Deb, and Amrit P Mathur. Multi–objective evolution- ary algorithms for the risk–return trade–off in bank loan management. <i>International Transactions</i> <i>in operational research</i> , 9(5):583–597, 2002.
624 625 626	Victor M Panaretos and Yoav Zemel. Statistical aspects of wasserstein distances. <i>Annual review of statistics and its application</i> , 6:405–431, 2019.
627 628	Dana Pessach and Erez Shmueli. A review on fairness in machine learning. <i>ACM Computing Surveys</i> ( <i>CSUR</i> ), 55(3):1–44, 2022.
629 630 631 632 633	Amirarsalan Rajabi and Ozlem Ozmen Garibay. Tabfairgan: Fair tabular data generation with generative adversarial networks. <i>Machine Learning and Knowledge Extraction</i> , 4(2):488–501, 2022. ISSN 2504-4990. doi: 10.3390/make4020022. URL https://www.mdpi.com/ 2504-4990/4/2/22.
634 635 636	Filippo Santambrogio. Optimal Transport for Applied Mathematicians. Springer International Publishing, 2015. doi: 10.1007/978-3-319-20828-2. URL https://doi.org/10.1007% 2F978-3-319-20828-2.
637 638 639	Max Sommerfeld and Axel Munk. Inference for empirical wasserstein distances on finite spaces. Journal of the Royal Statistical Society Series B: Statistical Methodology, 80(1):219–238, 2018.
640 641 642	Carla Tameling, Max Sommerfeld, and Axel Munk. Empirical optimal transport on countable metric spaces: Distributional limits and statistical applications. <i>Annals of applied probability</i> , 29(5): 2744–2781, 2019.
643 644 645	Cuong Tran, Keyu Zhu, Ferdinando Fioretto, and Pascal Van Hentenryck. Sf-pate: scalable, fair, and private aggregation of teacher ensembles. <i>arXiv preprint arXiv:2204.05157</i> , 2022.
646 647	Boris van Breugel, Trent Kyono, Jeroen Berrevoets, and Mihaela van der Schaar. Decaf: Generating fair synthetic data using causally-aware generative networks. <i>Advances in Neural Information Processing Systems</i> , 34:22221–22233, 2021.

648 649	Cédric Villani. Topics in optimal transportation, volume 58. American Mathematical Soc., 2021.
650	Sandra Wachter, Brent Mittelstadt, and Chris Russell. Why fairness cannot be automated: Bridging
651	the gap between EU non-discrimination law and AI. Computer Law & Security Review, 41:
652	105567, 2021.
653	$\mathbf{D}$ (1) $\mathbf{V}$ (1) $\mathbf{V}$ (1) $\mathbf{V}$ (1) $\mathbf{U}$ (2) $\mathbf{T}$ (1) $\mathbf{U}$ (1)
654	wasserstein barycenter arYiv preprint arYiv:2211.01528, 2022
655	wasserstein-barycenter. <i>urxiv preprint urxiv.2211.01526</i> , 2022.
656	Depeng Xu, Shuhan Yuan, Lu Zhang, and Xintao Wu. Fairgan: Fairness-aware generative adver-
657	sarial networks. In 2018 IEEE International Conference on Big Data (Big Data), pp. 570–575.
658	IEEE, 2018.
659	Depeng Xu, Yongkai Wu, Shuhan Yuan, Lu Zhang, and Xintao Wu. Achieving causal fairness
660	through generative adversarial networks. In <i>Proceedings of the Twenty-Eighth International Joint</i>
661	Conference on Artificial Intelligence, 2019.
662	
663	Rich Zemel, Yu Wu, Kevin Swersky, Ioni Pitassi, and Cynthia Dwork. Learning fair representations.
664	In International conjerence on machine learning, pp. 323–355. PMLR, 2015.
665	Lu Zhang, Yongkai Wu, and Xintao Wu. A causal framework for discovering and removing direct
666	and indirect discrimination. In Proceedings of the Twenty-Sixth International Joint Conference
667	on Artificial Intelligence, IJCAI-17, pp. 3929–3935, 2017a. doi: 10.24963/ijcai.2017/549. URL
668	https://doi.org/10.24963/ijcai.2017/549.
669	Lu Zhang, Yongkai Wu, and Xintao Wu. A causal framework for discovering and removing di-
670	rect and indirect discrimination. In <i>Proceedings of the 26th International Joint Conference on</i>
671	Artificial Intelligence, pp. 3929–3935, 2017b.
672	
673	Han Zhao and Geoffrey J Gordon. Inherent tradeoffs in learning fair representations. <i>The Journal</i>
674	of Machine Learning Research, 25(1).2521–2552, 2022.
675	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word
675 676 677	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i>
675 676 677 678	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 683 684 685 686	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 685 686 687	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 684 685 686 687 688	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 684 685 686 687 688 689	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 684 685 686 687 688 689 690	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 684 685 686 687 688 689 690 691	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 684 685 686 687 688 689 690 691 692	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 684 685 686 687 688 689 690 691 692 693	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 686 685 686 687 688 689 690 691 692 693 694	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 686 685 686 687 688 689 690 691 692 693 694 695	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 686 685 686 687 688 689 690 691 692 693 694 695 695	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 696	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 691 692 693 694 695 695 696	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 692 693 694 695 695 696 697 698 699	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.
675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 694 695 696 697 698 699 700	Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. Learning gender-neutral word embeddings. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> , pp. 4847–4853, 2018.

#### A DEFINITIONS AND NOTIONS

#### A.1 P-WASSERSTEIN DISTANCE

**Definition A.1** (p-Wasserstein distance). Let  $\mu$  and  $\nu$  be two probability measures on  $\mathbb{R}$ . For  $p \ge 1$ , the p-Wasserstein distance between  $\mu$  and  $\nu$  is defined as

$$W_p(\mu, \nu) = \left[\inf_{\gamma \in \Gamma(\mu, \nu)} \int |x - y|^p d\gamma(x, y)\right]^{rac{1}{p}} \, ,$$

where  $\Gamma(\mu, \nu)$  is the set of joint probability measures on  $\mathbb{R} \times \mathbb{R}$  with marginals are  $\mu$  and  $\nu$ . Namely, for all measurable sets  $A, B \subseteq \mathbb{R}$ , it holds that  $\gamma(A \times \mathbb{R}) = \mu(A)$  and  $\gamma(\mathbb{R} \times B) = \nu(B)$ .

#### 715 A.2 EVALUATION OF THE UNFAIRNESS MEASURE $\mathcal{UF}(\mathcal{P}_{\mathcal{D}})$

The closed-form computation of  $\mathcal{UF}(\mathcal{P}_{\mathcal{D}})$  for any  $\mathcal{P}_{\mathcal{D}}$  is not easy due to the complex computation of Wasserstein-2 distance in equation 1. In the following example, we present one example when  $\mathcal{UF}(\mathcal{P}_{\mathcal{D}})$  can be explicitly computed.

Consider the case when  $\mu_{Y|s} = \mathcal{N}(b_s, \sigma_s^2)$  is a normal distribution, then it is known the minimizer (Agueh & Carlier, 2011) of equation 1 is  $\nu = \mathcal{N}\left(\sum_{s=1}^{K} \omega_s b_s, \left(\sum_{s=1}^{K} \omega_s \sigma_s\right)^2\right)$ . Therefore, one can compute  $\mathcal{UF}(\mathcal{P}_{\mathcal{D}})$  (Dowson & Landau, 1982) as

$$\mathcal{UF}(\mathcal{P}_{\mathcal{D}}) = \sum_{s=1}^{K} \omega_s \left[ \left( b_s - \sum_{s=1}^{K} \omega_s b_s \right)^2 + \left( \sigma_s - \sum_{s=1}^{K} \omega_s \sigma_s \right)^2 \right]^{\frac{1}{2}} .$$

727 728 729

730

731

725 726

702

703 704

705 706

707

714

716

In practice, the unfairness measure can be estimated by replacing the Wasserstein-2 distance in equation 1 with its estimator. There are many well studied estimators of Wasserstein-2 distance, for example, the plug-in estimator (Sommerfeld & Munk, 2018; Tameling et al., 2019; Dvurechensky et al., 2018), the estimation based on Sinkhorn divergence (Chizat et al., 2020).

#### 732 733 734 735

746

753

754 755

## A.3 INTERPRETATION OF THE WEIGHTS IN EQUATION 1

The unfairness measure provides flexibility in choosing different weights  $(\omega_1, \dots, \omega_K)$  to accommodate various purposes, particularly when there are majority and/or minority groups with respect to the sensitive attribute S. For example, the majority group (with respect to  $s_{\text{majority}}$ ) can be identified when  $\mathbb{P}(S = s_{\text{majority}}) \gg \mathbb{P}(S = s)$  for all  $s \in [K] \setminus \{s_{\text{majority}}\}$ .

740In general, for any  $s \in [K]$ , the larger the  $\omega_s$ , the closer the  $\mu_{Y|X}$  is to the optimal  $\nu$  (the minimizer741in equation 1). This leads to some natural choices of  $(\omega_1, \cdots, \omega_K)$ , including  $\omega_s = \mathbb{P}(S = s)$ 742and  $\omega_s \propto 1/\mathbb{P}(S = s)$ . When using the former, the optimal  $\nu$  will be closer to the conditional743distribution of Y for the majority group; when using the latter, the optimal  $\nu$  will be closer to the744conditional distribution of Y for the minority group. Alternatively, one could use equal weights by745letting  $\omega_s = \frac{1}{K}$ , when all groups are similar in size.

#### 747 A.4 FAIRNESS NOTIONS 748

749 **Definition A.2** (Conditional Fairness). Let  $X = (\tilde{X}, F)$ ,  $\mathcal{P}_{\mathcal{D}}$  is said to satisfy conditional fairness 750 with respect to , if  $(Y|S = s_1, F = f) \stackrel{d}{=} (Y|S = s_2, F = f)$ , for any  $s_1, s_2 \in [K]$ . That is to say, 751  $\mathbb{P}(Y \in T|S = s_1, F = f) = \mathbb{P}(Y \in T|S = s_2, F = f)$  for any  $T \subseteq \mathbb{R}$ . 752

# B PROOFS FOR MAIN RESULTS IN SECTION 3

The detailed proofs of the results in the main paper are included in this appendix.

# 756 B.1 PROOF OF EQUATION EQUATION 6 757

*Proof.* As discussed in the main paper, given the estimates of the model parameters  $(\hat{\beta}, \hat{\eta}^2, \hat{\sigma}^2)$ , we draw the synthetic  $\hat{Y}_i$  from the following predictive distribution defined under  $\mathcal{M}_{\text{FDA}}$ :

$$p(Y_i|Z_i, X_i, S_i; \hat{\beta}, \hat{\eta}^2, \hat{\sigma}^2) \propto \underbrace{p(Y_i|X_i, S_i; \hat{\eta}^2, \hat{\beta})}_{\mathcal{M}_{\text{fair}}} \underbrace{p(Z_i|Y_i; \hat{\sigma}^2)}_{\mathcal{M}_{\text{fairhful}}}.$$
(13)

Under the fair model  $\mathcal{M}_{\text{fair}}$  defined in equation 2, we have

$$p(Y_i|X_i, S_i; \hat{\eta}^2, \hat{\beta}) = \mathcal{N}\left(f(X_i, S_i, \hat{\beta}), \hat{\eta}^2\right),$$

and under the faithful model  $\mathcal{M}_{\text{fair}}$  defined in equation 3, we have

$$p(Z_i|Y_i;\hat{\sigma}^2) = \prod_{m=1}^M \mathcal{N}(Y_i,\hat{\sigma}^2)$$

That is,

$$p(Y_i|X_i, S_i; \hat{\beta}, \hat{\eta}^2) \propto \exp\left\{-\frac{\left(Y_i - f\left(X_i, S_i, \hat{\beta}\right)\right)^2}{2\hat{\eta}^2}\right\}$$

775 and

$$p(Z_i|Y_i;\hat{\sigma}^2) \propto \prod_{m=1}^M \exp\left\{-\frac{(Y_i - Z_{im})^2}{2\hat{\sigma}^2}\right\}$$

Then, realizing that

$$\exp\left\{-\frac{\left(Y_i - f\left(X_i, S_i, \hat{\beta}\right)\right)^2}{2\hat{\eta}^2}\right\} \prod_{m=1}^M \exp\left\{-\frac{\left(Y_i - Z_{im}\right)^2}{2\hat{\sigma}^2}\right\}$$
$$\propto \exp\left\{-\frac{1}{2}\left(\frac{\frac{\hat{\sigma}^2}{M}\hat{\eta}^2}{\frac{\hat{\sigma}^2}{M} + \hat{\eta}^2}\right)^{-1}\left(\hat{Y}_i - \frac{\frac{\hat{\sigma}^2}{M}f\left(X_i, S_i, \hat{\beta}\right) + \frac{\sum_{m=1}^M Z_{im}}{M}\hat{\eta}^2}{\frac{\hat{\sigma}^2}{M} + \hat{\eta}^2}\right)^2\right\},$$

which corresponds to the kernel of the Gaussian distribution as defined in equation 6:

$$\mathcal{N}\left(\frac{\frac{\hat{\sigma}^2}{M}f\left(X_i, S_i, \hat{\beta}\right) + \frac{\sum_{m=1}^{M} Z_{im}}{M}\hat{\eta}^2}{\frac{\hat{\sigma}^2}{M} + \hat{\eta}^2}, \frac{\frac{\hat{\sigma}^2}{M}\hat{\eta}^2}{\frac{\hat{\sigma}^2}{M} + \hat{\eta}^2}\right)$$

793 This concludes the proof of equation 6.

For clarity and ease of reading, we present the proof of Theorem 3.6 before Theorem 3.5.

#### 797 B.2 PROOF OF THEOREM 3.6

*Proof.* To determine the Wasserstein distance between  $\mu_{\hat{Y}|X,S}$  and  $\mu_{Y|X,S}$ , it is convenient to rewrite the generating model as

$$\hat{Y} = \frac{\frac{\sigma^2}{M} f(X, S, \beta)}{\frac{\sigma^2}{M} + \eta^2} + \frac{\frac{\sum_{j=1}^M Z_j}{M} \eta^2}{\frac{\sigma^2}{M} + \eta^2} + \sqrt{\frac{\frac{\sigma^2}{M} \eta^2}{\frac{\sigma^2}{M} + \eta^2}} N_1 ,$$
  
$$Z_j = Y + \sigma N_2 , \quad \text{for } j = 1, \cdots, M ,$$

where  $N_1$  and  $N_2$  are independent standard normal random variables that is independent of Y,  $\{Z_j\}_{j=1}^M$  are M noisy copies of Y (we omit the individual index *i* for simplicity). Thus,

$$\hat{Y} = \frac{\frac{\sigma^2}{M}f(X,S,\beta)}{\frac{\sigma^2}{M} + \eta^2} + \frac{Y\eta^2}{\frac{\sigma^2}{M} + \eta^2} + \sqrt{\frac{\frac{\sigma^2}{M}\eta^2}{\frac{\sigma^2}{M} + \eta^2}}N_1 + \frac{\eta^2}{\frac{\sigma^2}{M} + \eta^2}\sqrt{\frac{\sigma^2}{M}}N_2 \,.$$

810  
811 Followed by 
$$\hat{Y} - Y = \frac{\frac{\sigma^2}{M}f(X,S,\beta)}{\frac{\sigma^2}{M} + \eta^2} - \frac{Y\frac{\sigma^2}{M}}{\frac{\sigma^2}{M} + \eta^2} + \sqrt{\frac{\frac{\sigma^2}{M}\eta^2}{\frac{\sigma^2}{M} + \eta^2}} N_1 + \frac{\eta^2}{\frac{\sigma^2}{M} + \eta^2} \sqrt{\frac{\sigma^2}{M}} N_2$$
, one has

$$W_2^2(\mu_{\hat{Y}|X,S},\mu_{Y|X,S}) = \inf_{\gamma \in \Gamma(\mu_{\hat{Y}|X,S},\mu_{Y|X,S})} \int (\hat{y} - y)^2 d\gamma(\hat{y}, y)$$

$$= \mathbb{E}\left[\left(\frac{\frac{\sigma^2}{M}f(X,S,\beta)}{\frac{\sigma^2}{M} + \eta^2} - \frac{Y\frac{\sigma^2}{M}}{\frac{\sigma^2}{M} + \eta^2} + \sqrt{\frac{\frac{\sigma^2}{M}\eta^2}{\frac{\sigma^2}{M} + \eta^2}}N_1 + \frac{\eta^2}{\frac{\sigma^2}{M} + \eta^2}\sqrt{\frac{\sigma^2}{M}}N_2\right)^2 \left|X,S\right]\right]$$
$$= \left(\frac{\sigma^2}{\sigma^2 + M\eta^2}\right)^2 \mathbb{E}\left[\left(Y - f(X,S,\beta)\right)^2 \left|X,S\right] + \frac{\sigma^2\eta^2}{\sigma^2 + M\eta^2} + \frac{M\eta^4\sigma^2}{(\sigma^2 + M\eta^2)^2},\right]$$

where the last equation is a direct computation by taking the expectation of the squared form. The proof is competed by the relationships  $\frac{\lambda\eta^2}{\lambda+\eta^2} + \frac{\eta^4\lambda}{(\lambda+\eta^2)^2} = (1-\alpha)\eta^2 + \alpha(1-\alpha)\eta^2 = (1-\alpha^2)\eta^2$ . To see the asymptotic result when  $M \to \infty$  or/and  $\sigma^2 \to 0$ , we shall check the moments of  $\hat{Y} - Y$ .

$$\mathbb{E}\left[(\hat{Y} - Y)^{2}|X, S\right]$$

$$= \mathbb{E}\left[\left(\frac{\frac{\sigma^{2}}{M}f(X, S, \beta)}{\frac{\sigma^{2}}{M} + \eta^{2}} - \frac{Y\frac{\sigma^{2}}{M}}{\frac{\sigma^{2}}{M} + \eta^{2}} + \sqrt{\frac{\frac{\sigma^{2}}{M}\eta^{2}}{\frac{\sigma^{2}}{M} + \eta^{2}}}N_{1} + \frac{\eta^{2}}{\frac{\sigma^{2}}{M} + \eta^{2}}\sqrt{\frac{\sigma^{2}}{M}}N_{2}\right)^{2} \left|X, S\right]$$

$$= \left(\frac{\sigma^{2}}{\sigma^{2} + M\eta^{2}}\right)^{2} \mathbb{E}\left[(Y - f(X, S, \beta))^{2} |X, S\right] + \frac{\sigma^{2}\eta^{2}}{\sigma^{2} + M\eta^{2}} + \frac{M\eta^{4}\sigma^{2}}{(\sigma^{2} + M\eta^{2})^{2}}.$$
(14)
(14)

It is trivial to see from equation 15 converges to 0 when  $M \to \infty$  or/and  $\sigma^2 \to 0$ .

By Markov's inequality, for any  $\epsilon > 0$ ,

 $\mathbb{P}(|\hat{Y} - Y| > \epsilon | X, S) \le \frac{\mathbb{E}\left[(\hat{Y} - Y)^2 | X, S\right]}{\epsilon^2},$ (16)

where the right hand side converges to 0 when  $M \to \infty$  or/and  $\sigma^2 \to 0$ . Therefore, we have  $\hat{Y}|X, S \to Y|X, S$  in probability.

#### B.3 PROOF OF THEOREM 3.7

The proof of Theorem 3.7 follows directly from the result in Theorem 3.6. The detail is given as follows.

*Proof.* If  $\tilde{\mathcal{D}}$  is only generated from the fair model  $\mathcal{M}_{\text{fair}}$ , we have

$$\tilde{Y} - Y = f(X, S, \beta) - Y + \eta N_1 \,,$$

where  $N_1$  is a standard normal random variable that is independent of Y. Thus,

$$W_2^2(\mu_{\tilde{Y}|X,S},\mu_{Y|X,S}) = \inf_{\gamma \in \Gamma(\mu_{\tilde{Y}},\mu_Y)} \int (\tilde{y}-y)^2 d\gamma(\tilde{y},y)$$
$$= \mathbb{E}\left[ \left( f(X,S,\beta) - Y + \eta N_1 \right)^2 |X,S] \right]$$
$$= \mathbb{E}\left[ \left( f(X,S,\beta) - Y \right)^2 |X,S] + \eta^2 \right].$$

To obtain the inequality  $W_2^2(\mu_{\hat{Y}|X,S},\mu_{Y|X,S}) \leq W_2^2(\mu_{\tilde{Y}|X,S},\mu_{Y|X,S})$ , one can use the following results.

For any  $\sigma^2$  and M > 0, we always have  $\left(\frac{\sigma^2}{\sigma^2 + M\eta^2}\right)^2 \le 1$ , where the quality holds only when  $M \to 0$  or/and  $\sigma^2 \to \infty$ .

On the other hand, it is trivial to obtain the following inequality.

$$\frac{\sigma^2 \eta^2}{\sigma^2 + M\eta^2} + \frac{M\eta^4 \sigma^2}{(\sigma^2 + M\eta^2)^2} = \eta^2 \left(1 - \frac{M^2 \eta^4}{(\sigma^2 + M\eta^2)^2}\right) \le \eta^2 \,,$$

where the equality of the last inequality holds only when  $M \to 0$  or/and  $\sigma^2 \to \infty$ .

873 Thus, the inequality in Remark 3.7 is obtained immediately.874

#### B.4 PROOF OF THEOREM 3.5

*Proof.* First, due to the scaling law of Wasserstein-2 distance and its corresponding barycenter (Santambrogio, 2015; Panaretos & Zemel, 2019; Chzhen & Schreuder, 2022; Villani, 2021), we have

$$\min_{\nu \in \mathcal{P}_2(\mathbb{R})} \sum_{s=1}^K \omega_s W_p(\mu_{cY|s}, \nu) = c \min_{\nu \in \mathcal{P}_2(\mathbb{R})} \sum_{s=1}^K \omega_s W_p(\mu_{Y|s}, \nu),$$
(17)

for any  $c \ge 0$ .

Secondly, by the translation invariant property of Wasserstein-2 distance (Santambrogio, 2015; Panaretos & Zemel, 2019; Villani, 2021), one has

$$W_2(\mu_{Y+Z+a},\mu_{X+Z+a}) = W_2(\mu_Y,\mu_X)$$

for any constant a and random variable Z that is independent of Y and X. Thus,

$$\min_{\nu \in \mathcal{P}_2(\mathbb{R})} \sum_{s=1}^K \omega_s W_2(\mu_{Y+Z+a|s}, \nu) = \min_{\nu \in \mathcal{P}_2(\mathbb{R})} \sum_{s=1}^K \omega_s W_2(\mu_{Y|s}, \nu) \,. \tag{18}$$

Based on the generating model of  $\hat{Y}$ , we have

$$\hat{Y} = \frac{\frac{\sigma^2}{M} f(X, S, \beta)}{\frac{\sigma^2}{M} + \eta^2} + \frac{Y\eta^2}{\frac{\sigma^2}{M} + \eta^2} + \sqrt{\frac{\frac{\sigma^2}{M}\eta^2}{\frac{\sigma^2}{M} + \eta^2}} N_1 + \frac{\eta^2}{\frac{\sigma^2}{M} + \eta^2} \sqrt{\frac{\sigma^2}{M}} N_2 ,$$

where  $N_1$  and  $N_2$  are independent standard normal random variables that are independent of Y. A direct application of equation 17 and equation 18 implies  $\mathcal{UF}(\mathcal{P}_{\hat{\mathcal{D}}}) = \frac{\eta^2}{\frac{\sigma^2}{M} + \eta^2} \mathcal{UF}(\mathcal{P}_{\mathcal{D}}) = \alpha \mathcal{UF}(\mathcal{P}_{\mathcal{D}})$ for any given weights  $(\omega_1, \cdots, \omega_K) \in \Delta^{K-1}$ .

B.5 PROOF OF PROPOSITION 3.8

*Proof.* Assume  $\nu_0 = \arg \min_{\nu \in \mathcal{P}_2(\mathbb{R})} \sum_{s=1}^K \omega_s W_2(\mu_{\hat{Y}|s}, \nu)$ , by Theorem 3.5 we have

$$\sum_{s=1}^{K} \omega_s W_2(\mu_{\hat{Y}|s}, \nu_0) = \alpha \mathcal{UF}(\mathcal{P}_{\mathcal{D}}).$$

By triangle inequality, we have

$$\sum_{s=1}^{K} \omega_s W_2(\mu_{g(X,S)|s},\nu_0) \le \sum_{s=1}^{K} \omega_s W_2(\mu_{\hat{Y}|s},\nu_0) + \sum_{s=1}^{K} \omega_s W_2(\mu_{g(X,S)|s},\mu_{\hat{Y}|s})$$
(19)

$$\alpha \mathcal{UF}(\mathcal{P}_{\mathcal{D}}) + \delta \,, \tag{20}$$

for any given  $(\omega_1, \cdots, \omega_K)$ . It follows that

$$\min_{\nu \in \mathcal{P}_2(\mathbb{R})} \sum_{s=1}^K \omega_s W_2(\mu_{g(X,S)|s}, \nu) \le \sum_{s=1}^K \omega_s W_2(\mu_{g(X,S)|s}, \nu_0) \le \alpha \mathcal{UF}(\mathcal{P}_{\mathcal{D}}) + \delta.$$

 $\leq$ 

## C ADDITIONAL EXPERIMENTS ON UCI ADULT DATASET

## C.1 COMPUTATION DETAILS OF OUR EXPERIMENTS

Data split: The UCI Adult dataset is randomly split into two sets: a training set with 63,000 data points and a testing set with data size 2,000 data points.

925 Model training: The model parameter  $\beta$  depends on the specific choice of  $\mathcal{M}_{\text{fair}}$ , for example, 926  $\hat{\beta} = \bar{Y}$  if  $\mathcal{M}_{\text{fair}}$  is CMM. Here, we estimate the parameters using their sample versions. That is, for 927 the CMM model,  $\hat{\beta} = \bar{Y}$  and  $\hat{\eta} = \frac{1}{n-1} \sum_{i=1}^{n} (Y_i - \bar{Y})^2$ . Additionally, since the hyperparameter  $\sigma^2$ 928 is tuned by the user, its estimation is not necessary and its true value is used. The average time spent 929 for the training process to obtain one synthetic dataset is around 1.3 seconds on the Adult dataset 930 with n = 63,000 in the training set. The simulations were performed using Python 3.8.8 on a PC 931 with a 12th Gen intel Core i5-12600K CPU with 32 GB of RAM running Windows 11.

Downstream model: To evaluate the fairness and utility of downstream models, we train Multi-layer Perceptron (MLP) models on the generated synthetic datasets based on different fair data generation models. For instance, MLP models are trained using MLPclassifier function from sklearn module with all default parameters.

#### C.2 ADDITIONAL FAIRNESS EVALUATION ON SYNTHETIC DATA BY FDA

In this section, we run experiments on the UCI Adult dataset  $\mathcal{D}$ . Synthetic dataset  $\hat{\mathcal{D}}$  is generated under FDA framework with different choices of  $\alpha$ . In addition to the unfairness measure  $\mathcal{UF}(\hat{\mathcal{D}})$  in Definition 2.2, we also assess the commonly used unfairness measure  $|\mathbb{P}(\hat{Y} = 1|S = 1) - \mathbb{P}(\hat{Y} = 1)|S = 1|S = 1$ 1|S=0| with respect to the bias reduction parameter  $\alpha$  and illustrate in Figure 4. Still, the result of FDA is compared with those of DECAF (van Breugel et al., 2021), FairGAN (Xu et al., 2018), OPPDP(Calmon et al., 2017), TabFairGAN(Rajabi & Garibay, 2022). Each experiment is repeated 10 times, the average performances are reported in solid lines with shadowed variation areas in Figure 4. It is clear that, FDA still shows a clear tuning mechanism on faithfulness and fairness with respect to the bias reduction parameter  $\alpha$ , and the tuning trend is the same as using the unfairness measurement  $\mathcal{UF}(\mathcal{P}_{\hat{\mathcal{D}}})$ . This further emphasises FDA does facilitate the trade-off between absolute fairness and perfect data faithfulness by varying the bias reduction factor  $\alpha$ . 



Figure 4: Fairness of the generated synthetic dataset by FDA, FairGAN, OPPDP, DECAF with respect to  $\alpha$ :  $|\mathbb{P}(\hat{Y} = 1|S = 1) - \mathbb{P}(\hat{Y} = 1|S = 0)| (\downarrow \text{ more fair})$ . The shadowed areas along each line represent the variations on 10 repeat ions of experiments.

Table 2: Data utility and fairness of the downstream MLP model trained on generated synthetic data by FDA-DECAF model with different tuning ratio  $\lambda$ .

FDA-DECAF		DATA UTILITY		FAIRNESS		
$\mathcal{M}_{\text{fair}}: \text{DECAF}$	PRECISION ↑	<b>R</b> ecall <sup>↑</sup>	AUROC↑	DP↓	FTU↓	
ORIGINAL DATASET	$0.879 \pm 0.012$	$0.933 \pm 0.012$	$0.773 \pm 0.021$	$0.182 \pm 0.019$	$0.028 \pm 0.013$	
$\lambda = 0.25$	$0.846 \pm 0.012$	$0.962 \pm 0.001$	$0.717 \pm 0.023$	$0.142 \pm 0.023$	$0.050\pm0.024$	
$\lambda = 0.5$	$0.846 \pm 0.001$	$0.958 \pm 0.011$	$0.716 \pm 0.019$	$0.141 \pm 0.025$	$0.065\pm0.028$	
$\lambda = 1$	$0.800\pm0.006$	$0.989 \pm 0.005$	$0.622 \pm 0.015$	$0.071\pm0.017$	$0.028 \pm 0.018$	
$\lambda = 2$	$0.754 \pm 0.001$	$1.000\pm0.000$	$0.507 \pm 0.002$	$0.002\pm0.002$	$0.003\pm0.002$	
$\lambda = 4$	$0.751 \pm 0.000$	$1.000\pm0.000$	$0.500\pm0.000$	$0.000\pm0.001$	$0.000\pm0.000$	
DECAF	$0.753 \pm 0.000$	$0.989 \pm 0.000$	$0.505 \pm 0.000$	$0.006\pm0.000$	$0.006\pm0.000$	

Table 3: Data utility and fairness of the downstream MLP model trained on generated synthetic data by FDA-CMM model with different tuning ratio  $\lambda$ .

FDA-CMM	DATA UTILITY			FAIRNESS		
$\mathcal{M}_{\text{fair}}:CMM$	PRECISION $\uparrow$	<b>R</b> ecall <sup>↑</sup>	AUROC↑	DP↓	FTU↓	
ORIGINAL DATASET	$0.877 \pm 0.009$	$0.934 \pm 0.009$	$0.768 \pm 0.016$	$0.169 \pm 0.022$	$0.031 \pm 0.026$	
$\lambda = 0.25$	$0.931 \pm 0.014$	$0.781 \pm 0.037$	$0.803 \pm 0.008$	$0.287 \pm 0.031$	$0.079 \pm 0.044$	
$\lambda = 0.5$	$0.861 \pm 0.007$	$0.928 \pm 0.019$	$0.738 \pm 0.009$	$0.155 \pm 0.041$	$0.079\pm0.039$	
$\lambda = 1$	$0.751 \pm 0.000$	$1.000\pm0.000$	$0.501 \pm 0.001$	$0.001\pm0.001$	$0.000\pm0.001$	
$\lambda = 2$	$0.751 \pm 0.000$	$1.000\pm0.000$	$0.500\pm0.000$	$0.000 \pm 0.000$	$0.000\pm0.000$	
$\lambda = 4$	$0.751 \pm 0.000$	$1.000\pm0.000$	$0.500\pm0.000$	$0.000 \pm 0.000$	$0.000\pm0.000$	

C.3 Improving data faithfulness and downstream utility using FDA based on tuning ratio  $\lambda$ 

We generate synthetic fair dataset by FDA-DECAF model and by FDA with different choices of  $\lambda$ , the evaluation of downstream faithfulness and fairness are presented and compared with DECAF in Table 2 and 3, respectively. As expected, It clearly shows the downstream faithfulness is decreasing when  $\lambda$  is increasing, while the fairness is increasing simultaneously. This scenario coincides with the theoretical findings in Section 3.1. In addition to DP fairness, we also evaluate a different fairness notion: Fairness Through Unawareness (FTU), which is the difference between the predicted variables of a downstream classifier for setting S = 1 and S = 0, respectively, while giving the same feature. FTU is evaluated by the metric as  $|\mathbb{P}_{S=1}(\hat{Y}|X) - \mathbb{P}_{S=0}(\hat{Y}|X)|$ .

Comparing results in Table 2 for FDA-DECAF and Table 3 for FDA-CMM, it is interesting to see when the tuning ratio  $\lambda$  is large (i.e.,  $\lambda \geq 2$ ) the data utility and fairness performances for FDA-DECAF and FDA-CMM coincide. That is to say, when high level of fairness is required, one can either use FDA-DECAF or FDA-CMM. However, it is clear that, not like FDA-DECAF (prior causal relationships knowledge is required), FDA-CMM is a very simple model with no prior knowledge requirement. FDA-CMM is very recommended due to its computation simplicity and less assumptions requirement. This scenario coincides the interpretation in Section 4.2 and its reason is when  $\lambda$  is large, the information from  $\mathcal{M}_{\text{faithful}}$  will significantly override the information from  $\mathcal{M}_{\text{fair}}$ , leading high fairness level of the synthetic dataset from FDA joint model. 

#### D ADDITIONAL EXPERIMENTS ON COMPAS DATASET

1020The proposed FDA framework is very general not only on the capability of using different fair model1021 $\mathcal{M}_{fair}$ , but also on the stability of its performance on various real data. To show the generalization1022of using FDA framework, we run experiments on COMPAS data (Angwin et al., 2016), which is a1023dataset contains information about defendants from Broward County, and contains attributes about1024defendants such as their ethnicity, language, sex, etc. ,and for each individual a Decile score showing1025the likelihood of recidivism (reoffending). It is known there is bias (Calmon et al., 2017; Rajabi<br/>& Garibay, 2022) between ethnicity and Decile score in the sense that Decile score for African-

1026 American group is more likely to be assigned a higher Decile score indicating higher likelihood 1027 of recidivism. Therefore, in this experiment, ethnicity is the sensitive feature, and we only keep 1028 individuals when the ethnicity is African American and Caucasian. Also, we drop features, such 1029 as FirstName, LastName, MiddleName, CASE ID, and DateOfBirth, as people usually do. We 1030 convert Decile score as binary variable: "Low Chance of recidivism" when Decile score is less than 1031 5; "High chance of recidivism" for the rest. In a word, sensitive attribute S = ethnicity and the 1032 outcome Y = Recidivism Chance.

In what follows, we repeat experiment on COMPAS dataset as we did for Adult dataset in Section 4.1 to show how FDA facilitates the trade-off between fairness (when  $\alpha = 0$ ) and perfect data faithfulness (when  $\alpha = 1$ ) by varying the bias reduction factor  $\alpha \in (0, 1)$ .

We generate fair synthetic data  $\hat{D}$  by using FDA with various choices of  $\alpha$ . For different levels of bias reduction factor  $\alpha \in [0, 1]$ , we report the average of

- the empirical estimates of the Wasserstein-2 distance between the synthetic and original data distributions 
   *Ŵ*<sub>2</sub>(μ<sub>ŷ</sub>, μ<sub>Y</sub>)) in Figure 5,
- (2) the empirical estimates of the unfairness measure  $\mathcal{UF}(\hat{\mathcal{D}})$  in the debiased synthetic data in in Figure 6,
- (3) the commonly used unfairness measure  $|\mathbb{P}(\hat{Y} = 1|S = 1) \mathbb{P}(\hat{Y} = 1|S = 0)|$  in Figure 7,
- 1046 across 10 repetitions of experiments.

1047 1048 As we have seen in Section 4.1, it is clear that FDA provides a tuning mechanism on faithfulness and 1048 fairness with respect to the bias reduction parameter  $\alpha$  comparing with other benchmark method. 1049 Furthermore, it is worth noting that the variation of  $\hat{W}_2(\mu_{\hat{Y}}, \mu_Y)$  and  $\widehat{\mathcal{UF}}(\mathcal{P}_{\hat{\mathcal{D}}})$  for FDA are very 1050 small, providing stability of FDA framework.



Figure 5: Faithfulness of the generated synthetic dataset by FDA with respect to  $\alpha$ :  $\hat{W}_2(\mu_{\hat{Y}}, \mu_Y)$ ( $\downarrow$  more faithful). The shadowed areas along each line represent the variations on 10 repeat ions of experiments.

- 1071 1072
- 1074

1070

1036

1039

1040

1041

1043

1044 1045

1052 1053

1054

1055

1056 1057 1058

1062

1064

- 1075
- 1076
- 1077
- 1078



1129 on 10 repeat ions of experiments.