
Repairing Regressors for Fair Binary Classification at Any Decision Threshold

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the problem of post-processing a supervised machine-learned regressor
2 to maximize fair binary classification at all decision thresholds. By decreasing
3 the statistical distance between each group’s score distributions, we show that
4 we can increase fair performance across all thresholds at once, and that we can
5 do so without a large decrease in accuracy. To this end, we introduce a formal
6 measure of *Distributional Parity*, which captures the degree of similarity in the
7 distributions of classifications for different protected groups. Our main result is to
8 put forward a novel post-processing algorithm based on optimal transport, which
9 provably maximizes Distributional Parity, thereby attaining common notions of
10 group fairness like Equalized Odds or Equal Opportunity at all thresholds. We
11 demonstrate on two fairness benchmarks that our technique works well empirically
12, while also outperforming and generalizing similar techniques from related work.

13 1 Introduction

14 A common approach to fair machine learning is to train a classifier with a chosen decision threshold
15 in order to attain a certain degree of accuracy, and then to post-process the classifier to correct for
16 unfairness according to a chosen fairness definition [4, 12, 20]. Despite the popularity of this approach,
17 it suffers from two major limitations. First, it is well-known that the specific choice of decision
18 threshold can influence both fairness and accuracy in practice [2] producing an undesirable trade-off
19 between the two objectives. Second, when deploying a classifier in the real world, practitioners
20 typically need to tinker with the threshold as they evaluate whether a model meets their domain-
21 specific needs [14, 5]. One strategy to address these limitations, is to develop a procedure
22 that produces regressors that guarantee a selected fairness notion at *all* possible thresholds, while
23 simultaneously preserving accuracy. If a regressor is fair at all thresholds, then a practitioner can
24 freely perform application-specific threshold tuning without ever needing to retrain.

25 Some prior work has investigated this strategy, by using optimal-transport methods to achieve a
26 single, often trivially satisfied, fairness notion – Demographic Parity – at all thresholds [13, 17, 6].
27 However, [12] and related impossibility results [15, 5] demonstrate that attaining fairness only at
28 Demographic Parity does *not* capture the nuances in unfairness arising from examining true positive
29 rates, false positive rates, and combinations thereof [2]. We therefore ask: *Can we train a regressor
30 once and obtain fair binary classifiers at all thresholds for more flexible group fairness notions?*

31 **Our Work.** Our key insight is to observe that parity in the distributions of a regressor’s output for
32 each sensitive group, prior to the application of a threshold, can be harnessed to attain fairness at all
33 thresholds simultaneously. This insight yields the following contributions: **(1)** We introduce a metric
34 called *Distributional Parity* (Definition 3.1) based on the Wasserstein-1 Distance, which enables
35 reasoning about fairness across all thresholds for a wide class of metrics. **(2)** We employ a technique

36 called Geometric Repair [10], which leverages an important connection between Wasserstein-2
 37 barycenters to post-process a regressor under a Distributional Parity constraint, attaining all-threshold
 38 fairness. **(3)** We prove that that Distributional Parity is convex on the set of models produced
 39 by Geometric Repair, thereby enabling efficient computation of our proposed post processing.
 40 Additionally, we show that the models produced by geometric repair are Pareto optimal in the multi-
 41 objective optimization of accuracy (via an ℓ_1 -type risk) and Distributional Parity. **(4)** Lastly, we
 42 synthesize these insights into a novel post-processing algorithm for a broad class of fairness metrics;
 43 our algorithm subsumes earlier work on all-threshold Demographic Parity, and we demonstrate its
 44 efficacy in experiments on common benchmarks.

45 2 Background

46 Let $X \subseteq \mathbb{R}^d$ be a feature space and $G = \{a, b\}$ be a set of binary protected attributes, for which
 47 a is the majority group and b is the minority group. We denote the label space as $Y = \{0, 1\}$,
 48 where 0 denotes the negative class and 1 the positive class. We assume elements in X , G , and Y
 49 are drawn from some underlying distribution, with corresponding random variables \mathbf{X} , \mathbf{G} , and
 50 \mathbf{Y} . The proportion of each group is represented $\rho_g = \Pr[\mathbf{G} = g]$. Let \mathcal{F} be a set of measurable
 51 *group-aware regressors* (from which binary classifiers are derived). Each regressor $f \in \mathcal{F}$ has
 52 signature $f : X \times G \rightarrow [0, 1]$ and outputs a *score* $s \in [0, 1]$ where $s = \Pr[\mathbf{Y} = 1 | \mathbf{X} = x, \mathbf{G} = g]$.
 53 For a fixed regressor $f \in \mathcal{F}$ and a decision threshold $\tau \in [0, 1]$, we can derive a binary classifier
 54 from f by computing $\mathbb{1}\{f \geq \tau\}$ for any $\tau \in [0, 1]$. For a group $g \in G$, the *group-conditional*
 55 *score distribution* is the distribution of scores produced by a regressor on that group. We denote this
 56 distribution $f(\mathbf{X}, \mathbf{G}) | \mathbf{G} = g$.

57 For $p \geq 1$, we define $\mathcal{P}_p([0, 1])$ be the space of probability measures on $[0, 1]$ with finite p -order-
 58 moments. We use $\mu_g \in \mathcal{P}_p([0, 1])$ to denote the probability measure associated with each group’s
 59 score distribution. We also make the following (standard) assumption on these measures.

60 **Assumption 2.1.** *Any measure $\mu \in \mathcal{P}_p([0, 1])$ with finite p -order moments is non-atomic and*
 61 *absolutely continuous with respect to the Lebesgue measure.*

62 This assumption provides two guarantees. First, it ensures the cumulative distribution function (CDF)
 63 of μ_g , denoted $F_{\mu_g}(\tau) = \mu_g([0, \tau])$, has a well defined inverse $F_{\mu_g}^{-1}$. Second, it ensures that certain
 64 optimal transport operations, upon which our contributions crucially rely, are well-defined.

65 2.1 Wasserstein Distance and Wasserstein Barycenters

66 Before introducing our solution, we present some necessary background on Wasserstein distance
 67 and Wasserstein barycenters.

68 **Wasserstein Distance.** Informally, the Wasserstein distance captures the difference between
 69 probability measures by measuring the *cost* of transforming one probability measure into the other.
 70 In the special case when distributions are univariate, the Wasserstein distance has a nice closed form.

71 **Definition 2.1** (Wasserstein Distance). *For two measures $\mu_1, \mu_2 \in \mathcal{P}_p([0, 1])$*

$$72 \mathcal{W}_p^p(\mu_1, \mu_2) = \int_{[0,1]} |F_{\mu_1}^{-1}(q) - F_{\mu_2}^{-1}(q)|^p dq. \quad (1)$$

73 We can also define the Wasserstein distance using transport plans; this is commonly referred
 74 to as Monge’s Formulation. A transport plan is a function $T \in \mathcal{T}$ where every function in \mathcal{T}
 75 satisfies standard pushforward constraints, i.e, $T_{\#}\mu_1 = \mu_2$ such that $\mu_2(B) = \mu_1(T^{-1}(B))$ for all
 76 measurable $B \subseteq [0, 1]$.

Definition 2.2 (Wasserstein Distance [Monge]).

$$77 \mathcal{W}_p^p(\mu_1, \mu_2) = \inf_{T \in \mathcal{T}} \int_{[0,1]} |q - T(q)|^p d\mu_1(q). \quad (2)$$

78 In our specific case where Assumption 2.1 is satisfied, we know that these transport plans which
 79 solve Monge’s formulation exist and we can define them in closed form.

80 **Remark 2.1.** *The transport plan from $\mu_1 \rightarrow \mu_2$ which minimizes Eq. (2) is defined $T_1^2(x) =$*
 81 *$F_{\mu_2}^{-1}(F_{\mu_1}(x))$ for all $p \geq 1$ [21, Remark 2.6]*

82 **Wasserstein Barycenter.** The Wasserstein barycenter is a weighted composition of two distribu-
 83 tions, much like a weighted average or midpoint in the Euclidean sense; it provides a principled way
 84 to compose two measures.

85 **Definition 2.3** (Wasserstein Barycenter). *For two measures $\mu_1, \mu_2 \in \mathcal{P}_p([0, 1])$ their α -weighted*
 86 *Wasserstein barycenter is denoted μ_α and is computed*

$$\mu_\alpha \leftarrow \arg \min_{\nu \in \mathcal{P}_p([0,1])} (1 - \alpha)\mathcal{W}_p^p(\mu_1, \nu) + \alpha\mathcal{W}_p^p(\mu_2, \nu), \quad (3)$$

87 *and in the special case when $\alpha = \rho_b$ we write μ^* .*

88 To complete the weighted-average analogy, α behaves like a tunable knob: As $\alpha \rightarrow 0$ then μ_α
 89 will appear more like μ_1 , and as $\alpha \rightarrow 1$ the more μ_α will appear like μ_2 . As a consequence of this
 90 definition and Remark 2.1, we can express the transport plan to a barycenter in closed form, as well:

91 **Corollary 2.1.** *Let μ_* be the ρ_b -weighted barycenter of μ_a, μ_b then the transport plan from $\mu_a \rightarrow \mu_*$*
 92 *(wlog) is computed $T_a^*(\omega) = (\rho_a F_{\mu_a}^{-1} + \rho_b F_{\mu_b}^{-1}) \circ F_{\mu_a}(\omega)$.*

93 **A Note on Our Use of \mathcal{W}_1 and \mathcal{W}_2 .** In this work, we make use of both \mathcal{W}_1 and \mathcal{W}_2 . Our use of
 94 \mathcal{W}_1 is restricted to Distributional Parity computations (see Section 3). This choice is motivated by
 95 the fact when $\gamma = \mathcal{U}_{PR}$, the Wasserstein-1 distance recovers \mathcal{U}_{PR} . We use \mathcal{W}_2 to compute Wasserstein
 96 barycenters. Given that \mathcal{W}_2 is known to be strictly convex, and provided that some μ_g is non-atomic,
 97 for $p = 2$ the barycenter that minimizes Eq. 2.3 is unique [1, Proposition 3.5].

98 3 A Distributional View of Fairness

99 Our goal is to post-process a regressor such that all binary classifiers derived from thresholding this
 100 regressor are group fair, i.e., attain fairness in the regressor at every threshold. To attain fairness at
 101 every threshold, we look to create parity in outcomes at the level of the regressor – before thresholds
 102 are applied – rather than at the level of the derived predictor. The intuition is simple: if a regressor
 103 outputs similar scores for two groups, then no matter what threshold is selected, the output derived
 104 predictor will be fair. Specifically, we show that fairness can be attained at all thresholds by enforcing
 105 parity in the *distribution* of scores output by a regressor on some groups.

106 At the core of our new distributional definition of fairness are familiar metrics, namely: Positive Rate
 107 (PR), True Positive Rate (TPR), and False Positive Rate (FPR). From these metrics, we can obtain
 108 popular fairness definitions, such as Demographic Parity (PR Parity) [4], Equal Opportunity (TPR
 109 Parity), and Equalized Odds (TPR and FPR Parity) [12].

110 Let the set of these metrics be $\Gamma = \{\text{PR}, \text{TPR}, \text{FPR}\}$ and any arbitrary metric be $\gamma \in \Gamma$. We write
 111 $\gamma_g(\tau; f)$ to denote the rate γ on group g at threshold τ for a score distribution produced by f . When
 112 obvious from context, we omit f from this γ notation, writing only $\gamma_g(\tau)$. Additionally, as we show
 113 via Corollary 5.1, we can combine these metrics additively, e.g., producing Equalized Odds which
 114 combines TPR and FPR.

115 At a single threshold, (un)fairness is commonly measured by taking the difference in some metric
 116 across groups — e.g., for the case of Demographic Parity where $\gamma = \text{PR}$, we can measure fairness by
 117 simply computing $|\text{PR}_a(\tau) - \text{PR}_b(\tau)|$. A natural way to leverage these single-threshold measurements
 118 into an all-threshold measurement is to take their average across every possible τ . We formalize
 119 this idea in the following definition of *Distributional Parity*.

120 **Definition 3.1** (Distributional parity). *Let $U([0, 1])$ be the uniform distribution on $[0, 1]$. For a fairness*
 121 *metric $\gamma \in \Gamma$, a regressor f satisfies Distributional Parity denoted $\mathcal{U}_\gamma(f) \triangleq \mathbb{E}_{\tau \sim U([0,1])} |\gamma_a(\tau) -$*
 122 *$\gamma_b(\tau)|$, when $\mathcal{U}_\gamma(f) = 0$.*

123 A useful property of this definition is that when $\gamma = \text{PR}$, this definition is closely related to the
 124 Wasserstein Distance, a distance which is frequently used to measure distance between probability
 125 distributions.

126 **Proposition 3.1.** *For $\mu_a, \mu_b \in \mathcal{P}_2([0, 1])$ which are the groupwise score distributions of f , then*
 127 *$\mathcal{W}_2(\mu_a, \mu_b) = 0$ if and only if $\mathcal{U}_{PR}(f) = 0$.*

128 It is from this property that distributional parity is named. At its core, distributional parity is a way to
 129 quantify differences between *outcome* distributions – specifically the groupwise score distributions of

130 f . This relationship between distributional parity – an all threshold fairness metric – and the Wasser-
 131 stein distance – a measure of statistical distance – anchors our proposed shift in focus from thresholds
 132 to distributions. Next, we introduce our proposed post-processing objective for computing fair regres-
 133 sors under a distributional parity constraint. We will also begin to outline how we efficiently compute
 134 this post-processing, and how our solution elegantly addresses fairness-accuracy trade-off concerns.

135 3.1 Distributionally Fair Post-Processing

136 Our goal is to post-process a learned regressor f , such that it becomes (distributionally) fair while
 137 remaining accurate. The risk of some other regressor \hat{f} (with respect to f) is computed

$$\mathcal{R}(\hat{f}) = \|\hat{f} - f\|_1 = \mathbb{E} |\hat{f}(\mathbf{X}, \mathbf{G}) - f(\mathbf{X}, \mathbf{G})|$$

138 Using this definition of risk, a simple fair post-processing objective can be written as follows,

$$\arg \inf_{\hat{f} \in \mathcal{F}} \mathcal{R}(\hat{f}) \quad \text{s.t.} \quad \mathcal{U}_\gamma(\hat{f}) \leq c, \quad \text{where } c \text{ is some small constant.} \quad (4)$$

139 **The special case of PR.** In the special case where $\gamma = \text{PR}$ and $c = 0$, the solution to Eq. 4 can
 140 be computed using a solution based on optimal transport [13]. In this solution, a learned regressor
 141 f is transformed into a new regressor we call f^* which provably minimizes risk (with respect to f)
 142 while attaining distributional parity for $\gamma = \text{PR}$, i.e, demographic parity at every threshold. This
 143 all threshold guarantee is attained with minimal impact to risk. It was shown in [17, 6] that f^*
 144 is the regressor which increases risk the *least* amongst all regressors which satisfy all threshold
 145 demographic parity constraints.
 146

$$f^* \leftarrow \arg \min \mathcal{R}(\cdot) \quad \text{s.t.} \quad \mathcal{U}_{\text{PR}}(f^*) = 0. \quad (5)$$

147 This solution is strict – it enforces exact demographic parity, which may not always be desired [7].
 148 We can address this concern by considering a relaxation of Eq. 4 which uses a parameter λ to balance
 149 the a trade-off between fairness and accuracy. Specifically, for every $\lambda \in [0, 1]$ there is some $f_\lambda \in \mathcal{F}$
 150 which attains λ -increase in the fairness, in exchange for a λ -reduction in risk. We prove the existence
 151 of f_λ which satisfies this property in the following lemma .

152 **Lemma 3.1.** *Let f be some learned regressor. For all $\lambda \in [0, 1]$ the set of optimally fair regressors
 153 for λ -relaxations of f with respect to risk and distributional parity for $\gamma = \text{PR}$ are given by*

$$f_\lambda \leftarrow \arg \min_{\hat{f} \in \mathcal{F}} \lambda \mathcal{R}(\hat{f}) \quad \text{s.t.} \quad \mathcal{U}_{\text{PR}}(\hat{f}) = (1 - \lambda) \mathcal{U}_{\text{PR}}(f) \quad (6)$$

154 The functions f_λ are Pareto-optimal: indeed, we show in Theorem 4.2 that all $\{f_\lambda\}_{\lambda \in [0, 1]}$ are Pareto
 155 optimal in the multi-objective minimization of \mathcal{R} and \mathcal{U}_{PR} . This means we can view these regressors
 156 as being optimally accuracy preserving, while also being fair. As a result, the above optimization
 157 (with $\gamma = \text{PR}$) can be rewritten simply as

$$\arg \min_{\lambda \in [0, 1]} \mathcal{U}_\gamma(f_\lambda), \quad (7)$$

158 replacing the risk minimization objective with an objective that enforces distributional parity, given
 159 the aforementioned accuracy-preserving properties of f_λ .

160 **Extending to other fairness metrics.** The above approach to achieving all-threshold fairness
 161 has two steps: firstly, a characterization of a solution that achieves ideal fairness, and then the
 162 construction of a space of functions f_λ that allow for an optimal fairness-accuracy tradeoff. In order
 163 to generalize this to other fairness measures, we need version of both steps. The exact result for
 164 PR however relies heavily on optimal transport in a way that does not naturally generalize to other
 165 fairness measures $\gamma \neq \text{PR}$.

166 To address this, we provide two key insights. Firstly, that the f_λ can be expressed explicitly using
 167 optimal transport ideas in terms of score distributions that are independent of the choice of fairness
 168 measure, and secondly, that the optimization described in equation (7) describes a convex function of
 169 λ independent of the choice of γ .

170 This means that for any choice of γ , we can find the optimal value of λ to minimize $\mathcal{U}_\gamma(f_\lambda)$. And we
 171 will show empirically that this choice yields an almost perfect minimization of distributional parity,
 172 thus achieving an all-threshold fairness result as desired.

173 **4 Maximizing Distributional Parity with Geometric Repair**

174 We now introduce the method and main theorem that we use to compute distributionally fair post
 175 processing. The method, defined in Section 4.1 is called Geometric Repair [10, 7], and is how we
 176 efficiently compute solutions to the objective stated in Equation 7. Our main theoretical result is
 177 stated in Theorem 4.1. Subsequently, we show in subsection 4.2 that we can make use of an elegant
 178 transformation to an optimal transport problem in order to achieve approximate distributional parity
 179 for all γ .

180 **4.1 Defining Geometric Repair**

181 Geometric repair is a technique for constructing a regressor that interpolates between the output of
 182 some learned regressor f (assumed to be accurate), and the output of a certain fair function f^* . Note,
 183 f^* must be specifically chosen in order to prove our results, but for ease of exposition, we defer
 184 formal definition of f^* to the following subsection.

185 **Definition 4.1** (Geometric Repair). *We call $\lambda \in [0, 1]$ the repair parameter and define a geometrically
 186 repaired regressor f_λ as $f_\lambda(x, g) \triangleq (1 - \lambda)f(x, g) + \lambda f^*(x, g)$.*

187 Geometric repair enumerates a well structured set of regressors which achieve λ -relaxations of R and
 188 \mathcal{U}_{PR} as described in Section 3.

189 **Proposition 4.1.** *For any $\lambda \in [0, 1]$, a repaired regressor f_λ satisfies $R(f_\lambda) = \lambda R(f^*)$ and
 190 $\mathcal{U}_{PR}(f_\lambda) = (1 - \lambda)\mathcal{U}_{PR}(f)$.*

191 This is the set of regressors used to maximize distributional parity. The key to computing such a
 192 maximization lies in the following theorem, which shows that distributional parity is convex, on the
 193 set of repaired regressors. This convexity guarantee certifies our ability to locate the f_λ , amongst the
 194 set of repaired regressors, which *best* minimizes Distributional Parity for *any* γ .

195 **Theorem 4.1.** *Fix $\gamma \in \Gamma$. Let $f : X \times G \rightarrow [0, 1]$ be a regressor, and f_λ be the geometrically
 196 repaired regressor for any $\lambda \in [0, 1]$. The map $\lambda \mapsto \mathcal{U}_\gamma(f_\lambda)$ is convex in λ .*

197 The proof of this theorem crucially depends on the connection between f^* and Wasserstein barycen-
 198 ters. In the next section we, leverage this connection to analytically compute the distributions of f_λ ,
 199 which is a crucial piece needed in proving the convexity of $\mathcal{U}_\gamma(f_\lambda)$.

200 **4.2 How f^* Enables Geometric Repair**

201 Here, we formalize the earlier definition of f^* from Section 4.2 and its connection between to
 202 Wasserstein barycenters in the context of geometric repair.

203 **Definition 4.2** (Fully Repaired Regressor). *The regressor f^* which satisfies distributional parity for
 204 $\gamma = PR$ while minimizing risk (with respect to f) is the computed*

$$f^* \leftarrow \arg \min_{f \in \mathcal{F}} \mathcal{R}(\cdot) \quad \text{s.t.} \quad \mathcal{U}_{PR}(f) = 0. \quad (8)$$

205 *We call this regressor fully repaired in that $f_{\lambda=1}$ is equivalent to f^* .*

206 The aforementioned property which relates f^* to \mathcal{W}_2 barycenters is the the fairness constraint in Eq.
 207 (8). To make this clear, recall Proposition 3.1 which states that removing the \mathcal{W}_2 distance between
 208 distributions is sufficient to satisfy distributional parity for $\gamma = PR$. The tool we will use to remove
 209 this distance is, indeed, Wasserstein barycenters. Prior work [16, 6] show that mapping μ_a, μ_b onto
 210 their ρ_b -weighted barycenter distribution, which we denote μ_* , removes the Wasserstein distance
 211 between μ_a, μ_b under this mapping, thereby satisfying $\mathcal{U}_{PR}(f^*) = 0$ and establishing that f^* is
 212 distributed like μ_* .

213 We can use this fact to rewrite the score distributions of each group under geometric repair. The
 214 following proposition formalizes this claim, by showing that the groupwise distributions of output
 215 by any f_λ can be computed as barycenters of μ_g and μ_* .

216 **Proposition 4.2.** *Let $\lambda \in [0, 1]$. Let $\mu_{g,\lambda}$ be the λ -weighted barycenter between μ_g and μ_* , i.e.,*

$$\mu_{g,\lambda} \leftarrow \arg \min_{\nu \in \mathcal{P}_2([0,1])} (1 - \lambda)\mathcal{W}_2^2(\mu_g, \nu) + \lambda\mathcal{W}_2^2(\mu_*, \nu), \quad \text{then } \mu_{g,\lambda} = \text{Law}(f_\lambda(\mathbf{X}, \mathbf{G}) | \mathbf{G} = g).$$

217

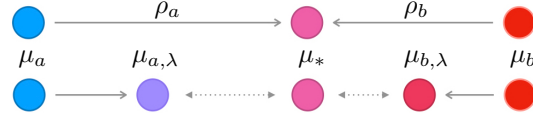


Figure 1: Let μ_a, μ_b be groupwise score distributions. We illustrate of the repaired score distributions $\mu_{g,\lambda}$ under geometric repair, where μ_* is the ρ_b -weighted barycenter.

218 This proposition shows us that the interpolation between f and f^* as parametrized by λ in geometric
 219 repair is replicated at the distributional level, i.e., λ also controls the interpolation from $\mu_{g,\lambda} \rightarrow \mu_*$;
 220 more importantly, the intermediate distributions of this interpolation have a special structure – they
 221 are barycenters. Note that under Assumption 2.1 and [1, Proposition 3.5], these $\mu_{g,\lambda}$ are unique
 222 and guaranteed to exist. For clarity, we visualize this interpolation (over distributions) in Figure 1.

223 Our proof of Theorem 4.1 computes distributional parity as a function of the score distributions of f_λ .
 224 With these established via Proposition 4.2, we are able to write a closed form expression for $\mathcal{U}_\gamma(f_\lambda)$.
 225 Using this expression, the proof proceeds computationally, showing that the second derivative of
 226 $\mathcal{U}_\gamma(f_\lambda)$ is non-negative to conclude convexity.

227 4.3 The Optimality of Geometric Repair in Balancing Fairness and Accuracy

228 Now, we will show that f_λ is optimal in the fairness-accuracy trade-off with respect to $\gamma = \text{PR}$.

229 **Definition 4.3** (Pareto Optimality). For $f, f' \in \mathcal{F}$ we say f Pareto dominates f' , denoted $f' \prec f$, if
 230 one of the following hold:

$$\mathcal{R}(f) \leq \mathcal{R}(f') \quad \mathcal{U}_{\text{PR}}(f) < \mathcal{U}_{\text{PR}}(f') \quad (9)$$

$$\mathcal{R}(f) < \mathcal{R}(f') \quad \mathcal{U}_{\text{PR}}(f) \leq \mathcal{U}_{\text{PR}}(f') \quad (10)$$

231 A regressor f is Pareto optimal if there is no other regressor f' that has improved risk without also
 232 having strictly more unfairness, or vice-versa.

233 The proof of Pareto optimality of f_λ follows from Proposition 4.1. The main idea of this result is
 234 the following: f^* is the lowest risk classifier where $\mathcal{U}_{\text{PR}}(\cdot) = 0$ meaning that it is Pareto optimal by
 235 construction. Since f_λ is a λ -relaxation of f^* with regards to both risk and unfairness, f_λ preserves
 236 the Pareto optimality of f^* .

237 **Theorem 4.2.** For all $\lambda \in [0, 1]$, the repaired regressor f_λ is pareto optimal in the multi-objective
 238 minimization of $\mathcal{R}(\cdot)$ and $\mathcal{U}_{\text{PR}}(\cdot)$.

239 5 Post-Processing Algorithms to Maximize Distributional Parity

240 Now that we have supported *why* we can use geometric repair to maximize distributional parity, we pro-
 241 vide some practical algorithms showing *how* to do so. First, we show how to estimate f_λ from samples.

242 **Plug-in Estimator for f_λ .** Indeed, computation of f^* , and therefore μ_* , requires exact knowledge
 243 of μ_a, μ_b . In practice we only have sample access to both score distributions, and so we must
 244 approximate these distributions, and consequently their barycenter and f_λ . We show a plug-in
 245 estimator w/ the following convergence guarantee (Theorem 5.1) to approximate f_λ in Algorithm 1.
 246 Our approach to approximating f_λ only requires a input regressor f and access to some unlabeled
 247 dataset $D = (x_1, g_1) \dots (x_n, x_g)$. Let n_g denote the number of samples from a group g .

248 **Theorem 5.1.** As $n_g \rightarrow \infty$ the empirical distribution of $\hat{f}_\lambda(x, g)$ converges to $\mu_{g,\lambda}$ in \mathcal{W}_2 almost
 249 surely.

250 **Post-Processing to Maximize Distributional Parity.** To actually compute the optimal λ_* for some
 251 metric, we propose the post-processing routine described in Algorithm 2. The algorithm consists
 252 of two main steps: approximating \hat{f}_λ in Step 1, and finding the optimal λ_* in Step 2. Note that our
 253 objective $\hat{\mathcal{U}}_\gamma(f_\lambda)$ is parametrized by the scalar λ , and so we find its minima using a univariate solver;
 254 we found success using Brent’s Method [3]. By the convexity of $\mathcal{U}_\gamma(\cdot)$ as proven in Theorem 4.1, we
 255 are guaranteed that the f_{λ_*} is optimal on the set of repaired regressors.

256 **Corollary 5.1.** Since convex functions are closed under addition, Theorem 4.1 also applies to
 257 additive combinations of metrics, meaning that the objective in Step (2) of Alg 2 can be replaced by
 258 $\mathcal{U}_{\gamma_1}(f_\lambda) + \mathcal{U}_{\gamma_2}(f_\lambda) + \dots + \mathcal{U}_{\gamma_m}(f_\lambda)$.

Algorithm 1 An Estimator for f_λ

Input: A regressor f , and an unlabeled dataset $D = (x_1, g_1) \dots (x_n, g_n)$

1. Let $n_g = \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{g_k=g}$. Use f to approximate the group-conditional distributions

$$\hat{\mu}_g = \frac{1}{n_g} \sum_{i=1}^n \delta_{f(x_i, g_i)} \mathbb{1}_{g_i=g}$$

2. Let $\hat{\rho}_g = \frac{n_g}{n}$ and compute the empirical optimal transport plans (see Remark 2.1)

$$\hat{T}_g^*(\omega) = (\hat{\rho}_a F_{\hat{\mu}_a}^{-1} + \hat{\rho}_b F_{\hat{\mu}_b}^{-1}) \circ F_{\hat{\mu}_a}(\omega)$$

3. For any $\lambda \in [0, 1]$, compute \hat{f}_λ where $\hat{f}_\lambda(x, g) = (1 - \lambda)f(x, g) + \lambda \hat{T}_g^*(f(x, g))$
-

Algorithm 2 Post-Processing for Distributional Parity

Input: A metric $\gamma \in \Gamma$, learned regressor f , and labeled dataset $E = (x_1, g_1, y_1) \dots (x_k, g_k, y_k)$

1. Using Algorithm (1) to approximate f_λ by computing \hat{T}_g such that for all $\lambda \in [0, 1]$ geometric repair is well defined, i.e., $\hat{f}_\lambda(x, g) = (1 - \lambda)f(x, g) + \lambda \hat{T}_g(f(x, g))$
2. Use Brent’s algorithm to find the optimal λ which minimizes $\lambda_* \leftarrow \text{Brent}_{\lambda \in [0, 1]} \hat{\mathcal{U}}_\gamma(f_\lambda)$ where $\hat{\mathcal{U}}(f_\lambda)$ is approximated for m randomly sampled $(\tau_1, \dots, \tau_m) \sim U([0, 1])$ via

$$\hat{\mathcal{U}}(f_\lambda) = \frac{1}{m} \sum_{\ell=1}^m |\gamma_a(\tau_\ell; f_\lambda) - \gamma_b(\tau_\ell; f_\lambda)|.$$

3. **Output:** $f_{\lambda_*}(x, g)$ such that $\hat{\mathcal{U}}_\gamma(f_{\lambda_*})$ is minimized (distributional parity is maximized)
-

259 6 Experiments

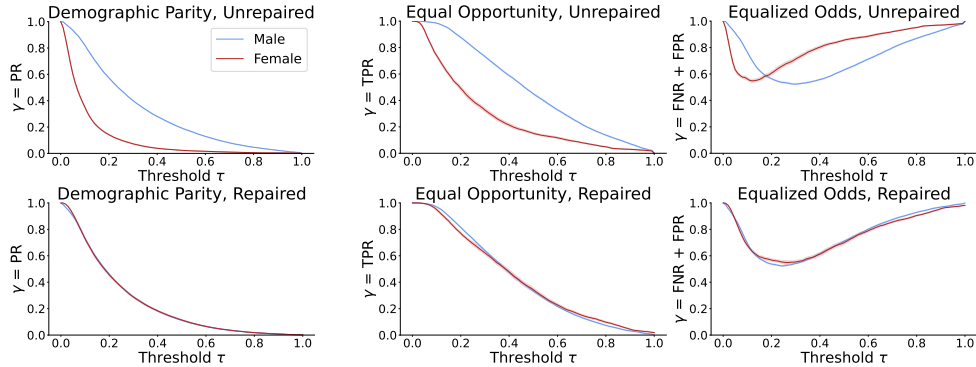
260 In this section we present experiments that demonstrate the effectiveness of our proposed algorithms
261 in Section 5. To that end, we provide two sets of results: 1) Figure 2 validates that Algorithm 2
262 achieves almost-exact distributional parity for Demographic Parity, Equal Opportunity, and Equalized
263 Odds; 2) Table 1 shows that Algorithm 2 outperforms related methods in maximizing Distributional
264 parity while preserving accuracy.

265 **Datasets.** We use two datasets: Adult Income-Sex from the the UCI repository [9], and Adult
266 Income-Race from the datasets produced in [8]. For both datasets, the task is to predict whether (1)
267 or not (0) an individual’s income exceeds \$50,000.

268 **Model Training.** To produce a model that we use in our experimentation, we implemented a
269 Logistic Regression (**LR**) with ℓ_2 regularization, and an Support Vector Machine (**SVM**) with an
270 Radial Basis Function kernel. Both were implemented using using scikit-learn with its default model
271 parameters and optimizers [18].

Table 1: Comparison of Geometric Repair (**GR**) against included baselines (abbreviations described under **baselines**). Results are averaged over ten trials, and the mean and standard deviation across all trials are reported for each metric.

		\mathcal{U}_{TPR}	TPR Optimized Worst case	AUC	\mathcal{U}_{EO}	EO Optimized Worst case	AUC
Income-Race (LR)	GR	0.025 ± 0.013	0.068 ± 0.034	<i>0.816 ± 0.004</i>	0.021 ± 0.007	0.055 ± 0.02	<i>0.816 ± 0.005</i>
	JIA	0.028 ± 0.007	0.129 ± 0.021	0.8 ± 0.005	0.049 ± 0.006	0.094 ± 0.014	0.8 ± 0.005
	FEL	0.042 ± 0.039	0.106 ± 0.051	0.81 ± 0.015	0.103 ± 0.039	0.213 ± 0.081	0.811 ± 0.014
	OG	0.219 ± 0.01	0.430 ± 0.024	0.834 ± 0.005	0.142 ± 0.007	0.299 ± 0.016	0.834 ± 0.005
Income-Sex (SVM)	GR	0.014 ± 0.005	0.042 ± 0.015	<i>0.882 ± 0.004</i>	0.032 ± 0.006	0.079 ± 0.015	<i>0.882 ± 0.004</i>
	JIA	0.052 ± 0.009	0.104 ± 0.013	0.878 ± 0.004	0.047 ± 0.006	0.110 ± 0.011	0.878 ± 0.004
	FEL	0.014 ± 0.007	0.08 ± 0.015	0.769 ± 0.007	0.026 ± 0.006	0.081 ± 0.013	0.769 ± 0.007
	OG	0.065 ± 0.01	0.114 ± 0.015	0.884 ± 0.003	0.035 ± 0.008	0.077 ± 0.013	0.884 ± 0.003



(a) Unrepaired vs. Full Repair

(b) Unrepaired vs. Optimal Repair

Figure 2: Performing geometric repair for $\gamma = \text{PR}$ (left), TPR(middle), EO (right) for Logistic Regression trained on Adult Income-Race. The top row depicts the rates for unrepaired regressors and the bottom row for the repaired regressor.

272 **Metrics.** We use the following measurements of model performance: (1) We approximate **Distribu-**
 273 **tional parity** $U_\gamma(\cdot)$ as per Step 2 of Algorithm 2 using $m = 100$ randomly sampled thresholds. We
 274 denote the Equalized Odds metric EO = FNR + FPR, i.e., the misclassification rate(2) We measure
 275 accuracy using the **Area Under the Curve (AUC)** given that AUC averages model performance
 276 across all thresholds similar to U_γ . (3) **Worst Case** refers to the worst disparity of the regressor at
 277 any threshold for the chosen γ , i.e., $\max_{\tau \in [0,1]} |\gamma_a(\tau) - \gamma_b(\tau)|$.

278 **Baselines.** We use the following algorithms as baselines : **(OG)** The learned classifier with no
 279 additional processing. **(JIA)** Post-processing algorithm proposed by Jiang et al. [13] which processes
 280 the output of a regressor such that model output is independent of protected group (shown to be equal
 281 to satisfying $U_{\text{PR}} = 0$, which is achieved by our method for $\lambda = 1$). **(FEL)** Pre-Processing of model
 282 inputs from Feldman et al. [10] which seeks to reduce disparate impact across all thresholds. The
 283 "amount" of pre-processing is parametrized by a λ similar to ours (just over inputs) – we search for
 284 the optimal λ for each metric we compare against. We abbreviate geometric repair with **(GR)**.

285 **Results.** In Table 1, as denoted by the bolded cells in the U_γ and *Worst Case* columns, our method
 286 outperforms almost all baselines on both the Adult Income-Sex and Adult Income-Race tasks
 287 datasets, for both TPR and EO. The one exception is for $\gamma = \text{EO}$ on the Income-Sex task, however
 288 our method still attains a reduction in all-threshold disparity, and preserves significant accuracy.
 289 For the AUC column, we italicize the cell which has AUC closest to that of the original regressor;
 290 for both metrics and datasets, our method was superior to the baselines in this aspect. We show
 291 illustrate the effect of geometric repair at every threshold in Figure 2. For the For $\gamma = \text{PR}$ (left) we
 292 show the *full* correction $\lambda = 1$. For $\gamma = \text{TPR}$ (middle) we the computed optimal repair parameter
 293 $\lambda_* \approx 0.73 \pm 0.04$, and for $\gamma = \text{EO}$ (right) we computed $\lambda_* \approx 0.75 \pm 0.03$.

294 7 Discussion and Related Work

295 In this work, we show that by interpolating between group-conditional score distributions we can
 296 achieve all-threshold fairness on fairness metrics like Equal Opportunity and Equalized Odds. To this
 297 end, we introduce Distributional parity to measures parity in a fairness metric at all thresholds, and
 298 provide a novel post-processing algorithm that 1) is theoretically-grounded by our convexity result,
 299 and 2) performs extremely well across benchmark datasets and tasks.

300 A number of prior works have demonstrated how to achieve exact distributional parity in the special
 301 case when $\gamma = \text{PR}$. Our work is most closely related to [13] who accomplish this using the \mathcal{W}_1
 302 distance, in both in/post processing settings. [6, 17] report a similar post-processing result to ours,
 303 deriving an optimal fair predictor (also limited to $\gamma = \text{PR}$) in a regression setting and using \mathcal{W}_2
 304 barycenters. We build on these approaches by extending them to a broader class of fairness metrics
 305 and definitions. Our technique is based on the *geometric repair* algorithm which was as originally
 306 introduced by [10] as a way to navigate the fairness-accuracy trade-off. Geometric repair was also
 307 studied by [11]. In the post-processing setting, the effect of geometric repair on classifier accuracy
 308 and $\gamma = \text{PR}$ fairness was studied in [7] – we extend these to all $\gamma \in \Gamma$ by showing convexity on the
 309 set of regressors enumerated by geometric repair.

References

- 310
- 311 [1] M. Agueh and G. Carlier. Barycenters in the Wasserstein Space. *SIAM Journal on Mathematical*
312 *Analysis*, 43(2):904–924, 2011.
- 313 [2] S. Barocas, M. Hardt, and A. Narayanan. *Fairness and Machine Learning*. fairmlbook.org,
314 2019. <http://www.fairmlbook.org>.
- 315 [3] R. P. Brent. *Algorithms for minimization without derivatives*. Courier Corporation, 2013.
- 316 [4] T. Calders, F. Kamiran, and M. Pechenizkiy. Building Classifiers with Independency Constraints.
317 In *2009 IEEE International Conference on Data Mining Workshops*, pages 13–18, 2009.
- 318 [5] A. Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction
319 instruments, 2016.
- 320 [6] E. Chzhen, C. Denis, M. Hebiri, L. Oneto, and M. Pontil. Fair regression with Wasserstein
321 barycenters. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances*
322 *in Neural Information Processing Systems*, volume 33, pages 7321–7331. Curran Associates,
323 Inc., 2020.
- 324 [7] E. Chzhen and N. Schreuder. A minimax framework for quantifying risk-fairness trade-off in
325 regression. *The Annals of Statistics*, 2022.
- 326 [8] F. Ding, M. Hardt, J. Miller, and L. Schmidt. Retiring adult: New datasets for fair machine
327 learning. *arXiv preprint arXiv:2108.04884*, 2021.
- 328 [9] D. Dua and C. Graff. UCI Machine Learning Repository, 2017.
- 329 [10] M. Feldman, S. A. Friedler, J. Moeller, C. Scheidegger, and S. Venkatasubramanian. Certifying
330 and Removing Disparate Impact. In *Proceedings of the 21th ACM SIGKDD International*
331 *Conference on Knowledge Discovery and Data Mining*, KDD '15, pages 259–268, 2015.
- 332 [11] P. Gordaliza, E. D. Barrio, G. Fabrice, and J.-M. Loubes. Obtaining Fairness using Optimal
333 Transport Theory. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th*
334 *International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning*
335 *Research*, pages 2357–2365. PMLR, 09–15 Jun 2019.
- 336 [12] M. Hardt, E. Price, and N. Srebro. Equality of Opportunity in Supervised Learning. In
337 *Proceedings of the 30th International Conference on Neural Information Processing Systems*,
338 *NeurIPS '16*, pages 3323–3331, 2016.
- 339 [13] R. Jiang, A. Pacchiano, T. Stepleton, H. Jiang, and S. Chiappa. Wasserstein Fair Classification.
340 In R. P. Adams and V. Gogate, editors, *Proceedings of The 35th Uncertainty in Artificial*
341 *Intelligence Conference*, volume 115 of *Proceedings of Machine Learning Research*, pages
342 862–872. PMLR, 22–25 Jul 2020.
- 343 [14] N. Kallus and A. Zhou. The fairness of risk scores beyond classification: Bipartite ranking and
344 the xauc metric, 2019.
- 345 [15] J. Kleinberg. Inherent Trade-Offs in Algorithmic Fairness. *SIGMETRICS Perform. Eval. Rev.*,
346 46(1):40, jun 2018.
- 347 [16] T. Le Gouic and J.-M. Loubes. Existence and consistency of wasserstein barycenters. *Probability*
348 *Theory and Related Fields*, 168(3):901–917, 2017.
- 349 [17] T. Le Gouic, J.-M. Loubes, and P. Rigollet. Projection to Fairness in Statistical Learning, 2020.
350 arXiv preprint.
- 351 [18] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
352 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
353 M. Perrot, and E. Duchesnay. Scikit-Learn: Machine Learning in Python. *J. Mach. Learn. Res.*,
354 12:2825–2830, nov 2011.
- 355 [19] G. Peyré and M. Cuturi. Computational optimal transport. 2018.

- 356 [20] G. Pleiss, M. Raghavan, F. Wu, J. Kleinberg, and K. Q. Weinberger. On fairness and calibration.
357 *Advances in Neural Information Processing Systems*, 30:5680–5689, 2017.
- 358 [21] F. Santambrogio. Optimal Transport for Applied Mathematicians. *Birkäuser, NY*, 55(58-63):94,
359 2015.
- 360 [22] V. S. Varadarajan. On the convergence of sample probability distributions. *Sankhyā: The Indian*
361 *Journal of Statistics (1933-1960)*, 19(1/2):23–26, 1958.
- 362 [23] C. Villani. Optimal transport: Old and new. 2008.

363 A Additional Background on Optimal Transport

364 In this section of the appendix, we present some additional background and theory from Optimal
365 Transport. These results are necessary to prove some of the results in the main paper body.

366 A.1 Wasserstein Geodesics

367 One key property of Wasserstein Barycenters that we exploit in this work, which is not refereed to
368 explicitly in the main paper body, is that Wasserstein Barycenters under Geometric Repair form a
369 curve in the space of probability measures called a constant speed geodesic.

370 **Definition A.1** ([21], pg. 182). *Let (X, d) be some metric space. A curve $\omega : [0, 1] \rightarrow X$ is a
371 constant speed geodesic between $\omega(0)$ and $\omega(1)$ if it satisfies*

$$d(\omega(t), \omega(s)) = |t - s|d(\omega(0), \omega(1)) \quad \forall t, s \in [0, 1]$$

372 The following result from [21, Theorem 5.27] proves that a specific interpolation of an optimal
373 transport plan forms a geodesic in the space of probability measures metricized by the Wasserstein
374 Distance.¹ We also remind the reader that in the following expression, $\#$ denotes the pushforward
375 operator on measures and id denotes the identity function.²

376 **Theorem A.1.** *Suppose that Ω is convex, take $\mu, \nu \in \mathcal{P}_p(\Omega)$, and $\gamma \in \Gamma(\mu, \nu)$ an optimal transport
377 plan for the cost $c(x, y) = |x - y|^p$ w/ ($p \geq 1$). Define $\pi_t : \Omega \times \Omega \rightarrow \Omega$ through $\pi_t(x, y) =$
378 $(1 - t)x + ty$. Then the curve $\mu_t(\pi_t)_{\#}\gamma$ is a constant speed geodesic connecting $\mu_0 = \mu$ to
379 $\mu_1 = \nu$. In the particular case where μ is absolutely continuous then this curve is obtained as
380 $((1 - t)\text{id} + tT)_{\#}\mu$*

381 For $p = 2$, this special form of the interpolation between measures given in the above theorem is
382 actually the exact same interpolation that is carried out by Wasserstein Barycenters.³

383 **Proposition A.1** ([1]). *Let $\mu, \nu \in \mathcal{P}_2([0, 1])$ satisfy Assumption 2.1 then α -weighted barycenters*

$$\mu_\alpha \leftarrow \arg \min_{\in \mathcal{P}_p([0, 1])} (1 - \alpha)\mathcal{W}_2^2(\mu, \cdot) + \alpha\mathcal{W}_2^2(\nu, \cdot),$$

384 *can be equivalently computed $((1 - t)\text{id} + tT)_{\#}\mu$ where T is the transport plan that solves transport
385 from $\mu \rightarrow \nu$.*

386 This means, under our mild assumptions, that barycenters both (a) follow the special form in Theorem
387 A.1 and (b) are constant speed geodesics. We use this fact to show that the distance between λ
388 repaired measures $\mu_{a,\lambda}, \mu_{b,\lambda}$ can be written as a $1 - \lambda$ weighted fraction of the Wasserstein distance
389 of the unrepaired measures μ_a, μ_b .

390 **Proposition A.2.** *Since $\mu_{g,\lambda}$ is a constant speed geodesic, the Wasserstein distance between repaired
391 measures is proportional to the repair amount, i.e., $\mathcal{W}_1(\mu_{a,\lambda}, \mu_{b,\lambda}) = (1 - \lambda)\mathcal{W}_1(\mu_a, \mu_b)$.*

392 *Proof.* Let $\mu_a, \mu_b \in \mathcal{P}_2$ and T_a^b be the optimal transport plan from $\mu_a \rightarrow \mu_b$. Suppose we parametrize
393 the interpolation from μ_a to μ_b with a function $w : [0, 1] \rightarrow \mathcal{P}_1([0, 1])$ where $w(\alpha) = ((1 - \alpha)\text{id} +$
394 $\alpha T_a^b)_{\#}\mu_a$. By Theorem A.1, this curve is a constant speed geodesic. Now, consider the geometric
395 repair score distributions $\mu_{a,\lambda}$ and $\mu_{b,\lambda}$. We see from Proposition A.2 that each distribution $\mu_{g,\lambda}$ is the
396 result of the λ -weighted interpolation of μ_g to the barycenter μ_* . These barycenters can alternatively
397 computed by interpolating from $\mu_a \rightarrow \mu_b$, i.e.,

$$\begin{aligned} \mu_{a,\lambda} &= ((1 - \lambda\rho_b)\text{id} + \lambda\rho_b T_a^b)_{\#}\mu_a \\ \mu_{b,\lambda} &= (\lambda\rho_a \text{id} + (1 - \lambda\rho_a)T_a^b)_{\#}\mu_a. \end{aligned}$$

398 From this, a reparametrization of the above interpolation under geometric using $w(\cdot)$ yields,

$$\mu_{a,\lambda} = w(\lambda\rho_b) \quad \text{and} \quad \mu_{b,\lambda} = w(1 - \lambda\rho_a).$$

¹We can metricize \mathcal{P}_p with \mathcal{W}_p under [23, Theorem 6.9]

²In this section, we use a sub-scripted μ_t to denote a measure that is the result of some interpolation when clear from context; this subscript notation should not be confused with the sub-scripts used on measures, e.g. μ_2 , in other places in the paper.

³This result is stated in [1] as the conclusion of Section 4 (see eq. 4.10) and in Section 6.2 of the same work

399 Then, the corollary follows from the definition of constant speed geodesics in Definition A.1, i.e.,

$$\mathcal{W}_1(\mu_{a,\lambda}, \mu_{b,\lambda}) = \mathcal{W}_1(w(\lambda\rho_b), w(1-\lambda\rho_a)) \quad (11)$$

$$= |\lambda\rho_b - (1-\lambda\rho_a)|\mathcal{W}_1(\mu_a, \mu_b) \quad (12)$$

$$= |\lambda(\underbrace{\rho_a + \rho_b}_{=1 \text{ by Def}}) - 1|\mathcal{W}_1(\mu_a, \mu_b) \quad (13)$$

$$= (1-\lambda)\mathcal{W}_1(\mu_a, \mu_b) \quad (14)$$

400

□

401 The last additional result we'll need to aid our effort to prove Theorem 4.1, is the following Lemma.
 402 Please note this lemma differs from the above corollary due to the specific optimal transport problems
 403 being solved. In the above, we are consider a parametrization of $\mu_{g,\lambda}$ along the interpolation from
 404 $\mu_a \rightarrow \mu_b$. In the below result we consider the interpolation of repaired distributions to their barycenter,
 405 i.e., $\mu_{a,\lambda} \rightarrow \mu_*$.

406 **Lemma A.1.** *Let $\mu_a, \mu_b \in \mathcal{P}_2([0, 1])$ satisfy Assumption 2.1 and let $\mu_{a,\lambda}$ be the λ -barycenter of μ_a
 407 and μ_* , and let $\mu_{b,\lambda}$ be the λ -barycenter of μ_b and μ_* then*

$$\mu_{a,\lambda} = \mu_{b, \frac{1-\rho_a\lambda}{1-\rho_a}}$$

$$\mu_{b,\lambda} = \mu_{a, \frac{1-\lambda}{\rho_a} + \lambda}$$

408 *Proof.* Let μ_* be the ρ_b barycenter of μ_a, μ_b . It is easy to show from their definitions that $\mu_{a,\lambda_1} =$
 409 $\mu_{\lambda_1(1-\rho_a)}$ and $\mu_{b,\lambda_2} = \mu_{1-\lambda_2\rho_a}$ (Figure 1 provides a nice illustration of this fact). To prove the
 410 Lemma, we let $\lambda_1(1-\rho_a) = 1-\lambda_2\rho_a$. Solving for λ_1 , yields the proposition, i.e., $\lambda_1 = \frac{1-\lambda_2\rho_a}{\rho_b}$
 411 and therefore $\mu_{a,\lambda_1} = \mu_{b, \frac{1-\rho_a\lambda_2}{\rho_b}}$. Letting $\lambda_1 = \lambda_2$, such that both μ_a, μ_b are controlled by the same
 412 repair parameter yields the first equality. Solving for λ_2 and making the same substitution ($\lambda_2 = \lambda_1$)
 413 yields the second equality. □

414 A.2 The Relationship Between Fair Risk Minimization and Barycenters

415 In this subsection we give an additional result relating the lowest risk $\gamma = \text{PR}$ regressor to the distance
 416 of that regressors groupwise score distributions, to their barycenter.

417 **Lemma A.2.** *Let $\mathcal{F}_{PR} \subset \mathcal{F}$ be a subset of regressors where $\mathcal{F}_{PR} = \{f \in \mathcal{F} : \mathcal{U}_{PR}(f) = 0\}$. The
 418 minimum risk in \mathcal{F}_{PR} is defined*

$$\min_{\hat{f} \in \mathcal{F}_{PR}} \mathcal{R}(\hat{f}) = \min_{\nu \in \mathcal{P}_1([0,1])} \sum_{g \in \mathcal{G}} p_g \mathcal{W}_1(\mu_g, \nu)$$

419 *Proof.* Suppose h is the regressor which minimizes the l.h.s. and let $\mu_h = \text{Law}(h(\mathbf{X}, \mathbf{G}))$. We can
 420 re-write

$$\begin{aligned} \sum_{g \in \mathcal{G}} p_g \mathcal{W}_1(\mu_g, \mu_h) &= \sum_{g \in \mathcal{G}} p_g \min_{T \in \mathcal{T}_g^h} \int_{[0,1]} |x - T(x)| d\mu_g \\ &= \sum_{g \in \mathcal{G}} p_g \min_{T \in \mathcal{T}_g^h} \int_{\mathbf{X}} |f(\mathbf{X}, g) - T(f(\mathbf{X}, g))| d\mu_{\mathbf{X}|g} \end{aligned}$$

421 Let $T_g^{\hat{h}} = F_{\mu_{\hat{h}}} \circ F_{\mu_g}^{-1}$ be the optimal transport maps which minimize the above, and let $\hat{h}(x, g) =$
 422 $T_g^{\hat{h}}(f(x, g))$. We can continue

$$\begin{aligned} \sum_{g \in \mathcal{G}} p_g \min_{T \in \mathcal{T}_g^{\hat{h}}} \int_{\mathbf{X}} |f(\mathbf{X}, g) - \hat{h}(\mathbf{X}, g)| d\mu_{\mathbf{X}|g} &= \mathbb{E}_{g \sim \mathbf{G}} \left[\mathbb{E}_{\mathbf{X}} [|f(\mathbf{X}, g) - \hat{h}(\mathbf{X}, g)| | \mathbf{G} = g] \right] \\ &= \mathbb{E}_{\mathbf{X}, \mathbf{G}} [|f(\mathbf{X}, \mathbf{G}) - \hat{h}(\mathbf{X}, \mathbf{G})|]. \end{aligned}$$

423 From the above equalities we've shown,

$$\sum_{g \in \mathcal{G}} p_g \mathcal{W}_1(\mu_a, \mu_h) = \mathbb{E}_{\mathbf{X}, \mathbf{G}} [|f(\mathbf{X}, \mathbf{G}) - \hat{h}(\mathbf{X}, \mathbf{G})|]. \quad (15)$$

424 and by the presumed optimality of h it follows,

$$\mathbb{E}_{\mathbf{X}, \mathbf{G}} [|f(\mathbf{X}, \mathbf{G}) - \hat{h}(\mathbf{X}, \mathbf{G})|] \geq \mathbb{E}_{\mathbf{X}, \mathbf{G}} [|f(\mathbf{X}, \mathbf{G}) - h(\mathbf{X}, \mathbf{G})|]. \quad (16)$$

425 On the other hand suppose T_g^h is an optimal transport plan such that $h(x, g) = T_g^h(f(x, g))$ then, by
426 the optimality of $T_*^{\hat{h}}$ it follows

$$\sum_{g \in \mathcal{G}} p_g \int_{\mathbf{X}} |f(\mathbf{X}, g) - T_g^{\hat{h}}(f(\mathbf{X}, g))| d\mu_{\mathbf{X}|g} \leq \sum_{g \in \mathcal{G}} p_g \int_{\mathbf{X}} |f(\mathbf{X}, g) - T_g^h(f(\mathbf{X}, g))| d\mu_{\mathbf{X}|g}.$$

427 Using similar properties as the above derivations we can re-write this relationship as

$$\mathbb{E}_{\mathbf{X}, \mathbf{G}} [|f(\mathbf{X}, \mathbf{G}) - \hat{h}(\mathbf{X}, \mathbf{G})|] \leq \mathbb{E}_{\mathbf{X}, \mathbf{G}} [|f(\mathbf{X}, \mathbf{G}) - h(\mathbf{X}, \mathbf{G})|]. \quad (17)$$

428 Therefore by Steps (16) and (17) we have

$$\mathbb{E}_{\mathbf{X}, \mathbf{G}} [|f(\mathbf{X}, \mathbf{G}) - \hat{h}(\mathbf{X}, \mathbf{G})|] = \mathbb{E}_{\mathbf{X}, \mathbf{G}} [|f(\mathbf{X}, \mathbf{G}) - h(\mathbf{X}, \mathbf{G})|],$$

429 and combining Step 15 with the above concludes

$$\min_{\nu \in \mathcal{P}_1([0,1])} \sum_{g \in \mathcal{G}} p_g \mathcal{W}_1(\mu_a, \nu) \leq \mathbb{E}_{\mathbf{X}, \mathbf{G}} [|f(\mathbf{X}, \mathbf{G}) - h(\mathbf{X}, \mathbf{G})|], \quad (18)$$

430 where $\mathcal{U}_{\text{PR}}(h) = 0$ by assumption. To prove the other direction, now let

$$\bar{\nu} \leftarrow \arg \min_{\nu \in \mathcal{P}_1([0,1])} \sum_{g \in \mathcal{G}} p_g \mathcal{W}_1(\mu_a, \nu)$$

431 and $T_g^{\bar{\nu}}$ be the optimal transport maps from $\mu_g \rightarrow \bar{\nu}$ and $\bar{h}(x, g) = T_g^{\bar{\nu}}(f(x, g))$. Now, if we consider

$$\sum_{g \in \mathcal{G}} p_g \mathcal{W}_1(\mu_a, \bar{\nu}) = \mathbb{E}_{\mathbf{X}, \mathbf{G}} [|f(\mathbf{X}, \mathbf{G}) - \bar{h}(\mathbf{X}, \mathbf{G})|]$$

432 then we can easily conclude by the assumed optimality of h that,

$$\min_{\nu \in \mathcal{P}_1([0,1])} \sum_{g \in \mathcal{G}} p_g \mathcal{W}_1(\mu_a, \nu) \geq \mathbb{E}_{\mathbf{X}, \mathbf{G}} [|f(\mathbf{X}, \mathbf{G}) - h(\mathbf{X}, \mathbf{G})|]. \quad (19)$$

433 Finally, recalling that \bar{h} satisfies $\mathcal{U}_{\text{PR}}(\bar{h}) = 0$ since \bar{h} is a Barycenter (Corollary 3.1). Combining
434 Steps 18 and 19 to yield the proof. \square

435 B Proofs

436 B.1 Proof of Proposition 3.1

437 **Proposition 3.1.** For $\mu_a, \mu_b \in \mathcal{P}_2([0, 1])$ which are the groupwise score distributions of f , then
438 $\mathcal{W}_2(\mu_a, \mu_b) = 0$ if and only if $\mathcal{U}_{\text{PR}}(f) = 0$.

439 *Proof.* Let μ_a, μ_b be the groupwise score distributions of some regressor f . Since W_p is a metric
440 on $\mathcal{P}_p([0, 1])$ (according to Proposition 2.3 in [19]) if $\mathcal{W}_2(\mu_a, \mu_b) = 0$ then $\mu_a = \mu_b$. Similarly, by the
441 same property we know that $\mathcal{W}_2(\mu_a, \mu_b) = \mathcal{W}_1(\mu_a, \mu_b) = 0$. Showing that $\mathcal{W}_1(\mu_a, \mu_b) = \mathcal{U}_{\text{PR}}(f)$
442 completes the proof. To show this equality, recall by definition that

$$\gamma_g(\tau) = \Pr[f(\mathbf{X}, \mathbf{G}) \geq \tau | \mathbf{G} = g] \quad (20)$$

$$= 1 - \Pr[f(\mathbf{X}, \mathbf{G}) \leq \tau | \mathbf{G} = g] \quad (21)$$

$$= 1 - F_g(\tau) \quad (22)$$

443 Plugging this into the expression for \mathcal{U}_{PR}

$$\mathcal{U}_{\text{PR}}(f) = \mathbb{E}_{\tau \in U([0,1])} |\gamma_a(\tau) - \gamma_b(\tau)| = \int_{[0,1]} |\gamma_a(\tau) - \gamma_b(\tau)|^p d\tau \quad (23)$$

$$= \int_{[0,1]} |(1 - F_a(\tau)) - (1 - F_b(\tau))| d\tau \quad (24)$$

$$= \int_{[0,1]} |F_a(\tau) - F_b(\tau)| d\tau \quad (25)$$

$$= \int_0^1 |F_a^{-1}(t) - F_b^{-1}(t)| dt = \mathcal{W}_1(\mu_a, \mu_b) \quad (26)$$

444 where the second to last equality was proven in Lemma 6 from [13]. \square

445 B.2 Proof of Lemma 3.1

446 **Lemma 3.1.** Let f be some learned regressor. For all $\lambda \in [0, 1]$ the set of optimally fair regressors
447 for λ -relaxations of f with respect to risk and distributional parity for $\gamma = \text{PR}$ are given by

$$f_\lambda \leftarrow \arg \min_{\hat{f} \in \mathcal{F}} \lambda R(\hat{f}) \quad \text{s.t.} \quad \mathcal{U}_{\text{PR}}(\hat{f}) = (1 - \lambda)\mathcal{U}_{\text{PR}}(f) \quad (27)$$

448 *Proof.* By Definition of f^* we know that $\mathcal{R}(f^*)$ is

$$\min_{\hat{f} \in \mathcal{F}} \mathcal{R}(\hat{f}) \quad \text{s.t.} \quad \mathcal{U}_{\text{PR}}(\hat{f}) = 0.$$

449 It follows that

$$\lambda(\min_{\hat{f} \in \mathcal{F}} \mathcal{R}(\hat{f})) = \min_{\hat{f} \in \mathcal{F}} \lambda R(\hat{f}) = \lambda \mathcal{R}(f^*). \quad (28)$$

450 By definition of f_λ it is straightforward to show that $\mathcal{R}(f_\lambda) = \lambda \mathcal{R}(f^*)$. Under Proposition A.2, it
451 is straightforward to show that $\mathcal{U}_{\text{PR}}(f_\lambda) = (1 - \lambda)\mathcal{U}_{\text{PR}}(f)$. Combining these two facts proves the
452 result. \square

453 B.3 Proof of Corollary 2.1

454 **Corollary 2.1.** Let μ_* be the ρ_b -weighted barycenter of μ_a, μ_b then the transport plan from $\mu_a \rightarrow \mu_*$
455 (wlog) is computed

$$T_a^*(\omega) = (\rho_a F_{\mu_a} + \rho_b F_{\mu_b}) \circ F_{\mu_a}^{-1}(\omega)$$

456 *Proof.* Observe that by Theorem A.1 we can express barycenter from μ_a to μ_* (wlog)

$$\mu_* = (\rho_a \text{id} + \rho_b T_a^b) \# \mu_a = (\rho_a F_{\mu_a}^{-1} \circ F_{\mu_a} + \rho_b F_{\mu_b}^{-1} \circ F_{\mu_a}) \# \mu_a$$

457 The second equality follows from Remark 2.1. From this expression, we can define $T_a^* = (\rho_a F_{\mu_a}^{-1} \circ$
458 $+ \rho_b F_{\mu_b}^{-1}) \circ F_{\mu_a}$ as the function which computes the transport from $\mu_a \rightarrow \mu_*$. \square

459 B.4 Proof of Proposition 4.1

460 **Proposition 4.1 .** For any $\lambda \in [0, 1]$, a repaired regressor f_λ satisfies the following

$$R(f_\lambda) = \lambda R(f^*) \quad \text{and} \quad \mathcal{U}_{\text{PR}}(f_\lambda) = (1 - \lambda)\mathcal{U}_{\text{PR}}(f)$$

461 *Proof.* The first equality follows from the definition of R and linearity of expectation. It is easy to
462 show that

$$\begin{aligned} R(f_\lambda) &= R((1 - \lambda)f + \lambda f^*) \\ &= (1 - \lambda)R(f) + \lambda R(f^*) = \lambda R(f^*) \end{aligned}$$

463 where the last equality follows by noting that $R(f) = 0$ by definition. The proof that $\mathcal{U}_{\text{PR}}(f_\lambda) =$
464 $(1 - \lambda)\mathcal{U}_{\text{PR}}(f)$ follows from Proposition A.2. \square

465 **B.5 Proof of Theorem 4.1.**

466 In the proof of Theorem 4.1, we make use of the fact that this transport plans are bijective, under
 467 Assumption 2.1. In order to show that these plans are bijective we show that they are strictly monotone
 468 via the following Remark.

469 **Remark B.1** ([21], p. 55). *For two measures $\mu, \nu \in \mathcal{P}_p([0, 1])$, if ν is non-atomic, then the transport
 470 plan T from $\mu \rightarrow \nu$ is strictly monotone on a closed domain like $[0, 1]$.*

471 It is well known that strictly monotone functions on a closed domain are bijective, and therefore we
 472 claim bijectivity as a corollary of the above result.

473 **Corollary B.1.** *A transport plan T that is strictly monotone, on a closed domain, is also bijective.*

474 Now we begin the proof of Theorem A.1.

475 **Theorem 4.1.** Fix $\gamma \in \Gamma$. Let $f : X \times G \rightarrow [0, 1]$ be a regressor, and f_λ be the geometrically
 476 repaired regressor for any $\lambda \in [0, 1]$. The map $\lambda \mapsto \mathcal{U}_\gamma(f_\lambda)$ is convex in λ .

477 *Proof.* Let $\gamma \in \Gamma$. To prove convexity, we show that $\frac{d^2}{d\lambda^2} \mathcal{U}_\gamma(f_\lambda)$ is non-negative everywhere. First,
 478 we remind readers the definition of $\mathcal{U}_\gamma(f_\lambda)$ (distributional parity):

$$\mathcal{U}_\gamma(f_\lambda) \triangleq \mathbb{E}_{\tau \sim \mathcal{U}([0,1])} |\gamma_a(\tau) - \gamma_b(\tau)|.$$

479 where γ_g is a fairness metric on the score distributions of f_λ for group $g \in G$.

480 Recall the definition of $\gamma_g(\tau; f_\lambda)$

$$\gamma_g(\tau; f_\lambda) = \Pr[f_\lambda(\mathbf{X}, \mathbf{G}) \geq \tau | \mathbf{G} = g].$$

481 by Proposition 4.2 we know that $\mu_{g,\lambda}$ is the score distribution associated with $f_\lambda(\cdot, g)$ and so we
 482 re-write this expression as a conditional expectation

$$\Pr[f_\lambda(\mathbf{X}, \mathbf{G}) \geq \tau | \mathbf{G} = g] = \int_{[0,1]} \mathbb{1}_{[\tau,1]} d\mu_{g,\lambda} \quad (29)$$

483 In order to take this derivative, we need to invoke several change of variables to convert this Lebesgue
 484 integral to a Riemann integral. We'll proceed for $a \in G$ without loss of generality. Also note
 485 for brevity, we present the proof for $\mu_{g,\lambda}$, i.e., the measure associated with $\gamma = \text{PR}$. Similarly, if
 486 we condition the l.h.s. of Eq. 29 on \mathbf{Y} , our results follow similarly for corresponding probability
 487 measures associated with this conditional probability, .e.g, we would let $\mu_{g|\mathbf{Y},\lambda}$ be the measure
 488 associated with the conditional probability $\Pr[f_\lambda(\mathbf{X}, \mathbf{G}) | \mathbf{G} = g, \mathbf{Y} \geq \tau]$ for which setting \mathbf{Y}
 489 computes TPR and FPR respectively.

490 Following Claim A.1 can re-write $\mu_{a,\lambda} := ((1 - \lambda)id + \lambda T_a^*) \# \mu_a$. For notational ease, define
 491 $\pi_{a,\lambda} := (1 - \lambda)id + \lambda T_a^*$. Using these substitutions, we have that $\mu_{a,\lambda} = (\pi_{a,\lambda}) \# \mu_a$, so γ_a can be
 492 equivalently written

$$\gamma_a(\tau) = \int_{[0,1]} \mathbb{1}_{[\tau,1]} d(\pi_{a,\lambda} \# \mu_a).$$

493 By definition of the push-forward operator

$$\int_{[0,1]} \mathbb{1}_{[\tau,1]} d(\pi_{a,\lambda} \# \mu_a) = \int_{\pi_{a,\lambda}^{-1}([0,1])} \mathbb{1}_{[\tau,1]}(\pi_{a,\lambda}) d\mu_a = \int_{[0,1]} \mathbb{1}_{[\tau,1]}(\pi_{a,\lambda}) d\mu_a.$$

494 We note that the domain of integration is unchanged in the last equality because π is a bijective
 495 mapping from $[0, 1] \rightarrow [0, 1]$ by Corollary B.1, and so $\pi_{a,\lambda}^{-1}([0, 1]) = [0, 1]$.

496 For the last change of variables, Let ℓ be the Lebesgue measure. By Assumption 2.1 μ_a is absolutely
 497 continuous with respect to ℓ meaning that by the Radon Nikodym-Theorem

$$\int_{[0,1]} \mathbb{1}_{[\tau,1]}(\pi_{g,\lambda}) d\mu_a = \int_{[0,1]} \sigma_a \mathbb{1}_{[\tau,1]}(\pi_{g,\lambda}) d\ell$$

498 where σ_a is the Radon-Nikodym Derivative, i.e., the probability density function associated with μ_a .
 499 We'll also need to define γ_b similarly. To do this we invoke Lemma A.1 which yields that $\mu_{b,\lambda} =$
 500 $\mu_a, \frac{1-\lambda}{\rho_a} + \lambda$. Using this substitution we get

$$\gamma_b(\tau) = \int_{[0,1]} \rho_{\mu_a} \mathbb{1}_{[\tau,1]}(\pi_{b, \frac{1-\lambda}{\rho_a} + \lambda}) d\ell$$

501 Next, let $h_{a,\tau}(\lambda)$ be the mapping $\lambda \mapsto \mathbb{1}_{[\tau,1]}(\mu_{a,\lambda})$ and $h_{b,\tau}(\lambda)$ be $\lambda \mapsto \mathbb{1}_{[\tau,1]}(\mu_a, \frac{1-\lambda}{\rho_a} + \lambda)$. Taking
 502 the first derivative of this difference, we get

$$\begin{aligned} \frac{d}{d\lambda} [h_{a,\tau}(\lambda) - h_{b,\tau}(\lambda)] &= \\ \frac{d}{d\lambda} \int_{[0,1]} \rho_{\mu_a} \cdot [\mathbb{1}_{[\tau,1]}(\pi_{a,\lambda}) - \mathbb{1}_{[\tau,1]}(\pi_{b, \frac{1-\lambda}{\rho_a} + \lambda})] d\ell &= \\ \int_{[0,1]} \rho_{\mu_a} \cdot \left[\frac{d}{d\lambda} (\mathbb{1}_{[\tau,1]}(\pi_{a,\lambda}) - \mathbb{1}_{[\tau,1]}(\pi_{b, \frac{1-\lambda}{\rho_a} + \lambda})) \right] d\ell \end{aligned}$$

503 where the second equality follows from Leibniz Rule. To finish the derivative, we remind the reader
 504 that the derivative of $\frac{d}{d\lambda} \mathbb{1}_{[\tau,1]}(\pi_{g,\lambda})$ is the Dirac delta function $\delta(\pi_{g,\lambda} - \tau)$. It follows that

$$\begin{aligned} \int_{[0,1]} \rho_{\mu_a} \cdot \left[\frac{d}{d\lambda} (\mathbb{1}_{[\tau,1]}(\pi_{a,\lambda}) - \mathbb{1}_{[\tau,1]}(\pi_{b, \frac{1-\lambda}{\rho_a} + \lambda})) \right] d\ell &= \int_{[0,1]} \rho_{\mu_a} \cdot \left[T_a^* (\delta(\pi_{a,\lambda} - \tau)) + \left(\frac{1-\rho_a}{\rho_a} \right) \delta(\pi_{b, \frac{1-\lambda}{\rho_a} + \lambda} - \tau) T_b^* \right. \\ &\quad \left. - \left(\frac{\rho_a - 1}{\rho_a} \right) \delta(\pi_{b, \frac{1-\lambda}{\rho_a} + \lambda} - \tau) \text{id} - \delta(\pi_{a,\lambda} - \tau) \text{id} \right] \end{aligned}$$

505 and by definition of δ of the delta function, we at last obtain

$$\begin{aligned} \int_{[0,1]} \rho_{\mu_a} \cdot \left[T_a^* (\delta(\pi_{a,\lambda} - \tau)) + \left(\frac{1-\rho_a}{\rho_a} \right) \delta(\pi_{b, \frac{1-\lambda}{\rho_a} + \lambda} - \tau) T_b^* - \left(\frac{\rho_a - 1}{\rho_a} \right) \delta(\pi_{b, \frac{1-\lambda}{\rho_a} + \lambda} - \tau) \text{id} - \delta(\pi_{a,\lambda} - \tau) \text{id} \right] \\ = \left[T_a^* + \left(\frac{1-\rho_a}{\rho_a} \right) T_b^* - \left(\frac{\rho_a - 1}{\rho_a} \right) \text{id} - \text{id} \right] \circ \tau. \end{aligned}$$

506 To summarize, we have just shown that

$$\frac{d}{d\lambda} [h_{a,\tau}(\lambda) - h_{b,\tau}(\lambda)] = \left[T_a^* + \left(\frac{1-\rho_a}{\rho_a} \right) T_b^* - \left(\frac{\rho_a - 1}{\rho_a} \right) \text{id} - \text{id} \right] \circ \tau.$$

507 To prove convexity we must also compute the second derivative of the above. Since the above does
 508 not depend on λ , taking another derivative yields

$$\frac{d^2}{d\lambda^2} [h_{a,\tau}(\lambda) - h_{b,\tau}(\lambda)] = 0. \quad (30)$$

509 Now, to prove the convexity of $\mathcal{U}_\gamma(f_\lambda)$ we take the second derivative of the absolute value of this
 510 difference, i.e.,

$$\frac{d}{d^2\lambda} |h_{a,\tau} - h_{b,\tau}| = \text{sign}(h_{a,\tau} - h_{b,\tau}) \underbrace{\frac{d^2}{d\lambda^2} [h_{a,\tau}(\lambda) - h_{b,\tau}(\lambda)]}_{=0} \quad (31)$$

$$+ 2 \underbrace{\delta(h_{a,\tau} - h_{b,\tau})}_{\simeq 0 \text{ or } 1} \underbrace{\left(\frac{d}{d\lambda} [h_{a,\tau}(\lambda) - h_{b,\tau}(\lambda)] \right)^2}_{\geq 0}. \quad (32)$$

511 The first term on the r.h.s., we've already shown is zero, and the second term is also non-negative.
 512 Another application of Leibniz' Rule allows that

$$\underbrace{\frac{d}{d^2\lambda} \mathbb{E}_{\tau \sim U([0,1])} |h_{a,\tau} - h_{b,\tau}|}_{\mathcal{U}_\gamma(f_\lambda)} = \mathbb{E}_{\tau \sim U([0,1])} \left| \underbrace{\frac{d}{d^2\lambda} [h_{a,\tau} - h_{b,\tau}]}_{\geq 0 \text{ by (31)}} \right|.$$

513 This indicates that $\mathcal{U}_\gamma(f_\lambda)$ is convex (i.e., we have shown that the second derivative is non-negative).
 514 \square

515 **B.6 Proof of Proposition 4.2**

516 **Proposition 4.2.** Let $\lambda \in [0, 1]$. Let $\mu_{g,\lambda}$ be the λ -weighted barycenter between μ_g and μ_* , i.e.,

$$\mu_{g,\lambda} \leftarrow \arg \min_{\nu \in \mathcal{P}_2([0,1])} (1 - \lambda)\mathcal{W}_2^2(\mu_g, \nu) + \lambda\mathcal{W}_2^2(\mu_*, \nu), \quad \text{then } \mu_{g,\lambda} = \text{Law}(f_\lambda(\mathbf{X}, \mathbf{G}) | \mathbf{G} = g).$$

517

518 *Proof.* First we recall the definition of geometric repair

$$f_\lambda(x, g) \triangleq (1 - \lambda)f(x, g) + \lambda f^*(x, g).$$

519 It is easy to show that for T_g^* we have that (wlog) $f^*(x, a) = T_g^*(f(x, a))$

$$f^*(x, a) = (\rho_a \text{id} + \rho_b T_a^b) \circ F_{\mu_a}(f(x, a)) = \tag{33}$$

$$F_{\mu_*}^{-1}(F_{\mu_a}(f(x, a))) = \tag{34}$$

$$T_g^*(f(x, a)). \tag{35}$$

520 Where the first equality follows from Theorem 2.3. in [6], and the second equality is the definition of
521 μ_* . Using this equality in the definition of geometric repair we get

$$f_\lambda(x, g) = (1 - \lambda)f(x, g) + \lambda T_g^*(f(x, g)) \tag{36}$$

$$= ((1 - \lambda)\text{id} + \lambda T_g^*) \circ f(x, g) \tag{37}$$

522 If we let μ_g be the groupwise score distribution for group g then we know $\mu_g = \text{Law}(f(\mathbf{X}, \mathbf{G}) | \mathbf{G} =$
523 $g)$ by definition. If we pushforward μ_g using $((1 - \lambda)\text{id} + \lambda T_g^*)$, i.e.,

$$((1 - \lambda)\text{id} + \lambda T_g^*) \# \mu_g = \arg \min_{\nu \in \mathcal{P}_2([0,1])} (1 - \lambda)\mathcal{W}_2^2(\mu_g, \nu) + \lambda\mathcal{W}_2^2(\mu_*, \nu)$$

524 by Claim A.1 and the uniqueness of \mathcal{W}_2 barycenters. Noticing the $\mu_{g,\lambda}$ is the score distribution for
525 $f_\lambda(\mathbf{X}, \mathbf{G}) | \mathbf{G} = g$ completes the proof. \square

526 **B.7 Proof of Theorem 4.2**

527 **Theorem 4.2.** For all $\lambda \in [0, 1]$, the repaired regressor f_λ is pareto optimal in the multi-objective
528 minimization of $\mathcal{R}(\cdot)$ and $\mathcal{U}_{\text{PR}}(\cdot)$

529 *Proof.* It is clear from the definition of f_λ that $\{f_\lambda\}_{\lambda \in [0,1]}$ forms a pareto front. Indeed, recall that
530 for any level of unfairness, say $\lambda \mathcal{U}_{\text{PR}}(f^*)$, that f_λ is the regressor which minimizes risk, i.e.,

$$f_\lambda \leftarrow \arg \min_{\hat{f} \in \mathcal{F}} \lambda R(\hat{f}) \quad \text{s.t. } \mathcal{U}_{\text{PR}}(f_\lambda) = (1 - \lambda)\mathcal{U}_{\text{PR}}(f).$$

531 Due to the above, it is easy to see that no classifier can have risk less than f_λ , without decreasing
532 λ , which in turn increase $\mathcal{U}(\cdot)$, proving the pareto optimality of f_λ . Now, suppose for contradiction,
533 $\{f_\lambda\}_{\lambda \in [0,1]}$ did not form a pareto front, i.e., there exists some $h \notin \{f_\lambda\}_{\lambda \in [0,1]}$ such that $h \succ f_\lambda$
534 for some $\lambda \in [0, 1]$. Since $h \succ f_\lambda$ then clearly (WLOG) $\mathcal{R}(h) < \mathcal{R}(f_\lambda)$. However if we select
535 $\lambda_h = \frac{\mathcal{R}(h)}{\mathcal{R}(f_\lambda)}$ then $\mathcal{R}(f_{\lambda_h}) = \mathcal{R}(h)$ and subsequently $\mathcal{U}_{\text{PR}}(f_{\lambda_h}) = \mathcal{U}_{\text{PR}}(h)$, which by definition means
536 $h \in \{f_\lambda\}_{\lambda \in [0,1]}$. In the other case where $\mathcal{U}(h) < \mathcal{U}(f_\lambda)$ the proof follows identically. In both cases,
537 we arrive at a contradiction indicating that $\{f_\lambda\}_{\lambda \in [0,1]}$ is indeed a Pareto Frontier. \square

538 **B.8 Proof of Theorem 5.1**

539 **Theorem 5.1.** As $n_g \rightarrow \infty$ the empirical distribution of $\hat{f}_\lambda(x, g)$ converges to $\mu_{g,\lambda}$ in \mathcal{W}_2 almost
540 surely.

541 *Proof.* To complete this proof, it will be convenient to consider the following mixture distribution

$$\mathcal{P} = \sum_{g \in \mathcal{G}} \rho_g \delta_{\mu_g}$$

542 and its empirical variant $\hat{\mathcal{P}}$ using $\hat{\rho}_g$ and $\hat{\mu}_g$. Relying on the barycenters uniqueness under Assumption
543 2.1 in \mathcal{W}_2 (proven by [1]) and the consistency of the Wasserstein barycenter [16][Theorem 3], proving
544 that $\hat{\mathcal{P}} \rightarrow \mathcal{P}$ in the Wasserstein Distance is sufficient to prove the convergence $\hat{\mu}_{g,\lambda}$.

545 We now begin the proof. Recall that we can express $\mu_{g,\lambda}$ as λ -weighted barycenter between μ_g, μ_*
546 or as a $\lambda\rho_b$ weighted barycenter between μ_a and μ_b . Consider the latter formulation, i.e.

$$\mathcal{P}_\lambda = (1 - \lambda)\rho_b\delta_{\mu_a} + \lambda\rho_b\delta_{\mu_b}$$

547 Thus via the consistency of Wasserstein barycenters, as stated above, we must only show that $\hat{\rho}_g$
548 converges to ρ_g , and that $\hat{\mu}_g \rightarrow \mu_g$ in \mathcal{W}_2 . The convergence of $\hat{\rho}_g$ follows by the law of large
549 numbers. The convergence of $\hat{\mu}_g$ follows from the well known facts that the Wasserstein Distance
550 metrizes the weak convergence of probability measures [23, Theorem 6.9], and that an empirical
551 measure $\hat{\mu}_k \rightarrow \mu$ almost surely, [22]. From these facts it follows that $\mathcal{W}_2(\hat{\mu}_g, \mu_g) \rightarrow 0$ almost surely,
552 completing the proof. \square