



**HAL**  
open science

# Demographic parity in regression and classification within the unawareness framework

Vincent Divol, Solenne Gaucher

► **To cite this version:**

Vincent Divol, Solenne Gaucher. Demographic parity in regression and classification within the unawareness framework. 2024. hal-04684789

**HAL Id: hal-04684789**

**<https://hal.science/hal-04684789v1>**

Preprint submitted on 3 Sep 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Demographic parity in regression and classification within the unawareness framework

Vincent Divol <sup>\*1</sup> and Solenne Gaucher <sup>†1</sup>

<sup>1</sup>CREST, ENSAE, IP Paris

September 3, 2024

## Abstract

This paper explores the theoretical foundations of fair regression under the constraint of demographic parity within the unawareness framework, where disparate treatment is prohibited, extending existing results where such treatment is permitted. Specifically, we aim to characterize the optimal fair regression function when minimizing the quadratic loss. Our results reveal that this function is given by the solution to a barycenter problem with optimal transport costs. Additionally, we study the connection between optimal fair cost-sensitive classification, and optimal fair regression. We demonstrate that nestedness of the decision sets of the classifiers is both necessary and sufficient to establish a form of equivalence between classification and regression. Under this nestedness assumption, the optimal classifiers can be derived by applying thresholds to the optimal fair regression function; conversely, the optimal fair regression function is characterized by the family of cost-sensitive classifiers.

**Keywords**— Statistical fairness, demographic parity, optimal transport, unawareness framework

## 1 Introduction

### 1.1 Motivation

Recent breakthroughs in artificial intelligence have led to the widespread adoption of machine learning algorithms, exerting an increasingly influential and insidious impact on our lives. Essentially, these algorithms learn to detect and reproduce patterns using massive datasets. It is now widely recognized that these predictions carry the risk of perpetuating, or even exacerbating, the social discriminations and biases often present in these datasets [ALMK16, BHN23]. Algorithmic fairness seeks to measure and mitigate the unfair impact of algorithms; we refer the reader to the reviews by [BHN23, dBGL20, OC20] for an introduction.

Different approaches have been developed to mitigate algorithmic unfairness. One approach focuses on *individual fairness*, ensuring that similar individuals are treated similarly, regardless of potentially discriminatory factors. Another approach targets *group fairness*, aiming to prevent algorithmic predictions from discriminating against groups of individuals. *Statistical fairness* falls under the latter approach and relies on the formalism of supervised learning to impose fairness criteria while minimizing a risk measure. In this work, we study risk minimization under the demographic parity criterion, which requires that predictions be statistically independent of sensitive attributes. Although this criterion, introduced by [CKP09, ADW19], has some known limitations [HPPS16, ZVGRG19], it finds application in a wide range of scenarios [MZP21, DEHH24]. Its simplicity arguably makes it the most extensively studied criterion.

---

\*vincent.divol@ensae.fr

†solenne.gaucher@ensae.fr

The statistical fairness literature can be broadly divided into two currents, depending on whether the direct use of the protected attribute in predictions is permitted or not. A first line of works, studying the *awareness framework*, considers regression functions that make explicit use of discriminating attributes, thus treating individuals differently based on discriminating factors. For this reason, this approach is also often referred to as *disparate treatment*. In this work, we adopt the *unawareness framework*, in which disparate treatment is prohibited and the regression function cannot directly use the sensitive attribute. Empirical evidence from simulations [LMC18] indicates that within the unawareness framework, predictions often result in suboptimal trade-offs between fairness and accuracy and may induce within-group discrimination. Moreover, the authors conjecture that while the unawareness framework aims to prevent discrimination based on sensitive attributes, predictions in this setting implicitly rely on estimates of these attributes—a phenomenon later proven in [GSC23] for classification problems. Nevertheless, this framework remains crucial in practice, as the direct use of sensitive attributes may be legally prohibited or simply unavailable at prediction time.

In this paper, we investigate the problem of fair regression under demographic parity constraints within the unawareness framework. A key difficulty in overcoming algorithmic unfairness is the limited understanding of how fair algorithms make predictions. Therefore, we focus on providing a simple mathematical characterization of the optimal regression function in the presence of fairness constraints.

## 1.2 Problem statement

Let  $(X, S, Y)$  be a tuple in  $\mathcal{X} \times \mathcal{S} \times \mathbb{R}$  with distribution  $\mathbb{P}$ , where  $X$  corresponds to a non-sensitive feature in a feature space  $\mathcal{X}$ ,  $S$  is a sensitive attribute in a finite set  $\mathcal{S}$ , and  $Y$  is a response variable that we want to predict, which has a finite second moment. To illustrate this problem with an example, assume, as in [CS20b], that  $X$  represents a candidate’s skill,  $S$  is an attribute indicating groups of populations, and  $Y$  is the current market salary of the candidate. Due to historical biases, the distribution of the salary may be unbalanced between the groups. Our aim is to make predictions that are fair, and as close as possible to the current market value  $Y$ . In the unawareness framework, we cannot make explicit use of the sensitive attribute to make our predictions. Therefore, we consider regression functions of the form  $f : \mathcal{X} \rightarrow \mathbb{R}$  in the set of score functions  $\mathcal{F}$ . We want to ensure that our regression function satisfies the following demographic parity criterion.

**Definition 1** (Demographic parity). *The function  $f : \mathcal{X} \rightarrow \mathbb{R}$  verifies the Demographic Parity criterion if*

$$f(X) \perp S.$$

In essence, the demographic parity criterion requires that the distribution of predictions (in our example, the salary) be identical across all groups. We assess the quality of a regression function  $f$  through its quadratic risk

$$\mathcal{R}_{sq}(f) = \mathbb{E} \left[ (Y - f(X))^2 \right].$$

**Definition 2** (Fair regression). *An optimal fair regression function  $f^*$  satisfies*

$$f^* \in \arg \min_{f \in \mathcal{F}} \{ \mathcal{R}_{sq}(f) : f(X) \perp S \}, \tag{1}$$

where  $\mathcal{F}$  is the set of regression functions from  $\mathcal{X}$  to  $\mathbb{R}$ .

Classical results show that when no fairness constraints are imposed, the Bayes regression function  $f^{\text{Bayes}}$  minimizing the squared risk  $\mathcal{R}_{sq}$  is a.s. equal to the conditional expectation  $\eta$ , where

$$\eta(x) = \mathbb{E} [Y | X = x].$$

In this paper, we also investigate the relationship between classification and regression problem. When  $Y \in \{0, 1\}$  a.s., the quality of a classification function  $g : \mathcal{X} \rightarrow \{0, 1\}$  can be assessed through its expected weighted 0 – 1 loss  $\mathcal{R}_y(g)$ , where for  $y \in [0, 1]$ ,  $\mathcal{R}_y(g)$  is defined as

$$\mathcal{R}_y(g) = y \cdot \mathbb{P}[Y = 0, g(X) = 1] + (1 - y) \cdot \mathbb{P}[Y = 1, g(X) = 0].$$

For the choice  $y = 1/2$ , minimizing this risk measure corresponds to maximizing the classical accuracy measure.

**Definition 3** (Fair classification). *For a given value  $y \in [0, 1]$ , an optimal fair classification function  $g_y^*$  verifies*

$$g_y^* \in \arg \min_{g \in \mathcal{G}} \{\mathcal{R}_y(g) : g(X) \perp S\}, \quad (2)$$

where  $\mathcal{G}$  is the set of classification functions from  $\mathcal{X}$  to  $\{0, 1\}$ .

Let us again illustrate this problem with an example from recruitment. Assume that  $X$  represents a candidate’s skill,  $S$  is an attribute indicating different population groups, and  $Y$  denotes whether a human recruiter would consider the candidate qualified for a given position. Due to historical biases, the distribution of the binary response  $Y$  may be unbalanced across the groups. We aim to make a prediction, or equivalently take the decision to accept or reject a candidate. Our goal is to make predictions for the value of  $Y$ , or equivalently, to decide whether to accept or reject a candidate, in a way that is both accurate and fair. Specifically, under demographic parity, we aim to ensure that the probability of acceptance is the same across all groups.

Classical results show that when no fairness constraints are imposed, the classifier  $g_y^{\text{Bayes}}(x) = \mathbb{1}\{f^{\text{Bayes}}(x) \geq y\}$  is a Bayes classifier that minimizes  $\mathcal{R}_y(g)$ . This relationship is at the heart of the design and study of plug-in classifiers [Yan99, MN06, AT07, BDL08]. Interestingly, it was recently shown that a similar relationship holds under demographic parity constraints in the awareness framework [GSC23]. Extending this result to the unawareness framework has remained an open problem, which we address in this paper.

**Notation** We first set some notation. Recall that we are given a tuple  $(X, S, Y)$  in  $\mathcal{X} \times \mathcal{S} \times \mathbb{R}$  with distribution  $\mathbb{P}$ , where  $\mathcal{X}$  is any measurable space (the space of features) and  $\mathcal{S}$  is a finite set (the set of labels). For  $s \in \mathcal{S}$ , we denote by  $p_s$  the probability  $\mathbb{P}(S = s)$  and by  $\mu_s$  the conditional law of  $X|S = s$ . We let  $\mu = \sum_{s \in \mathcal{S}} p_s \mu_s$  be the marginal distribution of  $X$ . We let  $\mathcal{P}(\mathcal{X})$  be the set of probability measures on the measurable space  $\mathcal{X}$ . Moreover, we let  $L^1(\nu)$  be the space of functions integrable with respect to the probability measure  $\nu$ . Finally,  $\overset{\circ}{C}$  denotes the interior of the set  $C$ .

### 1.3 Related work

**Fair classification** Research on optimal prediction under demographic parity constraints has primarily focused on classification, where the goal is to predict a binary response in  $\{0, 1\}$ , as this problem is intrinsically linked to the issue of fair candidate selection, central in algorithmic fairness. This problem is well understood in the awareness setting from an algorithmic point of view [FFM+15, MW18, YCK20, SC21, ZDC22b, DEHH24]. On the theoretical side, [GSC23] recently proved that the optimal classifier for the risk  $\mathcal{R}_y$  can be obtained as the indicator that the optimal fair prediction function for the squared loss  $f^*$  is above the threshold  $y$ , a result that was later extended in [XYZ23] to multi-class classification.

Less is known about fair classification in the unawareness framework. On the algorithmic side, several works have proposed various relaxations of the demographic parity constraint, leading to tractable algorithms for computing classifiers [GCGF16, ZVGRG19, ODP20]. On the theoretical side, [LMC18] provided empirical evidence suggesting that fair classifiers may base their decisions on non-relevant features correlated with the sensitive attribute, potentially disrupting within-group ordering. This hypothesis was further confirmed by [GSC23], who characterized the optimal fair classifier in the unawareness framework. They showed that it is given by the indicator that the conditional expectation  $\eta(X)$  is above a threshold, which depends on the probabilities that the individual described by  $X$  belongs to the different groups. Notably, the question of whether this classifier can be obtained by thresholding the optimal fair prediction function for the squared loss remains an open problem.

**Fair regression** In the awareness framework, fair regression is well understood from both the algorithmic and theoretical points of view [CS20b, GLR20, CDH+20a]. On the theoretical front, it has been shown that the problem of fair regression under demographic parity can be rephrased as the problem of finding the

weighted barycenter of the distributions of  $\eta(X, S) = \mathbb{E}[Y|X, S]$  across different groups, with costs given by optimal transport problems.

**Theorem 1** ([CS20b, GLR20]). *Assume that for all  $s \in \mathcal{S}$ , the distribution  $\nu_s$  of  $\eta(X, S)$  for  $S = s$  has no atoms, and let  $p_s = \mathbb{P}(S = s)$ . Then,*

$$\min_{f \text{ is fair}} \mathcal{R}_{sq}(f) = \min_{\nu \in \mathcal{P}(\mathbb{R})} \sum_{s \in \mathcal{S}} p_s \mathcal{W}_2^2(\nu_s, \nu)$$

where  $\mathcal{W}_2^2(\nu_s, \nu)$  is the squared Wasserstein distance between  $\nu_s$  and  $\nu$ . Moreover, if  $f^*$  and  $\nu$  solve the left-hand side and the right-hand side problems respectively, then  $\nu$  is equal to the distribution of  $f^*(X, S)$ , and

$$f^*(x, s) = \left( \sum_{s' \in \mathcal{S}} p_{s'} \mathcal{Q}_{s'} \right) \circ F_s(\eta(x, s)),$$

where  $\mathcal{Q}_s$  and  $F_s$  are respectively the quantile function and the c.d.f. of  $\nu_s$ .

This result relates the problem of fair regression in the awareness framework to a more general optimal transport problem. Interestingly, this problem has an explicit solution, given by the quantiles and c.d.f.s of the conditional expectation  $\eta(X, S)$  across the different groups. This explicit formulation yields, as an immediate consequence, that the optimal fair regression function preserves order, a property introduced in [CDH<sup>+</sup>20a, CS20b] within the awareness framework. Recall that the Bayes regression function in the awareness framework is  $\eta$ . A prediction function  $f$  is said to *preserve order* if for any two candidates  $(x, x') \in \mathcal{X}^2$  in the same group  $s \in \mathcal{S}$ ,  $\eta(x, s) \leq \eta(x', s)$  implies  $f(x, s) \leq f(x', s)$ . Thus, this property implies that the fairness correction does not alter the ordering of the predictions within a group.

In contrast, the problem of fair regression within the unawareness framework has been seldom studied, particularly from a theoretical perspective. One reason for this is that the demographic parity constraint is more challenging to implement without disparate treatment. While algorithms complying with these constraints have been proposed by [CS20a] and [ZM23], the authors do not claim that the estimators obtained are optimal in terms of risk. [ADW19] propose an algorithm based on a discretization of the problem, followed by a reduction to cost-sensitive classification. However, their algorithm requires calling an oracle cost-sensitive classifier, which may not be available in practice. Additionally, their results are limited to a class of regression functions with bounded Rademacher complexity.

## 1.4 Outline and contribution

In this paper, we focus on the theoretical aspects of the problem of fair regression in the unawareness framework, specifically on characterizing and studying the optimal regression function. We extend results presented earlier in the awareness framework to this setting, albeit under the assumption that the sensitive attribute takes only two values; henceforth, we assume that  $\mathcal{S} = \{1, 2\}$ . Although restrictive, this assumption is not uncommon in the literature [LMC18] and covers the important case where one of the two groups includes protected individuals. Our results shed light on important phenomena, and we leave the extension to scenarios with more than two groups to future work.

Similarly to the awareness case characterized in Theorem 1, we show that the solution to the fair regression problem in the unawareness framework is given by the solution to a barycenter problem with optimal transport costs. We begin in Section 2 with a brief introduction to optimal transport theory and to the main tools used in the proofs of our results. In Section 3, we characterize the optimal fair regression function. First, we prove in Proposition 1 that *in general, the optimal fair regression function  $f^*$  does not preserve order*. Next, we demonstrate the following result, which relates fair regression in the unawareness framework to an optimal transport problem.

---

<sup>1</sup>This result is formalized in Theorem 4.

**Theorem 2** (Informal<sup>1</sup>). *Under mild assumptions, the optimal fair regression function  $f^*$  is given by the solution to a barycenter problem with optimal transport costs. In particular, there exists a function  $\mathbf{f}^*$  such that*

$$f^*(x) = \mathbf{f}^*(\eta(x), \Delta(x)),$$

where  $\Delta(x) \propto \frac{\mathbb{P}(S=1|X=x)}{p_1} - \frac{\mathbb{P}(S=2|X=x)}{p_2}$ .

Comparing this result to the one provided in Theorem 1 within the awareness framework, we note that there is no explicit formula for the optimal fair regression function within the unawareness framework. Moreover, Theorem 2 underscores that the optimal fair regression function effectively relies on an estimate  $\Delta(X)$  of the unobserved sensitive attribute  $S$  to make predictions, thereby indirectly implementing disparate treatment. This result provides a theoretical explanation for the empirical phenomenon observed by [LMC18]. As noted in their work, this behavior is problematic as it can lead to basing predictions on factors that are not relevant to predicting the outcome  $Y$ , simply because they are informative for predicting the sensitive attribute  $S$ .

In Section 4, we investigate the relationship between fair regression and fair classification when  $Y \in \{0, 1\}$ . We demonstrate the existence of a dichotomy based on a *nestedness* criterion. Recall that as the threshold  $y$  increases, the Bayes classifier  $g_y^{\text{Bayes}}$  predicts 1 for a decreasing proportion of candidates; we show that this also holds for the optimal fair classifier  $g_y^*$ . We say that the fair classification problem is *nested* if, almost surely with respect to the measure  $\mu$  of  $X$ , the prediction  $g_y^*(X)$  for the candidate  $X$  decreases as  $y$  increases. In other words, candidates rejected (i.e., with prediction 0) at low values of  $y$  cannot be accepted at higher values of  $y$ , when the proportion of accepted candidates is lower. For example, the Bayes classifier defined by  $g_y^{\text{Bayes}}(x) = \mathbf{1}\{f^{\text{Bayes}}(x) \geq y\}$  satisfies this condition. When the nestedness criterion holds, the decision boundaries for the optimal fair classifier for different risk  $\mathcal{R}_y$  form nested sets. The following informal result summarizes our findings.

**Theorem 3** (Informal<sup>2</sup>). *Under mild assumptions, if the fair classification problem is nested, then the regression function*

$$f^*(x) = \sup \{y \in \mathbb{R} : g_y^*(x) = 1\}$$

*is optimal for the fair regression problem (1); equivalently, the classifier*

$$g_y(x) = \mathbf{1}\{f^*(x) \geq y\}$$

*is optimal for the fair classification problem (2) for the risk  $\mathcal{R}_y$ . Conversely, if the classification problem is not nested and if  $f^*$  is the optimal fair regression function, then there exists  $y \in (0, 1)$  such that*

$$g_y(x) = \mathbf{1}\{f^*(x) \geq y\}$$

*is sub-optimal for the fair classification problem with risk  $\mathcal{R}_y$ .*

While nestedness may initially appear to be a natural assumption, it does not always hold. In Section 5, we show how to design examples of problems where this assumption is either met or violated.

## 2 A short introduction to optimal transport

In this section, we provide a brief introduction to optimal transport. We present the main tools that will be used in the proofs of the theorems in Sections 3 and 4. We begin by providing an overview of optimal transport in Section 2.1, before discussing the multi-to-one dimensional transport problem in Section 2.2.

---

<sup>2</sup>This result is formalized in Proposition 6 and in Corollary 1

## 2.1 The optimal transport problem

Optimal transport provides a mathematical framework to compare probability distributions. Consider a Borel probability measure  $\mu$  on a Polish space  $\mathcal{X}$  and a Borel probability measure  $\nu$  on some other Polish space  $\mathcal{Y}$ . We are given a continuous cost function  $c : \mathcal{X} \times \mathcal{Y} \rightarrow [0, +\infty]$ , where  $c(x, y)$  represents the cost of moving a unit of mass from  $x \in \mathcal{X}$  to  $y \in \mathcal{Y}$ . The optimal transport problem consists in finding the optimal way of moving the distribution of mass  $\mu$  to  $\nu$  by minimizing the total displacement cost. Formally, a transport map is a measurable map  $T : \mathcal{X} \rightarrow \mathcal{Y}$  such that the pushforward measure  $T\#\mu$  of  $\mu$  by  $T$  is equal to  $\nu$ , where the pushforward measure is defined for all measurable sets  $B \subset \mathcal{Y}$  by

$$T\#\mu(B) = \mu(T^{-1}(B)).$$

The optimal transport problem is then the following

$$\begin{aligned} & \text{minimize} && \int c(x, T(x)) d\mu(x) \\ & \text{under the constraint} && T\#\mu = \nu. \end{aligned} \tag{3}$$

The existence of minimizers of the optimization problem (3) is a delicate problem that depends on both the regularity of the cost function  $c$  and the properties of  $\mu$  and  $\nu$ . For instance, when  $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$  and  $c(x, y) = \|x - y\|^2$ , a solution exists whenever  $\mu$  gives zero mass to sets of dimensions smaller than  $d - 1$ ; otherwise, a solution may not exist, see [Vil09, Chapter 10]. When  $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$  and  $c(x, y) = \|x - y\|^2$ , the corresponding minimum is known as the (squared) Wasserstein distance between  $\mu$  and  $\nu$ , denoted by  $\mathcal{W}_2^2(\mu, \nu)$ . More generally, an optimal transport map exists whenever  $\mu$  gives zero mass to sets of dimensions smaller than  $d - 1$  and the cost function  $c(x, y) = \|x - y\|^2$  is replaced by any smooth cost function  $c$  satisfying the so-called *twist condition*, which states that the determinant  $\det(\frac{\partial^2 c}{\partial y_j \partial x_i})$  never vanishes.

The optimal transport problem also admits a relaxed version in terms of transport plans, which is often more convenient to work with. A transport plan is a probability measure  $\pi$  on the product space  $\mathcal{X} \times \mathcal{Y}$  which has first marginal equal to  $\mu$  and second marginal equal to  $\nu$ : for all measurable sets  $A \subset \mathcal{X}$  and  $B \subset \mathcal{Y}$ ,

$$\pi(A \times \mathcal{Y}) = \mu(A), \quad \pi(\mathcal{X} \times B) = \nu(B),$$

or, in probabilistic terms, if  $(X, Y) \sim \pi$ , then  $X \sim \mu$  and  $Y \sim \nu$ . Informally, for  $x \in \mathcal{X}$ , the conditional law of  $Y|X = x$  describes the different locations where the mass initially at  $x$  will be sent. The cost of a transport plan  $\pi$  is given by

$$\iint c(x, y) d\pi(x, y).$$

Note that a transport map  $T$  induces a transport plan by considering the law  $\pi$  of  $(X, T(X))$  (formally,  $\pi = (\text{id}, T)\#\mu$ ). The optimal transport cost is defined by the following minimization problem:

$$\text{OT}_c(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \int c(x, y) d\pi(x, y), \tag{4}$$

where  $\Pi(\mu, \nu)$  is the set of transport plans between  $\mu$  and  $\nu$ . Optimal transport plans always exist, whereas optimal transport maps may fail to do so. When optimal transport maps exist and the source measure  $\mu$  has no atoms, the minimization problem (3) gives the same value as the optimal transport cost defined in (4), see [Pra07].

Our proofs will rely heavily on the dual formulation of the optimal transport problem, which we now introduce. The  $c$ -transform of a function  $\varphi : \mathcal{Y} \rightarrow \mathbb{R} \cup \{+\infty\}$  is defined as

$$\forall x \in \mathcal{X}, \quad \varphi^c(x) = \sup_{y \in \mathcal{Y}} (\varphi(y) - c(x, y)).$$

The subdifferential of  $\varphi$  is defined as

$$\partial_c \varphi = \{(x, y) \in \mathcal{X} \times \mathcal{Y} : \varphi(y) - \varphi^c(x) = c(x, y)\}.$$

For the quadratic cost, these notions are closely related to the usual notions of convexity, with  $c$ -transforms being analogous to the concept of convex conjugates.

Kantorovich duality [Vil09, Theorem 5.10] states that

$$\text{OT}_c(\mu, \nu) = \sup_{\varphi \in L^1(\nu)} \left( \int \varphi(y) d\nu(y) - \int \varphi^c(x) d\mu(x) \right).$$

Moreover, under the mild assumption that there exist two functions  $a \in L^1(\mu)$  and  $b \in L^1(\nu)$  such that  $c(x, y) \leq a(x) + b(y)$  for all  $x \in \mathcal{X}$ ,  $y \in \mathcal{Y}$ , the previous supremum is attained by a function  $\varphi$ , which we call a Kantorovich potential. In that case, any optimal transport  $\pi$  is supported on the subdifferential of the  $c$ -convex function  $\varphi$ , meaning that

$$\pi(\partial_c \varphi) = 1.$$

This last condition imposes significant constraints on the structure of optimal transport plans. For the quadratic cost, this fact is the key ingredient in proving that optimal transport plans are induced by optimal transport maps.

## 2.2 Multi-to-one dimensional optimal transport

In the next section, we demonstrate that the fair regression problem within the unawareness framework can be reduced to a barycenter problem of the form:

$$\min_{\nu \in \mathcal{P}(\mathbb{R})} p_1 \text{OT}_c(\mu_1, \nu) + p_2 \text{OT}_c(\mu_2, \nu), \quad (5)$$

where  $\mu_1, \mu_2$  are *two-dimensional* probability measures and  $c : \mathbb{R}^2 \times \mathbb{R} \rightarrow [0, +\infty)$  is a cost function. This reduction raises the question of whether the solutions to the barycenter problem (5) can be characterized by transport maps.

Proving that the optimal transport problem  $\text{OT}_c(\mu_s, \nu)$  is solved by a transport map is nontrivial. Complications arise because the measures  $\mu_s$  and  $\nu$  are defined on spaces of different dimensions. Optimal transport problems involving spaces of different dimensions have not been as extensively studied and exhibit distinct properties compared to the standard case where both measures are defined on the same space, see [Pas12, CMP16, CMP17, MP20]. For instance, the classical twist condition  $\det(\frac{\partial^2 c}{\partial y_j \partial x_i}) \neq 0$  does not make sense in this setting: the Hessian matrix of  $c$  is not squared, so that the determinant is not even well-defined.

[CMP17] focus on the optimal transport problem between a measure  $\mu$  supported on a domain  $\mathcal{X} \subset \mathbb{R}^m$  (with  $m > 1$ ) and a measure  $\nu$  on an interval  $\mathcal{Y} \subset \mathbb{R}$  for some cost function  $c : \mathcal{X} \times \mathcal{Y} \rightarrow [0, +\infty)$ . They demonstrate that an optimal transport map  $T : \mathcal{X} \rightarrow \mathcal{Y}$  between  $\mu$  and  $\nu$  exists under a natural condition on  $(c, \mu, \nu)$  known as *nestedness*. For  $y \in \mathcal{Y}$ ,  $k \in \mathbb{R}$ , let

$$\mathcal{X}_{\leq}(y, k) = \{x \in \mathcal{X} : \partial_y c(x, y) \leq k\}.$$

Kantorovich duality implies that an optimal transport plan between  $\mu$  and  $\nu$  will match an interval  $(-\infty, y]$  to a set  $\mathcal{X}_{\leq}(y, k)$ , where  $k = k(y)$  is a solution of the equation  $\nu((-\infty, y]) = \mu(\mathcal{X}_{\leq}(y, k))$ . The triplet  $(c, \mu, \nu)$  is called nested if the collection of sets  $(\mathcal{X}_{\leq}(y, k(y)))_y$  increases with  $y$ . [CMP17] prove that an optimal transport map  $T$  between  $\mu$  and  $\nu$  exists when the problem is nested: informally, the monotonicity of  $(\mathcal{X}_{\leq}(y, k(y)))_y$  ensures that a given  $x_0 \in \mathcal{X}$  belongs to the boundary of a single set  $\mathcal{X}_{\leq}(y_0, k(y_0))$ , with  $y_0$  being equal to  $T(x_0)$ .

This nestedness condition will be crucial in Section 4, where it will be used to establish the equivalence between regression and classification problems. However, in Section 3, we will be able to show the existence of optimal transport maps for the barycenter problem (5) (and consequently of optimal fair regression functions) without any nestedness condition.

## 3 Fair regression and the barycenter problem

In this section, we characterize the solution to the fair regression problem. We begin by showing in Section 3.1 that, under mild assumptions, the fair regression function does not preserve order. Then, in Section 3.2,



we show that the fair regression problem can be reduced to a barycenter problem with optimal transport costs. Using the tools introduced in Section 2, we prove the existence of a fair optimal prediction function and study some of its properties.

### 3.1 Fair regression functions do not preserve order

Before analyzing in detail the fair regression problem in the unawareness framework, we establish a simple yet important property of fair regression functions. We begin by extending the definition of order preservation [CDH<sup>+</sup>20a, CS20b] to the unawareness framework. Recall that in this case, the Bayes prediction for a candidate  $x$  is given by  $\eta(x) = \mathbb{E}[Y|X = x]$ . A prediction function  $f$  is said to preserve order if, for any two candidates  $(x, x') \in \mathcal{X}^2$  in the same group  $s \in \mathcal{S}$ ,  $\eta(x) \leq \eta(x')$  implies  $f(x) \leq f(x')$ . This definition is formalized below.

**Definition 4** (Order preservation in regression - unawareness framework). *A prediction function  $f$  preserves order if  $\mathbb{P} \otimes \mathbb{P}$ -almost surely,*

$$\{\eta(X) < \eta(X') \text{ and } S = S'\} \implies f(X) < f(X').$$

This property implies that the fairness correction does not alter the ordering of the predictions of the Bayes prediction function within a group. It is related to the concept of “rational ordering” introduced by [LMC18] in the context of classification, where the authors require that within a group, the most able candidates are the ones accepted.

**Proposition 1.** *Let  $f : \mathcal{X} \rightarrow \mathbb{R}$  be a regression function with  $\mathbb{E}[f(X)^2] < \infty$  satisfying the demographic parity constraint. Assume that the Bayes regression function  $\eta$  does not satisfy the demographic parity constraint and that  $\mathbb{P}(S = s|X = x) \in (0, 1)$  for all  $s \in \mathcal{S}$ ,  $x \in \mathcal{X}$ . Then,  $f$  does not preserve order.*

*Proof.* We prove the contrapositive: if  $f$  is a regression function satisfying the demographic parity constraint and preserving order, then the Bayes regression function also satisfies the demographic parity constraint. For a fixed group  $s \in \mathcal{S}$ , consider the joint law  $\pi_s$  of  $(\eta(X), f(X))$ , where  $X \sim \mu_s$ . As  $f$  is envy-free, the support of the measure  $\pi_s$  is monotone, in the sense that

$$\forall (y_1, z_1), (y_2, z_2) \in \text{supp}(\pi_s), y_1 < y_2 \implies z_1 < z_2. \quad (6)$$

According to [San15, Lemma 2.8], this implies that  $\pi_s$  is actually the optimal transport plan for the quadratic cost between the first marginal of  $\pi_s$ , equal to  $\eta\#\mu_s =: \nu_s$  and the second marginal of  $\pi_s$ , equal to  $f\#\mu_s =: \nu$  (the second measure does not depend on  $s$  because of demographic parity). We claim that strict monotonicity implies that the transport plan  $\pi_s$  takes the form of a transport map  $T_s$  transporting  $\nu$  towards  $\nu_s$ , that is  $\pi_s = (T_s, \text{id})\#\nu$  (see a proof below). So, if  $X \sim \mu_s$ , we have  $(\eta(X), f(X)) = (T_s(f(X)), f(X))$  almost surely. To put it another way, we have for every  $s$ ,

$$\eta(x) = T_s \circ f(x) \quad \mu_s\text{-almost everywhere.}$$

As  $\mathbb{P}(S = s|X = x) > 0$  for all  $x \in \mathcal{X}$ , this equality is also satisfied  $\mu$ -almost everywhere. Hence, for  $\mu$ -almost all  $x$ , the quantity  $T_s \circ f(x)$  does not depend on  $s$  (it is equal to  $\eta(x)$ ). This defines a function  $T$  with  $\eta\#\mu_s = T\#\mu_s = T\#\nu$ . As this measure does not depend on  $s$ , this proves that  $\eta$  satisfies the demographic parity constraint.

To conclude our proof, it remains to prove our claim. Decompose  $\nu$  as  $\nu_1 + \nu_2$  where  $\nu_2$  is atomless and  $\nu_1 = \sum_j p_j \delta_{z_j}$ . If  $f(X) = z_j$ , then we have  $\eta(X) = y_j$  for some value  $y_j$ : this value  $y_j$  has to be unique, for otherwise it would contradict the monotonicity assumption. Therefore,  $\nu_s$  can be written as  $\nu_s = \nu_{1s} + \nu_{2s}$ , where  $\nu_{1s} = \sum_j p_j \delta_{y_j}$ . Consider the plan  $\pi_1 = \sum_j p_j \delta_{(y_j, z_j)}$ . Then  $\pi - \pi_1$  is a plan between  $\nu_{2s}$  and  $\nu_2$ . By [San15, Lemma 2.8], as  $\nu_2$  is atomless, the monotonicity condition implies that it is induced by a transport map  $\tilde{T}_s$ . In total, we can define  $T_s$  by  $T_s(z_j) = y_j$  and by  $T_s = \tilde{T}_s$  on the complementary set of the atoms.  $\square$

Proposition 1 implies, in particular, that in many instances, the optimal fair regression function does not preserve order. Consequently, highly qualified individuals who belong to protected groups could potentially suffer from fairness corrections due to the demographic parity constraint.

The condition that  $\mathbb{P}(S = s|X = x) \in (0, 1)$  for all  $s \in \mathcal{S}$ ,  $x \in \mathcal{X}$  ensures that the sensitive attribute  $S$  cannot be determined from the observation of  $X$ . When this condition is not satisfied, the distinction between the unawareness and awareness frameworks becomes blurred: if  $S$  can be inferred from  $X$  alone, it becomes meaningless to differentiate between a regression function that depends on both  $X$  and  $S$ , and one that depends solely on  $X$ . Furthermore, it is important to note that in the awareness framework, there do exist regression functions that satisfy the demographic parity constraint and preserve order, with the optimal fair regression function described in Theorem 1 being one such example.

## 3.2 Reduction to an optimal transport problem

In the following, we let  $\mathcal{S} = \{1, 2\}$ . We assume that  $\mu_1 \neq \mu_2$  (otherwise the Bayes regression function  $\eta$  already solves the fair regression problem). We now show how to transform the fair regression problem into a barycenter problem using optimal transport costs. To do so, we first leverage a reformulation of the demographic parity constraint due to [CS20a], which is based on the Jordan decomposition of the signed measure  $\mu_1 - \mu_2$ . Then, we show how to rephrase the regression problem as a barycenter problem, using this new constraint. Finally, we show that, under mild assumptions, the barycenter problem admits a unique solution, which is given by a transport map.

### 3.2.1 Reformulation of the demographic parity constraint

Let  $|\mu_1 - \mu_2|$  be the variation of  $\mu_1 - \mu_2$  and define

$$\begin{cases} (\mu_1 - \mu_2)_+ = \frac{1}{2}(|\mu_1 - \mu_2| + \mu_1 - \mu_2), \\ (\mu_1 - \mu_2)_- = \frac{1}{2}(|\mu_1 - \mu_2| - \mu_1 + \mu_2) \end{cases}$$

the Jordan decomposition of  $\mu_1 - \mu_2$ . The two measures  $(\mu_1 - \mu_2)_+$  and  $(\mu_1 - \mu_2)_-$  have the same mass, which we denote by  $m$ . We define the scaled Jordan decomposition of  $\mu_1 - \mu_2$  as the pair of probability measures

$$\mu_+ = (\mu_1 - \mu_2)_+/m \quad \text{and} \quad \mu_- = (\mu_1 - \mu_2)_-/m.$$

Let  $\frac{d\mu_+}{d\mu}$  (resp.  $\frac{d\mu_-}{d\mu}$ ) be the density of  $\mu_+$  (resp.  $\mu_-$ ) with respect to  $\mu$  (that are defined uniquely  $\mu$ -almost everywhere). As  $\mu_+$  and  $\mu_-$  are mutually singular measures, we can always find versions of  $\frac{d\mu_+}{d\mu}$  and  $\frac{d\mu_-}{d\mu}$  such that the sets

$$\begin{cases} \mathcal{X}_+ = \{x \in \mathcal{X} : \frac{d\mu_+}{d\mu}(x) > 0\}, \\ \mathcal{X}_- = \{x \in \mathcal{X} : \frac{d\mu_-}{d\mu}(x) > 0\}, \\ \mathcal{X}_= = \mathcal{X} \setminus (\mathcal{X}_+ \sqcup \mathcal{X}_-). \end{cases}$$

form a partition of  $\mathcal{X}$ , with  $\mu_+$  giving mass 1 to  $\mathcal{X}_+$  and  $\mu_-$  giving mass 1 to  $\mathcal{X}_-$ . Then, for any three functions  $f_+$ ,  $f_-$ , and  $f_=$  from  $\mathcal{X}$  to  $\mathbb{R}$ , we can define the associated function  $\mathcal{F}(f_+, f_-, f_=)$  equal to  $f_+$  on  $\mathcal{X}_+$ ,  $f_-$  on  $\mathcal{X}_-$ , and  $f_=$  on  $\mathcal{X}_=$ :

$$\mathcal{F}(f_+, f_-, f_=)(x) = \begin{cases} f_+(x) & \text{if } x \in \mathcal{X}_+ \\ f_-(x) & \text{if } x \in \mathcal{X}_- \\ f_=(x) & \text{if } x \in \mathcal{X}_= \end{cases}$$

Conversely, for any function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , there exist functions  $f_+$ ,  $f_-$ , and  $f_=$  corresponding respectively to the restriction of  $f$  on  $\mathcal{X}_+$ ,  $\mathcal{X}_-$ , and  $\mathcal{X}_=$ , i.e., such that  $f = \mathcal{F}(f_+, f_-, f_=)$ . The following lemma, due to [CS20a], rephrases the demographic parity constraint in terms of  $\mu_+$  and  $\mu_-$ .

**Lemma 1.** *A regression function  $f : \mathcal{X} \rightarrow \mathbb{R}$  verifies the demographic parity constraint if and only if*

$$f \# \mu_+ = f \# \mu_-.$$

Lemma 1 reveals that for any functions  $f$ , and  $f_+$ ,  $f_-$ ,  $f_=$  such that  $f = \mathcal{F}(f_+, f_-, f_=)$ , the regression function  $f$  satisfies the demographic parity constraint if and only if  $f_+ \# \mu_+ = f_- \# \mu_-$ . The two functions  $f_+$  and  $f_-$  can be chosen with disjoint support (in  $\mathcal{X}_+$  and  $\mathcal{X}_-$ , respectively). Thus, the demographic parity constraint essentially reduces to the equality of the pushforward measures of two distinct probabilities ( $\mu_+$  and  $\mu_-$ ) by two distinct functions ( $f_+$  and  $f_-$ ).

### 3.2.2 A barycenter problem

In order to rephrase the regression problem as a barycenter problem, we introduce further notation. We define

$$\begin{cases} \Delta(x) = \frac{d\mu_+}{d\mu}(x) & \text{if } x \in \mathcal{X}_+, \\ \Delta(x) = -\frac{d\mu_-}{d\mu}(x) & \text{if } x \in \mathcal{X}_-, \\ \Delta(x) = 0 & \text{if } x \in \mathcal{X}_=. \end{cases} \quad (7)$$

Equivalently,  $\Delta(x)$  is proportional to  $\frac{d\mu_+}{d\mu}(x) - \frac{d\mu_-}{d\mu}(x)$ . We also define the cost  $c : \mathcal{X} \times \mathbb{R} \rightarrow [0, +\infty]$  given by  $c(x, y) = \frac{(\eta(x) - y)^2}{|\Delta(x)|}$  for all  $x \in \mathcal{X}$  and all  $y \in \mathbb{R}$ . When  $X \sim \mu_{\pm}$ , the variables  $(\eta(X), \Delta(X))$  belong to  $\Omega := \{(h, d) \in \mathbb{R}^2 : d \neq 0\}$ . In the following, we use bold notation to denote functions related to these two-dimensional variables. For example, we denote by  $\boldsymbol{\mu}_+$  (resp.  $\boldsymbol{\mu}_-$ ) the distributions of  $(\eta(X), \Delta(X))$  when  $X$  follows the distribution  $\mu_+$  (resp.  $\mu_-$ ). Note that the support of  $\boldsymbol{\mu}_+$  is included in the upper half-plane  $\{d > 0\}$  while  $\boldsymbol{\mu}_-$  is included in the lower half-plane  $\{d < 0\}$ . We define the two-to-one dimensional cost  $\mathbf{c}$ , given by  $\mathbf{c}(\mathbf{x}, y) = \frac{(h-y)^2}{|d|}$  for all  $\mathbf{x} = (h, d) \in \Omega$  and  $y \in \mathbb{R}$ .

Consider the barycenter problem

$$\text{minimize } \text{OT}_{\mathbf{c}}(\boldsymbol{\mu}_+, \nu) + \text{OT}_{\mathbf{c}}(\boldsymbol{\mu}_-, \nu) \text{ over } \nu \in \mathcal{P}(\mathbb{R}) \quad (8)$$

where we recall that  $\text{OT}_{\mathbf{c}}(\boldsymbol{\mu}_{\pm}, \nu)$  is the optimal transport cost for sending  $\boldsymbol{\mu}_{\pm}$  to  $\nu$  with cost function  $\mathbf{c}$ , defined in Equation (4). We say that a solution  $\nu^{\text{bar}}$  of the barycenter problem is solved by optimal transport maps if

$$\text{OT}_{\mathbf{c}}(\boldsymbol{\mu}_{\pm}, \nu^{\text{bar}}) = \int \mathbf{c}(\mathbf{x}, \mathbf{f}_{\pm}(\mathbf{x})) d\boldsymbol{\mu}_{\pm}(\mathbf{x})$$

for some transport maps  $\mathbf{f}_{\pm} : \Omega \rightarrow \mathbb{R}$  from  $\boldsymbol{\mu}_{\pm}$  to  $\nu^{\text{bar}}$ .

**Lemma 2.** *There is a one-to-one correspondence between the set of solutions to the barycenter problem (8) solved by optimal transport maps and the set of optimal fair regression functions solving (1). This correspondence associates a barycenter  $\nu^{\text{bar}}$  with optimal transport maps  $\mathbf{f}_{\pm}$  to the optimal fair regression function  $f = \mathcal{F}(f_+, f_-, \eta)$ , where  $f_{\pm}(x) = \mathbf{f}_{\pm}(\eta(x), \Delta(x))$  for  $x \in \mathcal{X}$ .*

*Proof.* Classical computations show that  $\mathcal{R}_{sq}(f) = \mathbb{E}[(\eta(X) - f(X))^2] + \mathbb{E}[(\eta(X) - Y)^2]$ . Thus, minimizing the risk is equivalent to minimizing  $\mathbb{E}[(\eta(X) - f(X))^2]$ . Now,

$$\begin{aligned} \mathbb{E}[(\eta(X) - f(X))^2] &= \int (\eta(x) - f(x))^2 d\mu(x) \\ &= \int_{\mathcal{X}_+} (\eta(x) - f(x))^2 \frac{d\mu}{d\mu_+}(x) d\mu_+(x) + \int_{\mathcal{X}_-} (\eta(x) - f(x))^2 \frac{d\mu}{d\mu_-}(x) d\mu_-(x) \\ &\quad + \int_{\mathcal{X}_=} (\eta(x) - f(x))^2 d\mu(x). \end{aligned} \quad (9)$$

Using the definition of  $\Delta$  along with Lemma 1, we see that any solution to the fair regression problem can be written as  $f = \mathcal{F}(f_+, f_-, \eta)$ , where  $(f_+, f_-)$  is solution to the problem

$$\begin{aligned} \text{minimize } & \int_{\mathcal{X}_+} \frac{(\eta(x) - f_+(x))^2}{|\Delta(x)|} d\mu_+(x) + \int_{\mathcal{X}_-} \frac{(\eta(x) - f_-(x))^2}{|\Delta(x)|} d\mu_-(x) \\ \text{such that } & f_+ \# \mu_+ = f_- \# \mu_-. \end{aligned}$$

The triplet  $(\eta(X), \Delta(X), f(X))$  defines a coupling  $\pi_{f_+}$  between  $\boldsymbol{\mu}_+$  and  $\nu_{f_+} = f_+ \# \mu_+$ . Likewise, we define a coupling  $\pi_{f_-}$  between  $\boldsymbol{\mu}_-$  and  $\nu_{f_-}$ . We can rewrite (9) as

$$\begin{aligned} \mathbb{E}[(\eta(X) - \mathcal{F}(f_+, f_-, \eta)(X))^2] &= \int \mathbf{c}(\mathbf{x}, y) d\pi_{f_+}(\mathbf{x}, y) + \int \mathbf{c}(\mathbf{x}, y) d\pi_{f_-}(\mathbf{x}, y) \\ &\geq \text{OT}_{\mathbf{c}}(\boldsymbol{\mu}_+, \nu_{f_+}) + \text{OT}_{\mathbf{c}}(\boldsymbol{\mu}_-, \nu_{f_-}). \end{aligned}$$

The constraint  $f_+ \# \mu_+ = f_- \# \mu_-$  implies that  $\nu_{f_+} = \nu_{f_-}$ . Hence,

$$\inf_{f \text{ fair}} \mathbb{E} [(\eta(X) - f(X))^2] \geq \inf_{\nu \in \mathcal{P}(\mathbb{R})} \text{OT}_c(\mu_+, \nu) + \text{OT}_c(\mu_-, \nu). \quad (10)$$

Reciprocally, assume that there exists  $\nu^{\text{bar}}$  solving the above barycenter problem, and that an optimal transport map between  $\mu_+$  and  $\nu^{\text{bar}}$  is given by an application  $\mathbf{f}_+ : \Omega \rightarrow \mathbb{R}$ , with  $(\mathbf{f}_+) \# \mu_+ = \nu^{\text{bar}}$ . Likewise, we assume that there exists an optimal transport map  $\mathbf{f}_-$  between  $\mu_-$  and  $\nu^{\text{bar}}$ . Then,  $\mathbf{f}_- \# \mu_- = \mathbf{f}_+ \# \mu_+ = \nu^{\text{bar}}$ . Defining  $f_{\pm}(x) = \mathbf{f}_{\pm}(\eta(x), \Delta(x))$ , we have  $f_- \# \mu_- = f_+ \# \mu_+ = \nu^{\text{bar}}$ , and so  $\mathcal{F}(f_+, f_-, \eta)$  is a fair regression function. Also, we have by optimality that

$$\begin{aligned} \text{OT}_c(\mu_+, \nu^{\text{bar}}) + \text{OT}_c(\mu_-, \nu^{\text{bar}}) &= \int c(x, f_+(x)) d\mu_+(x) + \int c(x, f_-(x)) d\mu_-(x) \\ &= \mathbb{E}[(\eta(X) - \mathcal{F}(f_+, f_-, \eta)(X))^2]. \end{aligned}$$

Hence, by (10), the regression function  $\mathcal{F}(f_+, f_-, \eta)$  is optimal. This concludes the proof of Lemma 2.  $\square$

### 3.2.3 Transport maps for the barycenter problem

The rest of this section is devoted to proving that the barycenter problem indeed admits a solution given by transport maps, which will imply that there exists a solution to the fair regression problem. We show that this holds under the following mild regularity assumption.

**Assumption 1.** *The measures  $\mu_+$  and  $\mu_-$  give zero mass to graphs of functions in the sense that for any measurable function  $F : \mathbb{R} \setminus \{0\} \rightarrow \mathbb{R}$ ,  $\mu_{\pm}(\{(F(d), d) : d \neq 0\}) = 0$ .*

By Fubini's theorem, this assumption is trivially satisfied if  $\mu_+$  and  $\mu_-$  have a density with respect to the Lebesgue measure. Another interesting example is given by the awareness framework, seen as a particular instance of the unawareness framework.

**Remark 1** (Awareness as a special case of unawareness). *Consider a triplet of random variable  $(X, S, Y) \sim \mathbb{P}$ , where  $X \in \mathcal{X}$  is a feature,  $S \in \{1, 2\}$  is a sensitive attribute and  $Y \in \mathbb{R}$  is a response variable of interest. Let  $Z = (X, S)$  and let  $\mathbb{Q}$  be the law of the triplet  $(Z, S, Y)$ . Then, there is an equivalence between considering an aware regression function  $f(X, S)$  under law  $\mathbb{P}$  and an unaware regression function  $f(Z)$  under law  $\mathbb{Q}$ . Note that  $Z$  is a random variable on  $\tilde{\mathcal{X}} = \mathcal{X} \times \{1, 2\}$ . The laws  $\mu_1$  of  $Z|S = 1$  and  $\mu_2$  of  $Z|S = 2$  have disjoint support. It follows that  $\mathcal{X}_+ = \mathcal{X} \times \{1\}$  with  $\mu_+ = \mu_1$  and  $\mathcal{X}_- = \mathcal{X} \times \{2\}$  with  $\mu_- = \mu_2$ . Then,  $\Delta(x) = 1/p_1$  if  $x \in \mathcal{X}_+$  and  $\Delta(x) = -1/p_2$  if  $x \in \mathcal{X}_-$ . In particular, both measures  $\mu_+$  and  $\mu_-$  are supported on horizontal lines in  $\Omega$ .*

*In that case, Assumption 1 is equivalent to the fact that  $\mu_+$  and  $\mu_-$  have no atoms, which is exactly equivalent to the fact that the law of  $\eta(X)$  (for  $X \sim \mu$ ) has no atoms. This assumption is often considered to be a minimal assumption to ensure the existence of optimal fair regression functions in the awareness framework. Hence, Assumption 1 constitutes a generalization of this assumption to the unawareness framework.*

**Theorem 4.** *Assume that  $(X, Y, S) \sim \mathbb{P}$  is such that  $\mathbb{E}[Y^2] < \infty$ . Under Assumption 1, there is a unique minimizer  $\nu^{\text{bar}}$  of the barycenter problem*

$$\inf_{\nu} \text{OT}_c(\mu_+, \nu) + \text{OT}_c(\mu_-, \nu). \quad (11)$$

*Moreover, this problem is solved by optimal transport maps  $\mathbf{f}_{\pm}$ . In particular, there exists a unique solution  $f^*$  of the regression problem under the demographic parity constraint (1), which is given by*

$$\forall x \in \mathcal{X}, f^*(x) = \mathcal{F}(\mathbf{f}_+(\eta(x), \Delta(x)), \mathbf{f}_-(\eta(x), \Delta(x)), \eta(x)).$$

*Proof.* Using Lemma 2, it is enough to show that the barycenter problem admits a unique solution  $\nu^{\text{bar}}$  such that the corresponding transport problems  $\text{OT}_c(\mu_+, \nu^{\text{bar}})$  and  $\text{OT}_c(\mu_-, \nu^{\text{bar}})$  are solved by transport maps.

**Step 1: reduction to a standard transport problem.** We begin by reducing the barycenter problem (11) to a single two-to-two dimensional optimal transport problem  $\text{OT}_{\mathbf{C}}(\boldsymbol{\mu}_+, \boldsymbol{\mu}_-)$ . The multimarginal version of the barycenter problem reads

$$\inf_{\nu \in \mathcal{P}(\mathbb{R})} \text{OT}_{\mathbf{c}}(\boldsymbol{\mu}_+, \nu) + \text{OT}_{\mathbf{c}}(\boldsymbol{\mu}_-, \nu) = \inf_{\rho \in \Pi(\cdot, \boldsymbol{\mu}_+, \boldsymbol{\mu}_-)} \int (\mathbf{c}(\mathbf{x}_1, y) + \mathbf{c}(\mathbf{x}_2, y)) d\rho(y, \mathbf{x}_1, \mathbf{x}_2), \quad (12)$$

where  $\Pi(\cdot, \boldsymbol{\mu}_+, \boldsymbol{\mu}_-)$  stands for the set of measures on  $\mathbb{R} \times \Omega \times \Omega$  with second marginal  $\boldsymbol{\mu}_+$  and third marginal  $\boldsymbol{\mu}_-$ . Indeed, if  $\rho \in \Pi(\cdot, \boldsymbol{\mu}_+, \boldsymbol{\mu}_-)$ , then its two first marginals provide a transport plan between its first marginal  $\nu$  and  $\boldsymbol{\mu}_-$ , while the first and last marginals provide a transport plan between  $\nu$  and  $\boldsymbol{\mu}_+$ . This proves that the left-hand side of Equation (12) is smaller than the right-hand side. For the other inequality, consider  $\nu \in \mathcal{P}(\mathbb{R})$ , with associated optimal transport plans  $\pi_+ \in \Pi(\boldsymbol{\mu}_+, \nu)$  and  $\pi_- \in \Pi(\boldsymbol{\mu}_-, \nu)$ . By the gluing lemma (see, e.g., Lemma 5.5 in [San15]), there exists  $\rho \in \Pi(\cdot, \boldsymbol{\mu}_+, \boldsymbol{\mu}_-)$  such that the joint law of the first two marginals is equal to  $\pi_+$ , and the joint law of the first and last marginal is equal to  $\pi_-$ . Then,  $\text{OT}_{\mathbf{c}}(\boldsymbol{\mu}_+, \nu) + \text{OT}_{\mathbf{c}}(\boldsymbol{\mu}_-, \nu) = \int (\mathbf{c}(\mathbf{x}_1, y) + \mathbf{c}(\mathbf{x}_2, y)) d\rho(y, \mathbf{x}_1, \mathbf{x}_2)$ , proving that the right-hand side is smaller than the left-hand side in (12). This shows the validity of (12).

Furthermore, if  $\rho$  solves the right-hand side of (12), then its first marginal  $\nu$  is a barycenter. Actually, by optimality, for any  $(y, \mathbf{x}_1, \mathbf{x}_2)$  in the support of the optimal  $\rho$ , the point  $y$  necessarily minimizes the function  $z \mapsto \mathbf{c}(\mathbf{x}_1, z) + \mathbf{c}(\mathbf{x}_2, z)$ . Let us compute this minimizer. For  $\mathbf{x}_1 = (h_1, d_1)$  and  $\mathbf{x}_2 = (h_2, d_2)$ , we have

$$\mathbf{c}(\mathbf{x}_1, y) + \mathbf{c}(\mathbf{x}_2, y) = (h_1 - y)^2/|d_1| + (h_2 - y)^2/|d_2|.$$

This function is convex in  $y$ . The first order condition for optimality reads

$$(y - h_1)/|d_1| + (y - h_2)/|d_2| = 0 \iff y = m(\mathbf{x}_1, \mathbf{x}_2) := \frac{h_1/|d_1| + h_2/|d_2|}{1/|d_1| + 1/|d_2|}. \quad (13)$$

Moreover, the cost  $\mathbf{C}(\mathbf{x}_1, \mathbf{x}_2) := \inf_y \mathbf{c}(\mathbf{x}_1, y) + \mathbf{c}(\mathbf{x}_2, y)$  corresponding to this minimum is equal to

$$\mathbf{C}(\mathbf{x}_1, \mathbf{x}_2) = \frac{(h_2 - h_1)^2}{|d_1| + |d_2|}. \quad (14)$$

These considerations show that

$$\inf_{\nu \in \mathcal{P}(\mathbb{R})} \text{OT}_{\mathbf{c}}(\boldsymbol{\mu}_+, \nu) + \text{OT}_{\mathbf{c}}(\boldsymbol{\mu}_-, \nu) = \text{OT}_{\mathbf{C}}(\boldsymbol{\mu}_+, \boldsymbol{\mu}_-), \quad (15)$$

and that optimal transport plans  $\pi^* \in \Pi(\boldsymbol{\mu}_+, \boldsymbol{\mu}_-)$  are in correspondence with barycenters  $\nu$  through the formula  $\nu = m\#\pi^*$ . In particular, as there exists at least one optimal transport plan, the infimum in the barycenter problem is actually a minimum.

**Step 2: existence of a transport map.** Note that

$$\mathbf{C}(\mathbf{x}_1, \mathbf{x}_2) = \frac{(h_1 - h_2)^2}{|d_1| + |d_2|} \leq 2 \frac{h_1^2}{|d_1|} + 2 \frac{h_2^2}{|d_2|}. \quad (16)$$

This quantity is integrable against  $\boldsymbol{\mu}_+ \otimes \boldsymbol{\mu}_-$ . Indeed,

$$\int \frac{h_1^2}{|d_1|} d\boldsymbol{\mu}_+(h_1, d_1) = \int \frac{\eta(x)^2}{\Delta(x)} d\mu_+(x) = \int_{\mathcal{X}_+} \eta(x)^2 d\mu(x) \leq \mathbb{E}[\mathbb{E}[Y|X]^2] \leq \mathbb{E}[Y^2] < \infty.$$

In particular, the optimal cost  $\text{OT}_{\mathbf{C}}(\boldsymbol{\mu}_+, \boldsymbol{\mu}_-)$  is finite. Hence, by Kantorovich duality (see Section 2), there is a  $\mathbf{C}$ -convex function (called a Kantorovich potential)  $\varphi : \Omega \rightarrow \mathbb{R} \cup \{+\infty\}$  such that if we let

$$\Gamma = \{(\mathbf{x}_1, \mathbf{x}_2) : \varphi(\mathbf{x}_1) - \varphi(\mathbf{x}_2) = \mathbf{C}(\mathbf{x}_1, \mathbf{x}_2)\} \quad (17)$$

be the subdifferential of  $\varphi$ , then any optimal transport plan  $\pi$  satisfies  $\pi(\Gamma) = 1$ , see Section 2. We show in Appendix A.1 the following lemma.

**Lemma 3.** Let  $\varphi : \Omega \rightarrow \mathbb{R} \cup \{+\infty\}$  be a  $\mathbf{C}$ -convex function with  $\text{dom}(\varphi) := \{\mathbf{x} : \varphi(\mathbf{x}) < +\infty\}$ . Then, the set of points  $\mathbf{x} \in \text{dom}(\varphi)$  such that the partial derivative  $\partial_h \varphi(\mathbf{x})$  does not exist is included in a countable union of graphs of measurable functions  $F : d \in \mathbb{R} \setminus \{0\} \mapsto F(d) \in \mathbb{R}$ .

Let  $\Sigma$  be the countable union of graphs given by Lemma 3 for the Kantorovich potential  $\varphi$ . According to Assumption 1, if we let  $\Omega_0 = \Omega \setminus \Sigma$ , then  $\mu_+(\Omega_0) = 1$ .

Let  $\mathbf{x}_1 \in \Omega_0$  and let  $(\mathbf{x}_1, \mathbf{x}_2) \in \Gamma$ . Consider the function  $g_{\mathbf{x}_2} : \mathbf{x} \in \Omega \mapsto \varphi(\mathbf{x}) - \mathbf{C}(\mathbf{x}, \mathbf{x}_2)$ . As  $\varphi^{\mathbf{C}}(\mathbf{x}_2) = \varphi(\mathbf{x}_1) - \mathbf{C}(\mathbf{x}_1, \mathbf{x}_2)$ , by definition of the  $\mathbf{C}$ -transform, the function  $g_{\mathbf{x}_2}$  attains its maximum at  $\mathbf{x}_1$ . In particular, as  $\partial_{h_1} \varphi(\mathbf{x}_1)$  exists by assumption, we have

$$\partial_h \varphi(\mathbf{x}_1) = \partial_{h_1} \mathbf{C}(\mathbf{x}_1, \mathbf{x}_2) = \frac{2(h_1 - h_2)}{|d_1| + |d_2|}.$$

This implies that

$$h_2 = h_1 - \frac{|d_1| + |d_2|}{2} \partial_{h_1} \varphi(\mathbf{x}_1). \quad (18)$$

Using this expression, we find that

$$m(\mathbf{x}_1, \mathbf{x}_2) = h_1 - \frac{|d_1| \partial_{h_1} \varphi(\mathbf{x}_1)}{2}. \quad (19)$$

In particular,  $m(\mathbf{x}_1, \mathbf{x}_2)$  is uniquely determined by  $\mathbf{x}_1$ . This defines a measurable map  $\mathbf{x}_1 \in \Omega_0 \mapsto \mathbf{f}_+(\mathbf{x}_1)$ . We extend  $\mathbf{f}_+$  on  $\Omega$  by setting  $\mathbf{f}_+(\mathbf{x}_1) = 0$  if  $\mathbf{x}_1 \in \Omega \setminus \Omega_0$ . As explained in **Step 1**, for  $(\mathbf{x}_1, \mathbf{x}_2) \sim \pi^*$ , the law  $\nu$  of  $m(\mathbf{x}_1, \mathbf{x}_2)$  solves the barycenter problem  $\text{OT}_{\mathbf{C}}(\mu_+, \mu_-)$ . Hence,  $\nu = m \# \pi^* = (\text{id}, \mathbf{f}_+) \# \mu_+$  is a barycenter.

**Step 3: uniqueness of a transport map.** Likewise, we show the existence of a function  $\mathbf{f}_-$  such that  $\nu' = (\text{id}, \mathbf{f}_-) \# \mu_-$  is a barycenter. If we show that there is a unique barycenter, then  $\nu = \nu' = \nu^{\text{bar}}$ , and the theorem is proven.

We now show uniqueness of the barycenter. Let  $\nu$  be any measure that solves the barycenter problem. Let  $\pi_+$  (resp.  $\pi_-$ ) be an optimal transport plan for  $\text{OT}_{\mathbf{c}}(\mu_+, \nu)$  (resp.  $\text{OT}_{\mathbf{c}}(\mu_-, \nu)$ ). By the gluing lemma, there exists  $\rho \in \Pi(\cdot, \mu_+, \mu_-)$  whose joint law of the two first marginals is equal to  $\pi_+$ , and whose joint law of the first and last marginal is equal to  $\pi_-$ . The joint distribution  $\pi$  between the second and last marginal is a transport plan between  $\mu_+$  and  $\mu_-$ . Furthermore, as  $\nu$  is a barycenter and by definition of  $\mathbf{C}$ , we have

$$\begin{aligned} \text{OT}_{\mathbf{C}}(\mu_+, \mu_-) &= \text{OT}_{\mathbf{c}}(\mu_+, \nu) + \text{OT}_{\mathbf{c}}(\mu_-, \nu) = \int (\mathbf{c}(\mathbf{x}_1, y) + c(\mathbf{x}_2, y)) d\rho(y, \mathbf{x}_1, \mathbf{x}_2) \\ &\geq \int \mathbf{C}(\mathbf{x}_1, \mathbf{x}_2) d\pi(\mathbf{x}_1, \mathbf{x}_2), \end{aligned}$$

so that  $\pi$  is an optimal transport plan between  $\mu_+$  and  $\mu_-$ , with  $\nu = m \# \pi$ . But then, recall that (17) holds for *any* optimal transport plan  $\pi$  (for the same potential  $\varphi$ ). Hence, by the same arguments as before, we have  $\nu = (\mathbf{f}_+) \# \mu_+$  for the map  $\mathbf{f}_+ : \mathbf{x}_1 \mapsto h_1 - \frac{d_1 \partial_{h_1} \varphi(\mathbf{x}_1)}{2}$  (defined  $\mu_+$ -almost everywhere). In particular,  $\nu$  is uniquely determined by  $\mu_+$  and  $\mu_-$  through the potential  $\varphi$ .  $\square$

Theorem 4 is the counterpart of Theorem 1, established by [CDH+20b] and [GLR20] within the awareness framework. Both theorems demonstrate that the optimal fair regression function solves a barycenter problem with optimal transport costs. Remark 1 further indicates that Theorem 1 generalizes Theorem 2.3 in [CDH+20b], as the awareness framework can be considered as a special case of the unawareness framework. However, unlike in the awareness framework, there is no explicit formulation of the optimal fair prediction function in the unawareness framework, as the corresponding barycenter problem involves multi-to-one dimensional transport costs with no explicit solutions.

Theorem 4 reveals that the fair prediction  $f^*(x)$  only depends on the bi-dimensional feature  $(\eta(x), \Delta(x))$  of the candidate  $x$ . By definition,  $\Delta(x) \propto \frac{d\mu_1}{d\mu}(x) - \frac{d\mu_2}{d\mu}(x)$ . Moreover, we have  $\mathbb{P}(S=1|X=x) = p_1 \frac{d\mu_1}{d\mu}(x)$  and  $\mathbb{P}(S=2|X=x) = p_2 \frac{d\mu_2}{d\mu}(x)$ . Thus,  $\Delta(x) \propto \frac{\mathbb{P}(S=1|X=x)}{p_1} - \frac{\mathbb{P}(S=2|X=x)}{p_2}$ . In other words,  $\Delta(x)$  reflects the probability that  $x$  belongs to the different groups. Hence, in the unawareness framework, the optimal fair regression function effectively relies on estimates of  $S$  to make its prediction. This result provides a theoretical justification for the empirical observations of [LMC18]. As noted by these authors, this phenomenon may be undesirable, as it means that the predictions can rely on features not relevant to predict the response  $Y$ , simply because they are predictive of the group  $S$ .

## 4 Links between classification and regression problems

We now turn to the study of the relationship between fair regression and fair classification problems within the unawareness framework. When  $Y \in \{0, 1\}$ , classical results show that the Bayes classifier  $g_y^{\text{Bayes}}$  minimizing the risk

$$\mathcal{R}_y(g) = y \cdot \mathbb{P}[Y = 0, g(X) = 1] + (1 - y) \cdot \mathbb{P}[Y = 1, g(X) = 0].$$

is given by  $g_y^{\text{Bayes}}(x) = \mathbf{1}\{f^{\text{Bayes}}(x) \geq y\}$ , where  $f^{\text{Bayes}}$  is the Bayes regression function minimizing  $\mathcal{R}_{sq}$ . Similarly, recent results by [GSC23] demonstrate that in the awareness framework, the optimal fair classifier  $g_y^*$  minimizing the risk  $\mathcal{R}_y$  is given by  $g_y^*(x, s) = \mathbf{1}\{f^*(x, s) \geq y\}$ , where  $f^*$  is the optimal fair regression function minimizing  $\mathcal{R}_{sq}$ . These results can be leveraged to obtain plug-in classifiers  $\hat{g}$  using estimates  $\hat{f}$  of the regression function.

Somewhat less explored is the converse relationship: given a family of optimal classifiers  $(g_y)_{y \in [0,1]}$  for the risks  $(\mathcal{R}_y)_{y \in [0,1]}$ , one could define a regression function  $f$  of the form  $f(x) = \sup\{y : g_y(x) = 1\}$ . For example, this formulation yields the Bayes regression function when using Bayes classifiers and the optimal fair regression function when using optimal fair classifiers in the awareness framework. In both examples, this relationship may not be particularly useful since there already exists an explicit characterization of the optimal regression function. However, if this relationship were to hold in the unawareness framework, it would be significantly more valuable. Indeed, Theorem 4 rephrases the fair regression problem as a barycenter problem with optimal transport costs but does not provide an explicit solution.

[ADW19] proposed leveraging this relationship to address the problem of fair regression using cost-sensitive classifiers. The authors demonstrate an equivalence between minimizing a discretized version of the risk  $\mathcal{R}_{sq}$  and minimizing the average of the cost-sensitive risks  $(\mathcal{R}_y(g_{f,y}))_{y \in \mathcal{Z}}$  for a finite set  $\mathcal{Z}$ , where  $g_{f,y}$  is defined as  $g_{f,y}(x) = \mathbf{1}\{f(x) \geq y\}$ . To obtain the optimal fair regression function for this discretized risk, the authors assume access to an oracle that returns the regression function  $f$  such that  $g_{f,y}$  minimizes the average of the risks  $(\mathcal{R}_y(g_{f,y}))_{y \in \mathcal{Z}}$ . We emphasize that minimizing the average of the risks  $(\mathcal{R}_y(g_{f,y}))_{y \in \mathcal{Z}}$  remains an open and challenging problem.

In contrast, to define a regression function of the form  $f(x) = \sup\{y : g_y(x) = 1\}$ , one only needs to solve independent cost-sensitive classification problems. Recent results by [GSC23] offer an explicit characterization of these classifiers. This raises the intriguing possibility of constructing the optimal fair regression function in the unawareness framework using these fair classifiers. In this section, we demonstrate that such a construction is not always possible. To do so, we begin by providing some reminders on fair classification in the unawareness framework.

### 4.1 Fair classification

In this section, we assume that  $Y \in \{0, 1\}$  almost surely. We consider the problem of minimizing a family of risk measures  $\mathcal{R}_y$  under the demographic parity constraint. We show that the optimal fair classifier for the risk  $\mathcal{R}_y$  is of the form  $g_y^\kappa$  for some  $\kappa \in \mathbb{R}$ , where  $g_y^\kappa$  is given by

$$\forall x \in \mathcal{X}, \quad g_y^\kappa(x) = \mathbf{1}\{\eta(x) \geq y + \kappa \Delta(x)\}. \quad (20)$$

The following proposition extends Proposition 5.3 in [GSC23], and characterizes the optimal fair classifier.

**Proposition 2.** *Let  $y \in \mathbb{R}$ , and let  $\kappa^* \in \mathbb{R}$  verify*

$$\mu_+(\eta(X) \geq y + \kappa^* \Delta(X)) = \mu_-(\eta(X) \geq y + \kappa^* \Delta(X)).$$

*Under Assumption 1,  $g_y^{\kappa^*}$  solves the fair classification problem*

$$\begin{cases} \text{minimize} & \mathcal{R}_y(g) \\ \text{such that} & \mathbb{E}[g(X)|S = 1] = \mathbb{E}[g(X)|S = 2]. \end{cases} \quad (C_y)$$

*Moreover, all solutions to  $(C_y)$  are a.s. equal to  $g_y^{\kappa^*}$  on  $\mathcal{X} \setminus \{x \in \mathcal{X} : \eta(x) = y \text{ and } \Delta(x) = 0\}$ .*

The optimal classifier is uniquely defined outside of  $\{\eta(x) = y\}$ . While the set  $\{\eta(x) = y \text{ and } \Delta(x) \neq 0\}$  has null measure under Assumption 1, the set  $\{\eta(x) = y \text{ and } \Delta(x) = 0\}$  may have positive measure. On this set, the classifiers  $\mathbf{1}\{\eta(x) \geq y\}$  and  $\mathbf{1}\{\eta(x) > y\}$  differ, yet they are both optimal for the risk  $\mathcal{R}_y(g)$ .

The proof of this proposition is postponed to Appendix A.2. As discussed in the previous section, Assumption 1 encompasses as a special case the awareness framework. In this case, the optimal classifier presented in Proposition 2 reduces to the optimal fair classifier in the awareness framework given by

$$\forall(x, s) \in \mathcal{X} \times \{1, 2\}, \quad g_y^{\text{aware}}(x, s) = \begin{cases} \mathbf{1}\{\eta(x) \geq y + \frac{\kappa^*}{p_1}\} & \text{if } s = 1 \\ \mathbf{1}\{\eta(x) \geq y - \frac{\kappa^*}{p_2}\} & \text{if } s = 2, \end{cases} \quad (21)$$

as described in [SC21, ZDC22a].

In the unawareness framework, the optimal classifier relies on the probability that the observation  $X$  belongs to the different groups  $\Delta(X)$ . This behavior is similar to that of the optimal fair regression function, as established in Theorem 4. Next, we investigate whether the optimal classifier is envy-free.

**Definition 5** (Envy-free classifiers). *We say that a classifier  $g : \mathcal{X} \rightarrow \{0, 1\}$  is envy-free within group if  $\mathbb{P} \otimes \mathbb{P}$ -a.s., for  $(X, S)$  and  $(X', S')$  such that  $S = S'$  and  $g^{\text{Bayes}}(X) > g^{\text{Bayes}}(X')$ , we have  $g(X) \geq g(X')$ .*

In essence, this property ensures that no candidate who would have been accepted by the Bayes classifier but is rejected after fairness correction envies another candidate *from the same group* who would have been rejected by the Bayes classifier but is accepted after fairness correction. Note that this property is weaker than order preservation, as a classifier that preserves order is envy-free within groups.

Proposition 3 reveals that in the unawareness framework, the optimal fair classifier is generally not envy free. This behavior contrasts with that of optimal fair classifiers in the awareness framework: indeed, since these classifiers preserve order, they are also envy-free.

**Proposition 3.** *Let  $y \in \mathbb{R}$ . Under Assumption 1, if  $\mathbb{P}(S = s|X = x) \in (0, 1)$  for all  $s \in \mathcal{S}$ ,  $x \in \mathcal{X}$ , then one of the following cases hold:*

1.  $g_y^{\text{Bayes}}(x) = 1 \implies g_y^{\kappa^*}(x) = 1$   $\mu$ -a.s.
2.  $g_y^{\kappa^*}(x) = 1 \implies g_y^{\text{Bayes}}(x) = 1$   $\mu$ -a.s.
3. the classifier  $g_y^{\kappa^*}$  is not envy-free within group.

*Proof.* Assume that 1. and 2. do not hold. Then, we have  $\kappa^* \neq 0$ , and we can assume without loss of generality that  $\kappa^* > 0$ . Since 1. does not hold, we have that  $\mu(g_y^{\text{Bayes}}(X) = 1 \text{ and } g_y^{\kappa^*}(X) = 0) > 0$ . This implies in turn that  $\mu_+(g_y^{\text{Bayes}}(X) = 1 \text{ and } g_y^{\kappa^*}(X) = 0) > 0$ , since  $g_y^{\text{Bayes}}$  and  $g_y^{\kappa^*}$  coincide on  $\mathcal{X}_-$ , and since by definition when  $\kappa^* > 0$ , we have  $\mu_-(g_y^{\text{Bayes}}(X) = 1 \text{ and } g_y^{\kappa^*}(X) = 0) = 0$ . Similarly, we can show that since 2. does not hold,  $\mu_-(g_y^{\text{Bayes}}(X) = 0 \text{ and } g_y^{\kappa^*}(X) = 1) > 0$ . Now,  $\mathbb{P}(S = s|X = x) \in (0, 1)$  for all  $s \in \mathcal{S}$ , so  $\mu_1 \gg \mu_+$  and  $\mu_1 \gg \mu_-$ . Therefore,  $\mu_1(g_y^{\text{Bayes}}(X) = 1 \text{ and } g_y^{\kappa^*}(X) = 0) > 0$ , and  $\mu_1(g_y^{\text{Bayes}}(X) = 0 \text{ and } g_y^{\kappa^*}(X) = 1) > 0$ . This implies

$$\mathbb{P}\left(g_y^{\text{Bayes}}(X) > g_y^{\text{Bayes}}(X') \text{ and } g_y^{\kappa^*}(X) < g_y^{\kappa^*}(X') \mid S = S' = 1\right) > 0$$

which concludes the proof.  $\square$

**Extending cost-sensitive classification to  $\mathcal{Y} = \mathbb{R}$**  In the following, we consider the more general case where  $\mathcal{Y} = \mathbb{R}$ . Although the interpretation in terms of optimal classification is no longer applicable, we can still analyze the family of functions  $g_y^{\kappa^*}$  defined in Equation (20). The following proposition characterizes the values of the parameter  $\kappa^*$  such that  $g_y^{\kappa^*}$  satisfies demographic parity. These values partition the feature space equally, see Figure 1.



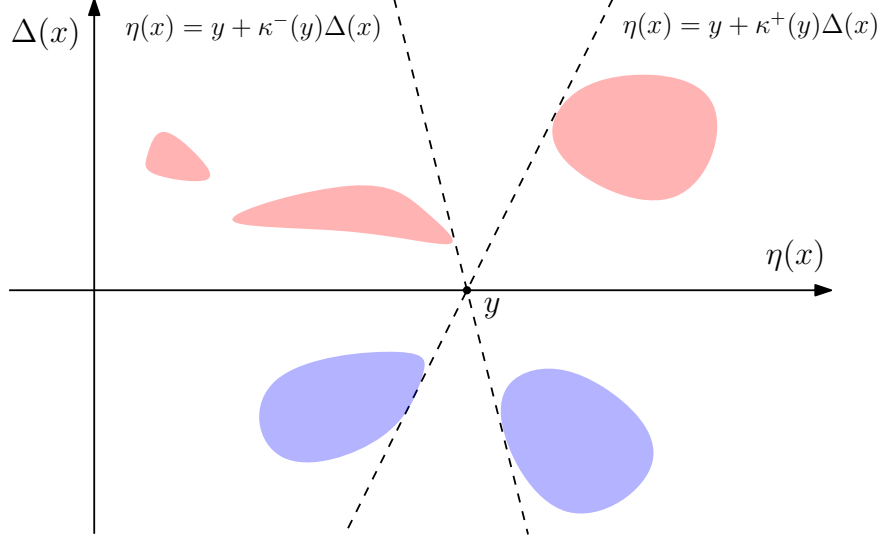


Figure 1: The measure  $\mu_+$  is displayed in red and the measure  $\mu_-$  is displayed in blue. By definition of  $\kappa^+(y)$  and  $\kappa^-(y)$ , the red region and the blue region to the right of the two dotted lines have equal masses. The region in between the two lines contains no mass.

**Proposition 4.** *Let  $y \in \mathbb{R}$ . Under Assumption 1, the set of numbers  $\kappa \in \mathbb{R}$  such that*

$$\mu_+(\eta(X) \geq y + \kappa\Delta(X)) = \mu_-(\eta(X) \geq y + \kappa\Delta(X)) \quad (22)$$

*is a nonempty closed interval  $I(y) = [\kappa^-(y), \kappa^+(y)]$ . The function  $y \mapsto \kappa^+(y)$  is upper semicontinuous and the function  $y \mapsto \kappa^-(y)$  is lower semicontinuous. Moreover, it holds that for  $\mu$ -almost all  $x$ , for all  $y \in \mathbb{R}$  and all  $\kappa, \kappa' \in I(y)$ ,  $g_y^\kappa(x) = g_y^{\kappa'}(x)$ .*

*Proof.* Introduce the function

$$G : (\kappa, y) \mapsto \mu_+(\eta(X) \geq y + \kappa\Delta(X)) - \mu_-(\eta(X) \geq y + \kappa\Delta(X)).$$

Under Assumption 1, the measures  $\mu_+$  and  $\mu_-$  give zero mass to non-horizontal lines, implying that the function  $G$  is continuous. Furthermore, for  $y \in \mathbb{R}$ , the function  $G(\cdot, y)$  is nonincreasing (recall that  $\Delta(X) < 0$  for  $X \sim \mu_-$ ). Hence, its zeroes form a closed interval  $I(y)$ . For  $\kappa \in \mathbb{R}$ , the set  $\{y \in \mathbb{R} : \kappa^-(y) > \kappa\}$  is equal to the set  $\{y \in \mathbb{R} : G(\kappa, y) > 0\}$ , which is an open set by continuity of  $G$ . This proves that  $\kappa^-$  is lower semicontinuous. We prove similarly that  $\kappa^+$  is upper semicontinuous.

It remains to prove the last statement. Fix  $y \in \mathbb{R}$ . First, we may assume without loss of generality that  $\kappa = \kappa^-(y)$  and that  $\kappa' = \kappa^+(y)$ . We have  $\mu_+(\eta(X) \geq y + \kappa^+(y)\Delta(X)) = \mu_+(\eta(X) \geq y + \kappa^-(y)\Delta(X))$  (and likewise for  $\mu_-$ ). Thus, we have

$$\begin{aligned} \mu_+(g_y^{\kappa'}(X) \neq g_y^\kappa(X)) &= \mu_+(g_y^\kappa(X) = 1, g_y^{\kappa'}(X) = 0) \\ &= \mu_+\left(\frac{\eta(X) - y}{\Delta(X)} \in [\kappa, \kappa']\right) = 0. \end{aligned}$$

The same equality holds for  $\mu_-$ . Also, the equality  $g_y^{\kappa^+(y)}(x) = g_y^{\kappa^-(y)}(x)$  holds on  $\mathcal{X}_=$  (as  $\Delta(x) = 0$  on  $\mathcal{X}_=$ ). Hence, for a fixed  $y$ , the equality  $g_y^{\kappa^+(y)}(x) = g_y^{\kappa^-(y)}(x)$  holds for  $\mu$ -almost all  $x$ .

However, the set of points  $x$  (of full measure) where this equality is satisfied depends on  $y$ , so that it is not trivial to show that this equality holds simultaneously for all  $y \in \mathbb{R}$ , almost surely.

To do so, we show that the set  $\{(h, d) \in \Omega : \exists y \in \mathbb{R}, y + \kappa^-(y)d \leq h \leq y + \kappa^+(y)d\}$  has mass 0 under  $\mu_+$  and  $\mu_-$ . For  $y \in \mathbb{R}$ , let  $C_y = \{(h, d) \in \Omega : y + \kappa^-(y)d \leq h \leq y + \kappa^+(y)d\}$ . We have previously shown that for any given  $y \in \mathbb{R}$ ,  $g_y^{\kappa^-(y)}(x) = g_y^{\kappa^+(y)}(x)$  for  $\mu$ -almost all  $x$ , implying that  $\mu_\pm(C_y) = 0$ . Let  $C = \bigcup_{y \in \mathbb{R}} C_y$ .

Let us show that  $\mu_+(C) = 0$ . Let  $C_1 = \bigcup_{y \in \mathbb{R}} \mathring{C}_y$ . First, it holds that  $\mu_+(C_1) = 0$ . If it were not the case, as the measure  $\mu_+$  is inner regular, there would exist a compact set  $K \subset C_1$  with  $\mu_+(K) > 0$ . But then, the compact set  $K$  is covered by the family of open sets  $(\mathring{C}_y)_{y \in \mathbb{R}}$ . By compactness, there exists a finite cover  $\mathring{C}_{y_1}, \dots, \mathring{C}_{y_N}$  covering  $K$ . As each  $\mathring{C}_{y_i}$  has zero mass, we obtain a contradiction with the positivity of  $\mu_+(K)$ . Furthermore, if  $(h, d) \in C \setminus C_1$ , then there exists  $y_0$  with either  $h = y_0 + \kappa^-(y_0)d$  or  $h = y_0 + \kappa^+(y_0)d$ , with also  $y + \kappa^-(y)d \leq h \leq y + \kappa^+(y)d$  for all  $y \in \mathbb{R}$ . This implies that  $C \setminus C_1$  is included in the union of the graphs of the two functions  $d \mapsto \sup_y (y + \kappa^-(y)d)$  and  $d \mapsto \inf_y (y + \kappa^+(y)d)$ . These two functions can easily be seen to be measurable because of the semicontinuity of  $\kappa^-$  and  $\kappa^+$ . Thus, by Assumption 1,  $\mu_+(C \setminus C_1) = 0$ . In conclusion, we have proven that  $\mu_+(C) = 0$ . We show likewise that  $\mu_-(C) = 0$ .  $\square$

## 4.2 The nestedness assumption

Recall that our goal is to determine whether the optimal fair regression function can be expressed as  $f^*(x) = \sup\{y : g_y^{\kappa(y)}(x) = 1\}$  for a certain choice  $\kappa(y) \in I(y)$ . In this section, we introduce an assumption regarding the family of classifiers  $g_y^{\kappa(y)}$  and demonstrate that this assumption is necessary for the relationship to hold. Specifically, we wish to use the decision boundaries of optimal classifiers at different risk levels to define the regression function. For this to be possible, the function  $y \mapsto g_y^{\kappa(y)}(x)$  must be nonincreasing for any choice of  $x$ : in other words, the rejection regions  $\{x : g_y^{\kappa(y)}(x) < y\}$  must be nested. We formalize this assumption in the following definition.

**Definition 6** (Nestedness). *We say that the problem corresponding to  $(X, Y, S) \sim \mathbb{P}$  is nested if there exists a choice of  $\kappa(y) \in I(y)$  for all  $y \in \mathbb{R}$  such that*

$$\text{for } \mu\text{-almost all } x \in \mathcal{X}, \text{ the map } y \in \mathbb{R} \mapsto g_y^{\kappa(y)}(x) \text{ is nonincreasing.} \quad (\mathbf{Nested})$$

A straightforward (but key) property implied by nestedness is the fact that the sets

$$\forall y \in \mathbb{R}, A(y) = \{x \in \mathcal{X} : \eta(x) < y + \kappa(y) \cdot \Delta(x)\} \quad (23)$$

are ‘‘almost’’ nested, in the sense that there exists a set  $\tilde{\mathcal{X}}$  of full  $\mu$ -measure such that for all  $y' \leq y$ , it holds that  $A(y') \cap \tilde{\mathcal{X}} \subseteq A(y) \cap \tilde{\mathcal{X}}$ .

**Lemma 4.** *Assume that Assumption 1 holds. Then, the problem is nested with choice  $\kappa(y) \in I(y)$  for all  $y \in \mathbb{R}$  if and only if for all  $y < y'$ ,*

$$\mu(g_y^{\kappa(y)}(X) = 0 \text{ and } g_{y'}^{\kappa(y')}(X) = 1) = 0. \quad (24)$$

*Furthermore, if the problem is nested, one can always choose  $\kappa(y) = \kappa^+(y)$  for all  $y \in \mathbb{R}$  in the definition of  $g_y^{\kappa(y)}$ .*

*Proof.* The direct implication is clear. For the converse one, assume that the nestedness assumption does not hold. By definition, there exists a measurable set  $\mathcal{X}_0$  of positive  $\mu$ -mass such that  $y \mapsto g_y^{\kappa(y)}(x)$  is not nonincreasing for all  $x \in \mathcal{X}_0$ . It holds that either  $\mu_+(\mathcal{X}_0) > 0$  or that  $\mu_-(\mathcal{X}_0) > 0$ . Assume without loss of generality that the first condition is satisfied, and let  $\tilde{\mathcal{X}}$  be the set of points  $x$  in  $\mathcal{X}_0 \cap \mathcal{X}_+$  that satisfy

$$\forall y \in \mathbb{R}, \frac{\eta(x) - y}{\Delta(x)} \notin [\kappa^-(y), \kappa^+(y)] \quad (25)$$

According to Proposition 4,  $\mu_+(\tilde{\mathcal{X}}) = \mu_+(\mathcal{X}_0 \cap \mathcal{X}_+) > 0$ . For  $x \in \tilde{\mathcal{X}}$ , there exists  $y < y'$  with  $g_y^{\kappa(y)}(x) = 0$  and  $g_{y'}^{\kappa(y')}(x) = 1$ . As  $x$  satisfies (25), we have that

$$\begin{cases} \eta(x) < y + \kappa^-(y)\Delta(x) \\ \eta(x) > y' + \kappa^+(y')\Delta(x). \end{cases} \quad (26)$$

Because the function  $\kappa^-$  is lower semicontinuous, for  $\tilde{y}$  close enough to  $y$ , we also have  $\eta(x) < \tilde{y} + \kappa^-(\tilde{y})\Delta(x)$ . Likewise, there exists  $\tilde{y}'$  close enough to  $y'$  with  $\eta(x) > \tilde{y}' + \kappa^+(\tilde{y}')\Delta(x)$ . In conclusion, we have shown that

$$\tilde{\mathcal{X}} \subset \bigcup_{\substack{y, y' \in \mathbb{Q} \\ y < y'}} \{x : \eta(x) < y + \kappa^-(y)\Delta(x) \text{ and } \eta(x) > y' + \kappa^+(y')\Delta(x)\} \quad (27)$$

In particular, as  $\mu_+(\tilde{\mathcal{X}}) > 0$ , there exists  $y < y' \in \mathbb{Q}$  with

$$\mu_+(\eta(X) < y + \kappa^-(y)\Delta(X) \text{ and } \eta(X) > y' + \kappa^+(y')\Delta(X)) > 0.$$

According to Proposition 4 and Assumption 1, as the equality  $\eta(X) = y' + \kappa^+(y')\Delta(X)$  happens with zero  $\mu_+$ -probability, we have

$$\mu_+(g_y^{\kappa(y)}(X) = 0 \text{ and } g_{y'}^{\kappa(y')}(X) = 1) > 0,$$

proving the first claim.

The second claim follows from the characterization of nestedness that we have just established. Indeed, let  $y < y'$ . By Proposition 4, for  $\mu$ -almost all  $x$ ,  $g_y^{\kappa(y)}(x) = g_y^{\kappa^+(y)}(x)$  and  $g_{y'}^{\kappa(y')}(x) = g_{y'}^{\kappa^+(y')}(x)$ . Thus, if (24) holds for  $\kappa(y)$  and  $\kappa(y')$ , it also holds for  $\kappa^+(y)$  and  $\kappa^+(y')$ .  $\square$

As a warm-up, we begin by showing that the nestedness assumption is always verified in the awareness setting.

**Proposition 5.** *Assume that  $S$  is  $X$ -measurable. Then, under Assumption 1, the classification problem is nested.*

*Proof.* We prove Proposition 5 by contradiction. Assume that the problem is not nested. Using Lemma 4, there exist  $y < y'$  and a set  $\mathcal{X}_0$  of positive  $\mu$  probability such that for all  $x \in \mathcal{X}_0$ ,  $g_y^{\kappa(y)}(x) = 0$  and  $g_{y'}^{\kappa(y')}(x) = 1$ . Using Proposition 4, we can also assume without loss of generality that  $\eta(x) < y + \kappa^-(y)\Delta(x)$  and  $\eta(x) > y' + \kappa^+(y')\Delta(x)$  for  $x \in \mathcal{X}_0$ . Letting  $x \in \mathcal{X}_0$ , that we assume without loss of generality is in  $\mathcal{X}_+$ , the previous inequalities become

$$y' + \frac{\kappa^+(y')}{p_1} < \eta(x) < y + \frac{\kappa^-(y)}{p_1}.$$

In words, the threshold for admission is lower at level  $y'$  than at level  $y$ . This implies in particular that  $\mu_+(g_{y'}^{\kappa^+(y')}(X) = 1) \geq \mu_+(g_y^{\kappa^-(y)}(X) = 1)$ . On the other hand, since  $y < y'$ , it also implies that  $\kappa^+(y') < \kappa^-(y)$ . Therefore,  $y - \frac{\kappa^-(y)}{p_2} < y' - \frac{\kappa^+(y')}{p_2}$ , so  $\mu_-(g_{y'}^{\kappa^+(y')}(X) = 1) \leq \mu_-(g_y^{\kappa^-(y)}(X) = 1)$ . Using  $\mu_+(g_{y'}^{\kappa^+(y')}(X) = 1) = \mu_-(g_{y'}^{\kappa^+(y')}(X) = 1)$  and  $\mu_+(g_y^{\kappa^-(y)}(X) = 1) = \mu_-(g_y^{\kappa^-(y)}(X) = 1)$ , we find that  $\mu_+(g_{y'}^{\kappa^+(y')}(X) = 1) = \mu_+(g_y^{\kappa^-(y)}(X) = 1)$ . It implies that

$$\mu_+ \left( \eta(X) \in \left[ y' + \frac{\kappa^+(y')}{p_1}, y + \frac{\kappa^-(y)}{p_1} \right] \right) = 0.$$

Likewise,

$$\mu_- \left( \eta(X) \in \left[ y - \frac{\kappa^-(y)}{p_2}, y' - \frac{\kappa^+(y')}{p_2} \right] \right) = 0.$$

In particular, for  $\kappa = \kappa^+(y') + p_1(y' - y)$ , we see that  $y + \frac{\kappa}{p_1} = y' + \frac{\kappa^+(y')}{p_1} < y + \frac{\kappa^-(y)}{p_1}$ . Thus,  $\kappa < \kappa^-(y)$ , and  $y - \frac{\kappa}{p_2} \geq y - \frac{\kappa^-(y)}{p_2}$ . Moreover,  $\kappa > \kappa^+(y')$ , and  $y < y'$ , so  $y - \frac{\kappa}{p_2} < y' - \frac{\kappa^+(y')}{p_2}$ . This implies that

$$y' - \frac{\kappa^+(y')}{p_2} - \left( y - \frac{\kappa}{p_2} \right) = y' - y + \frac{1}{p_2} (\kappa^+(y') + p_1(y' - y) - \kappa^+(y')) > 0.$$

so  $y - \frac{\kappa}{p_2} \in \left[ y - \frac{\kappa^-(y)}{p_2}, y' - \frac{\kappa^+(y')}{p_2} \right]$ . Thus,

$$\mu_+ \left( \eta(X) \geq y + \frac{\kappa}{p_1} \right) = \mu_- \left( \eta(X) \geq y - \frac{\kappa}{p_2} \right)$$

and  $\kappa \in I(y)$ . Since  $\kappa < \kappa^-(y)$ , this yields a contradiction.  $\square$

Somewhat surprisingly, although the nestedness assumption may appear intuitive, it is not always verified. In Section 5, we present examples where this assumption holds and others where it does not.

Before proving in the next section that under the nestedness assumption, the optimal fair classification functions  $g_y^{\kappa(y)}$  can be recovered by thresholding the optimal fair regression function  $f^*$ , we prove the converse: if the problem is not nested, there exists a value of  $y \in \mathbb{R}$  where the classifier  $\mathbf{1}\{f^*(x) \geq y\}$  is suboptimal for the fair classification problem  $(C_y)$ .

**Proposition 6.** *Assume that  $\mathcal{Y} = \{0, 1\}$ , that Assumption 1 holds and that the classification problem is not nested. Let  $f^*$  be the optimal fair regression function in the unawareness framework. Then, there exists  $y \in \mathbb{R}$  such that the classifier  $x \mapsto \mathbf{1}\{f^*(x) \geq y\}$  is not the optimal fair classifier for the risk  $\mathcal{R}_y$ .*

*Proof.* According to Lemma 4, there exists  $y < y'$  with

$$\mu(g_y^{\kappa(y)}(X) = 0 \text{ and } g_{y'}^{\kappa(y')}(X) = 1) > 0.$$

Let  $\mathcal{X}_0$  be the set corresponding to this event. Let us consider a classifier of the form  $g_y(x) = \mathbf{1}\{f(x) \geq y\}$ . On the one hand, if  $\mu(X \in \mathcal{X}_0 \text{ and } f(X) < y) > 0$ , then the probability  $\mu(X \in \mathcal{X}_0 \text{ and } f(X) < y')$  is also positive, so  $g_{y'}(x)$  and  $g_{y'}^{\kappa(y')}(x)$  disagree on a set of positive probability. Now, Proposition 2 implies that all optimal classifiers are a.s. equal, so  $g_{y'}$  is sub-optimal. On the other hand, if  $\mu(X \in \mathcal{X}_0 \text{ and } f(X) < y) = 0$ , then  $g_y(x) = 1$  for almost all  $x \in \mathcal{X}_0$ . This implies that  $g_y(x)$  and  $g_y^{\kappa(y)}(x)$  disagree on a set of positive probability, so  $g_y$  is sub-optimal.  $\square$

### 4.3 Constructing a regression function using nested classifiers

In the previous section, we proved that under mild assumptions, nestedness is a necessary condition for the relationship  $g_y^*(x) = \mathbf{1}\{f^*(x) \geq y\}$  between the optimal fair classification and regression functions to hold. We now conclude by showing that nestedness is also a sufficient condition for this relationship to hold.

We begin by defining the function  $f^* : \mathcal{X} \rightarrow \mathbb{R}$  as

$$\forall x \in \mathcal{X}, f^*(x) = \sup\{y : g_y^{\kappa(y)}(x) = 1\} \quad (28)$$

where  $g_y^{\kappa(y)}$  is given by Equation (20). We assume without loss of generality (using Lemma 4) that  $\kappa(y) = \kappa^+(y)$  for all  $y \in \mathbb{R}$ . Remark that  $f^*$  is then almost measurable because of the upper semicontinuity of  $\kappa^+$ , in the sense that its restriction to some set of full measure is measurable (here given by the set of full measure where  $y \mapsto g_y^{\kappa(y)}(x)$  is nonincreasing).

**Theorem 5.** *Assume that the classification problem is nested and that Assumption 1 is satisfied. Then, the regression function  $f^*$  is optimal for the fair regression problem (1).*

Before proving Theorem 5, we state the following corollary.

**Corollary 1.** *Assume that the classification problem is nested, that Assumption 1 is satisfied, and that  $\mathcal{Y} = \{0, 1\}$ . Then, the classification function  $g_y : y \mapsto \mathbf{1}\{f^*(x) \geq y\}$  is optimal for the fair classification problem with cost  $\mathcal{R}_y$ , where  $f^*$  is the solution to the fair regression problem (1).*

The proof of Corollary 1 follows immediately by noticing that by Theorem 4,  $f^*$  is uniquely defined, and that the nestedness assumption and Theorem 5 imply that  $g_y(x) = g_y^{\kappa(y)}(x)$  a.s.

The rest of the section is devoted to proving Theorem 5. To do so, we begin by proving that  $f^*$  is a fair regression function, and by defining  $F$ , the c.d.f. of the predictions under  $\mu_+$  and  $\mu_-$ .

**Lemma 5.** *Assume that the problem is nested and that Assumption 1 is satisfied. Let  $F : \mathbb{R} \rightarrow \mathbb{R}$  be defined by*

$$\forall y \in \mathbb{R}, F(y) = \mu_+(\eta(X) \leq y + \kappa(y)\Delta(X)) = \mu_-(\eta(X) \leq y + \kappa(y)\Delta(X)). \quad (29)$$

*Then, there exists a probability measure  $\nu^*$  with continuous c.d.f.  $F$  and finite second moment such that  $f^* \# \mu_+ = f^* \# \mu_- = \nu^*$ . In particular,  $f^*$  is a fair regression function.*

*Proof.* The ‘‘almost’’ nestedness of the sets  $(A(y))_y$  implies that  $F$  is nondecreasing. Let us show that the function  $F$  is the c.d.f. of some continuous random variable, i.e., that it goes to 0 in  $-\infty$ , that it goes to 1 in  $+\infty$ , and that it is continuous.

First, recall that  $\Delta(X) > 0$  for  $X \sim \mu_+$  and that  $\Delta(X) < 0$  for  $X \sim \mu_-$ . Thus, if  $\kappa(y) \leq 0$ , then  $F(y) \leq \mu_+(\eta(X) - y \leq 0)$  and if  $\kappa(y) \geq 0$ , then  $F(y) \leq \mu_-(\eta(X) - y \leq 0)$ . Hence,

$$F(y) \leq \max\{\mu_+(\eta(X) - y \leq 0), \mu_-(\eta(X) - y \leq 0)\}, \quad (30)$$

and  $F$  converges to 0 in  $-\infty$ . Similarly,  $F(y) \rightarrow 1$  when  $y$  converges to  $+\infty$ .

Next, let us show that  $F$  is continuous. Let  $y_0, y_1 \in \mathbb{R}$  be such that  $y_0 < y_1$ . Now, if  $\kappa(y_0) \geq \kappa(y_1)$ , then

$$\begin{aligned} F(y_1) - F(y_0) &= \mu_+(\eta(X) \leq y_1 + \kappa(y_1)\Delta(X)) - \mu_+(\eta(X) \leq y_0 + \kappa(y_0)\Delta(X)) \\ &\leq \mu_+(\eta(X) \leq y_1 + \kappa(y_0)\Delta(X)) - \mu_+(\eta(X) \leq y_0 + \kappa(y_0)\Delta(X)) \end{aligned}$$

Similarly, if  $\kappa(y_0) \leq \kappa(y_1)$ , then, (recalling that  $\Delta(X) < 0$  for  $X \sim \mu_-$ )

$$\begin{aligned} F(y_1) - F(y_0) &= \mu_-(\eta(X) \leq y_1 + \kappa(y_1)\Delta(X)) - \mu_-(\eta(X) \leq y_0 + \kappa(y_0)\Delta(X)) \\ &\leq \mu_-(\eta(X) \leq y_1 + \kappa(y_0)\Delta(X)) - \mu_-(\eta(X) \leq y_0 + \kappa(y_0)\Delta(X)). \end{aligned}$$

Thus,

$$F(y_1) - F(y_0) \leq \mu_+(\eta(X) - \kappa(y_0)\Delta(X) \in [y_0, y_1]) + \mu_-(\eta(X) - \kappa(y_0)\Delta(X) \in [y_0, y_1])$$

We have shown that  $F$  is non-decreasing, so  $F(y_1) - F(y_0) \geq 0$ . Under Assumption 1,  $\mu_+$  and  $\mu_-$  give zero mass to the sets  $\{\eta(X) = y_0 + \kappa(y_0)\Delta(X)\}$ , so  $F(y_1) - F(y_0) \rightarrow 0$  as  $y_1 \rightarrow y_0^+$ . This proves that  $F$  is right-continuous. To show that  $F$  is left-continuous, we note that if  $\kappa(y_0) \geq \kappa(y_1)$ , then

$$F(y_1) - F(y_0) \leq \mu_+(\eta(X) \leq y_1 + \kappa(y_1)\Delta(X)) - \mu_+(\eta(X) \leq y_0 + \kappa(y_1)\Delta(X)).$$

Similarly, if  $\kappa(y_0) \leq \kappa(y_1)$ , then,

$$F(y_1) - F(y_0) \leq \mu_-(\eta(X) \leq y_1 + \kappa(y_1)\Delta(X)) - \mu_-(\eta(X) \leq y_0 + \kappa(y_1)\Delta(X)).$$

Thus,

$$F(y_1) - F(y_0) \leq \mu_+(\eta(X) - \kappa(y_1)\Delta(X) \in [y_0, y_1]) + \mu_-(\eta(X) - \kappa(y_1)\Delta(X) \in [y_0, y_1])$$

and  $F$  is also left-continuous.

Then, let us show that  $\nu^*$  has finite second moment. Let  $Z \sim \nu^*$ . We have

$$\mathbb{E}[Y^2] = \int_0^{+\infty} \mathbb{P}(Z^2 > t) dt = \int_0^{+\infty} (F(\sqrt{t}) + (1 - F(-\sqrt{t}))) dt$$

We use (30) to obtain that for  $y \in \mathbb{R}$ ,  $F(y) \leq \max(\mu_+(\eta(X) \leq y), \mu_-(\eta(X) \leq y))$ . But, as  $\mathbb{E}[Y^2] < +\infty$ , the random variable  $\eta(X)$  has a finite second moment under the law of either  $\mu_+$  or  $\mu_-$ . In particular,  $\int_0^{+\infty} F(\sqrt{t}) dt$  is finite. Similarly,  $\int_0^{+\infty} (1 - F(-\sqrt{t})) dt$  is finite.

Finally, we prove the statement  $f^* \# \mu_+ = \nu^*$ . Indeed, for all  $y_0 \in \mathbb{R}$ , we have using that upper semicontinuity of  $\kappa(y) = \kappa^+(y)$  that

$$\begin{aligned} f^* \# \mu_+((-\infty, y_0]) &= \mu_+(\sup\{y : g_y^{\kappa(y)}(X) = 1\} \leq y_0) \\ &= \mu_+(g_{y_0}^{\kappa(y_0)}(X) = 0) = \mu_+(\eta(X) < y_0 + \kappa(y_0)\Delta(X)) = F(y_0). \end{aligned}$$

where the second line follows from the nestedness assumption and the fact that the line  $\{\eta(X) = y_0 + \kappa(y_0)\Delta(X)\}$  has zero mass. We show similarly that  $f^* \# \mu_- = \nu^*$ , thus concluding the proof of Lemma 5.  $\square$

The function  $f^*$  depends only on  $x$  through the pair  $(\eta(x), \Delta(x))$ . Let  $\mathbf{f}^* : \Omega \rightarrow \mathbb{R}$  be defined by the relation  $f^*(x) = \mathbf{f}^*(\eta(x), \Delta(x))$  for  $x \in \mathcal{X}_\pm$ . We show that  $\mathbf{f}^*$  defines an optimal transport map between  $\mu_+$  and  $\nu^*$  with respect to the cost  $\mathbf{c}$ .

**Lemma 6.** *Assume that the problem is nested. Then,  $\mathbf{f}^*$  is an optimal transport map between  $\mu_+$  and  $\nu^*$  for the cost  $\mathbf{c}$ , with Kantorovich potential between  $\nu^*$  and  $\mu_+$  given by  $v : y \mapsto -2 \int_0^y \kappa(t) dt$ . The same holds for  $\mu_-$ , with Kantorovich potential given by  $-v$ .*

The proof of Lemma 6 relies on the following technical lemma, whose proof is postponed to Appendix A.4.

**Lemma 7.** *The function  $y \mapsto \kappa(y)$  satisfies  $|\kappa(y)| \leq C(1 + |y|)$  for some  $C > 0$ .*

*Proof.* To prove Lemma 6, we begin by remarking that the potential  $v$  is in  $L^1(\nu^*)$  because of Lemma 5 and Lemma 7. Let us now show that for almost all  $x \in \mathcal{X}_+$ ,

$$v^c(x) := \sup_{y \in \mathbb{R}} (v(y) - c(x, y)) = v(f^*(x)) - c(x, f^*(x)). \quad (31)$$

Let  $x \in \mathcal{X}_+$  be a point such that  $y \mapsto g_y^{\kappa(y)}(x)$  is nonincreasing (almost all points satisfy this condition by nestedness). We remark that  $\partial_y c(x, y) = \frac{2(y - \eta(x))}{\Delta(x)}$ . Hence,

$$\begin{aligned} c(x, y) - c(x, f^*(x)) &= \int_{f^*(x)}^y \frac{2(t - \eta(x))}{\Delta(x)} dt \\ &= -2 \int_{f^*(x)}^y \left( \frac{\eta(x) - t}{\Delta(x)} - \kappa(t) \right) dt - 2 \int_{f^*(x)}^y \kappa(t) dt \\ &= -2 \int_{f^*(x)}^y \left( \frac{\eta(x) - t}{\Delta(x)} - \kappa(t) \right) dt + v(y) - v(f^*(x)). \end{aligned}$$

Assume that  $y \geq f^*(x)$ . For  $t \in (f^*(x), y]$ , by definition of  $f^*$  and by nestedness,  $g_t^{\kappa(t)}(x) = 0$ . Thus,

$$\frac{\eta(x) - t}{\Delta(x)} - \kappa(t) < 0.$$

This implies that

$$-2 \int_{f^*(x)}^y \left( \frac{\eta(x) - t}{\Delta(x)} - \kappa(t) \right) dt \geq 0.$$

We obtain that

$$c(x, y) - c(x, f^*(x)) \geq v(y) - v(f^*(x)).$$

The same result holds when  $y < f^*(x)$ . Indeed, in that case, for all  $t \in [y, f^*(x))$ ,

$$\frac{\eta(x) - t}{\Delta(x)} - \kappa(t) \geq 0.$$

Hence,

$$-2 \int_{f^*(x)}^y \left( \frac{\eta(x) - t}{\Delta(x)} - \kappa(t) \right) dt = 2 \int_y^{f^*(x)} \left( \frac{\eta(x) - t}{\Delta(x)} - \kappa(t) \right) dt \geq 0.$$

This proves (31). This relation implies that the  $\mathbf{c}$ -transform of the function  $v \in L^1(\nu^*)$  is a function  $w : \Omega \rightarrow \mathbb{R}$  satisfying for  $\mu_+$ -almost all  $x \in \mathcal{X}_+$  (with  $\mathbf{x} = (\eta(x), \Delta(x))$ )

$$w(\mathbf{x}) = v(f^*(x)) - c(x, f^*(x)) = v(\mathbf{f}^*(\mathbf{x})) - \mathbf{c}(\mathbf{x}, \mathbf{f}^*(\mathbf{x})).$$

As  $v \in L^1(\nu^*)$ , Kantorovich duality implies

$$\begin{aligned} \text{OT}_{\mathbf{c}}(\mu_+, \nu^*) &\geq \int v(y) d\nu^*(y) - \int w(\mathbf{x}) d\mu_+(\mathbf{x}) = \int v(\mathbf{f}^*(\mathbf{x})) d\mu_+(\mathbf{x}) - \int w(\mathbf{x}) d\mu_+(\mathbf{x}) \\ &= \int \mathbf{c}(\mathbf{x}, \mathbf{f}^*(\mathbf{x})) d\mu_+(\mathbf{x}), \end{aligned}$$

see Section 2. This shows that  $\mathbf{f}^*$  is the optimal transport map between  $\mu_+$  and  $\nu^*$ . The same holds for  $\mu_-$ , where we use the potential  $-v$  instead of  $v$ : precisely, we can show that we have for almost all  $x \in \mathcal{X}_-$

$$(-v)^c(x) := \sup_{y \in \mathbb{R}} (-v(y) - c(x, y)) = -v(\mathbf{f}^*(x)) - c(x, \mathbf{f}^*(x)). \quad (32)$$

This concludes the proof of Lemma 6.  $\square$

Lemma 6 shows that  $\mathbf{f}^*$  defines an optimal transport map from  $\mu_+$  to  $\nu^*$ , and from  $\mu_-$  to  $\nu^*$ . To conclude the proof of Theorem 1, it remains to show that  $\nu^*$  is solution to the barycenter problem described in Lemma 2.

**Lemma 8.** *The distribution  $\nu^*$  is solution to the barycenter problem described in Lemma 2.*

*Proof.* Let  $\varphi : \mathbf{x}_1 \in \Omega \mapsto \mathbf{c}(\mathbf{x}_1, \mathbf{f}^*(\mathbf{x}_1)) - v(\mathbf{f}^*(\mathbf{x}_1))$  and let  $\psi : \mathbf{x}_2 \in \Omega \mapsto \mathbf{c}(\mathbf{x}_2, \mathbf{f}^*(\mathbf{x}_2)) + v(\mathbf{f}^*(\mathbf{x}_2))$ . Using (31) and (32), we see that for all  $y \in \mathbb{R}$ , for  $\mu_+$ -almost all  $\mathbf{x}_1$  and  $\mu_-$ -almost all  $\mathbf{x}_2$ , it holds that

$$\begin{aligned} \varphi(\mathbf{x}_1) + \psi(\mathbf{x}_2) &= \mathbf{c}(\mathbf{x}_1, \mathbf{f}^*(\mathbf{x}_1)) - v(\mathbf{f}^*(\mathbf{x}_1)) + \mathbf{c}(\mathbf{x}_2, \mathbf{f}^*(\mathbf{x}_2)) + v(\mathbf{f}^*(\mathbf{x}_2)) \\ &\leq \mathbf{c}(\mathbf{x}_1, y) - v(y) + \mathbf{c}(\mathbf{x}_2, y) + v(y) \\ &= \mathbf{c}(\mathbf{x}_1, y) + \mathbf{c}(\mathbf{x}_2, y). \end{aligned}$$

By taking the value  $y$  that minimizes this last term, we obtain that

$$\varphi(\mathbf{x}_1) + \psi(\mathbf{x}_2) \leq \mathbf{C}(\mathbf{x}_1, \mathbf{x}_2),$$

where  $\mathbf{C}$  is the cost function defined in (14). In particular,  $-\varphi(\mathbf{x}_1) \geq \psi^{\mathbf{C}}(\mathbf{x}_1)$ . Furthermore, remark that

$$-v(\mathbf{f}^*(\mathbf{x}_1)) \leq \varphi(\mathbf{x}_1) \leq c(\mathbf{x}_1, 0).$$

Thus, as  $v \in L^1(\nu^*)$  and  $\int \frac{h^2}{d} d\mu_+(h, d) < +\infty$ , it holds that  $\varphi \in L^1(\mu_+)$ . Likewise,  $\psi \in L^1(\mu_-)$ . By Kantorovich duality, it holds that

$$\begin{aligned} \text{OT}_{\mathbf{C}}(\mu_+, \mu_-) &\geq \int \psi(\mathbf{x}_2) d\mu_-(\mathbf{x}_2) - \int \varphi(\mathbf{x}_1) d\mu_+(\mathbf{x}_1) \\ &\geq \int \psi(\mathbf{x}_2) d\mu_-(\mathbf{x}_2) + \int \varphi(\mathbf{x}_1) d\mu_+(\mathbf{x}_1) \\ &= \int \mathbf{c}(\mathbf{x}_1, \mathbf{f}^*(\mathbf{x}_1)) d\mu_+(\mathbf{x}_1) + \int \mathbf{c}(\mathbf{x}_2, \mathbf{f}^*(\mathbf{x}_2)) d\mu_-(\mathbf{x}_2) \\ &\quad - \int v(\mathbf{f}^*(\mathbf{x}_1)) d\mu_+(\mathbf{x}_1) + \int v(\mathbf{f}^*(\mathbf{x}_2)) d\mu_-(\mathbf{x}_2) \\ &= \text{OT}_{\mathbf{c}}(\mu_+, \nu^*) + \text{OT}_{\mathbf{c}}(\mu_-, \nu^*) + \int v(y) d\nu^*(y) - \int v(y) d\nu^*(y) \\ &= \text{OT}_{\mathbf{c}}(\mu_+, \nu^*) + \text{OT}_{\mathbf{c}}(\mu_-, \nu^*) \geq \text{OT}_{\mathbf{C}}(\mu_+, \mu_-). \end{aligned}$$

This proves that  $\nu^*$  is the solution to the barycenter problem, and that  $\mathbf{f}^*$  is an optimal regression function.  $\square$

## 5 Building examples and counterexamples

In the previous section, we proved that under mild assumptions, the relationship  $g_y^*(x) = \mathbf{1}\{f^*(x) \geq y\}$  only holds under the nestedness assumption. In this section, we now explain how to build large classes of triplets  $(X, Y, S) \in \mathcal{X} \times \mathbb{R} \times \{1, 2\}$  whose distributions  $\mathbb{P}$  either satisfy or do not satisfy this criterion. The starting point of our approach consisted in associating to each distribution  $\mathbb{P}$  a pair of distributions  $(\mu_+, \mu_-) = (\mu_+(\mathbb{P}), \mu_-(\mathbb{P}))$  on  $\Omega$ , where we recall that  $\mu_{\pm}(\mathbb{P})$  is the distribution of  $(\eta(X), \Delta(X))$  when  $X \sim \mu_{\pm}$ . Then, both the optimal fair regression function and the nestedness criterion are best understood in terms of the pair  $(\mu_+(\mathbb{P}), \mu_-(\mathbb{P}))$ .

However, given a pair of measure  $(\boldsymbol{\mu}_+, \boldsymbol{\mu}_-)$  on  $\Omega$ , it is not a priori clear whether there exists a triplet  $(X, Y, S) \sim \mathbb{P}$  with  $(\boldsymbol{\mu}_+, \boldsymbol{\mu}_-) = (\boldsymbol{\mu}_+(\mathbb{P}), \boldsymbol{\mu}_-(\mathbb{P}))$ . We give a definitive answer to this problem by providing a list of necessary and sufficient conditions for the existence of such a probability distribution  $\mathbb{P}$ . We then use this theoretical result to build probability distributions  $\mathbb{P}$  for which the associated fair classification problem is either nested or not nested.

Let  $\mathbb{P}$  be the distribution of a triplet  $(X, Y, S) \in \mathcal{X} \times \mathbb{R} \times \{1, 2\}$ , with  $\mathbb{E}[Y^2] < +\infty$ . Let  $(\boldsymbol{\mu}_+, \boldsymbol{\mu}_-) = (\boldsymbol{\mu}_+(\mathbb{P}), \boldsymbol{\mu}_-(\mathbb{P}))$  be the associated pair of measures on  $\Omega$ . Then, it always holds that

$$\int_{\Omega} |d|^{-1} d\boldsymbol{\mu}_+(h, d) = \int_{\mathcal{X}} \frac{1}{\Delta(x)} d\mu_+(x) = \int_{\mathcal{X}} \frac{d\mu}{d\mu_+}(x) d\mu_+(x) = \mu(\mathcal{X}_+),$$

while  $\int_{\Omega} |d|^{-1} d\boldsymbol{\mu}_-(h, d) = \mu(\mathcal{X}_-)$ . In particular,

$$0 < \int_{\Omega} |d|^{-1} d\boldsymbol{\mu}_+(h, d) + \int_{\Omega} |d|^{-1} d\boldsymbol{\mu}_-(h, d) \leq 1. \quad (33)$$

Also, note that  $\mu = p_1\mu_1 + p_2\mu_2$ , so that  $\Delta(x) = \frac{d\mu_+}{d\mu}(x) \leq \frac{1}{p_1 m}$  when  $x \in \mathcal{X}_+$ , whereas  $\Delta(x) \geq -\frac{1}{p_2 m}$  when  $x \in \mathcal{X}_-$  (recall that  $m$  is the mass of the measure  $(\mu_1 - \mu_2)_+$ ). In particular, the supports of  $\boldsymbol{\mu}_+$  and  $\boldsymbol{\mu}_-$  are located in an horizontal strip of the form  $\{(h, d) : -M \leq d \leq M\}$  for some  $M > 0$ . The next proposition states that these two conditions are actually sufficient for the existence of a probability measure  $\mathbb{P}$  with  $(\boldsymbol{\mu}_+, \boldsymbol{\mu}_-) = (\boldsymbol{\mu}_+(\mathbb{P}), \boldsymbol{\mu}_-(\mathbb{P}))$ .

**Proposition 7.** *Assume that  $\mathcal{X}$  is an uncountable standard Borel space (e.g.,  $\mathcal{X} = [0, 1]$ ). Then, the set of pairs of measures  $(\boldsymbol{\mu}_+(\mathbb{P}), \boldsymbol{\mu}_-(\mathbb{P}))$  that can be obtained from a distribution  $\mathbb{P}$  of a triplet  $(X, Y, S) \in \mathcal{X} \times \mathbb{R} \times \{1, 2\}$  with  $\mathbb{E}[Y^2] < \infty$  is exactly equal to the set of pairs  $(\boldsymbol{\mu}_+, \boldsymbol{\mu}_-)$  supported on bounded horizontal strips, satisfying Equation (33), with  $\boldsymbol{\mu}_+$  supported on  $\{d > 0\}$  and  $\boldsymbol{\mu}_-$  supported on  $\{d < 0\}$ .*

This proposition allows us to easily build examples where either nestedness or nonnestedness is satisfied: one does not need to build from scratch a joint distribution on  $\mathcal{X} \times \mathbb{R} \times \{1, 2\}$ , but can simply define a pair of measures  $(\boldsymbol{\mu}_+, \boldsymbol{\mu}_-)$  on  $\Omega$ . As long as this pair satisfies the conditions given in Proposition 7, the existence of a probability distribution  $\mathbb{P}$  such that  $(\boldsymbol{\mu}_+, \boldsymbol{\mu}_-) = (\boldsymbol{\mu}_+(\mathbb{P}), \boldsymbol{\mu}_-(\mathbb{P}))$  is ensured.

*Proof.* We have already established that the pairs of measures  $(\boldsymbol{\mu}_+(\mathbb{P}), \boldsymbol{\mu}_-(\mathbb{P}))$  satisfy the conditions stated in Proposition 7. Reciprocally, consider a pair  $(\boldsymbol{\mu}_+, \boldsymbol{\mu}_-)$  satisfying Equation (33), supported on bounded horizontal strips, with  $\boldsymbol{\mu}_+$  supported on  $\{d > 0\}$  and  $\boldsymbol{\mu}_-$  supported on  $\{d < 0\}$ . Let  $a_{\pm} = \int_{\Omega} |d|^{-1} d\boldsymbol{\mu}_{\pm}(h, d)$ .

Due to the Borel isomorphism theorem,  $\mathcal{X}$  is Borel isomorphic to  $\mathbb{R}^2$ , so we may assume without loss of generality that  $\mathcal{X} = \mathbb{R}^2$ . Let  $\mathcal{X}_+ = \{(h, d) \in \mathbb{R}^2 : d > 0\}$ ,  $\mathcal{X}_- = \{(h, d) \in \mathbb{R}^2 : d < 0\}$  and  $\mathcal{X}_= = \{(h, 0) : h \in \mathbb{R}\}$ . Let  $\mu_+ = \delta_0$ . We let  $\mu_+ = \boldsymbol{\mu}_+$  and  $\mu_- = \boldsymbol{\mu}_-$ .

Let

$$d\mu(h, d) = \frac{1}{|d|} d\boldsymbol{\mu}_+(h, d) + \frac{1}{|d|} d\boldsymbol{\mu}_-(h, d) + (1 - a_+ - a_-) d\mu_=(h, d). \quad (34)$$

Remark that  $\mu$  is a probability measure:

$$\int d\mu = \int \frac{1}{|d|} d\boldsymbol{\mu}_+(h, d) + \int \frac{1}{|d|} d\boldsymbol{\mu}_-(h, d) + (1 - a_+ - a_-) \int d\mu_+ = 1.$$

Consider  $m$  small enough so that the inequality  $m|d|/2 \leq 1$  holds on the support of  $\mu$  (this is possible because the  $d$  coordinate is bounded in the support of  $\boldsymbol{\mu}_+$  and  $\boldsymbol{\mu}_-$ ). We define

$$d\mu_1(h, d) = (1 + \frac{m}{2}d)d\mu(h, d) \quad \text{and} \quad d\mu_2(h, d) = (1 - \frac{m}{2}d)d\mu(h, d) \quad (35)$$

Note that  $\mu = \frac{1}{2}\mu_1 + \frac{1}{2}\mu_2$ . Also,  $\mu_1$  and  $\mu_2$  are probability measures, as

$$\int d\mu_1 = \int \frac{d}{|d|} d\boldsymbol{\mu}_+(h, d) + \int \frac{d}{|d|} d\boldsymbol{\mu}_-(h, d) = 1 - 1 = 0.$$

Let  $\eta(h, d) = h$ . We define the triplet  $(X, Y, S)$  by letting  $S$  be uniform on  $\{1, 2\}$ . If  $S = 1$ , we draw  $X \sim \mu_1$  and let  $Y = \eta(X)$ . If  $S = 2$ , we draw  $X \sim \mu_2$  and let  $Y = \eta(X)$ . Let  $\mathbb{P}$  be the distribution of  $(X, Y, S)$ . One can easily check that  $(\boldsymbol{\mu}_+(\mathbb{P}), \boldsymbol{\mu}_-(\mathbb{P})) = (\boldsymbol{\mu}_+, \boldsymbol{\mu}_-)$ , as desired.  $\square$



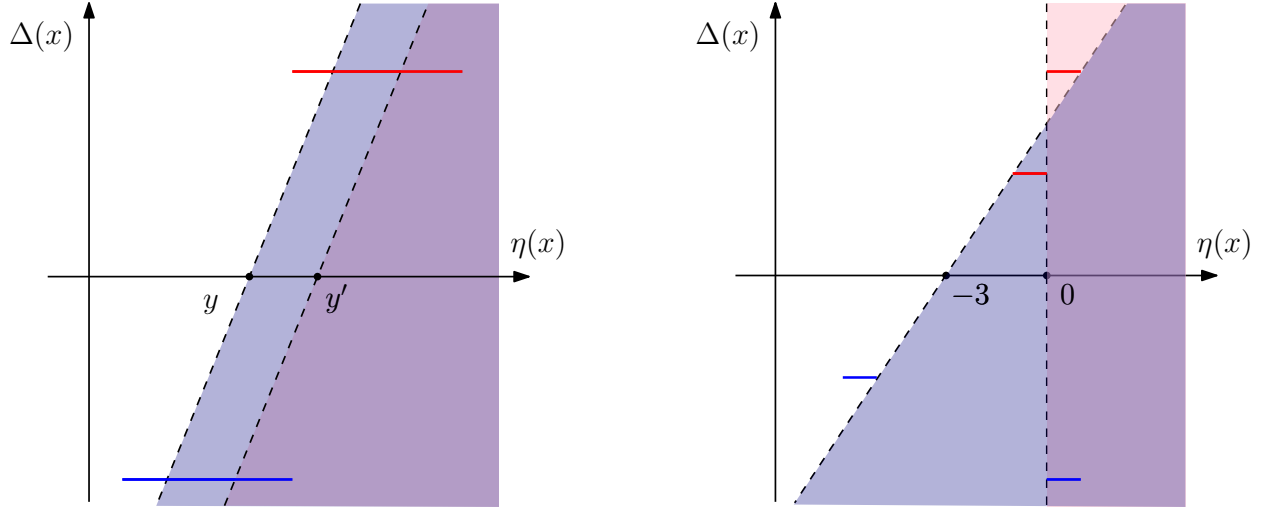


Figure 2: Left: example of a nested problem. The distributions of  $\mu_+$  and  $\mu_-$  are depicted in red and blue, corresponding to the distributions given in Example 1. The acceptance region for  $g_y^{\kappa(y)}$  and  $g_{y'}^{\kappa'(y')}$  are so that the masses of  $\mu_+$  and  $\mu_-$  to the right of the decision boundaries are equal. One can observe that these two regions are nested. Right: example of a non-nested problem. The distributions  $\mu_+$  and  $\mu_-$  are the ones described in Example 2. The region in pink is rejected for  $y = -3$  but accepted for  $y = 0$ , contradicting the nestedness assumption.

To build examples of nested and non-nested problems, we consider probability measures  $\mu_+$  and  $\mu_-$  supported on small horizontal segments:

$$\mu_{\pm} = \frac{1}{K} \sum_{i=1}^K \nu_{\pm}^{(i)} \quad (36)$$

where  $\nu_{\pm}^{(i)}$  is the uniform measure on  $[a_{\pm}^{(i)}, a_{\pm}^{(i)} + 1] \times \{d_{\pm}^{(i)}\}$ .

**Example 1** (A nested classification problem). Take  $K = 1$ ,  $d_+^{(1)} = d_-^{(1)} = 1$  and  $a_+^{(1)} = 0$ ,  $a_-^{(1)} = -1$ . Let  $\mathbb{P}$  be the probability associated with the pair  $(\mu_+, \mu_-)$  defined for this choice of parameters. Then, it holds that  $1/2 \in I(y)$  for all  $y \in \mathbb{R}$ . By choosing  $\kappa(y) = 1/2$  for all  $y \in \mathbb{R}$ , we see that the classification problem associated with  $\mathbb{P}$  is nested. See also Figure 2.

**Example 2** (A non-nested classification problem). Take  $K = 2$ . Let  $d_+^{(1)} = d_-^{(1)} = 1$  and  $a_+^{(1)} = a_-^{(1)} = 0$ . Let  $d_+^{(2)} = d_-^{(2)} = 1/2$ , and  $a_+^{(2)} = -1$ ,  $a_-^{(2)} = -6$ . Then, for  $y = 0$ ,  $I(y) = \{0\}$ , so the support of  $\nu_+^{(2)}$  is to the left of the classification threshold for  $y = 0$ . But for  $y = -3$ , we have  $I(y) = \{4\}$  and the support of  $\nu_+^{(2)}$  is to the right of the classification threshold. Hence, the classification problem is non-nested. See also Figure 2.

## 6 Conclusion and future work

This work presents the first theoretical characterization of the optimal fair regression function as the solution to a barycenter problem with an optimal transport cost. Our results also demonstrate that, under the nestedness assumption, the optimal fair regression function can be represented by the family of classifiers  $g_y^{\kappa(y)}$ . Although both approaches—whether based on optimal transport or cost-sensitive classifiers—depend on the underlying distribution  $\mathbb{P}$  which is generally unknown, they pave the way for developing new algorithms that estimate these unknown quantities from observed data. Designing such estimators, along with bounding their excess risk and potential unfairness, represents a critical step toward the development of fair algorithms.

While this work provides an initial characterization of the optimal fair regression function in the unawareness framework, it also has notable limitations. For instance, our results are currently limited to cases where the sensitive attribute is binary and apply only to univariate regression. Addressing these limitations and extending our findings to more general cases would be a valuable direction for future research.

## Acknowledgements

The authors gratefully acknowledge valuable and insightful discussions with Evgenii Chzhen and Nicolas Schreuder. S. G. gratefully acknowledges funding from the Fondation Mathématique Jacques Hadamard and from the ANR TopAI chair (ANR-19-CHIA-0001).

## References

- [AB99] Charalambos D. Aliprantis and Kim C. Border. *Measurable correspondences*, pages 557–586. Springer Berlin Heidelberg, Berlin, Heidelberg, 1999.
- [ADW19] Alekh Agarwal, Miroslav Dudik, and Zhiwei Steven Wu. Fair regression: Quantitative definitions and reduction-based algorithms. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 120–129. PMLR, 09–15 Jun 2019.
- [ALMK16] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Machine bias – there’s software used across the country to predict future criminals. and it’s biased against blacks. *ProPublica*, 2016.
- [AT07] Jean-Yves Audibert and Alexandre B. Tsybakov. Fast learning rates for plug-in classifiers. *The Annals of Statistics*, 35(2):608 – 633, 2007.
- [BDL08] Gerard Biau, Luc Devroye, and Gábor Lugosi. Consistency of random forests and other averaging classifiers. *Journal of Machine Learning Research*, 9:2015–2033, 09 2008.
- [BHN23] Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning: Limitations and Opportunities*. MIT Press, 2023.
- [CDH<sup>+</sup>20a] Evgenii Chzhen, Christophe Denis, Mohamed Hebiri, Luca Oneto, and Massimiliano Pontil. Fair regression via plug-in estimator and recalibration with statistical guarantees. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 19137–19148. Curran Associates, Inc., 2020.
- [CDH<sup>+</sup>20b] 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.
- [CKP09] Toon Calders, Faisal Kamiran, and Mykola Pechenizkiy. Building classifiers with independency constraints. In *2009 IEEE International Conference on Data Mining Workshops*, pages 13–18, 2009.
- [CMP16] Pierre-André Chiappori, Robert McCann, and Brendan Pass. Multidimensional matching. *arXiv preprint arXiv:1604.05771*, 2016.
- [CMP17] Pierre-André Chiappori, Robert J. McCann, and Brendan Pass. Multi-to one-dimensional optimal transport. *Communications on Pure and Applied Mathematics*, 70(12):2405–2444, 2017.
- [CS20a] Evgenii Chzhen and Nicolas Schreuder. An example of prediction which complies with demographic parity and equalizes group-wise risks in the context of regression. In *NeurIPS 2020 Workshop on Algorithmic Fairness through the Lens of Causality and Interpretability*, 2020.

- [CS20b] Evgenii Chzhen and Nicolas Schreuder. A minimax framework for quantifying risk-fairness trade-off in regression. *The Annals of Statistics*, 2020.
- [dBGL20] Eustasio del Barrio, Paula Gordaliza, and Jean-Michel Loubes. Review of Mathematical frameworks for Fairness in Machine Learning. arXiv admin note: substantial text overlap with arXiv:2001.07864, arXiv:1911.04322, arXiv:1906.05082 by other authors, October 2020.
- [DEHH24] Christophe Denis, Romuald Elie, Mohamed Hebiri, and François Hu. Fairness guarantees in multi-class classification with demographic parity. *Journal of Machine Learning Research*, 25(130):1–46, 2024.
- [FFM<sup>+</sup>15] Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and removing disparate impact. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '15, page 259–268, New York, NY, USA, 2015. Association for Computing Machinery.
- [GCGF16] Gabriel Goh, Andrew Cotter, Maya Gupta, and Michael P Friedlander. Satisfying real-world goals with dataset constraints. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [GLR20] Thibaut Le Gouic, Jean-Michel Loubes, and Philippe Rigollet. Projection to fairness in statistical learning. *arXiv preprint arXiv:2005.11720*, 2020.
- [GSC23] Solenne Gaucher, Nicolas Schreuder, and Evgenii Chzhen. Fair learning with Wasserstein barycenters for non-decomposable performance measures. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent, editors, *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pages 2436–2459. PMLR, 25–27 Apr 2023.
- [HPPS16] Moritz Hardt, Eric Price, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [LMC18] Zachary Lipton, Julian McAuley, and Alexandra Chouldechova. Does mitigating ml's impact disparity require treatment disparity? In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [MN06] Pascal Massart and Élodie Nédélec. Risk bounds for statistical learning. *The Annals of Statistics*, 34(5), October 2006.
- [MP20] Robert J McCann and Brendan Pass. Optimal transportation between unequal dimensions. *Archive for Rational Mechanics and Analysis*, 238(3):1475–1520, 2020.
- [MW18] Aditya Krishna Menon and Robert C. Williamson. The cost of fairness in binary classification. In Sorelle A. Friedler and Christo Wilson, editors, *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, volume 81 of *Proceedings of Machine Learning Research*, pages 107–118. PMLR, 23–24 Feb 2018.
- [MZP21] Karima Makhlof, Sami Zhioua, and Catuscia Palamidessi. On the applicability of machine learning fairness notions. *SIGKDD Explor. Newsl.*, 23(1):14–23, may 2021.
- [OC20] Luca Oneto and Silvia Chiappa. *Fairness in Machine Learning*, page 155–196. Springer International Publishing, 2020.
- [ODP20] Luca Oneto, Michele Donini, and Massimiliano Pontil. General fair empirical risk minimization. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2020.

- [Pas12] Brendan Pass. Regularity of optimal transportation between spaces with different dimensions. *Mathematical Research Letters*, 19(2):291–307, 2012.
- [Pra07] Aldo Pratelli. On the equality between Monge’s infimum and Kantorovich’s minimum in optimal mass transportation. In *Annales de l’Institut Henri Poincaré (B) Probability and Statistics*, volume 43, pages 1–13. Elsevier, 2007.
- [San15] Filippo Santambrogio. Optimal transport for applied mathematicians. *Birkäuser, NY*, 55(58-63):94, 2015.
- [SC21] Nicolas Schreuder and Evgenii Chzhen. Classification with abstention but without disparities. In Cassio de Campos and Marloes H. Maathuis, editors, *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*, volume 161 of *Proceedings of Machine Learning Research*, pages 1227–1236. PMLR, 27–30 Jul 2021.
- [Sri08] Sashi Mohan Srivastava. *A course on Borel sets*, volume 180. Springer Science & Business Media, 2008.
- [Vil09] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer, 2009.
- [XYZ23] Ruicheng Xian, Lang Yin, and Han Zhao. Fair and optimal classification via post-processing. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 37977–38012. PMLR, 23–29 Jul 2023.
- [Yan99] Yuhong Yang. Minimax nonparametric classification. I. Rates of convergence. *IEEE Transactions on Information Theory*, 45(7):2271–2284, 1999.
- [YCK20] Forest Yang, Mouhamadou Cisse, and Sanmi Koyejo. Fairness with overlapping groups; a probabilistic perspective. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 4067–4078. Curran Associates, Inc., 2020.
- [ZDC22a] Xianli Zeng, Edgar Dobriban, and Guang Cheng. Bayes-optimal classifiers under group fairness. *ArXiv*, abs/2202.09724, 2022.
- [ZDC22b] Xianli Zeng, Edgar Dobriban, and Guang Cheng. Fair bayes-optimal classifiers under predictive parity. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 27692–27705. Curran Associates, Inc., 2022.
- [ZM23] Quan Zhou and Jakub Marecek. Group-blind optimal transport to group parity and its constrained variants, 2023.
- [ZVGRG19] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez-Rodriguez, and Krishna P. Gummadi. Fairness constraints: A flexible approach for fair classification. *Journal of Machine Learning Research*, 20(75):1–42, 2019.

# A Additional proofs

## A.1 Proof of Lemma 3

Let  $f : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$  be a lower semicontinuous convex function. The domain of such a function is an interval  $\text{dom}(f)$ . Its right derivative  $f'_+$  is defined and finite everywhere on  $\text{dom}(f)$ , except on the right endpoint of the interval (should the right endpoint be included in  $\text{dom}(f)$ ) where it is equal to  $+\infty$ . Such a function is upper semicontinuous, with the representation:

$$\forall h \in \text{dom}(f), f'_+(h) = \inf_{u > h} \frac{f(u) - f(h)}{u - h}, \quad (37)$$

where the infimum can be restricted to a countable dense collection of values  $u$  if needed. Likewise, the right derivative  $f'_-$  can be defined on  $\text{dom}(f)$ , and is lower semicontinuous. Note also that the oscillation function  $\text{osc}(f) = f'_+ - f'_- \in [0, +\infty]$  can be defined on  $\text{dom}(f)$ , and is upper semicontinuous. Indeed, only one of  $f'_+$  and  $f'_-$  can be infinite on  $\text{dom}(f)$  (and only at one of the endpoints of the domain), so that the difference is well defined.

Recall that a function  $\varphi$  is **C**-convex if

$$\forall (h, d) \in \Omega, \varphi(h, d) = \sup_{(h', d') \in \Omega} \left( \varphi^{\mathbf{C}}(h', d') - \frac{(h - h')^2}{|d| + |d'|} \right). \quad (38)$$

The function  $\varphi$  is lower semicontinuous as a supremum of continuous functions. Furthermore, for any  $d \neq 0$ , the function

$$\varphi_d : h \mapsto \varphi(h, d) + \frac{h^2}{|d|} = \sup_{(h', d') \in \Omega} \left( \varphi^{\mathbf{C}}(h', d') + \frac{h^2}{|d|} - \frac{(h - h')^2}{|d| + |d'|} \right)$$

is convex as a supremum of lower semicontinuous convex functions. Let  $G : \text{dom}(\varphi) \rightarrow [0, +\infty]$  be defined as  $G(h, d) = \text{osc}(\varphi_d)(h)$  for  $(h, d) \in \text{dom}(\varphi)$ . Let  $L > 0$ , and consider the set  $\Sigma_{d,L}$  defined as the set of points  $h \in \text{dom}(\varphi_d)$  such that  $G(h, d) \geq L^{-1}$ ,  $(\varphi_d)'_+(h) \geq -L$  and  $(\varphi_d)'_-(h) \leq L$ . As the left and right derivatives of  $\varphi_d$  are nondecreasing, the set  $\Sigma_{d,L}$  is finite and its cardinality is bounded by a constant depending only on  $L$ . Let  $\Sigma_L = \bigcup_{d \neq 0} \Sigma_{d,L}$  and  $\Sigma = \bigcup_{L \in \mathbb{N}} \Sigma_L$ . The set of points  $\mathbf{x} \in \text{dom}(\varphi)$  such that  $\partial_h \varphi(\mathbf{x})$  does not exist is equal to  $\Sigma$ . Let us show that for any integer  $L$ , the set  $\Sigma_L$  is included in a countable union of graphs of measurable functions.

To do so, we use the following general result, see [AB99, Corollary 18.14]. A correspondence  $\Phi$  from a measurable set  $S$  to a topological space  $X$  assigns to each  $s \in S$  a subset  $\Phi(s)$  of  $X$ . We say that the correspondence is (weakly) measurable if for each open subset  $U \subset X$ , the set  $\Phi^\ell(U) = \{s \in S : \Phi(s) \cap U \neq \emptyset\}$  is measurable.

**Theorem 6** (Castaing's theorem). *If  $X$  is a Polish space and  $\Phi$  is a measurable correspondence with non-empty closed values between  $S$  and  $X$ , then there exists a sequence  $(f_n)_n$  of measurable functions  $S \rightarrow X$  such that for every  $s \in S$ ,  $\Phi(s) = \overline{\{f_1(s), f_2(s), \dots\}}$ .*

Let  $\Phi$  be the correspondence that assigns to each  $d \neq 0$  the subset  $\Sigma_{d,L} \cup \{0\}$  of  $\mathbb{R}$ . As each set  $\Sigma_{d,L}$  is finite, this correspondence takes non-empty closed values. If we show that this correspondence is measurable, then Castaing theorem asserts the existence of a sequence of measurable functions  $(f_n)_n$  such that for every  $d \neq 0$ ,  $\Sigma_{d,L} \cup \{0\} = \overline{\{f_1(d), f_2(d), \dots\}}$ . For each  $d$ , the set  $\Sigma_{d,L}$  is finite, so that  $\Sigma_{d,L} \cup \{0\} = \{f_1(d), f_2(d), \dots\}$ , implying that  $\Sigma_L$  is included in a countable union of graphs of measurable functions. It remains to show the measurability of  $\Phi$ .

**Lemma 9.** *The function  $G$  is measurable.*

*Proof.* The representation (37) implies that  $(h, d) \mapsto (\varphi_d)'_+(h)$  is given by a countable infimum of measurable functions, and is therefore measurable. Likewise,  $(h, d) \mapsto (\varphi_d)'_-(h)$  is measurable, so that  $G$  is also measurable.  $\square$

Let  $U \subset \mathbb{R}$  be an open set. If  $0 \in U$ , then  $\Phi^\ell(U) = \mathbb{R} \setminus \{0\}$  is measurable. If  $0 \notin U$ , we have

$$\Phi^\ell(U) = \{d \neq 0 : \exists h \in [-L, L] \cap U, G(h, d) \geq L^{-1}, (\varphi_d)'_+(h) \geq -L, (\varphi_d)'_-(h) \leq L\}.$$

This set is the projection on the  $d$ -axis of the measurable set

$$B = \{(h, d) \in \Omega : h \in [-L, L] \cap U, G(h, d) \geq L^{-1}, (\varphi_d)'_+(h) \geq -L, (\varphi_d)'_-(h) \leq L\}$$

Furthermore, for each  $d$ , the section  $\{h \in \mathbb{R} : (h, d) \in B\} = \Sigma_{d,L}$  is compact. By [Sri08, Theorem 4.7.11], this implies that  $\Phi^\ell(U)$  is measurable, concluding the proof of Lemma 3.

## A.2 Proof of Proposition 2

Classical manipulations show that the risk  $\mathcal{R}_y(g)$  of a classifier  $g$  can be expressed as

$$\begin{aligned} \mathcal{R}_y(g) &= y\mathbb{E}[(1-Y)g(X)] + (1-y)\mathbb{E}[Y(1-g(X))] \\ &= (1-y)\mathbb{E}[Y] + \mathbb{E}[g(X)(y-\eta(X))]. \end{aligned}$$

Using the definition of  $\mu_+$ ,  $\mu_-$  and  $\Delta$  given in Section 3, we find that

$$\begin{aligned} &\mathbb{E}[g(X)(y-\eta(X))] \\ &= \int_{\mathcal{X}_+} g(x)(y-\eta(x)) \frac{d\mu}{d\mu_+}(x) d\mu_+(x) + \int_{\mathcal{X}_-} g(x)(y-\eta(x)) \frac{d\mu}{d\mu_-}(x) d\mu_-(x) \\ &\quad + \int_{\mathcal{X}_=} g(x)(y-\eta(x)) d\mu(x) \\ &= \int_{\mathcal{X}_+} g(x) \frac{y-\eta(x)}{|\Delta(x)|} d\mu_+(x) + \int_{\mathcal{X}_-} g(x) \frac{y-\eta(x)}{|\Delta(x)|} d\mu_-(x) + \int_{\mathcal{X}_=} g(x)(y-\eta(x)) d\mu(x). \end{aligned}$$

Moreover, Lemma 1 implies that the demographic parity constraint is equivalent to the constraint  $\mathbb{E}_{X \sim \mu_+}[g(X)] = \mathbb{E}_{X \sim \mu_-}[g(X)]$ . Using the decomposition  $g = \mathcal{F}(g_+, g_-, g_-)$ , we see that the fair classification problem can be rephrased as follows

$$\begin{cases} \text{minimize} & \mathbb{E}_{\mu_+} \left[ g_+(X) \frac{y-\eta(X)}{\Delta(X)} \right] - \mathbb{E}_{\mu_-} \left[ g_-(X) \frac{y-\eta(X)}{\Delta(X)} \right] \\ & + \mathbb{E}_{\mu} [\mathbf{1}_{\mathcal{X}_=}(X) g_+(X)(y-\eta(X))] \\ \text{such that} & \mathbb{E}_{\mu_+}[g_+(X)] = \mathbb{E}_{\mu_-}[g_-(X)]. \end{cases} \quad (C'_y)$$

The following lemma characterizes the solutions to the problem  $(C'_y)$ .

**Lemma 10.** *Under Assumption 1, for any optimal classifier  $g$ , there exist  $\kappa^+$ ,  $\kappa^-$  such that  $g = \mathcal{F}(g^{\kappa^+}, g^{\kappa^-}, g_-)$ , with*

$$\begin{aligned} g_-(x) &= \mathbf{1}\{\eta(x) > y\} \quad \text{or} \quad g_-(x) = \mathbf{1}\{\eta(x) \geq y\}, \\ \text{and} \quad g^\kappa(x) &= \mathbf{1}\{\eta(x) \geq y + \kappa\Delta(x)\}. \end{aligned}$$

To conclude the proof of Proposition 2, it remains to prove that all optimal classifier are a.s. equal when  $\Delta(X) \neq 0$ , and that the optimal classifier can be chosen as  $g^* = \mathcal{F}(g^{\kappa^*}, g^{\kappa^*}, g_-)$  for some  $\kappa^*$ .

Denote by  $F_+$  the c.d.f. of the random variable  $Z_+ = \frac{\eta(X)-y}{\Delta(X)}$  when  $X \sim \mu_+$  and by  $F_-$  the c.d.f. of  $Z_- = \frac{y-\eta(X)}{\Delta(X)}$  when  $X \sim \mu_-$ . Let  $\mathcal{Q}_+$  (resp.  $\mathcal{Q}_-$ ) be the associated quantile function. To verify the demographic parity constraint, the classifier  $\mathcal{F}(g^{\kappa^+}, g^{\kappa^-}, g_-)$  must be such that

$$F_+(\kappa^+) = F_-(-\kappa^-)$$

(recall that  $\Delta(X) < 0$  when  $X \sim \mu_-$ , so that  $g^{\kappa^-}(X) = 1$  if and only if  $Z_- \geq -\kappa^-$ ). Denoting  $\beta = F_+(\kappa^+) = F_-(-\kappa^-)$  and using the definition of the quantile function, we see that the law of  $g^{\kappa^\pm}(X)$ , where  $X \sim \mu_\pm$  is equal to the law of

$$\mathbf{1}\{U \geq \beta\} = \mathbf{1}\{\mathcal{Q}_\pm(U) \geq \pm\kappa^\pm\},$$

where  $U$  is a uniform random variable on  $[0, 1]$ .

Then, minimizing the risk of the classifier is equivalent to maximizing

$$\mathbb{E}[(\mathcal{Q}_+(U) + \mathcal{Q}_-(U)) \mathbf{1}\{U \geq \beta\}] \quad (39)$$

Since  $F_+$  and  $F_-$  are continuous,  $u \rightarrow \mathcal{Q}_+(u) + \mathcal{Q}_-(u)$  is strictly increasing and left-continuous. Straightforward computations show that the expression in (39) has a unique maximum, which is attained for

$$\beta^* = \max\{\beta : \mathcal{Q}_+(\beta) + \mathcal{Q}_-(\beta) \leq 0\}.$$

Hence, it holds that  $F_+(\kappa^+) = F_-(-\kappa^-) = \beta^*$ . Let  $g_1^* = \mathcal{F}(g^{\kappa_1^+}, g^{\kappa_1^-}, g_-)$  and  $g_2^* = \mathcal{F}(g^{\kappa_2^+}, g^{\kappa_2^-}, g_-)$  be two optimal classifiers. For  $\Delta(X) > 0$ , they take different values only if  $Z_+ \in [\kappa_1^+, \kappa_2^+]$ . As  $F_+(\kappa_1^+) = F_+(\kappa_2^+) = \beta^*$ , this happens with zero probability. Likewise, the two classifiers are a.s. equal when  $\Delta(X) < 0$ .

It remains to show that we can pick  $\kappa^+ = \kappa^-$ . If  $\mathcal{Q}_+ + \mathcal{Q}_-$  is continuous at  $\beta^*$ , the proof is complete: in this case,  $\mathcal{Q}_+(\beta^*) + \mathcal{Q}_-(\beta^*) = 0$ , and the choice  $\kappa^* = \mathcal{Q}_+(\beta^*) = -\mathcal{Q}_-(\beta^*)$  satisfies  $F_+(\kappa^*) = F_-(-\kappa^*) = \beta^*$ .

Otherwise,  $\mathcal{Q}_+(\beta^*) + \mathcal{Q}_-(\beta^*) < 0$ . Defining

$$q_+ = \liminf_{\beta \rightarrow \beta_+^*} \mathcal{Q}_+(\beta) \quad \text{and} \quad q_- = \liminf_{\beta \rightarrow \beta_+^*} \mathcal{Q}_-(\beta),$$

we have  $q_+ + q_- > 0$ . Because  $q_+ > -q_-$  and  $\mathcal{Q}_+(\beta^*) < -\mathcal{Q}_-(\beta^*)$ , there exists  $\kappa^* \in [\mathcal{Q}_+(\beta^*), q_+] \cap [-q_-, -\mathcal{Q}_-(\beta^*)]$ . By construction,  $F_+(\kappa^*) = F_-(-\kappa^*) = \beta^*$ , concluding the proof.

### A.3 Proof of Lemma 10

Let  $g^*$  be a solution to the problem  $(C'_y)$ , and let  $g_+^*$ ,  $g_-^*$  and  $g_-^*$  be the restrictions of  $g^*$  to  $\mathcal{X}_+$ ,  $\mathcal{X}_-$  and  $\mathcal{X}_=$ , so that  $g^* = \mathcal{F}(g_+^*, g_-^*, g_-^*)$ . Straightforward computations show that we necessarily have  $g_-^*(X) = \mathbf{1}\{\eta(X) \geq y\} \mathbf{1}_{\mathcal{X}_=}(X)$  or  $g_-^*(X) = \mathbf{1}\{\eta(X) > y\} \mathbf{1}_{\mathcal{X}_=}(X)$ .

Now, let us assume (without loss of generality) that  $g_+^*$  is not of the form  $g^{\kappa^+}$ . More precisely, assume that for  $\kappa^+$  such that  $\mathbb{E}_{\mu_+}[g^{\kappa^+}(X)] = \mathbb{E}_{\mu_+}[g_+^*(X)]$ , we have  $g^{\kappa^+}(X) \neq g_+^*(X)$  with positive  $\mu_+$ -probability. Then, the classifier  $\mathcal{F}(g^{\kappa^+}, g_-^*, g_-^*)$  verifies the demographic parity constraint. Moreover, we have

$$\begin{aligned} & \mathbb{E}_{\mu_+} \left[ g^{\kappa^+}(X) \frac{y - \eta(X)}{\Delta(X)} \right] - \mathbb{E}_{\mu_+} \left[ g_+^*(X) \frac{y - \eta(X)}{\Delta(X)} \right] \\ &= \mathbb{E}_{\mu_+} \left[ \frac{y - \eta(X)}{\Delta(X)} g^{\kappa^+}(X) (1 - g_+^*(X)) - \frac{y - \eta(X)}{\Delta(X)} (1 - g^{\kappa^+}(X)) g_+^*(X) \right] \\ &= \mathbb{E}_{\mu_+} \left[ \left( \kappa^+ - \frac{\eta(X) - y}{\Delta(X)} \right) g^{\kappa^+}(X) (1 - g_+^*(X)) \right] - \kappa^+ \mathbb{E}_{\mu_+} \left[ g^{\kappa^+}(X) (1 - g_+^*(X)) \right] \\ & \quad - \mathbb{E}_{\mu_+} \left[ \left( \kappa^+ - \frac{\eta(X) - y}{\Delta(X)} \right) (1 - g^{\kappa^+}(X)) g_+^*(X) \right] + \kappa^+ \mathbb{E}_{\mu_+} \left[ (1 - g^{\kappa^+}(X)) g_+^*(X) \right]. \end{aligned}$$

Since  $\mathbb{E}_{\mu_+}[g^{\kappa^+}(X)] = \mathbb{E}_{\mu_+}[g_+^*(X)]$ , we obtain that

$$\mathbb{E}_{\mu_+} \left[ g^{\kappa^+}(X) (1 - g_+^*(X)) \right] = \mathbb{E}_{\mu_+} \left[ (1 - g^{\kappa^+}(X)) g_+^*(X) \right].$$

Then, the definition of  $g^{\kappa^+}$  implies

$$\begin{aligned} & \mathbb{E}_{\mu_+} \left[ g^{\kappa^+}(X) \frac{y - \eta(X)}{\Delta(X)} \right] - \mathbb{E}_{\mu_+} \left[ g_+^*(X) \frac{y - \eta(X)}{\Delta(X)} \right] \\ &= -\mathbb{E}_{\mu_+} \left[ \left( \kappa^+ - \frac{\eta(X) - y}{\Delta(X)} \right)_- (1 - g_+^*(X)) \right] - \mathbb{E}_{\mu_+} \left[ \left( \kappa^+ - \frac{\eta(X) - y}{\Delta(X)} \right)_+ g_+^*(X) \right] < 0. \end{aligned}$$

This implies that  $\mathcal{R}_y(\mathcal{F}(g^{\kappa^+}, g_-^*, g_-^*)) < \mathcal{R}_y(\mathcal{F}(g_+^*, g_-^*, g_-^*))$ , which is absurd. Using a similar argument for  $g_-^*$ , we arrive to the conclusion that any optimal classifier is of the form  $\mathcal{F}(g^{\kappa^+}, g^{\kappa^-}, g_-^*)$  with  $g_-^*(X) = \mathbf{1}\{\eta(X) \geq y\} \mathbf{1}_{\mathcal{X}_=}(X)$  or  $g_-^*(X) = \mathbf{1}\{\eta(X) > y\} \mathbf{1}_{\mathcal{X}_=}(X)$ .

## A.4 Proof of Lemma 7

Let us first show that  $\kappa$  is locally bounded. Let  $[y_0, y_1]$  be a bounded interval. Then, if  $\kappa(y) \geq M$  for some  $y \in [y_0, y_1]$ ,

$$F(y) = \mu_+(\eta(X) \leq y + \kappa(y)\Delta(X)) \geq \mu_+(\eta(X) \leq y_0 + M\Delta(X))$$

and

$$F(y) = \mu_-(\eta(X) \leq y + \kappa(y)\Delta(X)) \leq \mu_-(\eta(X) \leq y_1 + M\Delta(X)).$$

However, for  $M$  large enough,

$$\mu_-(\eta(X) \leq y_1 + M\Delta(X)) < \mu_+(\eta(X) \leq y_0 + M\Delta(X)),$$

and therefore it holds that  $\kappa(y) \leq M$  for all  $y \in [y_0, y_1]$ . Likewise, we show that for  $M$  large enough,  $\kappa(y) \geq -M$  for all  $y \in [y_0, y_1]$ . Hence,  $\kappa$  is locally bounded.

We refine this argument to obtain a control on  $\kappa$  for large values of  $y$ . Let  $L > 1$  and let  $y > 1$  be such that  $\kappa(y) \geq Ly$ . Then, because of the definition of  $\kappa(y)$  and as  $\Delta(X) < 0$  for  $X \sim \mu_-$ ,

$$\begin{aligned} \mu_+(\eta(X) \geq y) &\geq \mu_+(\eta(X) \geq y + \kappa(y)\Delta(X)) \\ &= \mu_-(\eta(X) \geq y + \kappa(y)\Delta(X)) \geq \mu_-(\eta(X) \geq y + Ly\Delta(X)). \end{aligned}$$

The complementary of the region  $\{(h, d) \in \Omega : d < 0, h \geq y + Lyd\}$  in the lower half-plane  $\{(h, d) \in \Omega : d < 0\}$  is contained in the set  $A_L$  given by the union of the horizontal strip  $\{d \geq -1/L\}$  with the region  $\{(h, d) \in \Omega : d < 0, h \leq 1 + Ld\}$ . For  $L$  large enough,  $\mu_-(A_L) < 1/2$ . Then it holds that  $\mu_+(\eta(X) \geq y) \geq 1/2$ . For  $y$  large enough, this is not possible. Thus, we have shown that there exist  $L > 0$  and  $C > 0$  such that for  $y > C$ , we have  $\kappa(y) \leq Ly$ . Likewise, we show that there exist constants  $L, C > 0$  such that for  $|y| > C$ ,  $|\kappa(y)| \leq L|y|$ . As  $\kappa$  is also locally bounded, the conclusion follows.