# EquiPy: Sequential Fairness using Optimal Transport in Python

Agathe Fernandes Machado[✉,1], Suzie Grondin[2], Philipp Ratz[1], Arthur Charpentier[1], and François Hu[3]

[1]Département de Mathématiques, Université du Québec à Montréal, Montréal, Québec, Canada
[2]ENSAE Institut Polytechnique, Paris, France
[3]Milliman, Paris, France

March 14, 2025

### Abstract

Algorithmic fairness has received considerable attention due to the failures of various predictive AI systems that have been found to be unfairly biased against subgroups of the population. Many approaches have been proposed to mitigate such biases in predictive systems, however, they often struggle to provide accurate estimates and transparent correction mechanisms in the case where multiple sensitive variables, such as a combination of gender *and* race, are involved. This paper introduces a new open source Python package, **EquiPy**, which provides a easy-to-use and model agnostic toolbox for efficiently achieving fairness across multiple sensitive variables. It also offers comprehensive graphic utilities to enable the user to interpret the influence of each sensitive variable within a global context. **EquiPy** makes use of theoretical results that allow the complexity arising from the use of multiple variables to be broken down into easier-to-solve sub-problems. We demonstrate the ease of use for both mitigation and interpretation on publicly available data derived from the US Census and provide sample code for its use.

## 1 Introduction

With the increased usage of machine learning (ML) models across different industries, the discovery of unjust biases in predictive models has also proliferated. Multiple cases have been prominently reported in the media, where examples include sexist recruiting algorithms, facial recognition systems that perform poorly for females with darker skin and challenges in recognizing specific subgroups in self-driving cars (Goodall, 2014; Nyholm and Smids, 2016; Dastin, 2022; Buolamwini and Gebru, 2018). Though none of these underlying models aimed to be biased in the first place, the predictions were found to replicate or even exaggerate existing biases in the training data. Simply not using a sensitive variable, such as gender or race, is not sufficient to avoid such biases, as more complex ML models can simply start to proxy for the omitted variables (Obermeyer et al., 2019). This has shown the need to develop metrics

and methods that can systematically measure and mitigate such biases in a more consistent manner.

In recent years, many approaches to address this issue have been developed. Broadly speaking, these methods can be categorized into pre-processing, in-processing, and post-processing methods. Pre-processing ensures fairness in the input data by removing biases within the data (Park et al., 2021; Qiang et al., 2022; Plečko and Meinshausen, 2020), in-processing methods incorporate fairness constraints during model training (Agarwal et al., 2018; Wang et al., 2020; Joo and Kärkkäinen, 2020), where fairness constraints are usually incorporated into the loss function, and post-processing methods (of which **EquiPy** makes use), that achieve fairness through modifications of the final scores (Karako and Manggala, 2018; Kim et al., 2019; Zeng et al., 2022). Fairness of the resulting corrected predictions is then evaluated with respect to a *group fairness metric*, such as independence, separation, or sufficiency (Barocas et al., 2018). Here, we focus on the independence criterion, which enforces Demographic Parity (DP) by requiring similar prediction distributions across sensitive groups.

Of particular importance for the development of this package is the literature using optimal transport, a mathematical framework for measuring distributional differences. Intuitively, achieving DP-fairness in predictions involves *transporting* unfair scores to fair ones, while minimizing the effects of this intervention to maintain predictive accuracy. In regression, methods like Chzhen et al. (2020) and Gouic et al. (2020) minimize the Wasserstein distance across groups to reduce discrimination. Similarly, in classification, Chiappa et al. (2020) and Gaucher et al. (2023) leverage optimal transport to achieve fair scores. Building upon this studies, Hu et al. (2023) achieved fairness in multi-task learning through a task-wise correction of the scores and extended the results to a sequential correction in Hu et al. (2024). **EquiPy** incorporates these insights and ensures DP-fairness in ML scores from binary classifiers or regression models, first for a single sensitive attribute using the Wasserstein barycenter-based mitigation of Chzhen et al. (2020) and then provides an extension to multiple discrete sensitive attributes via a sequential correction. This sequential correction approach facilitates not only an efficient correction, but also provides interpretable pathways of the correction methods which allows for a more in-depth analysis of its impacts. **EquiPy**'s foundation in optimal transport theory provides a strong theoretical basis, guarantees of optimality, and a natural extension to related applications such as causal inference and interpretability (Charpentier et al., 2023; Ratz et al., 2023).

## 1.1  Related Software

Tackling unfairness in ML algorithms has been broadly studied and tested on real datasets in recent years. To facilitate this task in practice, multiple packages in different programming languages have been introduced for this task, for example **FairLens** in Python, the **fairness** (Kozodoi and Varga, 2021) and **fairadapt** (Plečko et al., 2024) packages in R or the multi-platform package **AIF360** (Bellamy et al., 2019). More general in its applications is the **FairLearn** package (Weerts et al., 2023) in Python, which provides a broad introduction into the field of algorithmic fairness and a demonstration on how to use different metrics and mitigation techniques implemented within. Within **FairLearn**, the bulk of the mitigation techniques can be classified as either pre- or in-processing techniques[1] that focus on achieving fairness either before or during the training of a model. More recently, **OxonFair**, introduced by Delaney et al. (2024), enables the evaluation and enforcement of fairness in binary classification through group-specific thresholds, with support for NLP and Computer Vision applications.

---

[1]With the exception being the implementation of a post-processing method developed by Hardt et al. (2016), which crucially also requires the label (or regressand) to achieve fairness—a restriction that is not required in **EquiPy**.

**EquiPy**, being a post-processing optimal transport-based procedure, adjusts scores from any binary classification or regression model at a guaranteed minimal cost to predictive accuracy to achieve DP. Its lightweight, model-agnostic design ensures compatibility with models trained in languages beyond Python, making fairness considerations accessible even in industries reliant on commercial software for their predictions. With simple `fit` and `predict` methods, it lowers the barrier to entry for analysts and practitioners, aligning with recent high-level packages Pappalardo et al. (2022); Boudt et al. (2022); Tierney and Cook (2023). The custom-build visualization module also allows for the analysis of various policy-driven fairness scenarios, allowing non-technical decision-makers to understand and evaluate different approaches to achieving fairness and equity in predictive analytics, while illustrating iterative fairness-performance trade-offs arising from corrections across multiple sensitive attributes. Moreover, unlike **FairLearn** and **OxonFair**, the main method calibrates using only model predictions and sensitive variable(s), requiring no access to labels, facilitating seamless deployment in existing pipelines. Finally, given that most of the implementation of **EquiPy** does not rely on purpose-built software nor specialised hardware, but simple high-level mathematical operations, a future implementation in other programming languages such as R should therefore not pose significant compatibility issues.

A stable release can be installed via the Python Package index using

```
pip install equipy
```

The development version of the package is available on GitHub. A preliminary version can be installed directly from the GitHub repository using

```
pip install --upgrade git+https://github.com/equilibration/equipy.git
```

**Paper outline.**    This paper is structured as follows: In Section 2, we provide a brief introduction to algorithmic fairness and the optimal transport theory, establishing connections with the **EquiPy** package. Sections 3 and 4 introduce respectively the main mathematical propositions and functionalities of **EquiPy**, mainly derived from Hu et al. (2024). In Section 5, we present examples of these functions through a case study on an open-source dataset. Finally, Section 6 summarizes our contribution and offers concluding remarks.

## 2   Algorithmic Fairness and Optimal Transport

Algorithmic fairness has close links to the field of optimal transport, to which we provide a short formal discussion below. As an informal introduction to the topic, we focus on the concept of *Demographic Parity* (DP), also formally defined below, to establish fairness[2]. Intuitively speaking, DP requires the predictions of a model to be independent of a sensitive variable. As other variables used within a predictive model often correlate with the sensitive one, independence of the scores goes beyond isolating the direct effects of a sensitive variable. That is, DP-fair scores should be independent of the direct effects as well as the indirect effects of the variable. As an example, body height is generally strongly correlated with sex, a DP-fair prediction on the sensitive variable *sex* would therefore also need to correct for the discrepancies arising due to different distributions of height between the sexes. This also explains why it is not sufficient to simply omit the sensitive attribute, as a model can simply learn to proxy using

---

[2]DP is probably also the most commonly used fairness metric. Since **EquiPy** is intended for usage on both classification *and* regression tasks, we restrict the analysis on this metric, as other common metrics such as Equalized Odds are not directly applicable to the regression case.
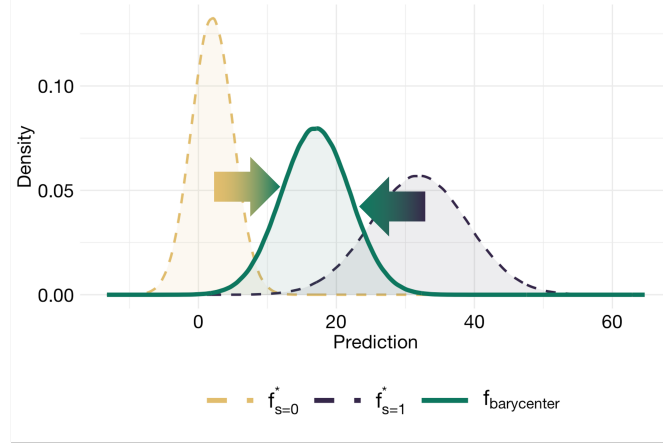
Figure 1: *Group-wise distribution of predictions and barycenter. Note that the $y$-axis represents the density and the barycenter contains the both of the group wise predictions.*

a mix of related variables, to which ML models with complex feature interactions are especially prone to.

To see how transporting scores could help achieve DP, consider Figure 1. Suppose that we have trained a model $f^*$ and analyze the distribution of the predictions given some sensitive variable $S \in \{0, 1\}$, resulting in the yellow and purple density curve. Clearly, DP would not be satisfied in this case. If instead the predictions were grouped in a single "middle-ground" distribution that builds a form of center for both group-wise prediction distributions, DP would be satisfied. This middle ground can be calculated using a Wasserstein barycenter, formally discussed below, and comes with a number of handy properties, such as being the distribution that satisfies DP but also minimizes the transportation cost. Intuitively, minimizing the transportation cost is a desirable property, as any change to the optimal predictions (that is, the initial *unfair* predictions) will necessarily lead to a degradation in predictive metrics. By keeping the changes to the original predictions as small as possible, this negative impact is in turn minimized.

## 2.1 Background on Optimal Transport

In this Section, we provide a brief formal introduction to the concept of the Wasserstein barycenter, emphasizing its characteristics within one-dimensional optimal transport theory and its correlation with algorithmic fairness. For those seeking further insights, excellent resources such as Santambrogio (2015); Villani (2021) are available for an overview of optimal transport theory and Chiappa et al. (2020); Chzhen et al. (2020) provide a more specific overview of its implications in algorithmic fairness.

Throughout this article, we refer to $\mathcal{V}$ as the space of univariate probability measures on $\mathcal{Y} \subset \mathbb{R}$ with finite variance (see Agueh and Carlier (2011) for technical assumptions in general cases). Let $\nu_1$ and $\nu_2$ represent two probability measures in $\mathcal{V}$. A commonly employed method for measuring the difference between distributions involves the *Wasserstein distance*, which computes the minimum "cost" required to transform one distribution into the other. In our case, we are interested in the $p$-Wasserstein distance (with $p \geq 1$) between $\nu_1$ and $\nu_2$, defined as:

$$\mathcal{W}_p^p(\nu_1, \nu_2) = \inf_{\pi \in \Pi(\nu_1, \nu_2)} \mathbb{E}_{(Z_1, Z_2) \sim \pi} |Z_2 - Z_1|^p,$$

where $\Pi(\nu_1, \nu_2)$ is the set of distributions on $\mathcal{Y} \times \mathcal{Y}$ having $\nu_1$ and $\nu_2$ as marginals. In the

univariate case, this expression can be rewritten as

$$\mathcal{W}_p^p(\nu_1, \nu_2) = \int_{u \in [0,1]} |Q_1(u) - Q_2(u)|^p \, du \ ,$$

where $Q_1$ and $Q_2$ are the associated quantile functions of $\nu_1$ and $\nu_2$ respectively. A coupling that achieves this infimum is called the optimal coupling between $\nu_1$ and $\nu_2$. An alternative way to look at this problem is to seek for a mapping $T$, solution of

$$\inf_{T: T_\# \nu_1 = \nu_2} \mathbb{E}_{Z_1 \sim \nu_1} |T(Z_1) - Z_1|^p$$

where $T_\# \nu_1 = \nu_2$ means that we consider push-forward measures of $\nu_1$ onto $\nu_2$. If $\nu_1$ is non-atomic and $p \geq 1$, there exists an optimal deterministic mapping $T^\star$. And if $\nu_1$ and $\nu_2$ are univariate measures in a subset of $\mathbb{R}$ (as considered here), and $p > 1$, it is unique, and given by $T^\star(z_1) = Q_{|2} \circ F_{|1}(z_1)$ where $F_{|1}$ is the cumulative distribution function associated with $\nu_1$ and $Q_{|2}$ is the quantile function associated with $\nu_2$. If $p = 1$, that mapping is not unique, but it is still optimal.

**Wasserstein Barycenter**     Throughout this article, we will frequently make use of *Wasserstein Barycenters* (Agueh and Carlier, 2011). The Wasserstein Barycenter finds a representative distribution that lies between multiple given distributions in the Wasserstein space. It is defined for a family of $K$ measures $(\nu_1, \ldots, \nu_K)$ in $\mathcal{V}$ and some positive weights $(w_1, \ldots, w_K) \in \mathbb{R}_+^K$. The Wasserstein barycenter, denoted as $\mathrm{Bar}\{(w_k, \nu_k)_{k=1}^K\}$ is the minimizer

$$\mathrm{Bar}\{(w_k, \nu_k)_{k=1}^K\} = \operatorname*{argmin}_{\nu \in \mathcal{V}} \sum_{k=1}^K w_k \cdot \mathcal{W}_2^2(\nu_k, \nu) \ . \tag{1}$$

The work in Agueh and Carlier (2011) demonstrates that in our setup, the barycenter exists, and a sufficient condition for uniqueness is that one of the measures $\nu_i$ has a density with respect to the Lebesgue measure. For the sake of simplicity and readability, we assume that this condition holds for all distributions introduced in this article. It is worth noting that intuitively, comparing distributions $f(\boldsymbol{X}, \boldsymbol{A})|\boldsymbol{A} = \boldsymbol{a}$ with $\boldsymbol{A}$ the sensitive attributes, where $\boldsymbol{a} \in \mathcal{A}$, and $f$ a given predictor seems natural. This concept becomes clearer in the subsequent sections, where we explore how the notion of unfairness is closely linked to these group-wise predictor response distributions.

**Notation**     Following the notation introduced in Hu et al. (2024), we examine a predictor $f : \mathcal{X} \times \mathcal{A} \to \mathcal{Y}$ where $(\boldsymbol{X}, \boldsymbol{A}, Y)$ is a random tuple drawn from the probability distribution $\mathbb{P}$. Here, $\boldsymbol{X} \in \mathcal{X} \subset \mathbb{R}^d$ represents the $d$ features, $\boldsymbol{A} = (A_1, A_2, \ldots, A_r) \in \mathcal{A} := \mathcal{A}_1 \times \mathcal{A}_2 \times \ldots \times \mathcal{A}_r \subset \mathbb{N}^r$ represents the $r$ sensitive features (or attributes), each of which is discrete and subject to fairness enforcement. For convenience, we define $A_{i:i+k} := (A_i, A_{i+1}, \cdots, A_{i+k})$ as the sequence of $k + 1$ sensitive attributes ranging from $i$ to $i + k$. Thus, $\boldsymbol{A} = A_{1:r}$. The task to be estimated is represented by $Y \in \mathcal{Y} \subset \mathbb{R}$ and may involve either regression or a probability score for classification. Further, additional quantities are introduced below. Given $\boldsymbol{A} = (A_1, \ldots, A_i, \ldots, A_r)$ and $\boldsymbol{a} = (a_1, \ldots, a_i, \ldots, a_r) \in \mathcal{A}$, we denote them as follows:

- $\nu_f$ the probability measure of $f(\boldsymbol{X}, \boldsymbol{A})$ in $\mathcal{V}$;

- $\nu_{f|\boldsymbol{a}}$ and $\nu_{f|a_i}$ (in $\mathcal{V}$) the probability measure of $f(\boldsymbol{X}, \boldsymbol{A})|\boldsymbol{A} = \boldsymbol{a}$ and $f(\boldsymbol{X}, \boldsymbol{A})|A_i = a_i$ respectively;

- $F_{f|\boldsymbol{a}}(u) := \mathbb{P}(f(\boldsymbol{X}, \boldsymbol{A}) \leq u | \boldsymbol{A} = \boldsymbol{a})$ and, with an abuse of notation, $F_{f|a_i}(u) := \mathbb{P}(f(\boldsymbol{X}, \boldsymbol{A}) \leq u | A_i = a_i)$ as their cumulative distribution function (CDF);

- $Q_{f|\boldsymbol{a}}(v) := \inf\{u \in \mathbb{R} : F_{f|\boldsymbol{a}}(u) \geq v\}$ and, with an abuse of notation, $Q_{f|a_i}(v) := \inf\{u \in \mathbb{R} : F_{f|a_i}(u) \geq v\}$ as their associated quantile functions.

## 2.2 Definitions of Unfairness

As described earlier Demographic Parity is used to determine the fairness of a predictor, and is applicable to both classification and regression tasks. **EquiPy** is able to handle the presence of multiple sensitive attributes, and hence as a special case also the situation where a single sensitive attribute is present. Since most of the existing literature does not consider the more general case, we extend from the standard case to define unfairness in the context of multiple sensitive attributes.

**Single sensitive attribute (SSA).** In algorithmic fairness literature regarding a sensitive attribute $A_i$, the DP-unfairness requires $f(\boldsymbol{X}, \boldsymbol{A}) \perp\!\!\!\perp A_i$ and its measure is typically defined with the following Total Variation,

$$\mathcal{U}_{TV}(f) = \max_{a_i \in \mathcal{A}_i} \sup_{I \subset \mathbb{R}} |\mathbb{P}(f(\boldsymbol{X}, \boldsymbol{A}) \in I | A_i = a_i) - \mathbb{P}(f(\boldsymbol{X}, \boldsymbol{A}) \in I)| \ ,$$

or defined based on the Kolmogorov-Smirnov test,

$$\mathcal{U}_{KS}(f) = \max_{a_i \in \mathcal{A}_i} \sup_{u \in \mathbb{R}} \left| F_f(u) - F_{f|a_i}(u) \right| \ .$$

In our study, the unfairness measure of the predictor on the feature $A_i$ is given by the corresponding 1-Wasserstein based measure:

$$\mathcal{U}_i(f) = \max_{a_i \in \mathcal{A}_i} \mathcal{W}_1(\nu_f, \nu_{f|a_i}) = \max_{a_i \in \mathcal{A}_i} \int_{u \in [0,1]} \left| Q_f(u) - Q_{f|a_i}(u) \right| du \ . \tag{2}$$

We say that the predictor $f$ is fair with respect to $A_i$ under Demographic Parity if and only if $\mathcal{U}_{\square}(f) = 0$, where $\square \in \{TV, KS, i\}$. This also ties in with the Wasserstein Distance, as it is clear that the DP condition is satisfied if and only if the Wasserstein distance between the group-wise distributions is zero. To facilitate both the attribution of each individual sensitive variable to a global measure of unfairness of a predictor, we opt for an additive measure of unfairness.

**Multiple sensitive attributes (MSA).** For the multiple sensitive attributes $A_i, \ldots, A_{i+k}$, their collective unfairness is assessed simply by summing the unfairness of each marginal sensitive attribute:

$$\mathcal{U}_{\{i,\ldots,i+k\}}(f) = \mathcal{U}_{i:i+k}(f) = \mathcal{U}_i(f) + \cdots + \mathcal{U}_{i+k}(f) \ . \tag{3}$$

Note here that from the definition of the unfairness measure in the MSA context, an ordered sequence of fairness variables is supposed. This will be important later on, as it allows to *sequentially* render fair a model.

Within **EquiPy** the measure for both SSA and MSA is implemented in the function `unfairness`, which expects a vector of scores of size $N \times 1$ (`predictions` in the example below) and a matrix of size $N \times r$, with $N$ the number of observations and $r$ the number of sensitive attributes (`sensitive_features` in the example), that can be either in the form of a dataframe **pandas** or an array using **numpy**. We recommend using **pandas** objects to specify the sensitive attribute names, facilitating graph reading.

```
import numpy as np
import pandas as pd
from equipy.metrics import unfairness

predictions = np.array([0.05, 0.08, 0.9, 0.5, 0.18, 0.92, 0.9, 0.5,
    0.16, 0.79])
sensitive_features = pd.DataFrame({'origin': [1, 0, 0, 1, 1, 1, 0, 0,
    0, 1], 'gender': [1, 1, 1, 0, 0, 1, 0, 0, 0, 1]})
unfairness(predictions, sensitive_features)
>>> 0.472
```

Under the hood, **EquiPy** can calculate the unfairness metric directly using the optimal transportation plan through the option `approximate=False`, but as this requires solving an $N \times N$ linear program, the default method approximates the transportation cost using a grid on the empirical quantile functions specified in Equation (2), which is computationally less demanding.

### 2.3   Optimal Fair Projection

The core problem within algorithmic fairness is to conduct the trade-off between model accuracy and minimizing unfairness. Having a constant predictor would ensure fairness across the predictions but is certainly undesirable from the view of predictive accuracy. To address this issue, the literature has developed tools to obtain optimal predictive accuracy under the fairness constraint. To both quantify and mitigate predictions from a model, we consider $f^*(\boldsymbol{X}, \boldsymbol{A}) := \mathbb{E}[Y | \boldsymbol{X}, \boldsymbol{A}]$ the Bayes rule that minimizes the squared risk

$$\mathcal{R}(f) := \mathbb{E}(Y - f(\boldsymbol{X}, \boldsymbol{A}))^2 \ , \tag{4}$$

as the *optimal* model,

$$f^\star := \underset{f \in \mathcal{F}}{\mathrm{argmin}} \big\{ \mathcal{R}(f) \big\}$$

In practical terms, this represents the model estimated on the data without any constraints. To achieve fair predictions, restrictions need to be imposed on the class of available models, denoted $\mathcal{F}$.

**SSA case.**   We define the class of DP-fair models for the $i$-th sensitive feature $A_i$,

$$\mathcal{F}_{\mathrm{fair},i} := \big\{ f \in \mathcal{F} \text{ s.t. } f(\boldsymbol{X}, \boldsymbol{A}) \perp\!\!\!\perp A_i \big\} = \big\{ f \in \mathcal{F} \text{ s.t. } \mathcal{U}_i(f) = 0 \big\} \subset \mathcal{F} \ .$$

Fairness for the feature $A_i$ is achieved by projection onto a fair subspace, resulting in a fair predictor $f_{B_i}$

$$f_{B_i} \in \underset{f \in \mathcal{F}_{\mathrm{fair},i}}{\mathrm{argmin}} \big\{ \mathcal{R}(f) \big\} \ .$$

Given a risk $\mathcal{R}$, a class $\mathcal{F}$ and the fair subclass $\mathcal{F}_{\mathrm{fair},i}$ for $A_i$, the price of fairness is quantified by

$$\mathcal{E}_{\mathrm{fair},i}(\mathcal{F}) = \min_{f \in \mathcal{F}_{\mathrm{fair},i}} \big\{ \mathcal{R}(f) \big\} - \min_{f \in \mathcal{F}} \big\{ \mathcal{R}(f) \big\} = \min_{f \in \mathcal{F}_{\mathrm{fair},i}} \big\{ \mathcal{R}(f) \big\} - \mathcal{R}(f^*) \ .$$

Chzhen et al. (2020) and Gouic et al. (2020) have both demonstrated that in the SSA case, under the DP notion of fairness, the price of fairness for $A_i$ corresponds to the following Wasserstein barycenter:

$$\mathcal{E}_{\mathrm{fair},i}(\mathcal{F}) = \inf_{f \in \mathcal{F}_{\mathrm{fair},i}} \sum_{a_i \in \mathcal{A}_i} p_{a_i} \cdot \mathcal{W}_2^2 \big( \nu_{f^*|a_i}, \nu_f \big) \ , \tag{5}$$

where $p_{a_i} := \mathbb{P}(A_i = a_i)$, and its minimum is achieved when using (see Equation (1)) $\mathrm{Bar}(p_{a_i}, \nu_{f^*|a_i})_{a_i \in \mathcal{A}_i}$—also denoted as $\mu_{\mathcal{A}_i}(\nu_{f^*})$, with $\nu_{f^*}$ representing the measure associated with the optimal (DP-unconstrained) predictor $f^*$.

**MSA case.** Hu et al. (2024) extends these formulations to the case of multiple sensitive attributes, where if we denote $\mathcal{F}_{\mathrm{fair}} := \big\{ f \in \mathcal{F} \mid \mathcal{U}(f) = 0 \big\} \subset \mathcal{F}$ as the class of DP-fair models for all sensitive attributes $\boldsymbol{A} = (A_1, \ldots, A_r)$ (or a subset), we have in summary the corresponding:

$$
\begin{cases}
\text{fair predictor:} & f_B \in \underset{f \in \mathcal{F}_{\mathrm{fair}}}{\mathrm{argmin}} \big\{ \mathcal{R}(f) \big\} \\
\text{price of fairness:} & \mathcal{E}_{\mathrm{fair}}(\mathcal{F}) := \underset{f \in \mathcal{F}_{\mathrm{fair}}}{\min} \big\{ \mathcal{R}(f) \big\} - \underset{f \in \mathcal{F}}{\min} \big\{ \mathcal{R}(f) \big\} = \mathcal{R}(f_B) - \mathcal{R}(f^*) \\
\text{optimal fair distribution:} & \mu_{\mathcal{A}}(\nu_{f^*}) := \mathrm{Bar}(p_{\boldsymbol{a}}, \nu_{f^*|\boldsymbol{a}})_{\boldsymbol{a} \in \mathcal{A}} \ .
\end{cases}
$$

In this context, $p_{\boldsymbol{a}} := \mathbb{P}(\boldsymbol{A} = \boldsymbol{a})$. Throughout this article, we consistently refer to $f_{B_i}$ as the aforementioned fair predictor in SSA and $f_B$ as the above overall fair predictor, both derived through the Wasserstein barycenter method. The associated measures are denoted respectively as $\nu_{f_{B_i}} = \mu_{\mathcal{A}_i}(\nu_{f^*})$ and $\nu_{f_B} = \mu_{\mathcal{A}}(\nu_{f^*})$.

Drawing upon Chzhen et al. (2020) for SSA and Hu et al. (2024) for MSA, the fair predictor $f_B$ is considered optimal as it minimizes the risk $\mathcal{R}$ among fair predictors. Its closed-form solution can be expressed as:

$$
f_B(\boldsymbol{x}, \boldsymbol{a}) = \left( \sum_{\boldsymbol{a}' \in \mathcal{A}} p_{\boldsymbol{a}'} Q_{f^*|\boldsymbol{a}'} \right) \circ F_{f^*|\boldsymbol{a}} \left( f^*(\boldsymbol{x}, \boldsymbol{a}) \right) \quad \text{for all } (\boldsymbol{x}, \boldsymbol{a}) \in \mathcal{X} \times \mathcal{A} \ . \tag{6}
$$

**In practice.** Within the context of SSA, DP-fairness can be achieved using a straightforward procedure. **EquiPy** offers the class `FairWasserstein` to calibrate DP-fairness by internally calculating the empirical counterparts of $p_{\boldsymbol{a}}$, $F_{f^*|\boldsymbol{a}}$, and $Q_{f^*|\boldsymbol{a}}$ in Equation (6), where $\boldsymbol{A} \in \mathcal{A}$ contains a unique sensitive attribute here. The class expects as input the sensitive variable as a vector, represented by `sensitive_feature_calib` either using **pandas** or **numpy**, and the predictions `predictions_calib` of an estimator of $f^*$. Again, the estimator can be any trained ML algorithm. DP-fair predictions can then be obtained by transforming given predictions from a test-set. The resulting empirical formulation is implemented in the **EquiPy** package and can be accessed as follows:

```python
# Single Sensitive Attribute (SSA) case
import numpy as np
import pandas as pd
from equipy.fairness import FairWasserstein

calibrator = FairWasserstein(sigma = 0.0001)

# calibration set
predictions_calib = np.array([0.05, 0.08, 0.9, 0.5, 0.18, 0.92, 0.9,
    0.5])
sensitive_feature_calib = pd.DataFrame({'origin': [1, 0, 0, 1, 1, 1,
    0, 0]})
calibrator.fit(predictions_calib, sensitive_feature_calib)
```

```
# studied set
predictions = np.array([0.16, 0.79])
sensitive_feature = pd.DataFrame({'origin': [0, 1]})

calibrator.transform(predictions, sensitive_feature)
>>> array([0.271, 0.752])
```

Note that setting `sigma = 0.0001`, which is the default value for `sigma`, corresponds to what is commonly referred to as the *jittering* process. This involves introducing a small perturbation into a given predictor by adding continuous random noise. The purpose of this perturbation is to prevent the occurrence of *atoms* in the predictions; for further details on regression, refer to (Chzhen et al., 2020), and for classification, refer to (Denis et al., 2021).

## 3 Sequential Fairness and Software

Whereas the application of the Wasserstein Barycenter is straightforward in the SSA context, multiple sensitive attributes can pose problems within both the estimation and interpretation of the effects of the projection. A simple possibility would be to recode multiple sensitive features into a single, multi-class sensitive feature. However, this introduces problems in the robustness of the estimation as it introduces exponentially many subgroups on a fixed size data set. Further, such a formulation hampers approximate fairness, discussed below. In line with the article by Hu et al. (2024), we present in this section a concise exposition of Equation (6) to improve transparency regarding the inter-correlations among sensitive features. This breakdown facilitates a clearer understanding. Following this, we show how the **EquiPy** Package addresses these aspects in practical applications.

### 3.1 A Sequentially Fair Mechanism

Recall that $f_{B_i}$ represents a predictor that ensures fairness with respect to the sensitive attribute $A_i$. To simplify notation, we express the composition $f_{B_i} \circ f_{B_j}$ by a slight abuse of notation as:

$$\left(f_{B_i} \circ f_{B_j}\right)(\boldsymbol{x}, \boldsymbol{a}) = \left(\sum_{a_i' \in \mathcal{A}_i} p_{a_i'} Q_{f_{B_j}|a_i'}\right) \circ F_{f_{B_j}|a_i}\left(f_{B_j|a_i}(\boldsymbol{x}, \boldsymbol{a})\right) \quad \text{for all } (\boldsymbol{x}, \boldsymbol{a}) \in \mathcal{X} \times \mathcal{A} .$$

(7)

Drawing on insights from the article by Hu et al. (2024), the overarching barycentric fair predictor $f_B$ within the space $\mathcal{A}$ can be articulated as a sequence of compositions of marginal "barycentric" fair predictors, advancing sequentially through $\mathcal{A}_1, \ldots, \mathcal{A}_r$.

$$f_B(\boldsymbol{x}, \boldsymbol{a}) = \left(f_{B_1} \circ f_{B_2} \circ \ldots \circ f_{B_r}\right)(\boldsymbol{x}, \boldsymbol{a}) \quad \text{for all } (\boldsymbol{x}, \boldsymbol{a}) \in \mathcal{X} \times \mathcal{A} ,$$

corresponding to the composition of measures: $\mu_{\mathcal{A}}(\nu_{f^*}) = \mu_{\mathcal{A}_1} \circ \mu_{\mathcal{A}_2} \circ \cdots \circ \mu_{\mathcal{A}_r}(\nu_{f^*})$. It's important to note that the fairness mitigation remains unaffected by the order in which the barycentric mitigation is executed. Hence, for $r = 2$ sensitive features, $f_{B_1} \circ f_{B_2} = f_{B_2} \circ f_{B_1}$. The proposition and its corresponding proof can be found in (Hu et al., 2024).

**In practice.** The sequential approach for multiple sensitive variables is implemented in the `MultiWasserstein` class, which internalizes the computations, enabling a straightforward adaption similar to the `FairWasserstein` class. It can be executed as follows:

```
import numpy as np
import pandas as pd
from equipy.fairness import MultiWasserstein

calibrator = MultiWasserstein(sigma = 0.0001)

# calibration set
predictions_calib = np.array([0.05, 0.08, 0.9, 0.5, 0.18, 0.92, 0.9,
    0.5])
sensitive_features_calib = pd.DataFrame({'origin': [1, 0, 0, 1, 1, 1,
    0, 0], 'gender': [1, 1, 1, 0, 0, 1, 0, 0]})
calibrator.fit(predictions_calib, sensitive_features_calib)

# studied set
predictions = np.array([0.16, 0.79])
sensitive_features = pd.DataFrame({'origin': [0, 1],
    'gender': [0, 1]})

calibrator.transform(predictions, sensitive_features)

>>> array([0.203, 0.361])
```

For further details, we refer to the paper (Hu et al., 2024) or the dedicated package documentation in https://equilibration.github.io/equipy/.

## 3.2 Extensions to Approximate Fairness

In our context, achieving *approximate fairness* involves improving fairness relatively and approximately with a preset level of relative fairness improvement, denoted as $\boldsymbol{\varepsilon} = \varepsilon_{1:r} = (\varepsilon_1, \cdots, \varepsilon_r)$ (abbreviated as $\boldsymbol{\varepsilon}$-fairness in this article). The SSA framework discussed in the article by Chzhen and Schreuder (2022) employs geodesic parameterization to generate predictors that approximate fairness. Specifically, for any $A_i \in \mathcal{A}_i$ without loss of generality, the predictor achieving $\varepsilon_i$-fairness takes the form:

$$f_{B_i}^{\varepsilon_i}(\boldsymbol{X}, \boldsymbol{A}) := (1 - \varepsilon_i) \cdot f_{B_i}(\boldsymbol{X}, \boldsymbol{A}) + \varepsilon_i \cdot f^*(\boldsymbol{X}, \boldsymbol{A}) \ .$$

This predictor achieves (optimally) the following risk-fairness trade-off:

$$f_{B_i}^{\varepsilon_i} \in \operatorname*{argmin}_{f \in \mathcal{F}} \{\mathcal{R}(f) : \mathcal{U}_i(f) \leq \varepsilon_i \cdot \mathcal{U}_i(f^*)\} \ .$$

Quoting (Hu et al., 2024), "introducing a sequential approach is pivotal for enhancing clarity". This introduced sequentially fair framework aids in comprehending intricate concepts such as the aforementioned approximate fairness within the MSA context. Indeed, as described in the article by Hu et al. (2024), if we permit

$$f_B^{\boldsymbol{\varepsilon}} \in \operatorname*{argmin}_{f \in \mathcal{F}} \left\{\mathcal{R}(f) : \mathcal{U}(f) \leq \sum_{i=1,\ldots,r} \varepsilon_i \cdot \mathcal{U}_i(f^*)\right\} \ ,$$

then $f_B^{\boldsymbol{\varepsilon}} = f_{B_1}^{\varepsilon_1} \circ \cdots \circ f_{B_r}^{\varepsilon_r}$. This allows us to break down how various components of the sequential fairness mechanism interact to achieve fairness goals. It also facilitates the interpretation of the inherent effects of fairness adjustments.

**In practice:** In **EquiPy**, the empirical counterpart of approximate fairness is inherently incorporated within both the `FairWasserstein` and `MultiWasserstein` classes of the `fairness` module. An illustrative example of its application is provided below:

```python
import numpy as np
import pandas as pd
from equipy.fairness import MultiWasserstein

calibrator = MultiWasserstein(sigma = 0.0001)

# calibration set
predictions_calib = np.array([0.05, 0.08, 0.9, 0.5, 0.18, 0.92, 0.9,
    0.5])
sensitive_features_calib = pd.DataFrame({'origin': [1, 0, 0, 1, 1, 1,
    0, 0], 'gender': [1, 1, 1, 0, 0, 1, 0, 0]})
calibrator.fit(predictions_calib, sensitive_features_calib)

# studied set
predictions = np.array([0.16, 0.79])
sensitive_features = pd.DataFrame({'origin': [0, 1],
    'gender': [0, 1]})

# approximate fairness
epsilon = [0.1, 0.2]

calibrator.transform(predictions, sensitive_features,
    epsilon = epsilon)

>>> array([0.210, 0.348])
```
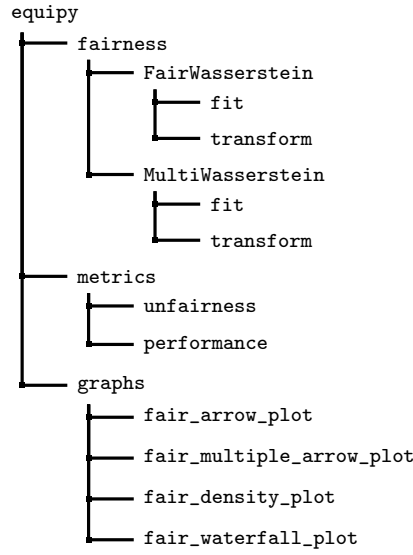
For additional details, we refer to the section below or the dedicated package documentation: https://equilibration.github.io/equipy/.

## 4 Description of the EquiPy Package Structure

The preceding sections offered an overview of the main modules and their functionalities illustrating the theoretical results. In this section, we will formally introduce the aforementioned Python classes and modules, along with all associated parameters (from version 0.0.11), before diving into practical illustrations demonstrating explicit applications and visual displays in Section 5.

### 4.1 Overview of the Package Structure

The package is structured around three core modules: `fairness`, `metrics` and `graphs`, where each module contains several sub-modules that help the user access the relevant functions. Refer to Figure 2 for an overview of the tree-structure. Additionally, Table 1 describes every sub-module in detail. Note that the the classes `FairWasserstein` and `MultiWasserstein` both contain a `fit` and `transform` method, similar to other packages used in ML, like **scikit-learn** in Python.

```
equipy
  ├── fairness
  │     ├── FairWasserstein
  │     │     ├── fit
  │     │     └── transform
  │     └── MultiWasserstein
  │           ├── fit
  │           └── transform
  ├── metrics
  │     ├── unfairness
  │     └── performance
  └── graphs
        ├── fair_arrow_plot
        ├── fair_multiple_arrow_plot
        ├── fair_density_plot
        └── fair_waterfall_plot
```

*Figure 2: Tree structure of the **EquiPy** package*

| Method | Description |
| --- | --- |
| `fit` | Fit distribution and quantile functions using the calibration data. Note that in practice, this calibration data could either be the training data or unlabeled data. |
| `transform` | Transform data to enforce fairness with Wasserstein distance using results of `fit`. In practical applications, we typically apply it to a hold-out (or test) data. |
| `unfairness` | Compute the unfairness measure (empirical version of Equation 3) for a given fair output and multiple sensitive attributes data. |
| `performance` | Compute the performance metric for predicted fair outputs compared to the true targets. |
| `fair_arrow_plot` | Create an arrow plot illustrating the progression of `fairness-performance` combinations of a predictive model, step by step according to the sensitive attributes. |
| `fair_multiple_arrow_plot` | This method uses `fair_arrow_plot` with respect to different permutations of sensitive attributes. |
| `fair_density_plot` | Visualize the distribution of predictions conditional on different sensitive features using kernel density estimates. We use a Beta kernel when the predictive task is binary classification. |
| `fair_waterfall_plot` | Generate a waterfall plot to visually represent how sequential fairness impacts a fairness metric associated with one or more sensitive attributes within a model. |

*Table 1: Overview of the core methods for the modules `fairness`, `metrics` and `graphs`*

## 4.2   Description of the Package Modules and Sub-Modules

The main modules follow the standard workflow in fairness mitigation. Mitigating biases through the `fairness` module, measuring both the predictive performance and unfairness

using `metric` to enable a quantitative impact assessment and plotting utilities in `graphs`, for further interpretation and reporting. This subsection provides details on the three main modules, including comprehensive information about their parameters and validation functions.

**Module fairness**    This module comprises two primary classes : `FairWasserstein`, ensuring fair predictions regarding a single sensitive attribute (Section 2.2), and `MultiWasserstein`, with a similar structure but addressing fairness with respect to multiple sensitive attributes (Section 2.2). Both classes implement fairness adjustment on model predictions related to one or multiple sensitive attributes, using Wasserstein distance for binary classification and regression tasks, as introduced in Section 3. In the case of binary classification, this class supports scores instead of predicted labels. A specific instance demonstrating the application of `FairWasserstein` and `MultiWasserstein` are presented respectively in Section 2.3 and Section 3. In addition, Figure 3 illustrates the application process of `FairWasserstein` and `MultiWasserstein`, highlighting their model-agnostic nature. Specifically, the data not used during the training phase of the ML model is divided into calibration and test sets. The `fit` method is trained on the calibration dataset to learn the mathematical quantities from Equation 7, after which the `transform` method is applied to the test set to obtain fair predictors using the functions learned during the calibration step. The resultant fair predictions, denoted as $\hat{y}_{\text{test}}^{\text{fair}}$, can subsequently be employed to calculate various methods within the `metrics` module. To enable approximate fairness, as described in Section 3.2, the user can specify values for the `epsilon` parameter, a vector of size $r$, in both `transform` functions of `FairWasserstein` and `MultiWasserstein`.
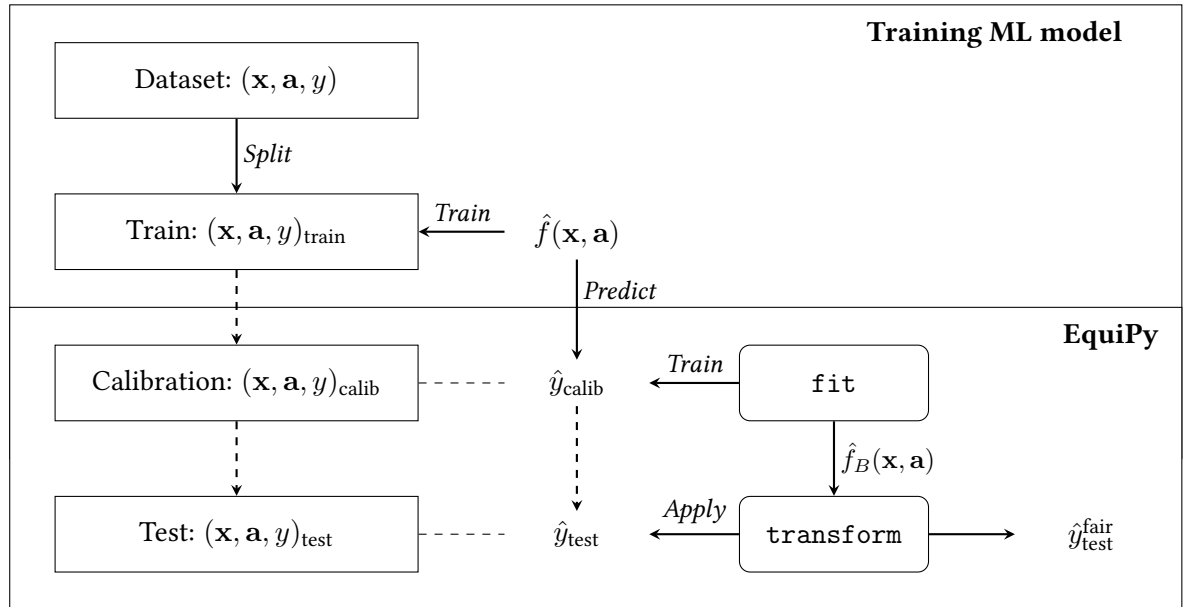


*Figure 3: Process of mitigating predictions using the methods of* `MultiWasserstein` *class*

**Module metrics**    The **EquiPy** package provides functionality for computing two types of metrics : performance and unfairness. Performance metrics, calculated through the method `performance`, refer to commonly used regression or classification metrics such as `accuracy_score` or `mean_squared_error` from `sklearn.metrics` module but also allow custom user specification. The `performance` method provides a wrapper around the specified functions which allows a seamless transition between the `metrics` and `graphs` modules. It

expects as inputs predictions and labels in array form, compatible with the underlying performance score. Hence for classification metrics, the predicted scores need to be transformed to label predictions first. The unfairness metric, as defined by Equation (3) in Section 2, can be computed using the function `unfairness`. Here, the sensitive variables need to be provided as a `pandas.DataFrame` or a `numpy.array` object. The unfairness metric calculation is either computed using the Wasserstein distance from the **POT** package or either determined as the maximum difference in quantiles defined on a grid on [0,1] between the two populations. The default approach is the latter (`approximate=True`), and we recommend using the **POT** package (setting `approximate=False` in the `unfairness` function) when the number of observations is insufficient to efficiently approximate the quantile functions by sensitive group. A practical example of the usage is provided in Section 2.2. As an extension, the module also includes hidden methods used in the `graphs` module to provide an overview and interpretations of variations in unfairness and performance across different permutations of the sensitive variables.

**Module graphs**    This module provides utilities for the visualization of the correction paths using three different graph types. The method `fair_arrow_plot` illustrates the fairness-performance relationship based on the ordered sequence of sensitive attributes specified by the user. The `fair_multiple_arrow_plot` provides an overview of the same relationships but for all possible permutations of the sensitive attributes, enabling a prioritization of different approaches. The `fair_density_plot` method presents the probability distributions of predictions relative to the values of the sensitive attributes. Lastly, `fair_waterfall_plot` demonstrates the sequential gain in fairness using waterfall plots. All of the `graphs` utilities also permit plots using approximate fairness, as explained in Section 3.2, if $\varepsilon$ is provided by the user. This is done by setting `epsilon` to specific vector values in all of the aforementioned `graphs` functions. The next Section 5 demonstrates all of these methods in a case study to illustrate their specific use.

All described methods are provided with doc-string documentation and example usage within the source code. For additional and in-depth details of their use we refer to the next section or the dedicated package documentation https://equilibration.github.io/equipy/.

## 5    Illustrations: Case Study on Census Data

### 5.1    Data and Basemodel

We demonstrate the functionalities of **EquiPy** on data derived from the US Census, made available within the **Folktables** package (Ding et al., 2021), which enables a straightforward replication of all results. Here, we specifically consider a regression task, where we predict the log of an individuals' total income (based on the *ACSIncome* task of **Folktables**). Note that the same logic also applies to binary classification tasks, with the condition that the model predictions for the *scores* are obtained and not the predicted class label.

We opt for a simple basemodel to construct our predictor, based on a gradient boosted machine of the **LightGBM** implementation (Ke et al., 2017), due to its proven track record for fast and accurate predictions and widespread use among the data science community. The data is split into three sets (`train`, `calib`, `test`), where for the `calib` set, no labels are needed. In practice, if there is no possibility to obtain a calibration set, the `train` data set may be used, although it is not recommended as described in (Denis et al., 2021). The tuning of the model itself is considered out of the scope of the package, and we largely omit a discussion thereof, but the interested reader can find all information in the *demo* section of the github

repository. We suppose that all relevant hyperparameters have been gathered in a dictionary called `optimized_parameters` and the base model is then simply fitted and predictions are obtained as:

```
train_data = lightgbm.Dataset(data = X_train, label = y_train)
model = lightgbm.train(train_set = train_data,
    params = optimized_parameters)

predictions_calib = model.predict(X_calib)
predictions_test = model.predict(X_test)
```

## 5.2   Unfairness Mitigation

We initially investigate the prediction for a sensitive ethnicity, indexed as 1 if a given observation is a member of this group and 0 otherwise. The left pane of Figure 4 depicts the densities of the initial predictions, as obtained directly from the model. Visually, there is a discrepancy between the predictions, and members of the sensitive group (indexed as 1) are predicted to have significantly lower incomes. We then proceed to measure the unfairness to quantify the bias, by using the sensitive variable `sens_ethn_test`:

```
from equipy.metrics import unfairness

unfairness(predictions_test, sens_ethn_test)
>>> 0.437
```

There seems to be a significant bias also according to the quantitative measure. The next step is to DP-calibrate the scores and post-process the test predictions. This is simply an application of the `FairWasserstein` function, because here we consider unfairness regarding a single sensitive attribute:

```
from equipy.fairness import FairWasserstein

calibrator = FairWasserstein()
calibrator.fit(predictions_calib, sens_ethn_calib)

predictions_test_fair = calibrator.transform(predictions_test,
    sens_ethn_test)

unfairness(predictions_test_fair, sens_ethn_test)
>>> 0.067
```

The simple approach reduced the unfairness significantly. To evaluate the impact on the predictive performance, we can then use the provided `performance` utility function:

```
from equipy.metrics import performance

performance(y_test, predictions_test)
>>> 0.544

performance(y_test, predictions_test_fair)
>>> 0.552
```

where we can detect a marginal increase in the mean squared error of the predictions, which is the default performance function. The changes can then easily be visualized using the plotting function `fair_density_plot` as depicted in Figure 4.

```
from equipy.graphs import fair_density_plot


fair_density_plot(sens_ethn_calib, sens_ethn_test, predictions_calib,
    predictions_test);
```
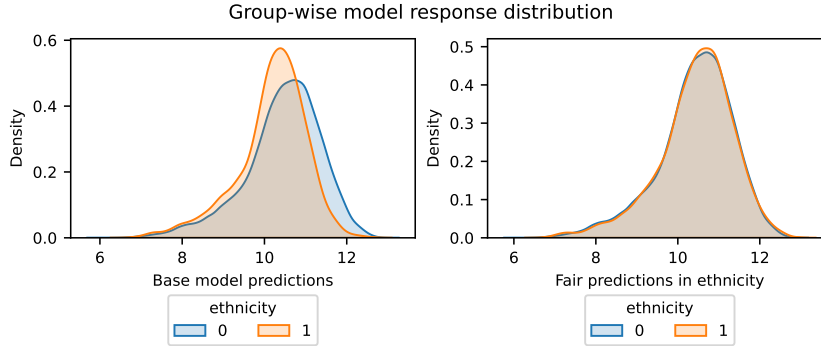


Figure 4: *Group-wise model response distribution. The fair model predictions exhibit no variations across different groups.*

As anticipated, based on predictions from the test set, the fair middle-ground distribution, derived from the Wasserstein barycenter, results in a perfect alignment of the conditional distributions. Consequently, the unfairness measure $\mathcal{U}_1$ is significantly reduced.

## 5.3 Unfairness Mitigation with MSA

The extension to the MSA case is also straightforward. As is often the case, predictions can be significantly biased with respect to more than one attribute. For example, the income distribution is often also skewed with respect to sex. Here, the predictions where fairness is imposed with respect to ethnicity can still exhibit a bias with respect to sex, or in some cases even exacerbate certain discrepancies, see (Hu et al., 2024) for an example thereof. **EquiPy** can handle such situations through the usage of the `MultiWasserstein` calibrator. Here, we specify the correction order, first addressing ethnicity and then sex. The application is similar to the SSA case:

```
from equipy.fairness import MultiWasserstein


sens_twovar_calib = pd.DataFrame({'ethnicity': sens_ethn_calib,
    'sex': sens_sex_calib})
sens_twovar_test = pd.DataFrame({'ethnicity': sens_ethn_test,
    'sex': sens_sex_test})


unfairness(predictions_test, sens_twovar_test)
>>> 0.783


calibrator_twovar = MultiWasserstein()
calibrator_twovar.fit(predictions_calib, sens_twovar_calib)
```
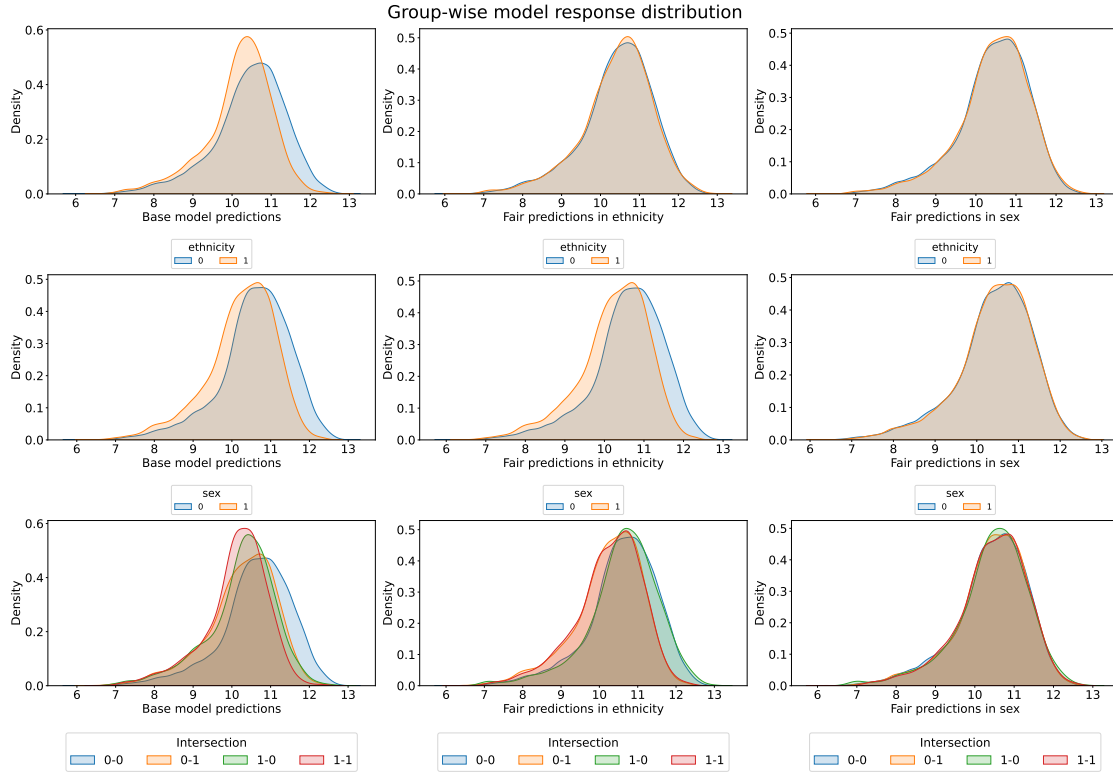
*Figure 5: Group-wise model response distribution for multiple sensitive attributes: ethnicity and sex. The fair model predictions exhibit no variations across different groups.*

```
predictions_twovar_test_fair = calibrator_twovar.transform(
    predictions_test, sens_twovar_test)

unfairness(predictions_twovar_test_fair, sens_twovar_test)
>>> 0.106
```

Again, the unfairness is significantly reduced. For multiple sensitive features the plotting utilities become paramount to compare the impact that each of the variables had on the final outcome. Further, the plotting utilities can help a decision maker to prioritize certain constraints over others, by providing a full impact analysis.

## 5.4 Visualization Using the Graphs Module

The simplest way to visualize the unfairness-performance trade-off is to use `fair_multiple_arrow_plot`, which plots the trade-off in terms of performance (that is, mean squared error for this application) and unfairness for all possible order of correction for the sensitive attributes. This enables verification of the algorithm's results, as the predictions should align after the final correction, regardless of the correction sequence. To be able to compare performance on test set before and after applying the mitigation technique, `y_test` must be specified, which contained the true labels. Figure 6 can be produced by calling the graphs module as:

```
from equipy.graphs import fair_multiple_arrow_plot
fair_multiple_arrow_plot(sens_twovar_calib, sens_twovar_test,
    predictions_calib, predictions_test, y_test)
```

The results show that correcting the scores for sex has a larger impact on the predictive performance as compared to correcting the scores with ethnicity first, while having a similar impact on the reduced unfairness, depicted with the higher slope for the former. This illustration enables to bring to light differences in correction strategies, in turn mitigating concerns about *fairwashing* (refer to the article Aïvodji et al. (2019) for a more in-depth discussion).
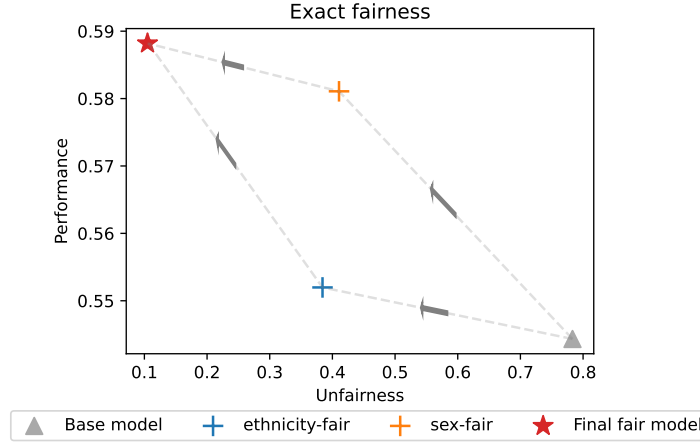


*Figure 6: Arrow graphic depicting the performance-unfairness trade-off given the two sensitive features.*

As discussed in Section 3.2, splitting the MSA problem into sub-problems has the additional advantage that each sensitive feature can be weighted to achieve approximate fairness. This can be visualized using the `fair_waterfall_plot` using a command such as:

```
from equipy.graphs import fair_waterfall_plot

fair_waterfall_plot(sens_twovar_calib, sens_twovar_test,
    predictions_calib, predictions_test, epsilon=[0.5, 0.25])
```

The goal of the waterfall plot is to visualize how the unfairness score is reduced at each step. Figure 7 depicts the gain in fairness while applying different sequential orders for ethnicity and sex with similar values for $\varepsilon$ in the `fair_waterfall_plot`. Here, correcting first for ethnicity with an epsilon factor of 50% enables the correction of 28% of the total unfairness (while reducing unfairness in ethnicity by 50%) whereas correcting secondly for ethnicity with the same epsilon factor enables the correction of 24% of the total unfairness. Given that sensitive variables can also correlate, an approximate correction might have influences in the mitigation step of the second variable. Nevertheless, the final unfairness value remains the same for both orders.

## 5.5 Contribution of Each Sensitive Variable to Unfairness

When applying the mitigation approach to MSA, it is also possible to decompose the unfairness attributed to each sensitive attribute along the sequential correction path for a given order (here, ethnicity then sex). This can be done by retrieving the dictionary `y_fair` from the `MultiWasserstein` module, which contains the initial predictions from the base model, the predictions after correcting for ethnicity in the first step, and the predictions after correcting for sex in the second and final step. The complete replication code for calculating unfairness in both ethnicity and sex at each correction step is provided in Appendix B. Table 2 presents the unfairness values for this illustrative example.
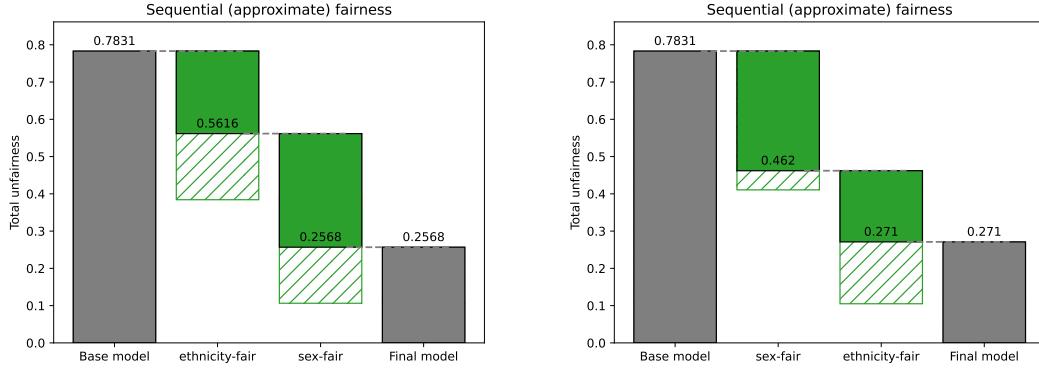
*Figure 7: Waterfall plots for the two correction approaches. Left pane, correcting first for ethnicity with a ε factor of 0.5 and then for sex with a ε factor of 0.25. Right pane, correcting first for sex with a ε factor of 0.25 and then for ethnicity with a ε factor of 0.5.*

|  | $\mathcal{U}_{\text{ethnicity}}$ | $\mathcal{U}_{\text{sex}}$ | $\mathcal{U}_{\text{sex,ethnicity}}$ |
|---|---|---|---|
| Base model $f$ | 0.4366 | 0.3465 | 0.7831 |
| 1. $f_{B_{\text{ethnicity}}}$ | 0.0466 | 0.3376 | 0.3842 |
| 2. $f_B = f_{B_{\text{ethnicity,sex}}}$ | 0.0726 | 0.0338 | 0.1064 |

*Table 2: Unfairness decomposition for sex and ethnicity calculated on Census Data. $f_{B_{ethnicity}}$ represents the model fair with respect to ethnicity, obtained after the first correction step in **EquiPy**, while $f_B$ corresponds to the model fair with respect to both ethnicity and sex, obtained after the second correction step in **EquiPy**.*

```
calibrator_twovar = MultiWasserstein()
calibrator_twovar.fit(predictions_calib, sens_twovar_calib)

predictions_twovar_test_fair = calibrator_twovar.transform(
    predictions_test, sens_twovar_test)
y_seq_fair = calibrator_msa.y_fair

# Unfairness in ethnicity before mitigation
 unfairness(y_seq_fair['Base model'],
    sens_twovar_test[['ethnicity']])
>>> 0.437

# Unfairness in sex before mitigation
 unfairness(y_seq_fair['Base model'], sens_twovar_test[['sex']])
>>> 0.347
```

Unfairness decomposition in Table 2 helps assess and interpret the contribution of each sensitive attribute to the total unfairness metric. Notably, as can be seen in the right-middle graph of Figure 5, correcting for sex in the second step *almost* does not increase unfairness in ethnicity (compared to the unfairness initial value), which is desirable, as the goal is to maintain fairness across all sensitive attributes throughout the sequential correction process. One might also want to evaluate the unfairness in ethnicity after mitigating only for sex, as these two variables can be correlated. Indeed, in practice, ethnicity is often unavailable to practitioners,

19

so it is important to ensure that ensuring fairness for one observed sensitive attribute (for example, sex) does not inadvertently compromise fairness for another unobserved attribute (for example, ethnicity).

## 6   Summary and Discussion

The **EquiPy** package provides simple and high level functionalities to achieve fairness as defined by the Demographic Parity condition. By relying on a post-processing approach, it remains truly agnostic to the underlying models predictions that should be corrected. Instead of addressing the issue arising from multiple sensitive attributes by regrouping features, **EquiPy** relies on the associativity of Wasserstein Barycenters to achieve a sequential fairness calibration. This in turn allows a more robust estimation and also an easier interpretation of the results. Next to the mitigation techniques, the package also provides extensive visualization utilities that enable non-technical users a broader understanding of the trade-off involved.

## A  More technical details

Consider two ($k = 2$) distributions, conditionnaly Gaussian, $\mathcal{N}(\mu_k, \sigma_k^2)$, Wasserstein barycenter is a Gaussian distribution $\mathcal{N}(\mu, \sigma^2)$ where

$$\begin{cases} \mu = w_1\mu_1 + w_2\mu_2 \\ \sigma^2 = (w_1\sigma_1 + w_2\sigma_2)^2 \end{cases}$$

Recall that the mixture has also mean $\mu = w_1\mu_1 + w_2\mu_2$, but has variance

$$v = w_1(\mu_1^2 + \sigma_1^2) + w_2(\mu_2^2 + \sigma_2^2) - \left(w_1\mu_1 + w_2\mu_2\right)^2 = w_1\sigma_1^2 + w_2\sigma_2^2 + w_1w_2(\mu_1 - \mu_2)^2.$$

Since the function $x \to x^2$ is convex, we can apply Jensen's inequality to $\sigma^2$,

$$(w_1\sigma_1 + w_2\sigma_2)^2 \leq w_1\sigma_1^2 + w_2\sigma_2^2$$

We recognize on the right the first part of the variance of the mixture, and the second part, $w_1w_2(\mu_1 - \mu_2)^2$ is positive, therefore

$$\sigma^2 = (w_1\sigma_1 + w_2\sigma_2)^2 \leq w_1\sigma_1^2 + w_2\sigma_2^2 + w_1w_2(\mu_1 - \mu_2)^2 = v$$

## B  Illustrations: Unfairness decomposition

Below is the replication code from the *demo* section of the github repository, which calculates unfairness for each sensitive attribute considered in the mitigation approach for multiple sensitive attributes using the Python package **EquiPy**, applied to the Census Data case study in Section 5.

```
calibrator_twovar = MultiWasserstein()
calibrator_twovar.fit(predictions_calib, sens_twovar_calib)

predictions_twovar_test_fair = calibrator_twovar.transform(
    predictions_test, sens_twovar_test)
y_seq_fair = calibrator_msa.y_fair

# Unfairness in ethnicity before mitigation
 unfairness(y_seq_fair['Base model'],
    sens_twovar_test[['ethnicity']])
>>> 0.437

# Unfairness in sex before mitigation
 unfairness(y_seq_fair['Base model'], sens_twovar_test[['sex']])
>>> 0.347

# Unfairness in ethnicity after mitigation w.r.t. ethnicity
unfairness(y_seq_fair['ethnicity'], sens_twovar_test[['ethnicity']])
>>> 0.047

# Unfairness in sex after mitigation w.r.t. ethnicity
unfairness(y_seq_fair['ethnicity'], sens_twovar_test[['sex']])
>>> 0.338
```

```
# Unfairness in ethnicity after mitigation w.r.t. ethnicity and sex
unfairness(y_seq_fair['ethnicity'], sens_twovar_test[['sex']])
>>> 0.073

# Unfairness in sex after mitigation w.r.t. ethnicity and sex
unfairness(y_seq_fair['sex'], sens_twovar_test[['sex']])
>>> 0.034
```

# References

Agarwal, A., Beygelzimer, A., Dudík, M., Langford, J. and Wallach, H. (2018). A Reductions Approach to Fair Classification. In *Proceedings of the 35th International Conference on Machine Learning*.

Agueh, M. and Carlier, G. (2011). Barycenters in the wasserstein space. *SIAM Journal on Mathematical Analysis* 43: 904–924.

Aïvodji, U., Arai, H., Fortineau, O., Gambs, S., Hara, S. and Tapp, A. (2019). Fairwashing: the risk of rationalization. In *International Conference on Machine Learning*. PMLR, 161–170.

Barocas, S., Hardt, M. and Narayanan, A. (2018). *Fairness and Machine Learning*. fairmlbook.org.

Bellamy, R. K., Dey, K., Hind, M., Hoffman, S. C., Houde, S., Kannan, K., Lohia, P., Martino, J., Mehta, S., Mojsilović, A. et al. (2019). Ai fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias. *IBM Journal of Research and Development* 63: 4–1.

Boudt, K., Kleen, O. and Sjørup, E. (2022). Analyzing intraday financial data in r: The highfrequency package. *Journal of Statistical Software* 104: 1–36.

Buolamwini, J. and Gebru, T. (2018). Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*. PMLR, 77–91.

Charpentier, A., Flachaire, E. and Gallic, E. (2023). Optimal transport for counterfactual estimation: A method for causal inference. In *Optimal Transport Statistics for Economics and Related Topics*. Springer, 45–89.

Chiappa, S., Jiang, R., Stepleton, T., Pacchiano, A., Jiang, H. and Aslanides, J. (2020). A general approach to fairness with optimal transport. *Proceedings of the AAAI Conference on Artificial Intelligence* 34: 3633–3640.

Chzhen, E., Denis, C., Hebiri, M., Oneto, L. and Pontil, M. (2020). Fair Regression with Wasserstein Barycenters. In *Advances in Neural Information Processing Systems*.

Chzhen, E. and Schreuder, N. (2022). A minimax framework for quantifying risk-fairness tradeoff in regression. *The Annals of Statistics* 50: 2416–2442.

Dastin, J. (2022). Amazon scraps secret AI recruiting tool that showed bias against women. In *Ethics of data and analytics*. Auerbach Publications, 296–299.

Delaney, E. D., Fu, Z., Wachter, S., Mittelstadt, B. and Russell, C. (2024). OxonFair: A Flexible Toolkit for Algorithmic Fairness. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

Denis, C., Elie, R., Hebiri, M. and Hu, F. (2021). Fairness guarantee in multi-class classification. *arXiv preprint arXiv:2109.13642* .

Ding, F., Hardt, M., Miller, J. and Schmidt, L. (2021). Retiring adult: New datasets for fair machine learning. *Advances in Neural Information Processing Systems* 34.

Gaucher, S., Schreuder, N. and Chzhen, E. (2023). Fair learning with Wasserstein barycenters for non-decomposable performance measures. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 2436–2459.

Goodall, N. J. (2014). Ethical decision making during automated vehicle crashes. *Transportation Research Record* 2424: 58–65.

Gouic, T. L., Loubes, J.-M. and Rigollet, P. (2020). Projection to fairness in statistical learning. *arXiv preprint arXiv:2005.11720* .

Hardt, M., Price, E. and Srebro, N. (2016). Equality of opportunity in supervised learning. In *Neural Information Processing Systems.*

Hu, F., Ratz, P. and Charpentier, A. (2023). Fairness in Multi-Task Learning via Wasserstein Barycenters. In Koutra, D., Plant, C., Gomez Rodriguez, M., Baralis, E. and Bonchi, F. (eds), *Machine Learning and Knowledge Discovery in Databases: Research Track.* Springer Nature Switzerland, 295–312.

Hu, F., Ratz, P. and Charpentier, A. (2024). A sequentially fair mechanism for multiple sensitive attributes. In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence.*

Joo, J. and Kärkkäinen, K. (2020). Gender slopes: Counterfactual fairness for computer vision models by attribute manipulation. In *Proceedings of the 2nd international workshop on fairness, accountability, transparency and ethics in multimedia*, 1–5.

Karako, C. and Manggala, P. (2018). Using image fairness representations in diversity-based re-ranking for recommendations. In *Adjunct Publication of the 26th Conference on User Modeling, Adaptation and Personalization*, 23–28.

Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q. and Liu, T.-Y. (2017). LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S. and Garnett, R. (eds), *Advances in Neural Information Processing Systems*, 30. Curran Associates, Inc.

Kim, M. P., Ghorbani, A. and Zou, J. (2019). Multiaccuracy: Black-box post-processing for fairness in classification. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, 247–254.

Kozodoi, N. and Varga, T. V. (2021). **fairness**: Algorithmic Fairness Metrics. R package version 1.2.2.

Nyholm, S. and Smids, J. (2016). The ethics of accident-algorithms for self-driving cars: An applied trolley problem? *Ethical theory and moral practice* 19: 1275–1289.

Obermeyer, Z., Powers, B., Vogeli, C. and Mullainathan, S. (2019). Dissecting racial bias in an algorithm used to manage the health of populations. *Science* 366: 447–453.

Pappalardo, L., Simini, F., Barlacchi, G. and Pellungrini, R. (2022). scikit-mobility: A python library for the analysis, generation, and risk assessment of mobility data. *Journal of Statistical Software* 103: 1–38.

Park, S., Hwang, S., Kim, D. and Byun, H. (2021). Learning disentangled representation for fair facial attribute classification via fairness-aware information alignment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 35, 2403–2411.

Plečko, D. and Meinshausen, N. (2020). Fair data adaptation with quantile preservation. *J. Mach. Learn. Res.* 21.

Plečko, D., Bennett, N. and Meinshausen, N. (2024). fairadapt: Causal reasoning for fair data preprocessing. *Journal of Statistical Software* 110: 1–35.

Qiang, Y., Li, C., Brocanelli, M. and Zhu, D. (2022). Counterfactual interpolation augmentation (cia): A unified approach to enhance fairness and explainability of dnn. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI*, 732–739.

Ratz, P., Hu, F. and Charpentier, A. (2023). Addressing fairness and explainability in image classification using optimal transport. *arXiv preprint arXiv:2308.11090* .

Santambrogio, F. (2015). Optimal transport for applied mathematicians. *Birkäuser, NY* 55: 94.

Tierney, N. and Cook, D. (2023). Expanding tidy data principles to facilitate missing data exploration, visualization and assessment of imputations. *Journal of Statistical Software* 105:

1–31.

Villani, C. (2021). *Topics in optimal transportation, 58.* American Mathematical Soc.

Wang, Z., Qinami, K., Karakozis, I. C., Genova, K., Nair, P., Hata, K. and Russakovsky, O. (2020). Towards fairness in visual recognition: Effective strategies for bias mitigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 8919–8928.

Weerts, H., Dudík, M., Edgar, R., Jalali, A., Lutz, R. and Madaio, M. (2023). Fairlearn: Assessing and improving fairness of AI systems. *Journal of Machine Learning Research* 24: 1–8.

Zeng, X., Dobriban, E. and Cheng, G. (2022). Fair bayes-optimal classifiers under predictive parity. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22. Red Hook, NY, USA: Curran Associates Inc.