

fairadapt: Causal Reasoning for Fair Data Pre-processing

Drago Plečko
ETH Zürich

Nicolas Bennett
ETH Zürich

Nicolai Meinshausen
ETH Zürich

Abstract

Machine learning algorithms are useful for various predictions tasks, but they can also learn how to discriminate, based on gender, race or other sensitive attributes. This realization gave rise to the field of fair machine learning, which aims to recognize, quantify and ultimately mitigate such algorithmic bias. This manuscript describes the R-package **fairadapt**, which implements a causal inference pre-processing method. By making use of a causal graphical model alongside the observed data, the method can be used to address hypothetical questions of the form “What would my salary have been, had I been of a different gender/race?”. Such individual level counterfactual reasoning can help eliminate discrimination and help justify fair decisions. We also discuss appropriate relaxations which assume that certain causal pathways from the sensitive attribute to the outcome are not discriminatory.

Keywords: algorithmic fairness, causal inference, machine learning.

1. Introduction

As society transitions to an economy driven by artificial intelligence (AI), an increasing number of prediction tasks is delegated to AI tools. Sometimes these tasks are within socially sensitive domains, such as determining credit-score ratings or predicting recidivism during parole. In the process, it has been recognized that machine learning algorithms are capable of learning societal biases, which we might not want them to learn, for example with respect to race (Larson, Mattu, Kirchner, and Angwin 2016b) or gender (Lambrecht and Tucker 2019; Blau and Kahn 2003). This realization seeded an important debate in the machine learning community about fairness of algorithms and their impact on decision-making.

1.1. Definitions of fairness

The first step towards understanding algorithmic fairness is about providing a formal definition of what fairness (or discrimination) means. In light of this, existing intuitive notions have been mathematically formalized, thereby also providing fairness metrics that can be used to quantify discrimination. However, various different notions of fairness exist, and these are sometimes mutually incompatible (Corbett-Davies and Goel 2018), meaning they cannot be satisfied at the same time for a given predictor \hat{Y} . In fact, there is currently no consensus on which notion of fairness is the correct one. Among the many proposals discussed in the literature, the most commonly considered ones include:

- (1) *Demographic parity* (Darlington 1971), which requires the protected attribute A (gender, race, religion etc.) to be independent of a constructed classifier or regressor \hat{Y} , written as

$$\hat{Y} \perp\!\!\!\perp A.$$

- (2) *Equality of odds* (Hardt, Price, Srebro *et al.* 2016), which requires equal false positive and false negative rates of classifier \hat{Y} between different groups (females and males for example) written as

$$\hat{Y} \perp\!\!\!\perp A \mid Y.$$

- (3) *Calibration* (Chouldechova 2017), which requires the protected attribute to be independent of the actual outcome given the prediction

$$Y \perp\!\!\!\perp A \mid \hat{Y}.$$

Intuitively, this means that the protected attribute A does not offer additional information about the outcome Y once we know what the prediction \hat{Y} is.

1.2. Fairness tasks

Apart from choosing a notion of fairness most appropriate to the setting that is analyzed, it is instructive to distinguish between two somewhat different tasks in fairness analysis. The first, usually simpler task is that of *bias quantification*, or bias measurement. In this case, we are interested in computing metrics from our dataset, in order to verify whether a definition is satisfied or not. Concretely, suppose we have a dataset in which \hat{Y} is the predicted salary of an employee and A is sex. If we are interested in demographic parity, $\hat{Y} \perp\!\!\!\perp A$, we can compute the average difference in salary between the male/female groups, that is

$$\mathbb{E}[\hat{Y} \mid A = \text{male}] - \mathbb{E}[\hat{Y} \mid A = \text{female}].$$

Based on this quantity (which is known as total variation, or TV for short), we might decide that either there is a disparity between sexes, or perhaps not. For other definitions of fairness, different metrics would be appropriate.

Sometimes performing the first step of bias measurement is not the end goal. Suppose that we found a large salary gap between sexes in the example above, raising issues about possible discrimination. We then might be interested in correcting this bias, by computing new, more fair predictions (in which, say, the salary gap is lower). This second task is that of *bias removal*, or bias mitigation, in which we want to remove the undesired bias from our predictions.

Different software tools can be used to perform the two fairness tasks described above. In Figure 1 we provide a graphical overview of some of the available software in `python` and `R`, which are the languages most commonly used for fairness analysis. For a given task and outcome/definition of fairness, we show the existing software packages that the reader might want to use. Since calibration is often satisfied by fitting an unconstrained model, our focus is on demographic parity and equality of odds.

The software landscape for fair ML in `python` is more mature. Currently three well developed packages exist, capable of computing various fairness metrics. Also, these tools support

		R	python	
Bias Removal	Bias Detection	DP & EO	<div> <div>fairness</div> <div>fairmodels</div> </div>	<div> <div>aif360</div> <div>fairlearn</div> <div>ethicML</div> </div> <div> <div>fairness indicators</div> </div>
	EO	<div> <div>fairml</div> </div>	<div> <div>aif360</div> <div>fairlearn</div> <div>ethicML</div> </div>	
	DP	<div> <div>fairmodels</div> <div>fairml</div> <div>fairadapt</div> </div> <div> <div>fairmodels</div> <div>fairml</div> <div>fairadapt</div> </div>	<div> <div>aif360</div> <div>fairlearn</div> <div>ethicML</div> </div>	<div> <div>aif360</div> <div>fairlearn</div> <div>ethicML</div> </div>
		Classification	Regression	Classification Regression

Figure 1: Software options for fair data analysis in R and python. In the left column, nodes correspond to R-packages and in the right column to python modules. DP stands for demographic parity, EO equality of odds.

training predictors that satisfy either equality of odds or demographic parity. These repositories are **aif360** (Bellamy, Dey, Hind, Hoffman, Houde, Kannan, Lohia, Martino, Mehta, Mojsilovic, Nagar, Ramamurthy, Richards, Saha, Sattigeri, Singh, Varshney, and Zhang 2018) maintained by IBM, **fairlearn** (Bird, Dudik, Edgar, Horn, Lutz, Milan, Sameki, Wallach, and Walker 2020) maintained by Microsoft and **EthicML** (Thomas, Kehrenberg, Bartlett, and Quadrianto 2018). A further package, **fairness indicators** (Shukla, Fang, and Jindal 2019), is narrower in scope but suitable for computing fairness metrics on very large datasets.

For R, i.e., distributed via CRAN, there are fewer packages that relate to fair machine learning (see Figure 1). Available packages include **fairml** (Scutari 2021), which implements the non-convex method of Komiyama, Takeda, Honda, and Shima (2018), as well as **fairness** (Kozodoi and V. Varga 2021) and **fairmodels** (Wiśniewski and Biecek 2022), which serve as diagnostic tools for measuring algorithmic bias and provide several pre- and post-processing methods for bias mitigation. The **fairadapt** package described in this manuscript is a bias removal method which is able to interpolate between demographic parity and calibration notions, and is applicable to both regression and classification settings. In particular, **fairadapt** is the only software in Figure 1 which is *causally aware*. This means that the bias removal performed by **fairadapt** can be explained by and related to the causal mechanisms that generated the discrimination in the first place.

1.3. A causal approach

The discussion on algorithmic fairness is, however, not restricted to the machine learning domain. There are many legal and philosophical aspects that are paramount. For example, the legal distinction between the disparate impact and disparate treatment doctrines (McGinley 2011; Barocas and Selbst 2016) is important for assessing fairness from a legal standpoint.

This in turn emphasizes the importance of the interpretation behind the decision-making process, which is often not the case with black-box machine learning algorithms. For this reason, research in fairness through a causal inference lens gained attention.

A possible approach to fairness is the use of counterfactual reasoning (Galles and Pearl 1998), which allows for arguing what might have happened under different circumstances that never actually materialized, thereby providing a tool for understanding and quantifying discrimination. For example, one might ask how a change in sex would affect the probability of a specific candidate being accepted for a given job opening. This approach has motivated another notion of fairness, termed *counterfactual fairness* (Kusner, Loftus, Russell, and Silva 2017), which states that the decision made, should remain fixed, even if, hypothetically, the protected attribute such as race or gender were to be changed (this can be written succinctly as $\hat{Y}_i(a) = \hat{Y}_i(a')$ in the potential outcomes notation). Causal inference can also be used for decomposing the total variation measure into its direct, mediated, and confounded contributions (Zhang and Bareinboim 2018), yielding further insights into demographic parity as a criterion. Furthermore, by introducing the notion of so-called resolving variables, Kilbertus, Carulla, Parascandolo, Hardt, Janzing, and Schölkopf (2017) described relaxations of demographic parity, which can possibly be a prohibitively strong notion.

The following sections describe an implementation of the fair data adaptation method outlined in Plecko and Meinshausen (2020), which combines the notions of counterfactual fairness and resolving variables, and explicitly computes counterfactual instances for individuals. The implementation is available as the R-package **fairadapt** from CRAN.

1.4. Novelty in the package

A first version of **fairadapt** was published with the original manuscript (Plecko and Meinshausen 2020). The software has since been developed further and novelty in the package presented in this manuscript includes the following:

- The methodology has been extended from the Markovian to the Semi-Markovian case (allowing noise variables to be correlated), which generalizes the scope of applications.
- Backdoor paths into the protected attribute A are now allowed, meaning that the attribute A does not need to be a root node of the causal graph.
- The user is provided with functionality for uncertainty quantification of the estimates.
- Introduction of S3 classes **fairadapt** and **fairadaptBoot**, alongside associated methods, provides a more formalized implementation.
- More flexibility is allowed in the quantile learning step, including different algorithms for quantile regression (linear, forest based, neural networks). Additionally, there is also a possibility of specifying a custom quantile learning function (utilizing S3 dispatch).
- The user is provided with functionality for assessing the quality of the quantile regression fit, which allows for tuning the hyperparameters of the more flexible algorithms.

The rest of the manuscript is organized as follows. In Section 2 we describe the methodology behind **fairadapt**, together with reviewing some important concepts of causal inference. In Section 3 we discuss implementation details and provide some general user guidance, followed by Section 4 in which we discuss how to perform uncertainty quantification with **fairadapt**. Section 5 illustrates the use of **fairadapt** through a large, real-world dataset and a hypothetical

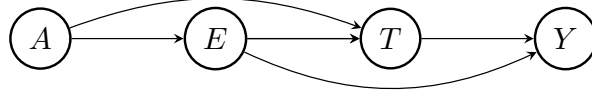


Figure 2: A visual representation of the university admission example. A denotes gender, E previous educational achievement, and T an admissions test score. The outcome Y represents the final score used for the admission decision. The protected attribute A has a discriminatory causal effect on variables E , T , and Y , which we wish to remove.

fairness application. Finally, in Section 6 we elaborate on some extensions, such as Semi-Markovian models and resolving variables.

2. Methodology

In this section, the intuition behind **fairadapt** is described using an example. This is followed by a more rigorous mathematical formulation, based on Markovian structural causal models (SCMs). Some relevant extensions, such as the Semi-Markovian case and the introduction of *resolving variables* are discussed in Section 6.

2.1. University admission example

Consider for example a dataset about students applying for university admission. Let variable A , gender, be the protected attribute ($A = a$ corresponding to females and $A = a'$ to males). Let E denote educational achievement (measured for example by grades achieved in school) and T the result of an admissions test for further education. Finally, let Y be the outcome of interest (final score) upon which admission to further education is decided. A visual representation of how the variables affect each other is given in Figure 2.

Attribute A , gender, has a causal effect on variables E , T , and Y . We consider this effect discriminatory and wish to eliminate it, in the following way. For each individual with observed values (a, e, t, y) we want to find a mapping

$$(a, e, t, y) \longrightarrow (a^{(fp)}, e^{(fp)}, t^{(fp)}, y^{(fp)}),$$

which represents the value the person would have obtained in an alternative world where everyone was female. To construct such a mapping, we adapt (transform) the variables in the dataset in order. Explicitly, to a male person with education value e , we assign the transformed value $e^{(fp)}$, chosen such that

$$\mathbb{P}(E \geq e \mid A = a') = \mathbb{P}(E \geq e^{(fp)} \mid A = a).$$

The key idea is that the *relative educational achievement within the subgroup* remains preserved if the protected attribute gender is changed. If, for example, a male person has a higher educational achievement value than 70% of males in the dataset, we assume that he would also be better than 70% of females, had he been female¹. After computing transformed

¹This assumption of course is not empirically testable, as it is impossible to observe both a female and a male version of the same individual.

educational achievement values corresponding to the *female* world ($E^{(fp)}$), the transformed test score values $T^{(fp)}$ can be calculated in a similar fashion, but conditional on educational achievement. That is, a male with values $(E, T) = (e, t)$ is assigned a test score $t^{(fp)}$ such that

$$\mathbb{P}(T \geq t \mid E = e, A = a') = \mathbb{P}(T \geq t^{(fp)} \mid E = e^{(fp)}, A = a),$$

where the value $e^{(fp)}$ was obtained in the previous step. This step can be visualized as shown in Figure ??.

In the final step, the outcome variable Y needs to be adjusted. The adaptation is based on the same principle as above, using transformed values of both education and the test score. The transformed value $y^{(fp)}$ of $Y = y$ is chosen to satisfy

$$\mathbb{P}(Y \geq y \mid E = e, T = t, A = a') = \mathbb{P}(Y \geq y^{(fp)} \mid E = e^{(fp)}, T = t^{(fp)}, A = a).$$

The form of counterfactual correction described above is known as *recursive substitution* (Pearl 2009, Chapter 7). We formalize this approach in the following sections. The reader who is satisfied with the intuitive notion provided by the above example is encouraged to go straight to Section 3.

2.2. Structural causal models

In order to describe the causal mechanisms of a system, a *structural causal model* (SCM) can be hypothesized, which fully encodes the assumed data-generating process. An SCM is represented by a 4-tuple $\langle V, U, \mathcal{F}, \mathbb{P}(u) \rangle$, where

- $V = \{V_1, \dots, V_n\}$ is the set of observed (endogenous) variables.
- $U = \{U_1, \dots, U_n\}$ are latent (exogenous) variables.
- $\mathcal{F} = \{f_1, \dots, f_n\}$ is the set of functions determining V , $v_i \leftarrow f_i(\text{pa}(v_i), u_i)$, where $\text{pa}(V_i) \subset V, U_i \subset U$ are the functional arguments of f_i and $\text{pa}(V_i)$ denotes the parent vertices of V_i .
- $\mathbb{P}(u)$ is a distribution over the exogenous variables U .

Any particular SCM is accompanied by a graphical model \mathcal{G} (a directed acyclic graph, DAG). The set of variables that are inputs of the mechanism f_{V_i} are the parents of V_i denoted by $\text{pa}(V_i)$. Therefore, the graph encodes how variables affect one another. Furthermore, we also write $\text{ch}(V_i)$, $\text{de}(V_i)$, $\text{an}(V_i)$ for the children, descendants, and ancestors of V_i in the graph \mathcal{G} . We assume throughout, without loss of generality, that

- $f_i(\text{pa}(v_i), u_i)$ is increasing in u_i for every fixed $\text{pa}(v_i)$.
- Exogenous variables U_i are uniformly distributed on $[0, 1]$.

Equipped with the notion of an SCM and the assumptions (i)-(ii), we can describe the adaptation procedure in the Markovian case, in which all exogenous variables U_i are mutually independent (for the Semi-Markovian case, where variables U_i are allowed to share information, see Section 6).

2.3. Markovian SCM formulation

Let Y be the outcome of interest, taking values in \mathbb{R} . Let A be the binary protected attribute taking two values a, a' . Denote by X the remaining covariates, and let $V = (A, X, Y)$ denote the observed variables. Our goal is to describe a pre-processing method which transforms the observed variables V into their fair version $V^{(fp)}$. This is achieved by computing the counterfactual values $V(A = a)$, which would have been observed if the protected attribute was fixed to a baseline value $A = a$ for the entire sample.

More formally, going back to the *university admission* example above, we want to align the distributions

$$V_i \mid \text{pa}(V_i), A = a \text{ and } V_i \mid \text{pa}(V_i), A = a',$$

meaning that the distribution of V_i conditional on $\text{pa}(V_i)$ should be indistinguishable between female and male groups (and this should hold for every variable V_i). Since each function f_i of the original SCM is reparametrized so that $f_i(\text{pa}(v_i), u_i)$ is increasing in u_i for every fixed $\text{pa}(v_i)$, and since variables U_i are uniformly distributed on $[0, 1]$, the U_i values can be interpreted as the latent *quantiles* associated with V_i . These latent quantiles are assumed to be preserved when performing the adaptation procedure.

The fair data adaption algorithms starts by fixing $A = a$ for all individuals. After this, the algorithm iterates over descendants of the protected attribute A , in any valid topological order (this topological order is inferred from the causal graph \mathcal{G} , which is also an input of the algorithm). For each V_i , the assignment function f_i and the corresponding quantiles U_i are inferred. Finally, transformed values $V_i^{(fp)}$ are obtained by evaluating f_i , using quantiles U_i and the transformed parents $\text{pa}(V_i)^{(fp)}$ (see Algorithm 1).

Algorithm 1: Fair Data Adaptation

Input: observed variables V , causal graph \mathcal{G}
 set $A \leftarrow a$ for everyone
for $V_i \in \text{de}(A)$ *in topological order* **do**
 learn function $V_i \leftarrow f_i(\text{pa}(V_i), U_i)$
 infer quantiles U_i associated with the variable V_i
 transform values as $V_i^{(fp)} \leftarrow f_i(\text{pa}(V_i)^{(fp)}, U_i)$
end
return $V^{(fp)}$

The assignment functions f_i of the SCM are always assumed to be unknown, but are inferred non-parametrically at each step. Algorithm 1 obtains the counterfactual values $V(A = a)$ under the $do(A = a)$ intervention for each individual, while keeping the latent quantiles U fixed. In the case of continuous variables, the latent quantiles U can be determined exactly, while for the discrete case, the situation is more subtle. A detailed discussion can be found in Plečko and Meinshausen (2020, Section 5).

3. Implementation

The main function for data adaption in the **fairadapt** package is `fairadapt()`. This function returns an S3 classed object of class **fairadapt**. The **fairadapt** class has associated implementations of the base R S3 generics `print()`, `summary()`, `plot()` and `predict()`. Furthermore, methods are available for the `autoplot()` generic exported from **ggplot2** (Wickham 2016), as well as **fairadapt**-specific implementations of S3 generics `visualizeGraph()`, `adaptedData()`, `quantFit()`, and `fairTwins()`.

The following sections describe the intended use of `fairadapt()`, together with the associated methods and their relations. The most important arguments of `fairadapt()` include:

- **formula**: Argument of type **formula**, specifying the dependent and explanatory variables.
- **adj.mat**: Argument of type **matrix**, encoding the adjacency matrix.
- **train.data** and **test.data**: Both of type **data.frame**, representing the respective datasets.
- **prot.attr**: Scalar-valued argument of type **character** identifying the protected attribute. Has to correspond to a column name in the **train.data** argument.

It is worth clarifying the possible data types that can be used in the **train.data** argument. We note the following:

- (i) *Attribute A*: the protected attribute is assumed to be binary. The **prot.attr** column in **train.data** can be of any data type coercible to a **factor**, but can only take two distinct values. Otherwise an error is thrown.
- (ii) *Outcome Y*: the dependent variable specified on the left hand side of the **formula** argument can be either a **numeric**, **logical**, **integer** (treated the same as an ordered factor), **factor**, or a **character**².
- (iii) *Remaining covariates X*: all other variables do not have limitations. Unordered **factor** or **character** inputs can also be used as covariates *X*.

As an example, we perform fair data adaption on the university admission dataset described in Section 2. We load the **uni_admission** dataset provided by **fairadapt** (inspired by the Berkeley admissions dataset (Bickel, Hammel, and O’Connell 1975)), consisting of synthetic university admission data of 1000 students. We subset the data, using the first **n_samp** rows as training data (**uni_trn**) and the following **n_samp** rows as testing data (**uni_tst**). Furthermore, we construct an adjacency matrix **uni_adj** with edges **gender** → **edu**, **gender** → **test**, **edu** → **test**, **edu** → **score**, and **test** → **score**, corresponding to the causal graph from Figure 2. We set **gender** as the protected attribute.

```
R> n_samp <- 500
R>
R> uni_dat <- data("uni_admission", package = "fairadapt")
R> uni_dat <- uni_admission[seq_len(2 * n_samp), ]
R>
R> head(uni_dat)
```

²Care needs to be taken when supplying unordered **factor** or **character** inputs since the adaptation procedure depends on the order of the levels of the outcome variable. For such inputs, the order of the levels will be chosen automatically, so using an ordered **factor** for the outcome variable is the recommended option.

	gender	edu	test	score
1	1	1.3499572	1.617739679	1.9501728
2	0	-1.9779234	-3.121796235	-2.3502495
3	1	0.6263626	0.530034686	0.6285619
4	1	0.8142112	0.004573003	0.7064857
5	1	1.8415242	1.193677123	0.3678313
6	1	-0.3252752	-2.004123561	-1.5993848

```
R> uni_trn <- head(uni_dat, n = n_samp)
R> uni_tst <- tail(uni_dat, n = n_samp)
R>
R> uni_dim <- c("gender", "edu", "test", "score")
R> uni_adj <- matrix(c(
+           0,      1,      1,      0,
+           0,      0,      1,      1,
+           0,      0,      0,      1,
+           0,      0,      0,      0),
+           ncol = length(uni_dim),
+           dimnames = rep(list(uni_dim), 2),
+           byrow = TRUE)
R>
R> set.seed(2022)
R> basic <- fairadapt(score ~ ., train.data = uni_trn,
+           test.data = uni_tst, adj.mat = uni_adj,
+           prot.attr = "gender")
R>
R> basic
```

Call:

```
fairadapt(formula = score ~ ., prot.attr = "gender", adj.mat = uni_adj,
          train.data = uni_trn, test.data = uni_tst)
```

Adapting variables:

score, edu, test

Based on protected attribute gender

AND

Based on causal graph:

	score	gender	edu	test
score	0	0	0	0
gender	0	0	1	1
edu	1	0	0	1
test	1	0	0	0

The implicitly called `print()` method in the previous code block displays some information about how `fairadapt()` was called. This information includes the variables that were adapted, the protected attribute, and the causal graph used for the adaptation (printed as an adjacency matrix). By additionally calling the `summary()` function, we can inspect the number of training and test samples, and the total variation before and after adaptation, written in our notation as

$$\mathbb{E}[Y \mid A = a] - \mathbb{E}[Y \mid A = a'] \text{ and } \mathbb{E}[Y^{(fp)} \mid A = a] - \mathbb{E}[Y^{(fp)} \mid A = a'],$$

respectively, shown below:

```
R> summary(basic)
```

Call:

```
fairadapt(formula = score ~ ., prot.attr = "gender", adj.mat = uni_adj,
  train.data = uni_trn, test.data = uni_tst)
```

```
Protected attribute:          gender
Protected attribute levels:    0, 1
Adapted variables:            edu, test, score

Number of training samples:    500
Number of test samples:        500
Quantile method:              rangerQuants

Total variation (before adaptation): -0.7045
Total variation (after adaptation):  -0.066
```

The adapted train and test data can be obtained using the `adaptedData()` function and passing the argument `train = TRUE` for the training data, and `train = FALSE` for the test data:

```
R> head(adaptedData(basic, train = FALSE))
```

```
      gender      edu      test
501      0 -2.2844949 -1.3101484
502      0 -0.2884019 -1.0954813
503      0  0.4161544  0.6885127
504      0 -0.6166185 -1.1286510
505      0 -0.1607580 -0.7841601
506      0 -0.2337741 -1.1327089
```

3.1. Specifying the graphical model

The algorithm used for fair data adaption in `fairadapt()` is based on graphical causal model \mathcal{G} (see Algorithm 1). To specify the causal graph, we pass the corresponding adjacency matrix

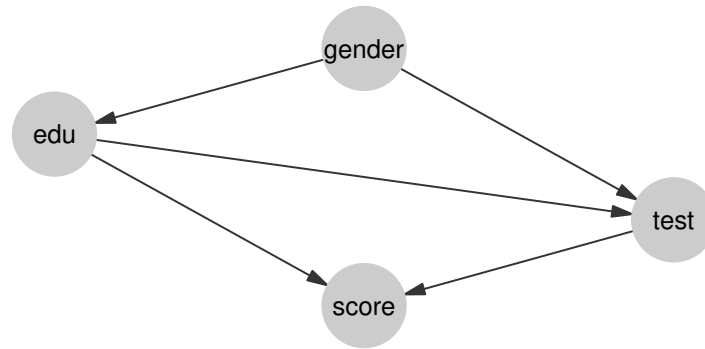


Figure 3: The underlying graphical model corresponding to the university admission example (also shown in Figure 2).

as the `adj.mat` argument. The convenience function `graphModel()` turns a graph specified as an adjacency matrix into an annotated graph using the **igraph** package (Csardi and Nepusz 2006). Alternatively, by calling the S3 generic `visualizeGraph()` on a **fairadapt** object, the user can also inspect the graphical model that was used for the data adaptation.

```
R> uni_graph <- graphModel(uni_adj)
```

Warning: Using the `size` aesthetic in this geom was deprecated in ggplot2 3.4.0.

i Please use `linewidth` in the `default_aes` field and elsewhere instead.

This warning is displayed once every 8 hours.

Call `lifecycle::last_lifecycle_warnings()` to see where this warning was generated.

A visualization of the **igraph** object returned by `graphModel()` is shown in Figure 3. The graph is the same as that in Figure 2. However, specifying the causal graph is not the only option to perform the data adaptation. A possible alternative is to specify a valid topological ordering over the observable variables V , and specify it as a **character** vector, using the `top.ord` argument.

3.2. Quantile learning step

The most common and recommended usage of `fairadapt()` follows a typical machine learning framework. We start by calling the `fairadapt()` function that performs the quantile learning step and counterfactual correction, based on the `train.data` argument. Following this, we can call the `predict()` function on the returned **fairadapt** S3 object, in order to perform data adaption on new test data. In such a workflow, the adaptation of the training and testing data is done separately. In specific situations, it might be desirable to input both `train.data` and `test.data` arguments directly to `fairadapt()`, which then transforms both the training and testing data jointly. This one-step procedure might be considered when the proportion of test samples compared to train samples is large, and when the `train.data` has a relatively small sample size. The benefit of this approach is that, even though the outcome Y is not available, other attributes X of `test.data` can be used for quantile learning step.

	Random forest	Neural network	Linear regression
R-package	ranger	qrnn	quantreg
<code>quant.method</code>	<code>rangerQuants</code>	<code>mcqrnnQuants</code>	<code>linearQuants</code>
Complexity	$O(np \log n)$	$O(np n_{\text{epochs}})$ 1 hidden layer	$O(p^2 n)$ "br" method of
Default parameters	$ntrees = 500$ $mtry = \sqrt{p}$	fully connected feed-forward network	Barrodale and Roberts used for fitting
$T_{\text{uni}}(200)$	0.4	37.7	0.2
$T_{\text{uni}}(500)$	0.9	96.8	0.4

Table 1: Summary table of different quantile regression methods. n is the number of samples, p number of covariates, n_{epochs} number of training epochs for the neural network. $T_{\text{uni}}(n)$ denotes the runtime of different methods on the university admission dataset, with n training and n testing samples. The runtimes were obtained on a system with Intel Core i7-8750H CPU @ 2.2GHz running MacOS Big Sur 11.6. The version of R was 4.2.0 "Vigorous Calisthenics" with **quantreg** version 5.93, **ranger** version 0.13.1, and **mcqrnn** version 2.0.5.

The data frames passed as `train.data` and `test.data` are required to have column names which also appear in the row and column names of the adjacency matrix. The protected attribute A , passed as scalar-valued character vector `prot.attr`, should also appear in the column names of `train.data` and `test.data`. The `test.data` argument defaults to NULL, with the intention that `test.data` is specified as an input to the `predict()` function at a later stage.

Quantile methods

The quantile learning step of Algorithm 1 can in principle be carried out by several methods, three of which are implemented in **fairadapt**:

- Quantile Regression Forests (Meinshausen 2006; Wright and Ziegler 2017).
- Non-Crossing Quantile Neural Networks (Cannon 2018, 2015).
- Linear Quantile Regression (Koenker and Hallock 2001; Koenker, Portnoy, Ng, Zeileis, Grosjean, and Ripley 2018).

Using linear quantile regression is the most efficient option in terms of runtime, while for non-parametric models and mixed data, the random forest approach is well-suited, at the expense of a slight increase in runtime. The neural network approach is substantially slower when compared to linear and random forest estimators and consequently does not scale well to large sample sizes. As default, the random forest based approach is used, due to its non-parametric nature and computational speed. However, for smaller sample sizes, the neural network approach can also demonstrate competitive performance. In Table 1 we provide a quick summary outlining some differences between the three natively supported methods, and also report runtimes on the `uni_admission` dataset with different sample sizes.

Influencing the fit

The quantile methods shown in Table 1 make calls to specific functions that perform quantile regression. These functions take varying arguments and for that reason, `fairadapt()` forwards arguments passed as `...` to the function specified as `quant.method`.

Computational speed An important consideration in choosing values for optional arguments of specific quantile regression functions is computational speed. For example, the function `rangerQuants()` internally calls the `ranger()` function (from **ranger**) and with respect to computational speed, an important argument is `num.trees`, the number of trees used when building the quantile regression forest. Clearly, choosing a smaller number of trees will be faster, but at the same time will result in a fit with larger variance.

Similarly, `mcqrnnQuants()` internally calls `mcqrnn.fit()` (from **qrnn**), which has a number of arguments that can be used for adapting the underlying neural network. In terms of computational speed, the most important arguments are `n.trials` (number of repeated initializations used to avoid local minima) and `iter.max` (maximum number of iterations of the optimization). Choosing smaller values will reduce the runtime. Lastly, function `linearQuants()` internally calls `rq()` (from **quantreg**). This function is less flexible, since the model is linear. However, its `method` argument can be used when the number of samples becomes large (using `method` equal to `"fn"` or `"pfn"` utilizes the Frisch-Newton interior point method, which may be preferable for large samples).

Fit quality Both `rangerQuants()` and `mcqrnnQuants()` expose flexible machine learning tools with several parameters that impact fit quality. In order to optimize the fitting procedure by tuning these parameters, we need a way of assessing the quality of our fit. In the context of quantile regression we can estimate the expected τ -quantile loss function,

$$\mathbb{E}[\rho_\tau(V_i, \mu_\tau(\text{pa}(V_i)))], \quad (1)$$

where $\mu_\tau(\text{pa}(V_i))$ is the function predicting the τ -quantile of variable V_i using the parents $\text{pa}(V_i)$ and ρ_τ is the asymmetric L1 loss function whose minimizer is the τ -quantile. The function ρ_τ is given by

$$\rho_\tau(x, y) = \begin{cases} \tau(x - y), & \text{for } x \geq y \\ (1 - \tau)(y - x), & \text{for } x < y. \end{cases}$$

A smaller empirical loss based on Equation 1 corresponds to better fit quality. For hyperparameter tuning we can perform cross-validation (fitting the quantile regression on separate folds), which is directly available within the `fairadapt()` function. The argument `eval.qfit` has a default value `NULL`, but if this argument is given a positive integer value, then it is used as the number of folds for performing cross-validation.

We compute the average empirical loss $\widehat{\mathbb{E}}[\rho_\tau(V_i, \mu_\tau(\text{pa}(V_i)))]$ for each variable V_i and $\tau = 0.25, 0.5, 0.75$ (corresponding to 25%, 50% and 75% quantiles). The average of these three values is reported at the end, and can be extracted from the resulting `fairadapt` object using the `quantFit()` method:

```
R> set.seed(22)
R> fit_qual <- fairadapt(score ~ ., train.data = uni_trn,
```

```

+               adj.mat = uni_adj, prot.attr = "gender",
+               eval.qfit = 3L)
R>
R> quantFit(fit_qual)

```

```

      edu      test      score
0.3405883 0.2803902 0.3457824

```

The function returns the quality of the quantile fit for each variable. A very reasonable objective to minimize is the average of these values, by iterating over a grid of possible values of the tuning parameters. The interesting parameters to optimize for the two methods include:

- (i) for `ranger()`: parameters `mtry` (number of candidate variables considered for splitting in each step), `min.node.size` (size of leaf nodes below which splitting is stopped) and `max.depth` (maximum depth of each tree),
- (ii) for `mcqrnn()`: parameters `n.hidden`, `n.hidden2` (number of nodes in the first and second hidden layers), `Th` (activation function), `method` (optimizer to be used), and a range of other parameters that are fed to the chosen optimizer via ellipsis.

The quantile methods included in **fairadapt** have reasonable default values, that serve as a good starting point. Optimizing the quantile fit should therefore be of interest mostly to advanced users (and hence we do not perform parameter tuning here explicitly).

Extending to custom methods

The above-mentioned set of methods is not exhaustive. Further options are conceivable and therefore **fairadapt** provides an extension mechanism to allow for custom quantile method specified by the user. The `fairadapt()` argument `quant.method` expects a function to be passed, a call to which will be constructed with three unnamed arguments:

1. A `data.frame` containing data to be used for quantile regression. This will either be the `data.frame` passed as `train.data`, or if `test.data` was specified, a concatenated, row-bound version of `train.data` and `test.data`.
2. A logical flag, indicating whether the protected attribute is the root node of the causal graph. If the attribute A is a root node, we know that

$$\mathbb{P}(X \mid \text{do}(A = a)) = \mathbb{P}(X \mid A = a).$$

Therefore, the interventional and conditional distributions are in this case the same, and we can leverage this knowledge in the quantile learning procedure, by splitting the data into $A = 0$ and $A = 1$ groups.

3. A logical vector of length `nrow(data)`, indicating which rows in the `data.frame` passed as `data` correspond to samples with baseline values of the protected attribute.

Arguments passed as `...` to `fairadapt()` will be forwarded to the function specified as `quant.method` and passed after the first three fixed arguments listed above. The return value of the function passed as `quant.method` is expected to be an S3-classed object. This

object should represent the conditional distribution $V_i \mid \text{pa}(V_i)$ (see function `rangerQuants()` for an example). Additionally, the object should have an implementation of the S3 generic function `computeQuants()` available. For each row $(v_i, \text{pa}(v_i))$ of the `data` argument, the `computeQuants()` function uses the S3 object to perform the following steps:

- (i) Infer the quantile of $v_i \mid \text{pa}(v_i)$.
- (ii) Compute the counterfactual value $v_i^{(fp)}$ under the change of protected attribute, using the counterfactual values of parents $\text{pa}(v_i^{(fp)})$ computed in previous steps (values $\text{pa}(v_i^{(fp)})$ are contained in the `newdata` argument).

For an example, see the `computeQuants.ranger()` method for a `ranger` object, which can be invoked by the `computeQuants()` generic.

3.3. Fair-twin inspection

We now turn to a useful property of **fairadapt**, which allows the user to explore counterfactual instances for different individuals in the dataset. The university admission example presented in Section 2 demonstrates how to compute counterfactual values for an individual while preserving their relative educational achievement. Setting candidate gender as the protected attribute and gender level *female* as baseline value, for a *male* student with values (a, e, t, y) , his *fair-twin* values $(a^{(fp)}, e^{(fp)}, t^{(fp)}, y^{(fp)})$, i.e., the values the student would have obtained, had he been *female*, are computed. These values can be retrieved from a **fairadapt** object by calling the S3-generic function `fairTwins()` as:

```
R> ft_basic <- fairTwins(basic, train.id = seq_len(n_samp))
R> head(ft_basic, n = 3)
```

	gender	score	score_adapted	edu	edu_adapted	test
1	1	1.9501728	0.4274580	1.3499572	0.8144842	1.6177397
2	0	-2.3502495	-2.3502495	-1.9779234	-1.9779234	-3.1217962
3	1	0.6285619	0.1898589	0.6263626	0.1383595	0.5300347
	test_adapted					
1	0.86888027					
2	-3.12179624					
3	-0.03061109					

In this example, we compute the values in a *female* world. Therefore, for *female* applicants, the values remain fixed, while for *male* applicants the values are adapted, as can be seen from the output. Having access to explicit counterfactual instances as above may help justify fair decisions in practice or help guide the choice of the assumed causal model and resolving variables (see Section 6 for resolving variables).

4. Uncertainty quantification

The user might naturally be interested in uncertainty quantification of the procedure performed in `fairadapt()`. In order to explain how this can be achieved, we give a visualization of the typical workflow when using `fairadapt()` (see Figure 4).

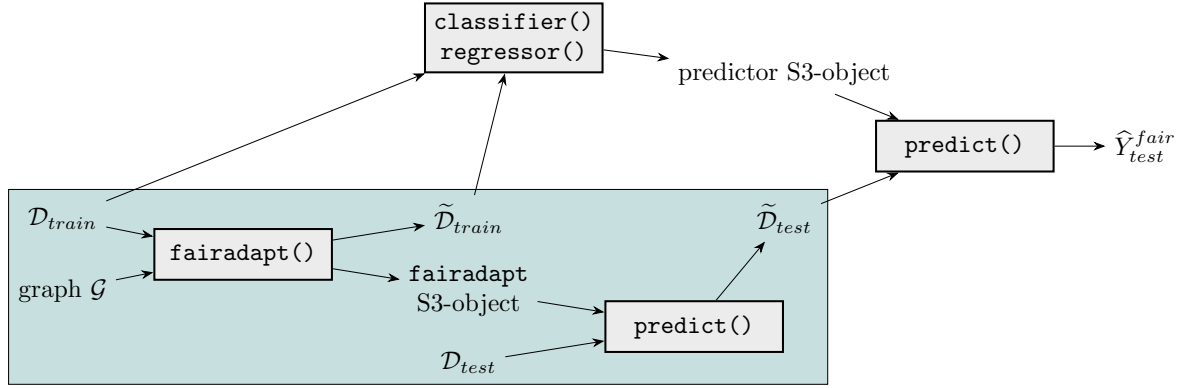


Figure 4: The typical workflow when using **fairadapt**. The shaded region represents the fair pre-processing which happens within the **fairadapt** package. Often, this is followed by applying a regressor or a classifier to the transformed data, in order to obtain fair predictions. The latter part is up to the user and not included in **fairadapt**.

Such a workflow can be described as follows. We start from the training data \mathcal{D}_{train} and the causal graph \mathcal{G} . These two arguments are used as inputs of the **fairadapt()** function, which returns a **fairadapt** S3 object (which also contains the transformed test data $\tilde{\mathcal{D}}_{train}$). The **fairadapt** object, together with the test data \mathcal{D}_{test} , is then used as a input to the **predict()** function. The **predict()** function returns the transformed test data $\tilde{\mathcal{D}}_{test}$. Often, the end goal is to obtain fair predictions on some new test data \mathcal{D}_{test} . To do so, we need to train a classifier/regressor. Either the training data \mathcal{D}_{train} or its transformed counterpart $\tilde{\mathcal{D}}_{train}$ can be used for building a predictor. The predictor then needs to be applied to the transformed train data $\tilde{\mathcal{D}}_{test}$. Building on the graphical visualization in Figure 4, which serves as a mental map of our workflow, we can now explain the distinct sources of uncertainty that can be considered:

- **Finite sample uncertainty:** The first, commonly encountered source of uncertainty is the one induced by finite sample size. The training data \mathcal{D}_{train} has a finite size, and for this reason inferences made using this data are imperfect. We wish to quantify the uncertainty in the predictions \hat{Y}_{test}^{fair} introduced by the finite sample size of \mathcal{D}_{train} . As Figure 4 shows, the training data \mathcal{D}_{train} affects the resulting fair predictions in two ways. Firstly, it affects the value of the transformed test data $\tilde{\mathcal{D}}_{test}$ (mediated by the **fairadapt** S3 object). Secondly, it affects the predictions \hat{Y}_{test}^{fair} through the predictor (since it is the input to the regressor/classifier). These finite sample uncertainties can be analyzed using *bootstrap* (Efron and Tibshirani 1994). This means that we repeat the procedure in Figure 4 many times, each time taking a different bootstrap sample of the training data. Below we will show how this can be done with the **fairadaptBoot()** function.
- **Inherent uncertainty in the quantiles:** A second source of uncertainty arises from the uncertainty in quantile estimation and is specific to **fairadapt**. As described in Section 2 (see also Figure ??), the **fairadapt()** procedure aims to preserve the relative quantile of the variable, when computing the $do(A = a)$ intervention. However, when we are working with variables that are not continuous, defining a quantile becomes more

difficult³. Therefore, in presence of discrete variables, due to the imperfect estimation of quantiles, the `fairadapt()` procedure has some inherent randomness. This randomness would still persist even if we had infinite training samples in \mathcal{D}_{train} . Importantly, to achieve fair predictions, taking an expectation over this randomness is not feasible. For a detailed discussion of why this is the case, refer to (Plecko and Meinshausen 2020, Section 5).

For quantifying uncertainty, we use the `fairadaptBoot()` function, where the most important arguments are:

- `formula`, `prot.attr`, `adj.mat`, and `train.data` arguments are the same as for the `fairadapt()` function (see Section 3).
- `test.data`, a `data.frame` containing the test data, defaults to `NULL`. Whenever the test data equals `NULL`, then `keep.object` must be `TRUE`.
- `keep.object`, a logical scalar, indicating whether all the `fairadapt` S3 objects built in bootstrap repetitions should be kept in working memory. Default value is `FALSE`.
- `rand.mode`, a character scalar, taking values `"finsamp"`, `"quant"`, or `"both"`, corresponding to considering finite sample uncertainty, quantile uncertainty, or both.

The function `fairadaptBoot()` returns an S3 object of class `fairadaptBoot`. Calling this function can be computationally expensive, both in terms of runtime and memory. Keeping the default value of `FALSE` for the `keep.object` argument reduces the memory consumption substantially, with the drawback that `test.data` has to be provided directly to `fairadaptBoot()`, and that the resulting `fairadaptBoot` object cannot be reused for making predictions at a later stage. Passing `TRUE` to `save.object`, on the other hand, might consume more memory, but then the object can be reused for transforming new test data. This can be done by using the `predict()` function for the `fairadaptBoot` S3 object, to which a `newdata` argument is available.

For illustration purposes, we now compute bootstrap repetitions for finite sample uncertainty and quantile uncertainty on the COMPAS dataset (Larson, Mattu, Kirchner, and Angwin 2016a). We begin by loading the COMPAS dataset and constructing its causal graph:

```
R> cmp_dat <- data("compas", package = "fairadapt")
R> cmp_dat <- get(cmp_dat)
R>
R> cmp_mat <- matrix(0, nrow = ncol(cmp_dat), ncol = ncol(cmp_dat),
+                   dimnames = list(names(cmp_dat), names(cmp_dat)))
R>
R> cmp_mat[c("race", "sex", "age"),
+          c("juv_fel_count", "juv_misd_count",
+            "juv_other_count", "priors_count",
+            "c_charge_degree", "two_year_recid")] <- 1
R> cmp_mat[c("juv_fel_count", "juv_misd_count", "juv_other_count"),
+          c("priors_count", "c_charge_degree", "two_year_recid")] <- 1
```

³For example in the case where we have a binary $X \in \{0, 1\}$, it is impossible to define what a 70% quantile is, as opposed to the continuous case (of a Gaussian variable X for example), where no such challenge exists.

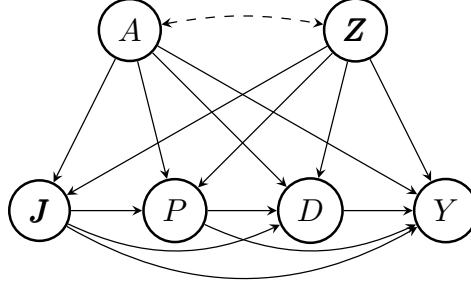


Figure 5: The causal graph for the COMPAS dataset. Z are demographic features, A is race, J juvenile offense counts, P priors count, D the degree of charge, and Y two year recidivism.

```
R> cmp_mat["priors_count", c("c_charge_degree", "two_year_recid")] <- 1
R> cmp_mat["c_charge_degree", "two_year_recid"] <- 1
R>
R> head(cmp_dat)
```

	sex	age	race	juv_fel_count	juv_misd_count	juv_other_count
1	Male	69	Non-White	0	0	0
2	Male	34	Non-White	0	0	0
3	Male	24	Non-White	0	0	1
4	Male	23	Non-White	0	1	0
5	Male	43	Non-White	0	0	0
6	Male	44	Non-White	0	0	0

	priors_count	c_charge_degree	two_year_recid
1	0	F	0
2	0	F	1
3	4	F	1
4	1	F	0
5	2	F	0
6	0	M	0

The COMPAS dataset contains information on 7214 individuals from Broward County, Florida, who were released on parole. The variables include race, sex, age, juvenile offense counts, priors count and the degree of criminal charge. The outcome of interest is recidivism within two years and the protected attribute is race (taking values Non-White and White). A possible causal graph for the COMPAS dataset is given in Figure 5.

After loading the dataset, we run the `fairadaptBoot()` function twice, with two different values of the `rand.mode` argument. First, we consider only the finite sample uncertainty.

```
R> cmp_trn <- tail(cmp_dat, n = 100L)
R> cmp_tst <- head(cmp_dat, n = 100L)
R>
R> n_itr <- 3L
R>
R> set.seed(2022)
```

```
R> fa_boot_fin <- fairadaptBoot(two_year_recid ~ ., "race", cmp_mat,
+                               cmp_trn, cmp_tst, rand.mode = "finsamp",
+                               n.boot = n_itr)
```

Then, we re-run the bootstrap procedure, but by considering only the inherent quantile uncertainty, by setting the `rand.mode` argument to `"quant"`.

```
R> set.seed(2022)
R> fa_boot_quant <- fairadaptBoot(two_year_recid ~ ., "race", cmp_mat,
+                                cmp_trn, cmp_tst, rand.mode = "quant",
+                                n.boot = n_itr)
```

The returned objects are of class `fairadaptBoot`. The object stores different replicates of the adapted test data (`n.boot` copies, the number of bootstrap repetitions) and some metadata. To obtain predictions, we train a random forest classifier on different bootstrap samples of `train.data`, and apply it to the transformed data bootstrap replicates. In doing so, we make use of the `boot.ind` list, contained in the `fairadaptBoot` object, representing row indices of all bootstrap repetitions.

```
R> fit_rf <- function(x) {
+   ranger(factor(two_year_recid) ~ ., cmp_trn[x, ], probability = TRUE)
+ }
R>
R> extract_pred <- function(x) x$predictions[, 2L]
R>
R> set.seed(2022)
R> cmp_rf <- lapply(fa_boot_fin$boot.ind, fit_rf)
R>
R> pred_fin <- Map(predict, cmp_rf, adaptedData(fa_boot_fin, train = FALSE))
R> pred_fin <- do.call(cbind, lapply(pred_fin, extract_pred))
R>
R> pred_quant <- Map(predict, cmp_rf, adaptedData(fa_boot_quant, train = FALSE))
R> pred_quant <- do.call(cbind, lapply(pred_quant, extract_pred))
```

4.1. Analyzing the uncertainty

In order to analyze the different sets of predictions \hat{Y}_{test}^{fair} , two slightly different perspectives can be taken, and we elaborate on both in the following sections.

Decision-maker analysis

The first way to analyze the uncertainty of the predictions is from the viewpoint of the decision-maker. By decision-maker, in this context we refer to the individual performing the analysis and obtaining a set of fair predictions. For a decision-maker, it is important to understand how sensitive the outcome of the classification is to uncertainties induced by both finite sample size and the inherent uncertainty induced by discrete variables. Let p_A be the

vector of predicted probabilities on the `test.data`, with length `nrow(test.data)`. Denote `nrow(test.data)` with n_{test} . Let p_B be another vector of predicted probabilities (under a different bootstrap repetition). We can consider the following four metrics of uncertainty:

1. For each threshold $t \in [0, 1]$, we compute the decision sets D_A, D_B , which are obtained by selecting all individuals for whom the value of $p_A \geq t$ (and p_B respectively). We can then compute the Jaccard similarity of decision sets D_A, D_B , for each threshold t . By comparing many pairs of bootstrap repetitions in this way, we can estimate what the average Jaccard similarity for each threshold t is. In practice, we may consider values of t within a subset of $[0, 1]$, i.e., between the 5% and 95% quantiles of the predicted probabilities in the repetition corresponding to full data:

```
R> jac_frm <- function(x, modes = "single") {
+   +
+   jac <- function(a, b) {
+     intersection <- length(intersect(a, b))
+     union <- length(a) + length(b) - intersection
+     intersection / union
+   }
+   +
+   res <- lapply(
+     seq(quantile(x[, 1L], 0.05), quantile(x[, 1L], 0.95), 0.01),
+     function(tsh) {
+       +
+       ret <- replicate(100L, {
+         col <- sample(ncol(x), 2L)
+         jac(which(x[, col[1L]] > tsh), which(x[, col[2L]] > tsh))
+       })
+       +
+       data.frame(tsh = tsh, y = mean(ret), sd = sd(ret),
+                 mode = modes)
+     }
+   )
+   +
+   do.call(rbind, res)
+ }
R>
R> jac_df <- rbind(jac_frm(pred_fin, "Finite Sample"),
+                 jac_frm(pred_quant, "Quantiles"))
```

2. Consider two indices i, j of the vectors p_A, p_B such that $i \neq j$, corresponding to two distinct individuals. For such pairs of individuals (i, j) we can analyze the probability $P((p_A)_i \geq (p_B)_j)$, where we can consider i, j to be drawn randomly, and p_A, p_B resulting from two random bootstrap repetitions. This probability tells us how likely it is that two randomly selected individuals appear in the same order in two repetitions.

```
R> ord_ind <- function(x, modes = "single") {
+   +
```

```

+   res <- replicate(5000L, {
+     row <- sample(nrow(x), 2)
+     ord <- mean(x[row[1], ] > x[row[2], ])
+     max(ord, 1 - ord)
+   })
+
+   data.frame(res = res, mode = modes)
+ }
R>
R> ord_df <- rbind(ord_ind(pred_fin, "Finite Sample"),
+                  ord_ind(pred_quant, "Quantiles"))

```

3. Another interesting metric is the inversion number. Notice that p_A, p_B define two permutations of the n_{test} individuals, when we consider a ranking of individuals according to their predicted probabilities. We can compute the inversion number of these two permutations π_A, π_B , which is the total number of pairs of individuals whose ordering is not the same in π_A and π_B . Notice that the maximum value of the inversion number is $\binom{n_{test}}{2}$. Hence, we normalize this quantity accordingly.

```

R> inv_frm <- function(x, modes = "single") {
+
+   gt <- function(x) x[1L] > x[2L]
+
+   res <- replicate(100L, {
+     col <- sample(ncol(x), 2L)
+     prm <- order(x[, col[2L]][order(x[, col[1L]])])
+     sum(combn(prm, 2L, gt)) / choose(length(prm), 2L)
+   })
+
+   data.frame(res = res, mode = modes)
+ }
R>
R> inv_df <- rbind(inv_frm(pred_fin, "Finite Sample"),
+                  inv_frm(pred_quant, "Quantiles"))

```

4. For each individual i , we can take the 5% and 95% quantiles of predicted probabilities in all of the `n.boot` bootstrap repetitions. We analyze the width of this interval across all individuals.

```

R> prb_frm <- function(x, modes = "single") {
+   qnt <- apply(x, 1L, quantile, probs = c(0.05, 0.95))
+   data.frame(width = qnt[2L, ] - qnt[1L, ], mode = modes)
+ }
R>
R> prb_df <- rbind(prb_frm(pred_fin, "Finite Sample"),
+                  prb_frm(pred_quant, "Quantiles"))

```

The results of these four metrics applied to different bootstrap repetitions on the COMPAS dataset are shown in Figure 6.

Warning: Removed 15 rows containing missing values (``geom_point()``).

Warning: Removed 15 rows containing missing values (``geom_line()``).

Warning: Removed 75 rows containing non-finite values (``stat_density()``).

Individual analysis

The other point of view we can take in quantifying uncertainty is that of the individual experiencing fair decisions. In this case, we are not so much interested in how much the overall decision set changes (as was the case above), but rather how much variation there is in the estimate for the specific individual. To this end, one might investigate the spread of the predicted values for a specific individual. The spread of the predictions for the first three individuals can be obtained as follows:

```
R> ind_prb <- data.frame(
+   prob = as.vector(t(pred_quant[seq_len(3), ])),
+   individual = rep(c(1, 2, 3), each = fa_boot_quant$n.boot)
+ )
```

This spread is visualized in Figure 7. While it might be tempting to take the mean of the different individual predictions, to have more stable results, this should not be done, as such an approach does not guarantee the fairness constraint `fairadapt()` aims to achieve. Hence, in order to achieve the desired fairness criteria overall, we have to accept some level of randomization at the level of individual predictions. It would be interesting for future work to quantify and optimize explicitly the trade-off between the desired fairness criteria and the necessary level of randomization.

5. Illustration

As a hypothetical real-world use of **fairadapt**, suppose that after a legislative change the US government has decided to adjust the salary of all of its female employees in order to remove both disparate treatment and disparate impact effects. To this end, the government wants to compute the counterfactual salary values of all male employees, that is the salaries that male employees would obtain, had they been female (i.e., the female group serves as the baseline). To do this, the government is using data from the 2018 American Community Survey by the US Census Bureau. This dataset is also available in pre-processed form as a package dataset from **fairadapt**. Columns are grouped into demographic (**dem**, including age, race, origin, citizenship, and economic region), familial (**fam**, including marital status, size of the family, and number of children), educational (**edu**, including number of years spent in schooling and the level of English proficiency) and occupational (**occ**, including the hours and week worked every year, occupation category, and industry of employment) categories and finally, salary is selected as response (**res**) and sex as the protected attribute (**prt**):

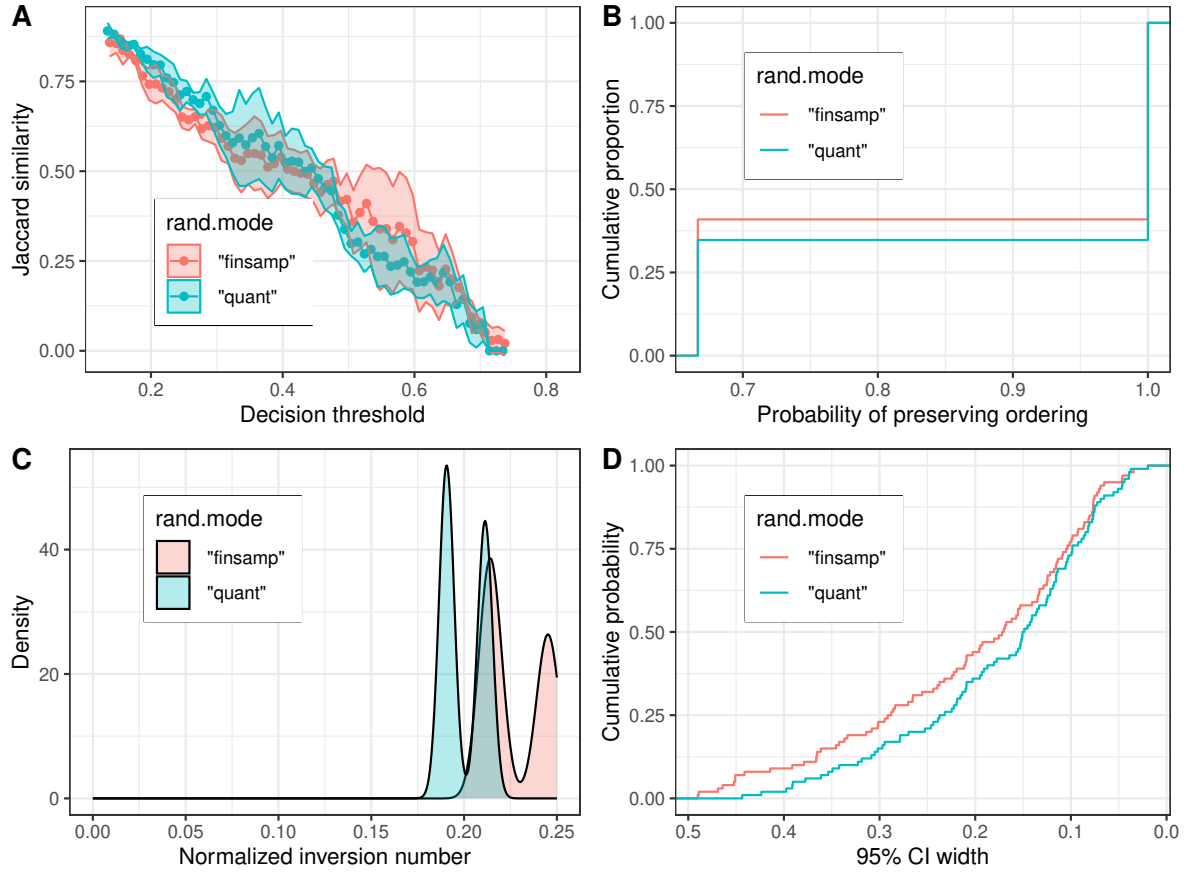


Figure 6: Analyzing uncertainty of predictions in the COMPAS dataset from decision-maker's point of view. Panel A shows how the Jaccard similarity of two repetitions varies depending on the decision threshold. Panel B shows the cumulative distribution of the random variable that indicates whether two randomly selected individuals preserve order (in terms of predicted probabilities) in bootstrap repetitions. Panel C shows the density of the normalized inversion number of between predicted probabilities in bootstrap repetitions. Panel D shows the cumulative distribution function of the 95% confidence interval (CI) width for the predicted probability of different individuals.

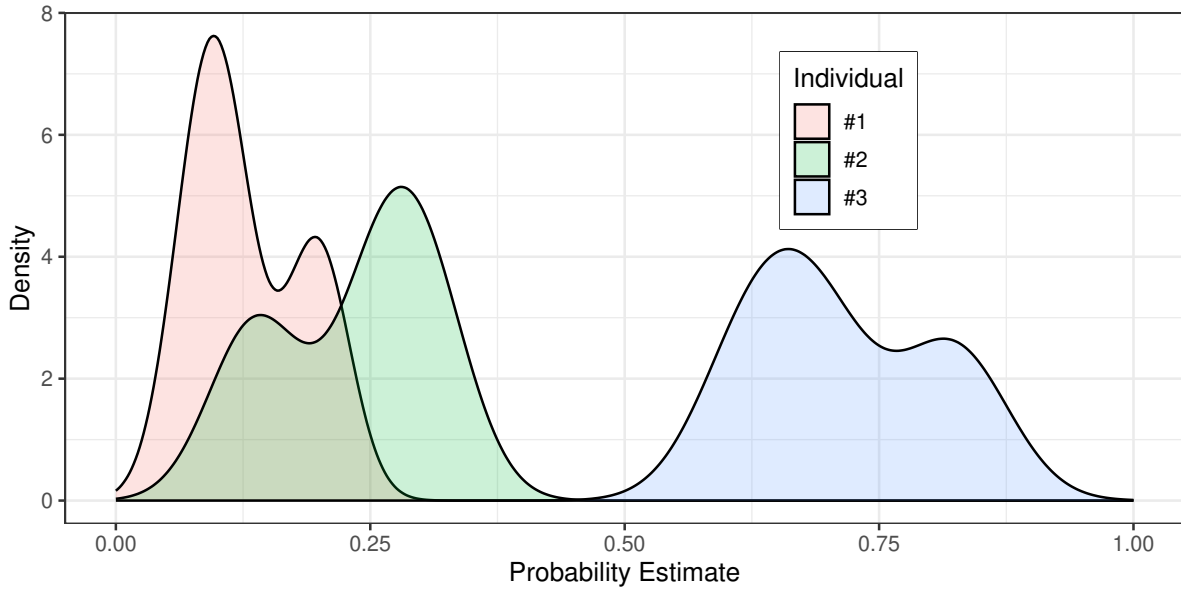


Figure 7: Analyzing the spread of individual predictions in the COMPAS dataset, resulting from different bootstrap repetitions.

```
R> gov_dat <- data("gov_census", package = "fairadapt")
R> gov_dat <- get(gov_dat)
R>
R> dem <- c("age", "race", "hispanic_origin", "citizenship",
+          "nativity", "economic_region")
R> fam <- c("marital", "family_size", "children")
R> edu <- c("education_level", "english_level")
R> occ <- c("hours_worked", "weeks_worked", "occupation",
+          "industry")
R>
R> prt <- "sex"
R> res <- "salary"
```

The hypothesized causal graph for the dataset is given in Figure 8. According to this, the causal graph can be specified as an adjacency matrix `gov_adj` and as confounding matrix `gov_cfd`:

```
R> cols <- c(dem, fam, edu, occ, prt, res)
R>
R> gov_adj <- matrix(0, nrow = length(cols), ncol = length(cols),
+                   dimnames = rep(list(cols), 2))
R> gov_cfd <- gov_adj
R>
R> gov_adj[dem, c(fam, edu, occ, res)] <- 1
R> gov_adj[fam, c(edu, occ, res)] <- 1
R> gov_adj[edu, c(occ, res)] <- 1
```

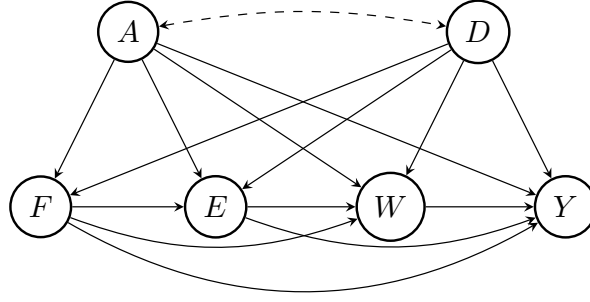


Figure 8: The causal graph for the government-census dataset. D are demographic features, A is gender, F represents marital and family information, E education, W work-related information and Y the salary, which is also the outcome of interest. The bidirected edge between A and D represents that the noise variables U_A and U_D share information.

```
R> gov_adj[occ,          res ] <- 1
R>
R> gov_adj[prt, c(fam, edu, occ, res)] <- 1
R>
R> gov_cfd[prt, dem] <- 1
R> gov_cfd[dem, prt] <- 1
R>
R> gov_grph <- graphModel(gov_adj, gov_cfd)
```

A visualization of the full graph using **igraph** is shown in Figure 9.

Before applying **fairadapt()**, we first log-transform the salaries to avoid dealing with a possibly heavy-tailed distribution for which quantile estimation may be more difficult. We then inspect the densities of variable **salary** by sex group, as shown in Figure 10A. There is a clear shift between the two distributions, indicating that **male** employees are better compensated than **female** employees. We perform the adaptation by using **n_samp** samples for training and **n_pred** samples for testing.

```
R> n_samp <- 750
R> n_pred <- 5

R> gov_dat$salary <- log(gov_dat$salary)
R>
R> gov_trn <- head(gov_dat, n = n_samp)
R> gov_prd <- tail(gov_dat, n = n_pred)
R>
R> set.seed(22)
R> gov_ada <- fairadapt(salary ~ ., train.data = gov_trn,
+                      adj.mat = gov_adj, cfd.mat = gov_cfd,
+                      prot.attr = prt)
```

After adapting the data, we investigate whether the salary gap has become smaller. This can be done by comparing distributions of variable **salary** using the **ggplot2**-exported S3 generic function **autoplot()** (Figure 10B).

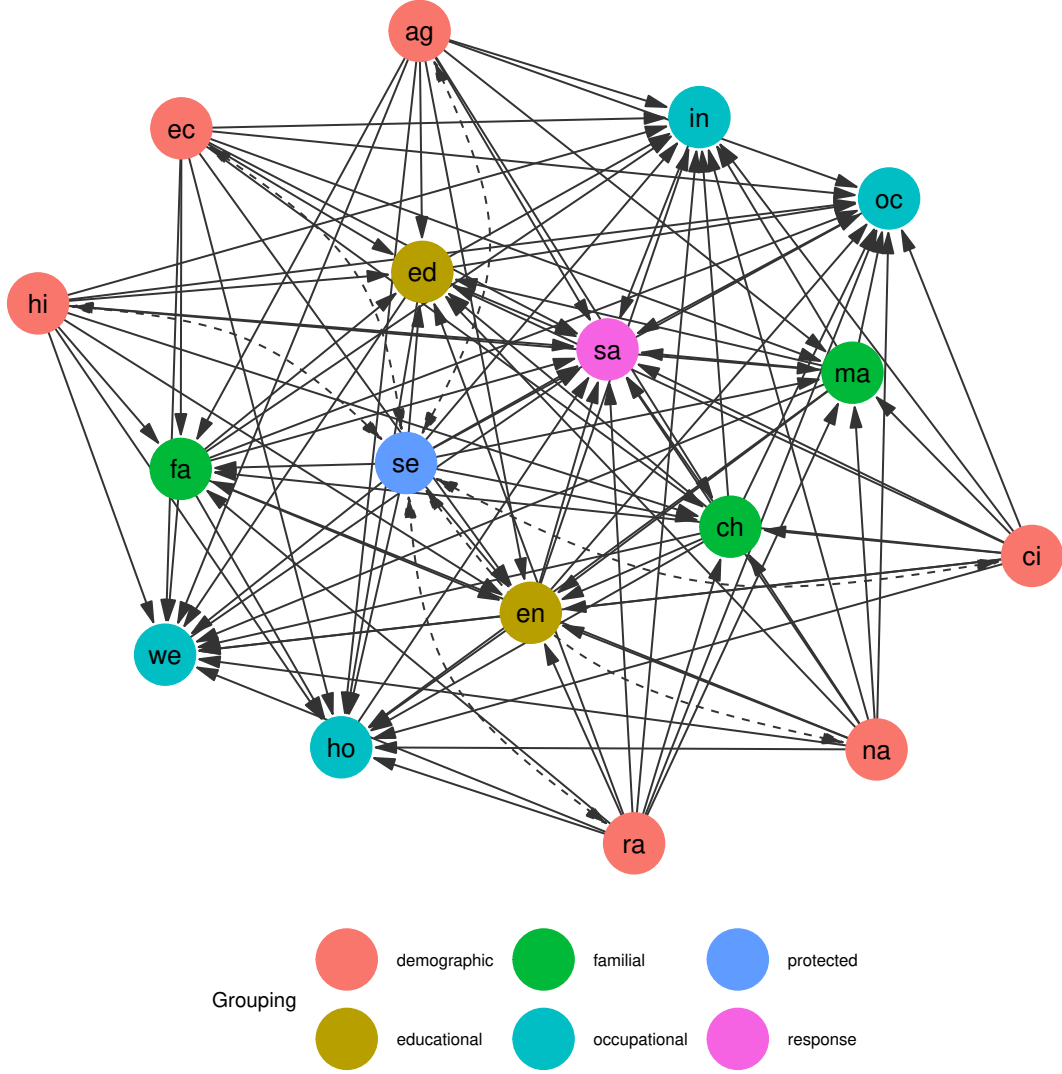


Figure 9: Full causal graph for the government census dataset, expanding the grouped view presented in Figure 8. *Demographic* features include age (**ag**), race (**ra**), whether an employee is of Hispanic origin (**hi**), is US citizen (**ci**), whether the citizenship is native (**na**), alongside the corresponding economic region (**ec**). *Familial* features are marital status (**ma**), family size (**fa**) and number of children (**ch**), *educational* features include education (**ed**) and English language levels (**en**), and *occupational* features, weekly working hours (**ho**), yearly working weeks (**we**), job (**oc**) and industry identifiers (**in**). Finally, the yearly salary (**sa**) is used as the *response* variable and employee sex (**se**) as the *protected* attribute variable.

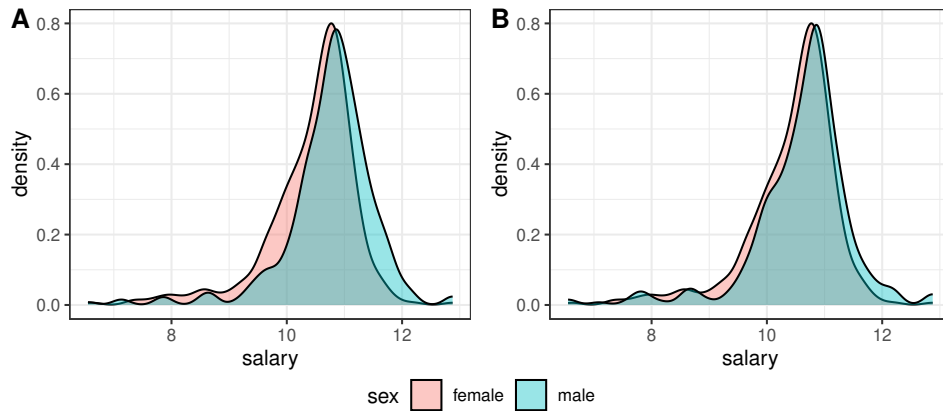


Figure 10: Visualization of salary densities grouped by employee sex, before (panel A) and after adaptation (panel B). Panel A indicates a shift towards higher values for male employees. In panel B, after the data is transformed, the gap between groups is reduced.

For adapting additional testing data, we use the base R S3 generic function `predict()` and output a selection of columns:

```
R> set.seed(2022)
R> gov_prd_ada <- predict(gov_ada, newdata = gov_prd)
R> gov_prd_ada[, c("sex", "age", "education_level", "salary")]
```

	sex	age	education_level	salary
1:	female	19	16	7.003065
2:	female	46	18	9.667765
3:	female	24	19	10.126631
4:	female	23	19	9.903488
5:	female	50	19	11.472103

Finally, we can do fair-twin inspection using the `fairTwins()` function of **fairadapt**, to retrieve counterfactual values of some features for different individuals:

```
R> fairTwins(gov_ada, train.id = 1:5,
+           cols = c("sex", "age", "salary"))
```

	sex	age	age_adapted	salary	salary_adapted
1	male	64	64	10.66896	10.51867
2	female	54	54	10.71442	10.71442
3	male	38	38	11.50288	11.50288
4	female	41	41	11.05089	11.05089
5	female	40	40	10.71885	10.71885

Note that values remain unchanged for female individuals (as *female* was used as baseline level). Variable `age`, which is not a descendant of the protected attribute `sex` (see Figure

9), also remains unchanged. However, variables `education_level` and `salary` do change for males, since these variables are descendants of the protected attribute `sex`.

We conclude the section with a remark. Notice that the variable `hours_worked` is a descendant of `A`. However, one might argue that this variable should *not* be adapted in the procedure, i.e., it should remain the same, irrespective of employee sex. In other words, one might argue it is acceptable to distinguish individuals based on this variable. This is the idea behind *resolving variables*, which are discussed in Section 6.1. It is worth emphasizing that we are not answering the question of how to choose which variables are resolving - this choice is left to social scientists familiar with the context of the dataset.

6. Extensions

Several extensions to the basic Markovian SCM formulation introduced in Section 2.3 exist, and these are outlined in the following sections.

6.1. Adding resolving variables

As we mentioned earlier, in some situations the protected attribute A might affect other variables in a non-discriminatory way. For instance, in the Berkeley admissions dataset (Bickel *et al.* 1975) we observe that females often apply for departments with lower admission rates and consequently have a lower admission probability. However, we perhaps would not wish to account for this difference in the adaptation procedure, if we were to argue that applying to a certain department is a choice everybody is free to make. Such examples motivated the idea of *resolving variables* by Kilbertus *et al.* (2017). A variable R is called resolving if

- (i) $R \in \text{de}(A)$, where $\text{de}(A)$ are the descendants of A in the causal graph \mathcal{G} .
- (ii) The causal effect of A on R is considered to be non-discriminatory.

In presence of resolving variables, computation of the counterfactual is carried out under the more involved intervention $\text{do}(A = a, R = R(a'))$. The potential outcome value $V(A = a, R = R(a'))$ is obtained by setting $A = a$ and computing the counterfactual while keeping the values of resolving variables to those they *attained naturally*. This is a nested counterfactual and the difference in Algorithm 1 is simply that resolving variables R are skipped in the for-loop. In order to perform fair data adaptation with the variable `test` being resolving in the `uni_admission` dataset used in Section 3, the string `"test"` can be passed as `res.vars` to `fairadapt()`.

```
R> res_basic <- fairadapt(score ~ ., train.data = uni_trn,
+                         test.data = uni_tst, adj.mat = uni_adj,
+                         prot.attr = "gender", res.vars = "test")
R> summary(res_basic)
```

Call:

```
fairadapt(formula = score ~ ., prot.attr = "gender", adj.mat = uni_adj,
          train.data = uni_trn, test.data = uni_tst, res.vars = "test")
```

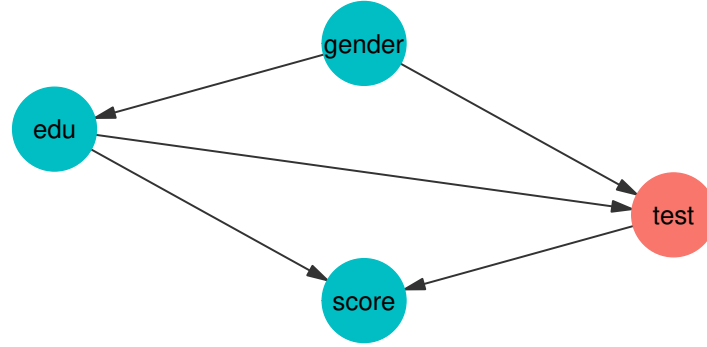


Figure 11: Visualization of the causal graph corresponding to the university admissions example introduced in Section 1 with the variable `test` chosen as a *resolving variable* and therefore highlighted in red.

Protected attribute:	<code>gender</code>
Protected attribute levels:	0, 1
Adapted variables:	<code>edu</code> , <code>score</code>
Resolving variables:	<code>test</code>
Number of training samples:	500
Number of test samples:	500
Quantile method:	<code>rangerQuants</code>
Total variation (before adaptation):	-0.7045
Total variation (after adaptation):	-0.3213

As can be seen from the respective model summary outputs, the total variation after adaptation, in this case, is larger than in the example from Section 3, with no resolving variables. The intuitive reasoning here is that resolving variables allow for some discrimination, so we expect to see a larger total variation between the groups.

```
R> uni_res <- graphModel(uni_adj, res.vars = "test")
```

A visualization of the corresponding graph is available from Figure 11, which highlights the resolving variable `test` in red, but the underlying graphical model remains the same.

6.2. Semi-Markovian and topological ordering variant

In Section 2 we focused on the Markovian case, which assumes that all exogenous variables U_i are mutually independent. However, in practice, this requirement is often not satisfied. If a mutual dependency structure between variables U_i exists, we are speaking about Semi-Markovian models. In the university admission example, we might have that $U_{\text{test}} \not\perp U_{\text{score}}$. That is, latent variables corresponding to variables `test` and `score` being correlated. Such dependencies between latent variables can be represented by dashed, bidirected arrows in the causal diagram, as shown in Figures 12 and 13.

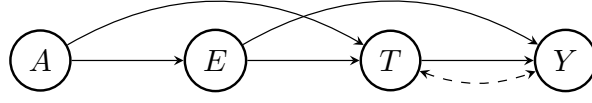


Figure 12: Causal graphical model corresponding to a Semi-Markovian variant of the university admissions example, introduced in Section 3. Here, we allow for the possibility of a mutual dependency between the latent variables corresponding to variables ‘test’ and ‘score’.

There is an important difference in the adaptation procedure for the Semi-Markovian case: when inferring the latent quantiles U_i of variable V_i , in the Markovian case, only the direct parents $\text{pa}(V_i)$ are needed. In the Semi-Markovian case, due to correlation of latent variables, using only the $\text{pa}(V_i)$ can lead to biased estimates of the U_i . Instead, the set of direct parents needs to be extended, as described in more detail by [Tian and Pearl \(2002\)](#). A brief sketch of the argument goes as follows: Let the *C-components* be a partition of the set V , such that each *C-component* contains a set of variables which are mutually connected by bidirected edges. Let $C(V_i)$ denote the entire *C-component* of variable V_i . We then define the set of extended parents as

$$\text{Pa}(V_i) := (C(V_i) \cup \text{pa}(C(V_i))) \cap \text{an}(V_i),$$

where $\text{an}(V_i)$ is the set of ancestors of V_i . The adaptation procedure in the Semi-Markovian case in principle remains the same as outlined in Algorithm 1, with the difference that the set of direct parents $\text{pa}(V_i)$ is replaced by $\text{Pa}(V_i)$ at each step.

To include the bidirected edges in the adaptation, we can pass a `matrix` as `cfd.mat` argument to `fairadapt()` such that:

- `cfd.mat` has the same dimension, column and row names as `adj.mat`.
- `cfd.mat` is symmetric.
- As with the adjacency matrix `adj.mat`, an entry `cfd.mat[i, j] == 1` indicates that there is a bidirected edge between variables `i` and `j`.

The following code performs fair data adaptation of the Semi-Markovian university admission variant with a mutual dependency between the variables representing test and final scores. For this, we create a matrix `uni_cfd` with the same attributes as the adjacency matrix `uni_adj` and set the entries representing the bidirected edge between vertices `test` and `score` to 1. Finally, we can pass this confounding matrix as `cfd.mat` to `fairadapt()`. A visualization of the resulting causal graph is available from Figure 13.

```

R> uni_cfd <- matrix(0, nrow = nrow(uni_adj), ncol = ncol(uni_adj),
+                   dimnames = dimnames(uni_adj))
R>
R> uni_cfd["test", "score"] <- 1
R> uni_cfd["score", "test"] <- 1
R>
R> semi <- fairadapt(score ~ ., train.data = uni_trn,
+                   test.data = uni_tst, adj.mat = uni_adj,
+                   cfd.mat = uni_cfd, prot.attr = "gender")

```

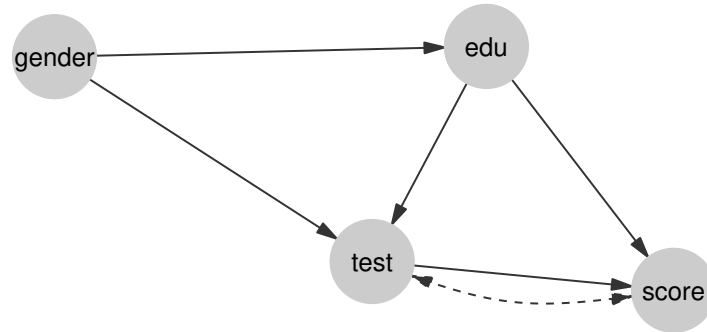


Figure 13: Visualization of the causal graphical model also shown in Figure 12, obtained when passing a confounding matrix indicating a bidirected edge between vertices `test` and `score` to `fairadapt()`. The resulting Semi-Markovian model can also be handled by `fairadapt()`, extending the basic Markovian formulation introduced in Section 2.3.

Alternatively, instead of using the extended parent set $\text{Pa}(V_i)$, we could also use the entire set of ancestors $\text{an}(V_i)$. This approach is implemented as well, and available by specifying a topological ordering. This is achieved by passing a `character` vector, containing the correct ordering of the names appearing in `names(train.data)` as `top.ord` argument to `fairadapt()`. The benefit of using this option is that the specific edges of the causal model \mathcal{G} need not be specified. However, in the linear case, specifying the edges of the graph, so that the quantiles are inferred using only the set of parents, will in principle have better performance. The topological variant can be invoked as follows:

```

R> set.seed(2022)
R> top_ord <- fairadapt(score ~ ., train.data = uni_trn, test.data = uni_tst,
+                       top.ord = c("gender", "edu", "test", "score"),
+                       prot.attr = "gender")
R>
R> summary(top_ord)

```

Call:

```

fairadapt(formula = score ~ ., prot.attr = "gender", train.data = uni_trn,
          test.data = uni_tst, top.ord = c("gender", "edu", "test",
          "score"))

```

```

Protected attribute:      gender
Protected attribute levels: 0, 1
Adapted variables:      edu, test, score

```

```

Number of training samples: 500
Number of test samples:    500
Quantile method:          rangerQuants

```

```

Total variation (before adaptation): -0.7045
Total variation (after adaptation):  -0.00599

```

Note that the topological variant for the university admissions dataset is the same as supplying the adjacency matrix, since the causal graph has no missing edges.

6.3. Questions of identifiability

So far we did not discuss whether it is always possible to carry out the counterfactual inference described in Section 2. In the causal literature, an intervention is termed *identifiable* if it can be computed uniquely using the data and the assumptions encoded in the graphical model \mathcal{G} . An important result by [Tian and Pearl \(2002\)](#) states that an intervention $\text{do}(X = x)$ on a variable X is identifiable if there is no bidirected path between X and $\text{ch}(X)$. Therefore, our intervention of interest is identifiable if one of the two following conditions are met:

- The model is Markovian.
- The model is Semi-Markovian and,
 - (i) there is no bidirected path between A and $\text{ch}(A)$ and,
 - (ii) there is no bidirected path between R_i and $\text{ch}(R_i)$ for any resolving variable R_i .

Based on this, the `fairadapt()` function might return an error, if the specified intervention is not possible to compute. An additional limitation is that **fairadapt** currently does not support *front-door identification* ([Pearl 2009](#), Chapter 3), meaning that certain special cases, which are in principle identifiable, are currently not handled.

6.4. Future avenues to be explored

We conclude with a brief look at the possible extensions of the **fairadapt** package, which we hope to consider in future work:

1. *Spurious pathways*: the **fairadapt** package allows for correcting discrimination along causal pathways from the protected attribute A to outcome Y . However, it is also possible that A and Y are associated through a confounding mechanism, and correcting for such spurious association poses an interesting methodological and practical challenge.
2. *General identification*: as discussed above, there are certain cases of causal graphs in which our $\text{do}(A = a, R = R(a'))$ intervention is identifiable, but the `fairadapt()` function currently does not support doing so. One of such examples is front-door identification, mentioned above. In a future version of **fairadapt**, we hope to cover all scenarios in which identification is possible.
3. *Path-specific effects*: when using resolving variables (Section 6.1), the user decides to label these variables as “non-discriminatory”, that is, the algorithm is free to distinguish between groups based on these variables. In full generality, a user might be interested in considering all path-specific effects ([Avin, Shpitser, and Pearl 2005](#)). Such an approach would offer even more flexibility in modeling, since for every attribute-outcome path $A \rightarrow \dots \rightarrow Y$, the user could decide whether it is fair or not.
4. *Selection bias*: a commonly considered problem in causal inference is that of selection bias ([Hernán, Hernández-Díaz, and Robins 2004](#)), when inclusion of individuals into the dataset depends on the observed variables in the dataset. In fairness applications,

the presence of selection bias could invalidate our conclusions about discrimination and make our fair predictions biased. Recovering from selection bias algorithmically would therefore be a desirable feature in the **fairadapt** package.

5. *Non-binary attribute A*: one assumption used currently is that the protected attribute *A* is binary. When considering possible protected attributes as socioeconomic status or education, this assumption may need to be relaxed. We hope to address this in future work.

References

- Avin C, Shpitser I, Pearl J (2005). “Identifiability of Path-Specific Effects.” *IJCAI’05*, p. 357–363.
- Barocas S, Selbst AD (2016). “Big Data’s Disparate Impact.” *Calif. L. Rev.*, **104**, 671.
- Bellamy RKE, Dey K, Hind M, Hoffman SC, Houde S, Kannan K, Lohia P, Martino J, Mehta S, Mojsilovic A, Nagar S, Ramamurthy KN, Richards J, Saha D, Sattigeri P, Singh M, Varshney KR, Zhang Y (2018). “AI Fairness 360: An Extensible Toolkit for Detecting, Understanding, and Mitigating Unwanted Algorithmic Bias.” URL <https://arxiv.org/abs/1810.01943>.
- Bickel PJ, Hammel EA, O’Connell JW (1975). “Sex Bias in Graduate Admissions: Data From Berkeley.” *Science*, **187**(4175), 398–404. doi:10.1126/science.187.4175.398. URL <https://doi.org/10.1126/science.187.4175.398>.
- Bird S, Dudik M, Edgar R, Horn B, Lutz R, Milan V, Sameki M, Wallach H, Walker K (2020). “Fairlearn: A Toolkit for Assessing and Improving Fairness in AI.” *Technical Report MSR-TR-2020-32*, Microsoft. URL <https://www.microsoft.com/en-us/research/publication/fairlearn-a-toolkit-for-assessing-and-improving-fairness-in-ai/>.
- Blau FD, Kahn LM (2003). “Understanding International Differences in the Gender Pay Gap.” *Journal of Labor Economics*, **21**(1), 106–144. doi:10.3386/w8200. URL <https://doi.org/10.3386/w8200>.
- Cannon AJ (2015). **qrnn**: *Quantile Regression Neural Network*. R package version 2.0.5, URL <https://cran.r-project.org/web/packages/qrnn>.
- Cannon AJ (2018). “Non-Crossing Nonlinear Regression Quantiles by Monotone Composite Quantile Regression Neural Network, With Application to Rainfall Extremes.” *Stochastic Environmental Research and Risk Assessment*, **32**(11), 3207–3225. doi:10.31223/osf.io/wg7sn. URL <https://doi.org/10.31223/osf.io/wg7sn>.
- Chouldechova A (2017). “Fair Prediction With Disparate Impact: A Study of Bias in Recidivism Prediction Instruments.” *Big data*, **5**(2), 153–163. doi:10.1089/big.2016.0047. URL <https://doi.org/10.1089/big.2016.0047>.

- Corbett-Davies S, Goel S (2018). “The Measure and Mismeasure of Fairness: A Critical Review of Fair Machine Learning.” *arXiv preprint arXiv:1808.00023*.
- Csardi G, Nepusz T (2006). “The **igraph** Software Package for Complex Network Research.” *InterJournal, Complex Systems*, 1695. URL <https://igraph.org>.
- Darlington RB (1971). “Another Look at Cultural Fairness.” *Journal of Educational Measurement*, 8(2), 71–82. doi:10.1111/j.1745-3984.1971.tb00908.x. URL <https://doi.org/10.1111/j.1745-3984.1971.tb00908.x>.
- Efron B, Tibshirani RJ (1994). *An Introduction to the Bootstrap*. CRC press. doi:10.1201/9780429246593. URL <https://doi.org/10.1201/9780429246593>.
- Galles D, Pearl J (1998). “An Axiomatic Characterization of Causal Counterfactuals.” *Foundations of Science*, 3(1), 151–182.
- Hardt M, Price E, Srebro N, *et al.* (2016). “Equality of Opportunity in Supervised Learning.” In *Advances in neural information processing systems*, pp. 3315–3323.
- Hernán MA, Hernández-Díaz S, Robins JM (2004). “A Structural Approach to Selection Bias.” *Epidemiology*, pp. 615–625. doi:10.1097/01.ede.0000135174.63482.43. URL <https://doi.org/10.1097/01.ede.0000135174.63482.43>.
- Kilbertus N, Carulla MR, Parascandolo G, Hardt M, Janzing D, Schölkopf B (2017). “Avoiding Discrimination Through Causal Reasoning.” In *Advances in Neural Information Processing Systems*, pp. 656–666.
- Koenker R, Hallock KF (2001). “Quantile Regression.” *Journal of Economic Perspectives*, 15(4), 143–156. doi:10.1257/jep.15.4.143. URL <https://doi.org/10.1257/jep.15.4.143>.
- Koenker R, Portnoy S, Ng PT, Zeileis A, Grosjean P, Ripley BD (2018). **quantreg**: *Quantile Regression*. R package version 5.86.
- Komiyama J, Takeda A, Honda J, Shima H (2018). “Nonconvex Optimization for Regression With Fairness Constraints.” In *International Conference on Machine Learning*, pp. 2737–2746. PMLR.
- Kozodoi N, V Varga T (2021). **fairness**: *Algorithmic Fairness Metrics*. R package version 1.2.2, URL <https://CRAN.R-project.org/package=fairness>.
- Kusner MJ, Loftus J, Russell C, Silva R (2017). “Counterfactual Fairness.” In *Advances in Neural Information Processing Systems*, pp. 4066–4076.
- Lambrecht A, Tucker C (2019). “Algorithmic Bias? An Empirical Study of Apparent Gender-Based Discrimination in the Display of STEM Career Ads.” *Management Science*, 65(7), 2966–2981. doi:10.2139/ssrn.2852260. URL <https://doi.org/10.2139/ssrn.2852260>.
- Larson J, Mattu S, Kirchner L, Angwin J (2016a). <https://github.com/propublica/compas-analysis>.

- Larson J, Mattu S, Kirchner L, Angwin J (2016b). “How We Analyzed the COMPAS Recidivism Algorithm.” *ProPublica* (5 2016), **9**.
- McGinley AC (2011). “Ricci v. DeStefano: Diluting Disparate Impact and Redefining Disparate Treatment.” *Scholarly Works*, **646**.
- Meinshausen N (2006). “Quantile Regression Forests.” *Journal of Machine Learning Research*, **7**(Jun), 983–999.
- Pearl J (2009). *Causality*. Cambridge University Press. doi:10.4288/kisoron.39.2_109. URL https://doi.org/10.4288/kisoron.39.2_109.
- Plecko D, Meinshausen N (2020). “Fair Data Adaptation With Quantile Preservation.” *Journal of Machine Learning Research*, **21**, 1–44.
- Scutari M (2021). **fairml**: *Fair Models in Machine Learning*. R package version 0.6, URL <https://CRAN.R-project.org/package=fairml>.
- Shukla K, Fang H, Jindal S (2019). “Tensorflow’s Fairness Evaluation and Visualization Toolkit.” URL <https://github.com/tensorflow/fairness-indicators>.
- Thomas O, Kehrenberg T, Bartlett M, Quadrianto N (2018). “EthicML: A Featureful Framework for Developing Fair Algorithms.” URL <https://github.com/predictive-analytics-lab/EthicML>.
- Tian J, Pearl J (2002). “A General Identification Condition for Causal Effects.” In *Eighteenth National Conference on Artificial Intelligence*, pp. 567–573. American Association for Artificial Intelligence, USA.
- Wickham H (2016). **ggplot2**: *Elegant Graphics for Data Analysis*. Springer-Verlag New York. ISBN 978-3-319-24277-4. URL <https://ggplot2.tidyverse.org>.
- Wiśniewski J, Biecek P (2022). “fairmodels: a Flexible Tool for Bias Detection, Visualization, and Mitigation in Binary Classification Models.” *The R Journal*, **14**, 227–243. doi:10.32614/RJ-2022-019. URL <https://rj.urbanek.nz/articles/RJ-2022-019/>.
- Wright MN, Ziegler A (2017). “**ranger**: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R.” *Journal of Statistical Software*, **77**(1). doi:10.18637/jss.v077.i01. URL <https://doi.org/10.18637/jss.v077.i01>.
- Zhang J, Bareinboim E (2018). “Fairness in Decision-Making: The Causal Explanation Formula.” In *Thirty-Second National Conference on Artificial Intelligence*. American Association for Artificial Intelligence, USA. doi:10.1609/aaai.v32i1.11564. URL <https://doi.org/10.1609/aaai.v32i1.11564>.

Affiliation:

Drago Plečko
ETH Zürich
Seminar for Statistics Rämistrasse 101 CH-8092 Zurich
E-mail: drago.plecko@stat.math.ethz.ch

Nicolas Bennett
ETH Zürich
Seminar for Statistics Rämistrasse 101 CH-8092 Zurich
E-mail: nicolas.bennett@stat.math.ethz.ch

Nicolai Meinshausen
ETH Zürich
Seminar for Statistics Rämistrasse 101 CH-8092 Zurich
E-mail: meinshausen@stat.math.ethz.ch