

---

# Explaining knock-on effects of bias mitigation

---

**Svetoslav Nizhnichenkov**  
IBM Research  
School of Computer Science  
University College Dublin  
Dublin, Ireland  
svetoslav.nizhnichenkov@ibm.com

**Rahul Nair**  
IBM Research  
Dublin, Ireland  
rahul.nair@ie.ibm.com

**Elizabeth Daly**  
IBM Research  
Dublin, Ireland  
elizabeth.daly@ie.ibm.com

**Brian Mac Namee**  
School of Computer Science  
University College Dublin  
Dublin, Ireland  
brian.macnamee@ucd.ie

## Abstract

In machine learning systems, bias mitigation approaches aim to make outcomes fairer across privileged and unprivileged groups. Bias mitigation methods work in different ways and have known “waterfall” effects, e.g., mitigating bias at one place may manifest bias elsewhere. In this paper, we aim to characterise impacted cohorts when mitigation interventions are applied. To do so, we treat intervention effects as a classification task and learn an explainable meta-classifier to identify cohorts that have altered outcomes. We examine a range of bias mitigation strategies that work at various stages of the model life cycle. We empirically demonstrate that our meta-classifier is able to uncover impacted cohorts. Further, we show that all tested mitigation strategies negatively impact a non-trivial fraction of cases, i.e., people who receive unfavourable outcomes solely on account of mitigation efforts. This is despite improvement in fairness metrics. We use these results as a basis to argue for more careful audits of static mitigation interventions that go beyond aggregate metrics.

## 1 Introduction

In the context of decision-making, **fairness** [17] is the absence of any prejudice or favouritism toward an individual or a group based on their inherent or acquired characteristics, while **bias** [17] occurs when an algorithm skews its decisions toward a particular group of individuals. Different types of biases can manifest themselves in many shapes and forms.

*Biases in data:* Data plays a significant role in the functionality of a wide range of AI systems. If data is biased in any way, the underlying algorithms within these systems will learn these biases and the generated predictions will reflect them. In automation settings, this can lead to biases being perpetuated at scale, along with the inequities they cause.

*Biases in algorithms:* Algorithms used to train AI systems can behave in a biased fashion. Often, this is due to specific design choices. For instance, a model trained to predict healthcare costs rather than illness will be biased against communities with low resources. The outcome of such systems may be used as input into real-life systems and affect the decisions of the people using them. This will also inevitably result in more data with underlying biases that will be used for training different algorithms further down the pipeline.

*Biases in user experience:* Users can promote biased behaviour when interacting with systems. For example, users of a search engine will likely interact with the results closer to the top of a list of results. Should this engagement signal be used as a proxy for popularity and the underlying system retrained to include this information, the results will exhibit a popularity bias where top results become even more popular [14].

To address these different biases, different algorithms have been proposed in the literature along with several notions of (statistical) fairness and associated metrics [3]. A fairness metric aims to assess how biased a model is and specific mitigation methods improve these metrics by altering data, algorithms, or outcomes. From the perspective of a machine learning practice, mitigation methods can be viewed as being able to work at three different stages of the machine learning pipeline: (1) pre-processing (on the data), (2) in-processing (during the model training phase), or (3) post-processing (adjusting the outcomes).

However, most mitigation methods are known to have a “waterfall” effect, i.e., bias addressed at one place can manifest biases elsewhere in the data [13]. Since mitigations are generally viewed as once-off static interventions, these knock-on waterfall effects have rarely been investigated beyond aggregate fairness measures. This paper aims to provide a characterisation of mitigation interventions in terms of impacted cohorts.

The paper makes two contributions. First, we develop a method that uses a supervised meta-classifier to describe impacted cohorts after bias mitigation. The method provides interpretable summaries of cohorts in the form of conjunctions. Second, we demonstrate over a wide range of mitigation strategies, fairness metrics and several datasets, that the meta-classifier is discriminative and highlight several important findings that call for a more careful audit of static mitigation interventions. Specifically, we empirically, show that there is always a negatively impacted cohort of individuals that are impacted solely on account of bias mitigation efforts, and some methods work consistently better than others.

The rest of this paper proceeds as follows. Section 2 introduces background information and related work. Section 3 describes the methodological approaches used in our experiments. Section 4 describes the characteristics of the datasets used in our experiments. Section 5 presents and discusses the results of our experiments. Section 6 summarises the contributions of the paper.

## 2 Background & Related Work

Our analysis is based on several mitigation methods from the literature. We use the AI Fairness 360 toolkit <sup>1</sup> implementations for these methods.

*Pre-processing* bias-mitigation algorithms work on the raw data by altering it in some way in order to reduce or eliminate the bias present in the data before this data is used for the training of some model. The pre-processing fairness-intervention techniques used in this study are Learning Fair Representations (LFR) [23] and Disparate Impact Remover (DIR) [4].

*In-processing* bias-mitigation algorithms have two parallel goals: accuracy and fairness. This means that the model should display sufficient accuracy while also being fair w.r.t. the fairness constraints. The idea is that such models directly take fairness into account and produce a classifier that will yield less bias compared to a model that is unaware of fairness. The in-processing approaches adopted in this study are Gerry Fair Classification (GF) [12] [11] and Prejudice Remover (PR) [10].

*Post-processing* bias-mitigation algorithms work on the predictions already made by a biased model and alter them in some way (e.g., change the target label of an instance from positive to negative or vice-versa) to reduce bias according to some fairness metric. The post-processing bias-mitigation approaches utilised in this study are Reject Option Classification (ROC) [9], Equalized Odds (EO) [6] and Calibrated Equalized Odds (CEO) [20].

Several other works have studied similar problems. In [13], the authors compared different fairness strategies to investigate their behaviour at the prediction level w.r.t. whether similar-performing techniques mitigate bias in the same way, impact a similar volume of people and if the same individuals are being affected. The results show that the fairness approaches do in fact impact a different number of individuals as well as even having different targeted cohorts. Furthermore, these observations hold true for multiple executions of the same fairness approaches.

---

<sup>1</sup><https://github.com/Trusted-AI/AIF360/>

Friedler et al. [5] compared various fairness approaches to uncover the differences w.r.t. performance, and whether fairness interventions have knock-on impacts. The findings indicate that, due to the fact that some portion of the fairness metrics correlate with one another, when a fairness approach optimizes one metric, it also performs well on all other correlated metrics. Moreover, the study shows that the fairness approaches have a tendency to be sensitive to variations in input (e.g., different train/test splits) which results in fluctuating fairness metrics for different splits. Finally, the outcomes of the different fairness approaches were observed to vary a lot even if a given fairness metric is being satisfied due to the underlying difference in the mechanisms of the fairness approaches.

Marchiori et al. [16] utilised an evaluation framework to perform a comparative study on the analysis of the effect of classification and explainability of two pre-processing bias-mitigation strategies, namely preferential sampling [8] and uniform sampling [7] on a set of known biased datasets. They provided a two-fold quantitative assessment, global high-level assessment and local-based explanations via SHAP [15] and LIME [21]. Results show the most variation w.r.t. the different classifier was displayed by the fairness metric output while the preferential sampling strategy yielded better results than uniform sampling at reducing the importance of the attributes with sensitive nature.

[1] conducted an empirical analysis on a binary classification task dataset [19] to compare the results of two bias-mitigation techniques, ROC and PR, w.r.t. group fairness and accuracy metrics while exploring different protected attribute settings (e.g., *gender*, *race*). Results showed that, overall, ROC outperformed PR based on the metric results. However, PR was found to improve fairness metrics better when changes to its hyperparameter settings were made, but this came at the cost of decreased accuracy. While this study looked at how bias-mitigation approaches compare when only aggregate metrics were taken into account, our study aims to dig deeper and uncover unintended consequences of bias-mitigation strategies that go beyond accuracy and fairness metrics.

### 3 Methodology

We focus on binary classification tasks with sensitive attributes that divide a population into privileged and unprivileged groups. The classification assigns each instance in a given group a favourable or an unfavourable label. While there are different notions of fairness, e.g., individual fairness (treating similar individuals the same way) and group fairness, we focus on the latter where the objective is to optimize a metric amongst the divided groups. We choose to optimize four different fairness metrics when it comes to group fairness, namely disparate impact, average odds, equal opportunity difference, and statistical parity difference. Finally, we use a batter of mitigation strategies that work at different stages of the model life cycle.

#### 3.1 Bias-Mitigation Pipeline

The bias-mitigation pipeline is a mapping  $g : D_x \rightarrow D_z$  that maps the input dataset,  $D_x$ , containing ground truth labels  $y$  and predictions  $y'$  generated from a model  $f$ , to  $D_z$ , an output dataset where predictions  $y'$  have been changed to bias-mitigated predictions  $y''$  as a result of applying some bias-mitigation model  $g$ .

Figure 1 provides a graphical overview. Given the dataset  $D_x$ , a fairness metric is used to measure the amount of bias present in it relative to a known privileged group. Following that, Multi-Dimensional Subset Scan [24] is performed to identify the cohorts in the data that show the most significant amount of bias. The resulting information regarding the biased cohorts is then passed to the bias-mitigation model,  $g$ , along with the dataset,  $D_x$ , and it produces a new set of bias-mitigated predictions,  $y''$ , which replace the predictions,  $y'$ , that came with the original dataset, thus outputting a dataset,  $D_z$ , containing bias-mitigated predictions, and a model,  $g$ , which can be used to debias future predictions.

#### 3.2 Meta-classifier Pipeline

This pipeline is a mapping  $D_z \rightarrow h$  where  $D_z$  is the output dataset from the bias-mitigation pipeline that contains bias-mitigated predictions  $y''$  and  $h$  is a meta classifier (in this case a decision tree) that has been trained on a mapping between the ground truth and the bias-mitigated predictions.

Figure 2 gives a high-level graphical overview of the meta-classifier pipeline, the aim of which is to learn a meta classifier in the form of an explainable decision tree,  $h$ , on a mapping between

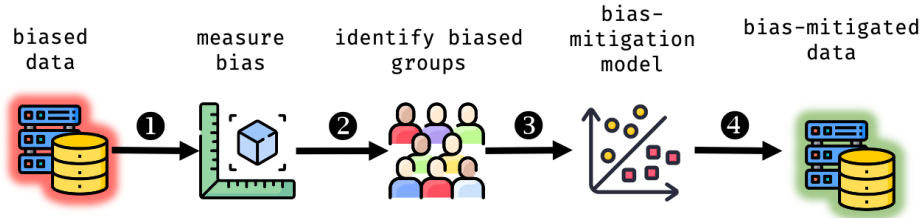


Figure 1: Bias-mitigation Pipeline. This pipeline depicts how a dataset with ground truth and predictions from a biased model is utilised to produce a new set of bias-mitigated predictions.

the ground truth and the bias-mitigated data, and generate explanations to reveal insights about the decision-making of the bias-mitigation model  $g$ . Given the dataset  $D_z$ , this pipeline applies a mapping between the ground truth,  $y$ , and the bias-mitigated predictions,  $y''$ , to produce a new output set,  $y'''$ , capturing observational changes in treatment. For instance, all individuals whose treatment didn't differ from the ground truth after applying bias mitigation will be assigned the value of 0, whereas the individuals whose treatment changed for the positive (e.g., an individual from the Adult dataset being predicted to earn more than 50,000 USD per annum post bias mitigation while their ground truth states otherwise) will be assigned the value of 1, and the people whose treatment changed for the negative (e.g., an individual from the Utrecht dataset being rejected for the job they applied post bias mitigation while their ground truth states they were approved) will be assigned the value of -1. As a result, this newly obtained set would capture changes in treatment between the ground truth and the bias-mitigated predictions.

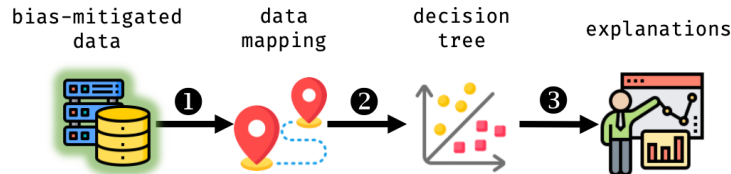


Figure 2: Meta-classifier Pipeline. This pipeline displays how an explainable model is used to learn a mapping of bias-mitigated predictions in order to generate explanations that can be used to infer the decision-making process of the bias-mitigation model.

Thus, this new output set,  $y'''$ , will be imposed as the new ground truth for the decision tree to learn from in order to capture the decision-making process of the bias-mitigation model,  $g$ , and be able to produce explanations in the form of decision trees or boolean decision rules that can be used to identify cohorts which will be affected in different ways (e.g., positively, negatively or no change in treatment).

## 4 Data

There are three different datasets utilised in this study. Most of them are well-known in the fairness research literature. The task associated with each of the datasets is binary classification and they all contain sensitive attributes of some nature (e.g., *age*, *race*, *gender* and *nationality*). The datasets used are:

*Utrecht Fairness Recruitment*: The Utrecht Fairness Recruitment dataset [22] is a synthetic dataset containing information about 4,000 individuals. Each individual is described by a number of sensitive attributes (e.g., *gender*, *age*, *nationality*) and general attributes (e.g., *highest degree achieved*, *degree score*, *programming experience*, *sport*, *international experience*). The target specifies whether that candidate was successful in receiving a job offer from the company they have applied to.

*Adult Census Income*: The Adult Census Income dataset [2] is composed of demographic, educational and work-related information about 32,561 entries. Sensitive features include *age*, *race*, *sex*, and more. The target indicates if individuals belonging to a certain entry in the dataset earn more than

50,000 US dollars per annum. Note that preprocessing was applied to this dataset whereby duplicates were removed and fields with missing values were replaced with *unknown* for better interpretability.

*Bank Marketing*: The Bank Marketing dataset [18] includes bank campaign-related and demographic information about 45,211 customers of a Portuguese banking institution. The sensitive feature for this data is *age*. The binary classification task is to predict whether an individual subscribes to a term deposit or not.

The datasets are split into 67% training and 33% testing subsets. The training data is used to train a model that does not take fairness into account and biased predictions are generated on the testing subset. The resulting testing subset (now containing ground truth and biased predictions) is then fed to the pipeline which trains bias-mitigation models and generates bias-mitigated predictions.

## 5 Results

We focus on the quantitative evaluation measures related to performance and fairness along with qualitative results describing the explanations. We summarise accuracy and fairness measures of the base (biased) model, i.e., a random forest classifier, along with all the tested mitigation methods in Table 1.

Table 1: Accuracy and fairness metrics across test sets for all datasets for biased and fairness models. Each metric result represents a tuple (x/y) where x is the mean while y is the standard deviation from a 5-fold cross-validation. (Note, the value in brackets for each fairness metric depicts its ideal value.)

	Model	Accuracy	Disparate Impact (1.0)	Average Odds (0.0)	Equal Opportunity (0.0)	Statistical Parity (0.0)
<b>Utrecht Fairness</b>	Biased	83.1%/±1.7	0.75/±.49	-0.13/±.19	-0.27/±.33	-0.08/±.13
	LFR	66.9%/±7.1	1.43/±1.17	0.01/±.18	-0.14/±0.15	0.08/±.22
	DIR	83.1%/±3.1	0.63/±.36	-0.13/±.19	-0.17/±.39	-0.12/±.13
	GF	74.0%/±2.7	0.68/±.80	-0.06/±.19	-0.06/±.29	-0.07/±.14
	PR	74.0%/±1.4	1.28/±.83	0.05/±.20	0.10/±.31	0.02/±.15
	ROC	81.4%/±3.8	0.80/±.23	-0.08/±.16	-0.09/±.32	-0.09/±.10
	EO	76.8%/±1.8	0.58/±.41	-0.15/±.22	-0.23/±.46	-0.13/±.14
	CEO	77.0%/±1.9	0.65/±.39	-0.12/±.19	-0.17/±.40	-0.11/±.13
<b>Adult Census</b>	Biased	84.3%/±1.0	0.55/±.12	-0.05/±.06	-0.04/±.12	-0.11/±.03
	LFR	80.7%/±1.3	0.57/±.12	-0.06/±.04	-0.08/±.09	-0.08/±.03
	DIR	83.3%/±1.3	0.55/±.07	-0.04/±.03	-0.01/±.07	-0.10/±.02
	GF	83.0%/±1.1	0.44/±.11	-0.09/±.05	-0.14/±.08	-0.10/±.03
	PR	84.3%/±1.1	0.39/±.06	-0.13/±.03	-0.20/±.05	-0.13/±.02
	ROC	80.8%/±2.0	1.13/±.26	0.11/±.07	0.14/±.08	0.04/±.07
	EO	79.5%/±1.1	0.60/±.07	-0.05/±.04	-0.06/±.12	-0.09/±.01
	CEO	79.5%/±1.0	0.43/±.06	-0.13/±.06	-0.19/±.10	-0.15/±.02
<b>Bank Marketing</b>	Biased	90.0%/±0.0	1.36/±.30	0.0/±.03	0.0/±.06	0.02/±.01
	LFR	86.9%/±2.4	1.05/±.45	0.0/±0.0	0.01/±.01	0.0/±0.0
	DIR	90.1%/±0.0	1.22/±.16	0.0/±.02	0.01/±.03	0.01/±.01
	GF	90.0%/±0.0	1.26/±.34	0.01/±.04	0.02/±.07	0.01±.01
	PR	90.2%/±0.0	1.43/±.20	0.04/±.03	0.07/±.05	0.02±.01
	ROC	90.1%/±0.0	1.24/±.14	0.0/±.04	0.0/±.07	0.02/±.01
	EO	90.0%/±0.0	1.36/±.30	0.01/±.03	0.01/±.06	0.01/±.01
	CEO	90.0%/±0.0	1.40/±.30	0.02/±.03	0.03/±.06	0.02/±.01

### 5.1 There is always a negatively impacted cohort

Table 2 highlights that bias mitigation interventions always induce knock-on effects that negatively impact a cohort (labelled Disagree(-) Ratio). Outcomes for this group that were previously favourable have been changed to be unfavourable solely on account of mitigation interventions. This result holds

Table 2: Cohort ratio distribution post bias mitigation and precision for each cohort based on the meta classifier’s predictions. Agreement shows the ratio of individuals whose treatment didn’t change post bias mitigation. Positive disagreement (+) shows the ratio of individuals whose treatment changed from negative to positive while negative disagreement (-) shows the ratio of individuals whose treatment changed from positive to negative. The precisions for the cohorts show the percentage of predictions made by the meta classifier that are correct for that cohort. Each result in this table is a tuple (x/y) where x is the mean while y is the standard deviation from a 5-fold cross-validation.

	Fairness Model	Agree Ratio	Disagree(+) Ratio	Disagree(-) Ratio	Meta Clf Accuracy	Agree Precision	Disagree(+) Precision	Disagree(-) Precision
Utrecht Fairness	LFR	66.9%/±7.1	12.0%/±4.7	21.1%/±5.4	69.0%/±6.0	96.6%/±6.5	0.0%/±0.0	14.2%/±30.5
	DIR	83.1%/±3.1	6.4%/±1.3	10.4%/±4.1	82.7%/±2.5	99.4%/±.06	0.0%/±0.0	0.0%/±0.0
	GF	74.0%/±2.7	6.2%/±2.1	19.8%/±3.5	74.6%/±4.1	90.0%/±3.0	19.9%/±13.6	33.9%/±13.0
	PR	73.9%/±1.4	7.4%/±.08	18.7%/±1.7	74.7%/±3.6	92.8%/±3.7	17.2%/±14.6	26.3%/±17.9
	ROC	81.5%/±3.8	12.1%/±1.3	6.4%/±3.6	81.3%/±3.8	99.4%/±.06	2.2%/±5.0	0.0%/±0.0
	EO	76.8%/±1.8	4.3%/±.08	18.9%/±2.4	77.2%/±3.4	95.8%/±1.4	0.0%/±0.0	20.2%/±13.9
	CEO	77.0%/±1.9	4.3%/±.07	18.7%/±2.5	76.8%/±3.6	95.0%/±1.5	0.0%/±0.0	20.7%/±13.6
Adult Census	LFR	80.7%/±1.3	5.7%/±.07	13.6%/±1.2	80.7%/±1.3	99.9%/±0.0	0.0%/±0.0	0.01%/±.02
	DIR	83.2%/±1.3	6.1%/±.08	10.7%/±.08	83.1%/±1.4	99.8%/±.02	0.0%/±0.0	0.02%/±.05
	GF	83.0%/±1.1	4.6%/±.07	12.4%/±.05	84.0%/±.07	99.1%/±.04	3.4%/±1.0	12.7%/±2.6
	PR	84.3%/±1.1	5.4%/±.07	10.3%/±.09	84.3%/±1.1	99.2%/±.02	12.4%/±5.3	0.07%/±.08
	ROC	80.8%/±2.0	12.8%/±1.9	6.4%/±1.0	80.8%/±1.9	99.9%/±0.0	0.07%/±.08	0.02%/±.05
	EO	79.6%/±1.0	3.7%/±.04	16.7%/±.08	80.0%/±1.2	97.6%/±1.1	0.0%/±0.0	13.6%/±5.4
	CEO	79.4%/±1.0	3.7%/±.07	16.9%/±.05	80.2%/±1.1	97.2%/±.09	0.0%/±0.0	17.8%/±5.7
Bank Marketing	LFR	86.9%/±2.4	1.7%/±2.9	11.4%/±.05	87.9%/±3.0	97.9%/±1.1	0.0%/±0.0	25.0%/±1.3
	DIR	90.0%/±.03	2.0%/±.03	8.0%/±.03	90.0%/±.03	100%/±0.0	0.0%/±0.0	0.0%/±0.0
	GF	90.0%/±.04	1.8%/±.03	8.2%/±.04	89.4%/±.04	97.8%/±1.4	25.9%/±11.0	11.6%/±8.3
	PR	90.2%/±.04	2.1%/±.04	7.7%/±.05	89.9%/±.03	98.6%/±.03	16.3%/±7.9	6.8%/±2.4
	ROC	90.1%/±.03	4.2%/±.06	5.7%/±.04	90.1%/±.03	100%/±0.0	0.0%/±0.0	0.0%/±0.0
	EO	89.9%/±.03	2.0%/±.03	8.1%/±.05	89.9%/±.03	100%/±0.0	0.0%/±0.0	0.01%/±.03
	CEO	90.0%/±.02	2.2%/±.02	7.8%/±.03	90.0%/±.02	99.9%/±0.0	0.0%/±0.0	0.02%/±.03

regardless of the bias-mitigation approach or dataset. Additionally, on average, across all datasets and bias-mitigation approaches, the negatively impacted cohorts are larger than the positively impacted ones. This suggests that mitigation interventions have a high likelihood of harming individuals outside of the targeted group.

## 5.2 Some fairness techniques do better than others in a consistent manner

Our experiments suggest that some mitigation methods consistently perform better than others. Table 1 shows that the Reject Option Classification (ROC) technique overall is better than the other methods tested. It demonstrates a good trade-off between accuracy and fairness measures (disparate impact in particular). On the other hand, Learning Fair Representations (LFR) consistently performs the worst in our battery of tests. To quantify the extent of knock-on effects, Table 2 displays the ratios of the people distributed in each of the affected cohorts. Here again, an ROC intervention impacts the smallest fraction of people, while LFR and Gerry Fair (GF) methods consistently induce the largest unintended effects.

A surprising observation for the Utrecht Fairness dataset in Table 1 is that the in-processing techniques Gerry Fair (GF) and Prejudice Remover (PR) perform worse than the other models given that they usually offer the best trade-off between accuracy and fairness metrics due to the fact that they have access to more information while removing bias as compared to pre and post-processing techniques. This, however, is only observed for the Utrecht Fairness dataset and can be attributed to its small size.

## 5.3 Explaining negatively affected groups

We now discuss qualitative results stemming from our meta-classifier which identifies cohorts that are negatively impacted. Table 2 shows relevant meta-classifier accuracy measures and class-specific precision. The lower performance of the model for the Utrecht dataset can be attributed to its smaller size. The precision of the decision tree for the affected cohorts follows the distribution of the cohort ratios as a result of the mapping after bias mitigation. That is, for large enough cohort ratios, the precision will be good as well. For example, the precision for the agreeable cohort for most datasets

and bias-mitigation approaches is around 100% and it follows the fact that the ratio of the agreeable cohort is composed of the majority of the population. Consequently, for the minority cohorts, the precision measurements can be quite low and sometimes even 0%. Therefore, this method has the potential to work well if there is a meaningful sample for minority classes.

As an example, the following two decision rule sets were taken from the Utrecht dataset where Gerry Fair (GF) was used to provide bias-mitigated predictions that exhibit 33.9% precision for the negatively affected cohort:

```
(1) ind-languages <= 0.50 & ind-university_grade <= 68.50 & sport_Rugby > 0.50 & company_C > 0.50 & ind-degree_bachelor > 0.50 & ind-university_grade > 55.00; [class: -1]
```

```
(2) ind-languages <= 0.50 & ind-university_grade > 68.50 & gender_female <= 0.50 & company_A <= 0.50 & sport_Swimming <= 0.50; [class: -1]
```

The first decision rule set can be interpreted as follows: "Individuals who don't speak any foreign languages, have a bachelor's degree with a grade between 55 and 68, play Rugby and have applied for a job for company C will receive a negative treatment."

The second decision rule set can be interpreted as follows: "Male individuals who don't speak any foreign languages, who don't practice swimming, who have achieved a university grade > 68.5 and have not applied to company A will receive a negative treatment."

Furthermore, the following decision rule set is taken from the Adult dataset where Calibrated Equalized Odds (CEO) was used to provide bias-mitigated predictions that exhibit 17.8% precision for the negatively impacted cohort:

```
(1) Marital-status_Married-civ-spouse <= 0.50 & capital-gain <= 7139.50 & hours-per-week > 43.50 & age <= 43.50 & (native-country_Poland > 0.50 or native-country_France > 0.50); [class: -1]
```

This rule set can be interpreted as follows: "Individuals who are not married, are below 44 years of age, are of either Polish or French nationality, work above 43.5 hours per week and have an annual capital gain of less than 7140 currency units will receive a negative treatment."

Thus, we can observe that we can identify the cohorts and the characteristics of individuals that will be impacted in a negative fashion with different success rates based on the ratio distribution among the different cohorts.

## 6 Conclusion

In this paper, we looked at discovering and explaining the knock-on effects of bias mitigation by conducting an empirical study involving a number of biased datasets, various bias-mitigation approaches and a set of fairness metrics. We utilised two pipelines, one that performs bias mitigation and gives us bias-mitigated data and a bias-mitigator, and a subsequent one that utilises the bias-mitigated data in order to build an explainable classifier with which we could narrow down on the affected cohorts. Using the results of these pipelines, we showed three main things, (1) regardless of dataset and bias-mitigation approach utilised, there are always individuals who will receive a negative treatment post bias mitigation, (2) some bias-mitigation techniques perform better in comparison to others in a consistent fashion across all datasets by showing that ROC harms the least amount of individuals while also exhibits good accuracy and satisfies most fairness metrics, and (3) there are methods to allow for the discovery of the affected cohorts, one of which can be a decision tree classifier that will generate decision rule sets, the characteristics of which can describe the affected individuals and tell who will be impacted in future predicting scenarios. We also note that this method is reliant on the cohort ratio distributions and it performs best when there is a balance between the different classes. Thus, with this study, we argue that there is more to bias mitigation than only observing accuracy and fairness metrics as negative impacts can also manifest themselves in other shapes and forms and we stress on the more careful audit of static bias-mitigation interventions that go beyond taking into account only aggregate measurements.

## 7 Acknowledgements

This work was funded by the European Union’s Horizon Europe research and innovation programme under grant agreement no. 101070568 (AutoFair) and Science Foundation Ireland under Grant number 18/CRT/6183.

## References

- [1] Tor H Aasheim and Knut T Hufthammer. Bias mitigation with AIF360: A comparative study.
- [2] Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996. DOI: <https://doi.org/10.24432/C5XW20>.
- [3] Simon Caton and Christian Haas. Fairness in Machine Learning: A Survey, October 2020.
- [4] Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and Removing Disparate Impact. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 259–268, Sydney NSW Australia, August 2015. ACM.
- [5] Sorelle A. Friedler, Carlos Scheidegger, Suresh Venkatasubramanian, Sonam Choudhary, Evan P. Hamilton, and Derek Roth. A comparative study of fairness-enhancing interventions in machine learning. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 329–338, Atlanta GA USA, January 2019. ACM.
- [6] Moritz Hardt, Eric Price, Eric Price, and Nati Srebro. Equality of Opportunity in Supervised Learning.
- [7] Faisal Kamiran and Toon Calders. Data preprocessing techniques for classification without discrimination. *Knowledge and Information Systems*, 33(1):1–33, October 2012.
- [8] Faisal Kamiran, Toon Calders, F Kamiran, Tue NI, T Calders, and Tue NI. Classification with No Discrimination by Preferential Sampling.
- [9] Faisal Kamiran, Asim Karim, and Xiangliang Zhang. Decision Theory for Discrimination-Aware Classification. In *2012 IEEE 12th International Conference on Data Mining*, pages 924–929, Brussels, Belgium, December 2012. IEEE.
- [10] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. Fairness-Aware Classifier with Prejudice Remover Regularizer. In David Hutchison, Takeo Kanade, Josef Kittler, Jon M. Kleinberg, Friedemann Mattern, John C. Mitchell, Moni Naor, Oscar Nierstrasz, C. Pandu Rangan, Bernhard Steffen, Madhu Sudan, Demetri Terzopoulos, Doug Tygar, Moshe Y. Vardi, Gerhard Weikum, Peter A. Flach, Tijl De Bie, and Nello Cristianini, editors, *Machine Learning and Knowledge Discovery in Databases*, volume 7524, pages 35–50. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012.
- [11] Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. Preventing Fairness Gerrymandering: Auditing and Learning for Subgroup Fairness.
- [12] Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. An Empirical Study of Rich Subgroup Fairness for Machine Learning. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 100–109, Atlanta GA USA, January 2019. ACM.
- [13] Natasa Krco, Thibault Laugel, Jean-Michel Loubes, and Marcin Detyniecki. When Mitigating Bias is Unfair: A Comprehensive Study on the Impact of Bias Mitigation Algorithms, February 2023.
- [14] Kristina Lerman and Tad Hogg. Leveraging Position Bias to Improve Peer Recommendation. *PLoS ONE*, 9(6):e98914, June 2014.
- [15] Scott M Lundberg and Su-In Lee. A Unified Approach to Interpreting Model Predictions.
- [16] Marta Marchiori Manerba and Riccardo Guidotti. Investigating Debiasing Effects on Classification and Explainability. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*, pages 468–478, Oxford United Kingdom, July 2022. ACM.
- [17] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A Survey on Bias and Fairness in Machine Learning, January 2022.
- [18] S. Moro, P. Rita, and P. Cortez. Bank Marketing. UCI Machine Learning Repository, 2012. DOI: <https://doi.org/10.24432/C5K306>.



- [19] U.S. Bureau of Labor Statistics. Demographics and Employment in the United States, 2013. URL: <https://www.kaggle.com/datasets/econdata/demographics-and-employment-in-the-united-states>.
- [20] Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q Weinberger. On Fairness and Calibration.
- [21] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why Should I Trust You?": Explaining the Predictions of Any Classifier, August 2016.
- [22] Sieuwert van Otterloo. Utrecht Fairness Recruitment Data, 2021. URL: <https://www.kaggle.com/datasets/ictinstitute/utrecht-fairness-recruitment-dataset>.
- [23] Richard Zemel. Learning Fair Representations.
- [24] Zhe Zhang and Daniel B. Neill. Identifying significant predictive bias in classifiers, 2017.