

The Impact of Data Preparation on the Fairness of Software Systems

Inês Valentim, Nuno Lourenço, Nuno Antunes
CISUC, Department of Informatics Engineering
University of Coimbra
Coimbra, Portugal
valentim@dei.uc.pt, naml@dei.uc.pt, nmsa@dei.uc.pt

Abstract—Machine learning models are widely adopted in scenarios that directly affect people. The development of software systems based on these models raises societal and legal concerns, as their decisions may lead to the unfair treatment of individuals based on attributes like race or gender. Data preparation is key in any machine learning pipeline, but its effect on fairness is yet to be studied in detail. In this paper, we evaluate how the fairness and effectiveness of the learned models are affected by the removal of the sensitive attribute, the encoding of the categorical attributes, and instance selection methods (including cross-validators and random undersampling). We used the Adult Income and the German Credit Data datasets, which are widely studied and known to have fairness concerns. We applied each data preparation technique individually to analyse the difference in predictive performance and fairness, using statistical parity difference, disparate impact, and the normalised prejudice index. The results show that fairness is affected by transformations made to the training data, particularly in imbalanced datasets. Removing the sensitive attribute is insufficient to eliminate all the unfairness in the predictions, as expected, but it is key to achieve fairer models. Additionally, the standard random undersampling with respect to the true labels is sometimes more prejudicial than performing no random undersampling.

Index Terms—Fairness, Data Preparation, Machine Learning.

I. INTRODUCTION

Software systems based on machine learning (ML) are being used at an increasingly higher rate and on a multitude of scenarios that have a significant impact on people’s lives. Their ubiquity raises several legal and societal concerns, as decisions based on the output of ML models may introduce or perpetuate historical bias against some individuals, based on their intrinsic characteristics, such as race, gender or age. The use of automated decision-making systems is often appealing due to the gains associated with it, and might even be perceived as a step towards the eradication of personal bias from the process. Nevertheless, many are the risks associated with a careless adoption of decisions supported by these systems.

In this context, fairness emerges as a key property in terms of the reliability and trustworthiness of software systems based on ML. These receive nowadays increased attention from regulatory institutions, with the recently approved European Union General Data Protection Regulation (GDPR) demanding organisations to handle personal data in a privacy-preserving, fair and transparent manner [1].

Techniques to assess fairness and build models capable of providing fairer predictions are of great help to organisations which intend to be GDPR compliant, but may lack

the resources or knowledge [2]. These organisations must be aware of the potential biases in their models at the design, implementation and deployment phases, and should make regular fairness evaluations of their systems [3]. Moreover, the assessment approaches may be used to audit non-compliant organisations, therefore providing valuable insight on violations of these fairness principles [2]. Individuals who rely on these organisations also benefit from the deployment of fairness-aware models and the adoption of such practices, since they provide an extra assurance that their data is not being used in ways that may negatively impact their daily lives.

In this paper, we assess the impact of widely adopted data preparation procedures on the fairness of systems based on ML. More precisely, we consider the removal of the sensitive attribute, the encoding of the categorical attributes, and instance selection methods, like cross-validation and random undersampling. Despite not being in the main scope of this work, we also consider the influence of the learning algorithm on fairness. From the many algorithms suitable for a supervised classification setting, we first focus on tree-based methods, like Decision Trees and Random Forests, partly due to the easier interpretability of the resulting models.

The obtained results show the importance of adopting the standard legal practices to mitigate discrimination, namely the removal of the sensitive attribute prior to training. However, we have also found that this procedure might not always lead to the expected behaviour, with the models’ predictions sometimes being more unfair than when the model has access to the sensitive attribute. We also report the drawbacks of using more complex learning algorithms, with Random Forests making more discriminatory predictions than Decision Trees. Furthermore, we emphasise that caution must be taken when dealing with datasets which show an imbalance with respect to both the true labels and the sensitive attribute. Standard sampling methods, such as random undersampling with respect to the true labels, may have undesired effects on fairness.

The remainder of this paper is organised as follows. Section II provides an overview of the key concepts of machine learning and fairness, and reviews related work on fairness of software systems based on machine learning models. Section III presents the research questions and details the experimental methodology. The obtained results are presented in Section IV and our findings are discussed in Section V. Section VI concludes the paper.

II. BACKGROUND AND RELATED WORK

In this section we overview machine learning concepts key to our work, after which we present the core concepts of fairness and related work on the topic, namely development of fairness-aware algorithms and fairness metrics.

A. Machine Learning

Machine learning (ML) aims at enabling software systems to learn from data, by modifying and adapting their actions towards the desired outcome [4]. Our focus is on **supervised classification problems**: we have access to the features (or attributes) of the training instances and to the *discrete outcome variable* (or true labels), which guides the learning process [5], [6]. The goal is to assign one of the possible classes to each instance. For example, we may want to determine whether someone has a high or a low risk of recidivism.

A system based on ML usually follows a pipeline as shown in Fig. 1. The **data collection** phase includes gathering representative data for the problem at hand, as well as labelling the training examples when in the presence of a supervised learning task. The **data preparation and pre-processing** steps may include handling missing data, encoding categorical features, discretisation, feature normalisation, feature selection and feature reduction techniques. It is crucial to apply these techniques for models to deliver the expected results, while helping to deal with overfitting. **Model selection** deals with the process of selecting the most appropriate model for the problem we are trying to solve, taking the complexity and flexibility of the models into account [7]. **Model assessment** deals with evaluating the performance of the chosen model by estimating its generalisation error on new unseen data [6], [7].

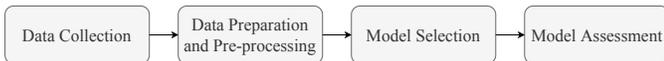


Fig. 1. Typical machine learning pipeline.

Decision Trees try to learn simple decision rules from the features of the training data [8]. A classification tree is built by following a recursive binary splitting process guided by the evaluation of the splits' quality using a criterion like the classification error rate, the Gini index, or cross-entropy [7]. Tree pruning can be used to avoid overfitting. Some well-known decision tree algorithms include ID3 and C4.5. **Random Forests** are collections of decision trees where the final prediction is given by a majority vote over the predictions of all the trees in the ensemble [9]. To reduce the correlation between the trees, the candidates for splitting are randomly selected from the full set of input features before each split [6]. This randomisation process also aims at reducing variance [7].

To choose a model we need to assess its generalisation performance (capability of making accurate predictions given new unseen instances) [6]. **Cross-validation** is a suitable approach when there is insufficient data to make a partition into training, validation, and test sets. Furthermore, the choice of a performance metric is dependent on the problem and the characteristics of the available data. A *confusion matrix*, as shown in Table I, summarises the results of a binary

classification problem, with four possible classification results. Several performance metrics, whose definitions can be found in [4], can be derived from this matrix.

TABLE I
CONFUSION MATRIX FOR A BINARY CLASSIFICATION PROBLEM.

		Predicted Class	
		Positive	Negative
True Class	Positive	True Positive (TP)	False Negative (FN)
	Negative	False Positive (FP)	True Negative (TN)

Despite the widespread use of **accuracy** to evaluate the performance of an algorithm, it may lead to misleading results in imbalanced scenarios and when incorrect classifications have a different cost. It is given by the ratio between correctly classified instances and the total number of instances.

Precision is given by the fraction of instances classified as positive that are correctly classified: $TP/(TP + FP)$.

Recall, also known as *true positive rate (TPR)* or *sensitivity*, is given by the fraction of positive instances that are correctly classified: $TP/(TP + FN)$.

The **F1-score** corresponds to the harmonic mean of precision and recall, also given by:

$$F1\text{-score} = \frac{2 \times TP}{2 \times TP + FN + FP} \quad (1)$$

Specificity is usually used alongside sensitivity and is given by the fraction of negative instances that are correctly classified: $TN/(TN + FP)$. *False positive rate (FPR)* is given by $1 - specificity$.

The *receiver operating characteristic (ROC) curve* depicts the trade-off between costs and benefits by plotting the recall against the FPR, as some threshold parameter of the classifier is varied. The *area under the curve (AUC)*, a single quantitative summary of a model's performance [6], can be computed from the ROC curve.

B. Fairness Concepts

Although the fairness of a software system is difficult to define due to its ambiguity [10], throughout this work we consider it to be the absence of bias or discrimination against people based on **protected or sensitive attributes**, such as race, gender, or age. We are particularly interested in the unfairness of the predictions made by ML models, and its relation to the data used to train them, which in turn can be biased if derived from discriminatory historical decisions [11].

Disparate treatment, a direct form of discrimination, results from a deliberated use of the sensitive attribute and can be avoided by removing it from the data prior to training the model [12]. Even when trained without the sensitive attribute, the predictions may still be discriminatory, leading to an unfair treatment of protected groups [12], [13]. This *red-lining effect* is due to the presence of features highly associated with the sensitive attribute [12], [13] and is linked to **disparate impact**. This indirect form of discrimination is not illegal in itself, as long as objective and reasonable justifications for it can be given [14], [15]. The rationale behind **disparate mistreatment**, proposed by [11], addresses

differences in the misclassification rates between the protected and the unprotected groups [11].

In a supervised classification problem, we are given a labelled dataset $\mathcal{D} = \{X, S, Y\}$ of n instances: X are the non-sensitive attributes, S denotes a sensitive attribute, and Y represents the true labels. The variable that represents the classifier’s predictions is referred to as \hat{Y} . A binary S partitions the dataset into the protected or unprivileged group (value of 0 for the sensitive attribute) and the unprotected or privileged group (value of 1 for the sensitive attribute).

We focus on fairness metrics which can be applied on different stages: prior to training a model, and afterwards, when the predictions are known. The following definitions can be applied to the datasets, if we use Y instead of \hat{Y} .

Statistical parity difference, or the Calders-Verwer score (CVS), considers the difference of the rate of favourable predictions between protected and unprotected groups [12]:

$$P(\hat{Y} = 1|S = 1) - P(\hat{Y} = 1|S = 0) \quad (2)$$

Measures of this metric lie in $[-1, 1]$, with 0 being optimal fairness. The sign of a measure indicates a skew in favour of either the protected or unprotected group [16].

Disparate impact (DI) is given by the ratio of the rate of favourable predictions for the protected group to that of the unprotected group [16]:

$$\frac{P(\hat{Y} = 1|S = 0)}{P(\hat{Y} = 1|S = 1)} \quad (3)$$

This is often referred to as the $p\%$ -rule and for a classifier to be fair, i.e. not to have DI, it should be greater than 80% but lower than 125% [15], [17]. The 80% rule is advocated by the US Equal Employment Opportunity Commission [15], and can be found in the Code of Federal Regulations, in the scope of labour regulations [18]. Measures of this metric lie in $[0, \infty]$, with a value different from 1, the optimal fairness, indicating a skew in favour of one of the groups.

The prejudice index (PI) corresponds to the mutual information between the predictions and the sensitive attribute [19]. The **normalised prejudice index** (NPI) results from the application of a normalisation technique for mutual information:

$$NPI = \frac{PI}{\sqrt{H(\hat{Y})H(S)}} \quad (4)$$

where $H(\cdot)$ is an information entropy function [19]. The NPI ranges between $[0, 1]$, with 0 being the optimal value.

C. Related Work

The approaches that have been proposed to enhance the fairness of ML models or mitigate the bias in their predictions can be grouped into three categories: pre-process, in-process and post-process. **Pre-process approaches** modify the training data to make it free of discrimination; **in-process approaches** change the models by adding constraints and regularisation terms to the objective functions; and **post-process approaches** directly change the predictions made by the models [12].

Pre-process approaches align with our work in that they explore ways to manipulate the data before it is used to

train the models. The Uniform Sampling and the Preferential Sampling methods proposed in [20] are similar to those used in our work. In addition to these sampling methods, the authors also propose suppression (removal of the sensitive attribute and those highly correlated with it), massaging (modification of the labelling of the training examples) and reweighing (assignment of weights to the training instances).

We also take inspiration from [19], where the authors compare models trained with and without the sensitive attribute, and propose the NPI as a fairness metric. They also consider different data preparation procedures, but no evaluation is performed with regards to using different versions of the same dataset to train the same learning algorithms.

The in-process approaches proposed by [21] and [22] focus on modifying tree-based methods so as to make them discrimination-aware, thus improving the fairness of the models’ predictions. This goal is accomplished by changing the evaluation of the splitting criterion and relabelling the leaves.

Some previous work, such as [23], [11] and [24], focused on the definition of new fairness notions and metrics capable of overcoming some of the known flaws of more traditional metrics, like statistical parity and the 80% rule.

More recently, [16] focused on defining a benchmark approach to evaluate fairness. A variety of fairness-enhancing methods are compared, and the relation between different fairness metrics is investigated. In contrast, we focus on assessing the impact of standard ML data preparation procedures rather than on fairness-aware methods. This work also alerts to the need to carefully specify the data pre-processing techniques applied to the training data, as they may have a significant impact on the fairness evaluation of a system.

III. METHODOLOGY

The main goal of this work is to **understand how the different procedures applied to a dataset during data preparation impact the fairness** of the predictions made by a model. To achieve this goal we aim at answering the following research questions:

- RQ1.** How does the removal of the sensitive attribute impact the fairness of the predictions made by an ML model?
- RQ2.** How does feature discretisation and the encoding of the categorical attributes impact the fairness of the predictions made by an ML model?
- RQ3.** What is the impact of different instance selection techniques on the fairness of the predictions made by an ML model? We consider cross-validators and sampling methods as instance selection techniques.

We also investigated the impact of the learning algorithm on the fairness of a system. The predictive performance of the models was also evaluated.

Fig. 2 shows the followed approach to assess the impact of data preparation procedures on the fairness of software systems, more precisely on the fairness of the predictions made by an ML model. The steps represented by green boxes are the main focus of this work, with dashed boxes representing optional steps: for instance, under some of the tested configurations, the sensitive attribute is not removed

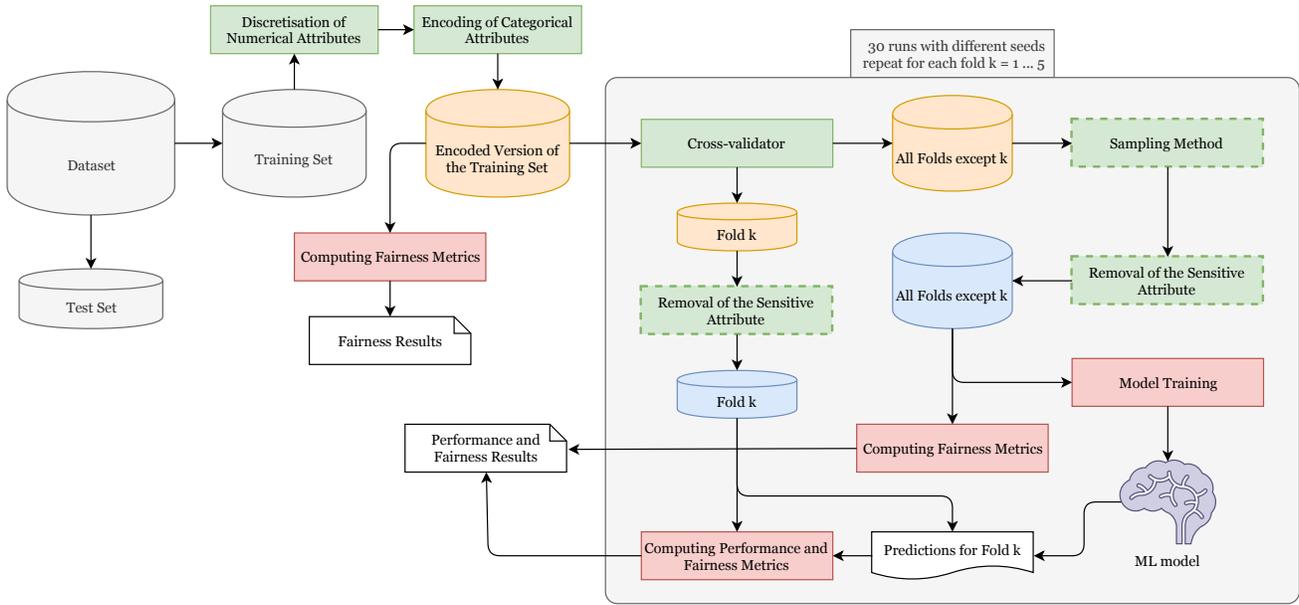


Fig. 2. Approach to assess the impact of data preparation on fairness.

prior to training. These steps are detailed in the remainder of this section. To get a better understanding of the impact of data preparation, we measure fairness at both the data level and the predictions. As depicted in Fig. 2, we only used the training set of both datasets in our experiments.

A. Datasets

We conducted our experiments with two datasets publicly available from the UCI Machine Learning Repository [25].

The *Adult* or *Census Income dataset* contains demographic data extracted from the 1994 US Census Bureau database, with each instance being described by 14 categorical and numerical attributes. There are 48,842 instances in the dataset and a split into training (32,561 instances) and test (16,281 instances) sets is provided. The main task is to predict whether a person earns over 50,000 dollars per year, therefore making a classification into high or low income. In our experiments, we followed previous work and used *sex* as the sensitive attribute with *female* being the unprivileged group.

The *German Credit Data dataset* contains financial information about 1,000 individuals, described by a set of 20 categorical and numerical attributes. The objective is to classify each person into good or bad credit risk. Similar to other studies, we considered *age* to be the sensitive attribute with *young* as the unprivileged group, based on the findings reported in [26]. The attribute *sex* can be derived from *personal-status-sex* of the original dataset. A pre-split into training and test data is not provided for this dataset. Thus, we performed a 70/30 stratified split and tried to maintain the distributions of the true labels and the sensitive attribute on each set.

B. Numerical / Categorical Attributes and Missing Values

Following the approaches of [13] and [19], the numerical attributes were discretised into 4 bins with the boundaries corresponding to those of the interquartile ranges. An additional transformation was performed for the *Adult Income dataset*, with bins which correspond to low frequency counts (less than

50 instances) being pooled together and the attribute values being replaced by the same `POOL` value. This additional transformation was only applied to originally categorical attributes. We refer to this as the integer encoded version of a dataset.

Furthermore, another version of each dataset was created with all features using a one-hot (or 1-of- K) encoding scheme [5] after being discretised, meaning that they are represented by binary dummy variables. We refer to this as the one-hot encoded version of a dataset.

Two exceptions to this approach occurred with *German Credit Data*. The *personal-status-sex* attribute was removed from both versions of the dataset after deriving *sex*. The *age* attribute was discretised into two bins defined by a value greater than or equal to 25, a threshold that was set based on the findings reported by [26].

Contrary to the *German Credit Data dataset*, *Adult Income* contains missing values. For the integer encoded version of this dataset, all instances containing at least one missing value were dropped prior to training or testing the models. As far as the one-hot encoded version of the dataset is concerned, all instances were kept regardless of the presence of missing values. In such cases, a missing value was represented by setting all the corresponding dummy variables to zero.

An overview of the *Adult Income dataset* is shown in Table II, where the data for the integer encoded version is shown between parenthesis. For the one-hot encoded version, the unprivileged group only represents around 33.08% of the dataset. Furthermore, only around 15.04% of the favourable classifications (*high-income*) correspond to *females*. For this version of *Adult Income*, the CVS is 0.1963, the NPI is 4.35×10^{-2} , and the DI is 35.80%.

In the integer encoded version of the dataset, *females* represent around 32.43% of the training data and around 14.81% of the favourable classifications, after removing the missing values. The CVS increases to 0.2002, the NPI suffers a slight increase to 4.36×10^{-2} , and the DI is now 36.22%.

None of the versions of the dataset can be considered fair under the 80% rule. The favourable classifications represent 24.08% and 24.89% of the training data, for the one-hot and integer encoded versions of Adult Income, respectively.

TABLE II
OVERVIEW OF THE ONE-HOT (INTEGER) ENCODED VERSION OF THE ADULT INCOME DATASET.

		Sensitive Attribute	
		Male	Female
True Label	High income	6,662 (6,396)	1,179 (1,112)
	Low income	15,128 (13,984)	9,592 (8,670)

Table III presents an overview of the other dataset. In this case, young individuals are represented by 15.00% of the dataset. The favourable classifications (good credit) represent 70.00% of the training data, with only 12.65% of them being assigned to the unprivileged group. For both versions of the dataset, the CVS is 0.1289, the NPI is 9.47×10^{-3} , and the DI is 82.09% when taking age as the sensitive attribute. According to the 80% rule, the training set of German Credit Data can be considered fair.

TABLE III
OVERVIEW OF THE GERMAN CREDIT DATA DATASET.

		Sensitive Attribute	
		Aged	Young
True Label	Good credit	428	62
	Bad credit	167	43

C. Imbalanced Data and Sampling Methods

The datasets have an imbalance with respect to not only the true labels, but also the sensitive attribute. In such scenarios, it is common to apply sampling methods, namely random undersampling, so as to train the models with an equal number of instances from each class. Besides the typical scenario in which random undersampling is performed w.r.t. the true labels (undersampling-label), we considered two additional configurations: in one it is applied w.r.t. the sensitive attribute (undersampling-protected), and in another it is applied w.r.t. a variable which combines the true labels and the sensitive attribute (undersampling-multivariate).

For each of these settings, we determine which group has fewer instances and keep them, while randomly removing instances from the remaining classes, until their number equals that minimum. After the application of undersampling-multivariate the training data can be considered perfectly fair under the CVS, the NPI, and DI.

We compared these sampling strategies to a baseline scenario in which no sampling method is applied (without-resampling).

D. Learning Algorithms

We performed our experiments with Decision Trees and Random Forests. Bearing in mind that we were dealing with categorical attributes, we looked for implementations of these methods which offered support for such attributes. We opted for the implementations provided by the Python API

of Apache SparkTM. We set the maximum depth of the trees to 30 (maximum supported by these implementations) and used the Gini index as the impurity criterion. Additionally, for the Random Forests, we considered ensembles of 10 trees and the squared root of the total number of features when looking for the best split. Our goal was not to fine tune the parameters, but to understand how the different combinations of data preparation techniques and classifiers impacted the system from a fairness point-of-view. Therefore, we used the default values for most of the remaining parameters.

Since the objective is also to analyse the impact of the removal of the sensitive attribute prior to training a model, we devised four possible scenarios: Decision Tree with and without the sensitive attribute (DT and DTns, respectively), and Random Forest with and without this attribute (RF and RFns, respectively).

E. Model Assessment

We performed five-fold cross-validation with the help of the methods provided by Scikit-learn [8]. In addition to the standard version of cross-validation (normal-cv), the experiments were repeated with stratification (stratified-cv) so as to maintain the class distributions of the original data. Furthermore, each configuration was run with 30 different seeds for the random generators.

We selected fairness metrics which can be applied to the datasets and to the predictions made by the models, so as to be able to compare the unfairness in the predictions to that originally found in the training data. The selected set of metrics includes statistical parity difference (CVS), disparate impact (DI), and the normalised prejudice index (NPI).

The F1-score is more suitable when dealing with imbalanced datasets. However, we also include accuracy in our analysis to facilitate the comparison with previous work.

IV. RESULTS

We first make a relative comparison between all the tested configurations, both in terms of fairness and predictive performance, and then analyse how the fairness in the predictions relates to the fairness in the training data.

A. Analysis of Fairness and Performance

We remind you that for CVS and the NPI the fairer results are closer to zero, while for DI they are closer to one.

Fig. 3 shows the average predictive performance and fairness for Adult Income, when training the models with a stratified-cv. For this dataset, the results with a normal-cv are omitted since the impact on fairness is almost negligible and we reach similar conclusions.

Removing the sensitive attribute prior to training results in models which make fairer predictions under all fairness metrics. Surprisingly, an exception to this behaviour is observed when applying undersampling-multivariate with both Decision Trees and Random Forests. It is also worth mentioning that none of the models can be considered fair according to the 80% rule.

The best sampling method depends on the fairness metric. Using the NPI, undersampling-label

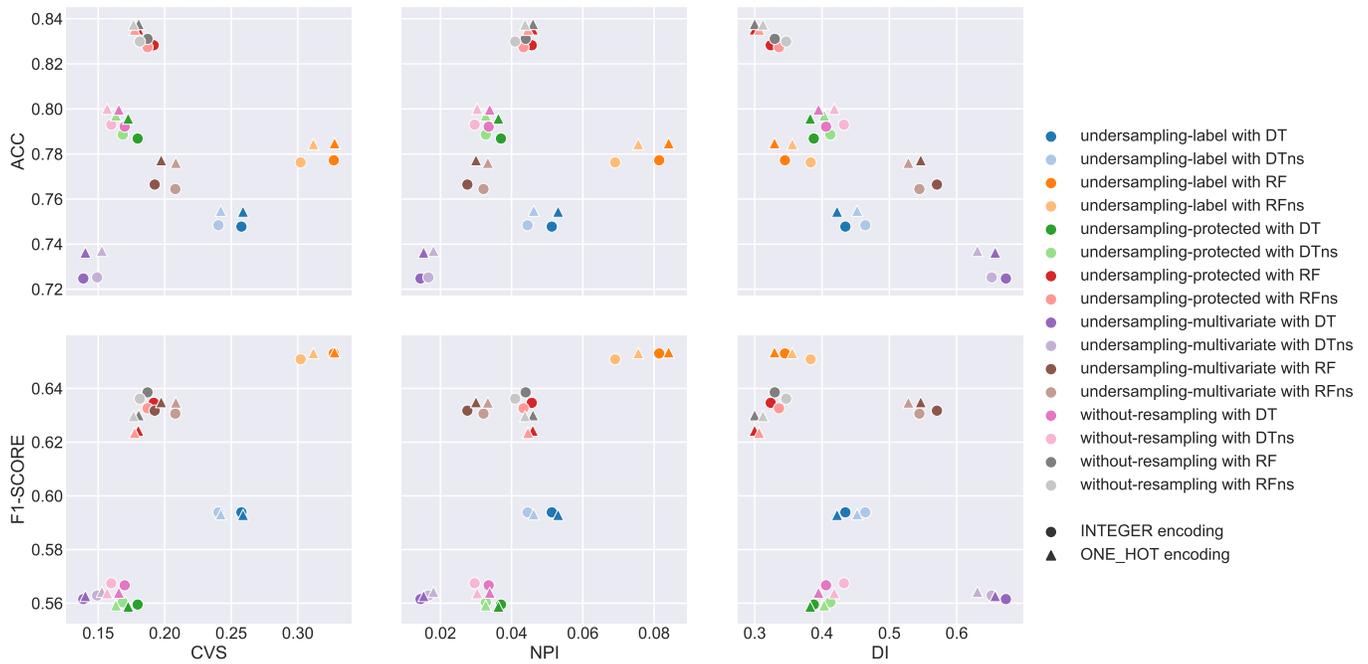


Fig. 3. Average predictive performance and fairness when using stratified folds with the Adult Income dataset. From left to right, fairness is given by CVS, the NPI, and DI. Predictive performance is given by accuracy (ACC) in the top row and by the F1-score in the bottom row.

and undersampling-multivariate seem to be the worst and the best option, respectively. Performing undersampling-protected seems to create models whose predictions are more unfair than without-resampling. Using DI, the sampling methods, from best to worst, are: undersampling-multivariate, undersampling-label, without-resampling, and undersampling-protected. Using CVS, the sampling methods to train Decision Trees, from best to worst, are: undersampling-multivariate, without-resampling, undersampling-protected, and undersampling-label. On the other hand, to train Random Forests, following the same order, we have: without-resampling, undersampling-protected, undersampling-multivariate, and undersampling-label.

The best encoding depends on the fairness metric and the learning algorithm. Using DI, the unfairness in the predictions made by models trained with integer encoded data tends to be lower than in the ones made by models trained with one-hot encoded data. An analysis based on the NPI suggests that all models, except for DT and DTns combined with undersampling-protected or without-resampling, may be able to make fairer predictions if trained with integer encoded. According to CVS, the unfairness in the predictions made by models trained with undersampling-protected and without-resampling in combination with integer encoded data appear to be higher than if those models are trained with one-hot encoded data. The opposite happens with RFns trained with undersampling-label. In fact, most of these differences seem negligible when using the NPI or

CVS, except for RFns with undersampling-label.

Measuring performance with F1-score, the results suggest that the second most important factor after the learning algorithm is the sampling method, with models trained with undersampling-label outperforming models trained with any of the other sampling methods that were tested. The encoding of the categorical attributes also seems to influence the models' performance. In most cases, a one-hot encoding leads to worse F1-scores than integer encoding. The exceptions to this behaviour occur when undersampling-multivariate is applied, regardless of the learning algorithm, and when undersampling-label is applied to the data used to train RF and RFns.

Regarding accuracy, the best options for the sampling method appear to be undersampling-protected or without-resampling, followed by undersampling-label, which seems likely to be a better choice than undersampling-multivariate. The results also suggest that, regardless of sampling strategy, training models with one-hot encoded data is preferable over an integer encoding.

Fig. 4 shows the average F1-score and fairness for German Credit Data. The analysis with accuracy leads to similar conclusions regarding the models' predictive performance. For this reason, these results are not shown here.

Removing the sensitive attribute prior to training seems to have a more pivotal role on improving fairness than the remaining factors under evaluation. The only exceptions regarding this procedure occur with undersampling-multivariate. The most relevant of which is observed when combining a normal-cv with one-hot encoded data. In this case, the predictions of RFns are more unfair than RF.

Moreover, models trained with the one-hot version of the dataset tend to produce fairer predictions than those trained

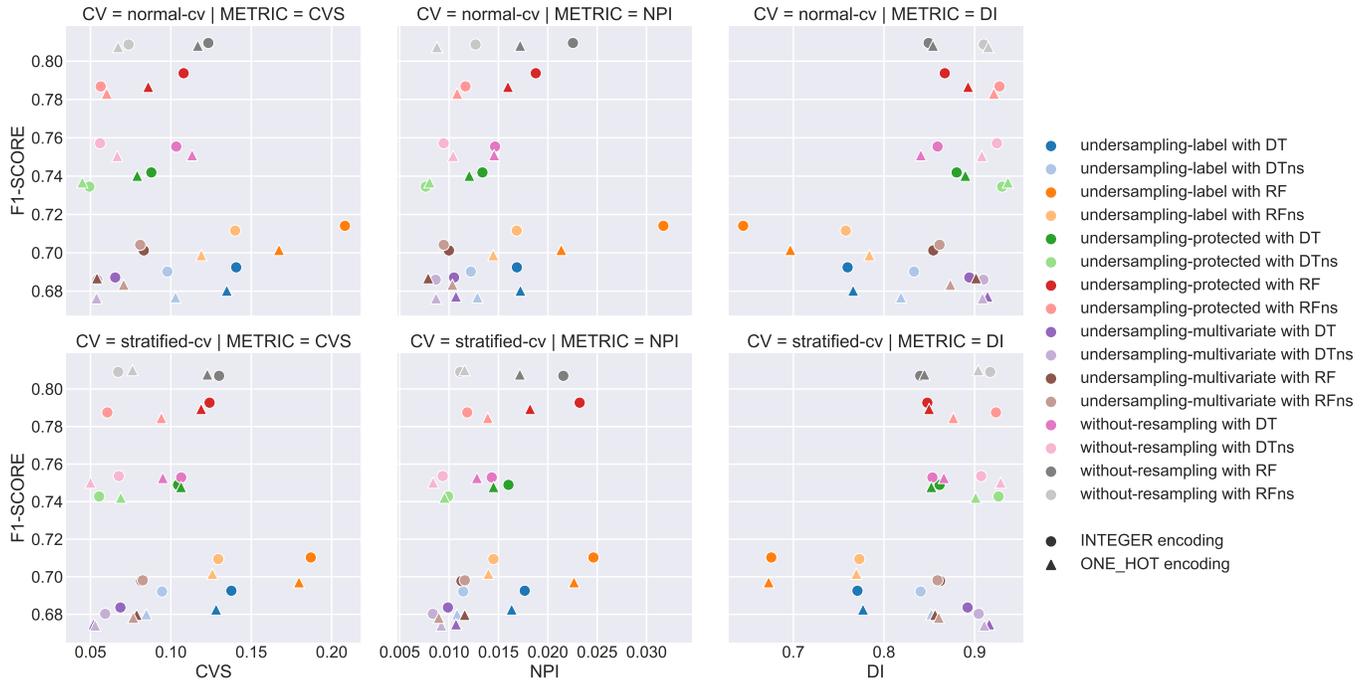


Fig. 4. Average F1-score and fairness of all tested configurations with the German Credit Data dataset. From left to right, fairness is given by CVS, the NPI, and DI. In the top row, models were trained with a `normal-cv` and in the bottom row with a `stratified-cv`.

with integer encoded data. A clear example of this behaviour occurs when combining RF and `undersampling-label` while using a `normal-cv`. Nevertheless, there are exceptions to this behaviour. When using a `stratified-cv` and combining RFns with `undersampling-protected` or `without-resampling`, the integer encoded version of the dataset actually seems to lead to fairer models.

Performing `undersampling-label` with Random Forests seems to be the worst of the tested configurations, with this sampling strategy often being the worst choice from a fairness point-of-view.

As far as predictive performance is concerned, models based on Random Forests tend to be better than models based on Decision Trees. Another interesting observation is the effect on predictive performance that the sampling method seems to introduce: applying no sampling method seems to be the best choice, followed by performing `undersampling-protected`. Performing `undersampling-label` or `undersampling-multivariate` leads to worse results, with a less significant improvement between Decision Trees and Random Forests. Nevertheless, moving from Decision Trees to Random Forests has a greater impact on performance than moving from `undersampling-protected` to `without-resampling`.

B. Fairness Comparison between Data and Predictions

Besides performing our analysis based on the fairness metrics mentioned in III-E, we computed the ratio between the CVS in the predictions and the CVS found in the data subset used to train the models (CVS Ratio), as well as a similar ratio regarding the NPI (NPI Ratio). These ratios give an indication of whether the unfairness in the training data

was increased or reduced under each configuration. A value of 1 indicates that the unfairness in the predictions is the same as in the training data, an absolute value greater than 1 means that the unfairness in the predictions is greater, and an absolute value lower than 1 means that the model makes fairer predictions than the procedure which produced the true labels of the training data. A DI Ratio was not computed since it would be difficult to interpret the results.

Caution must be taken when computing these ratios for models resulting from the application of `undersampling-multivariate`, since the subsets of data used to train these models have a CVS and an NPI equal to zero. In such cases, the value represented in the boxplots corresponds to the CVS or the NPI in the predictions instead of the invalid ratio.

The boxplots in Fig. 5 represent the distributions of the CVS Ratio and the NPI Ratio when a `stratified-cv` was used with Adult Income. Similar results were observed with a `normal-cv` and, for that reason, are not presented here.

The CVS Ratio suggests that performing `undersampling-protected` or not performing random undersampling at all (`without-resampling`) has similar effects on fairness, allowing for the creation of models that tend to reduce the unfairness in the training data. The opposite happens with models trained with `undersampling-multivariate` which always increase it. However, the average CVS of the training data with `undersampling-protected` and `without-resampling` is around 0.1915 and 0.2050, and so, the unfairness in the predictions is approximately the same between models trained with any of the three sampling strategies.

When applying `undersampling-label`, caution must be taken when choosing the learning algorithm, since DT and

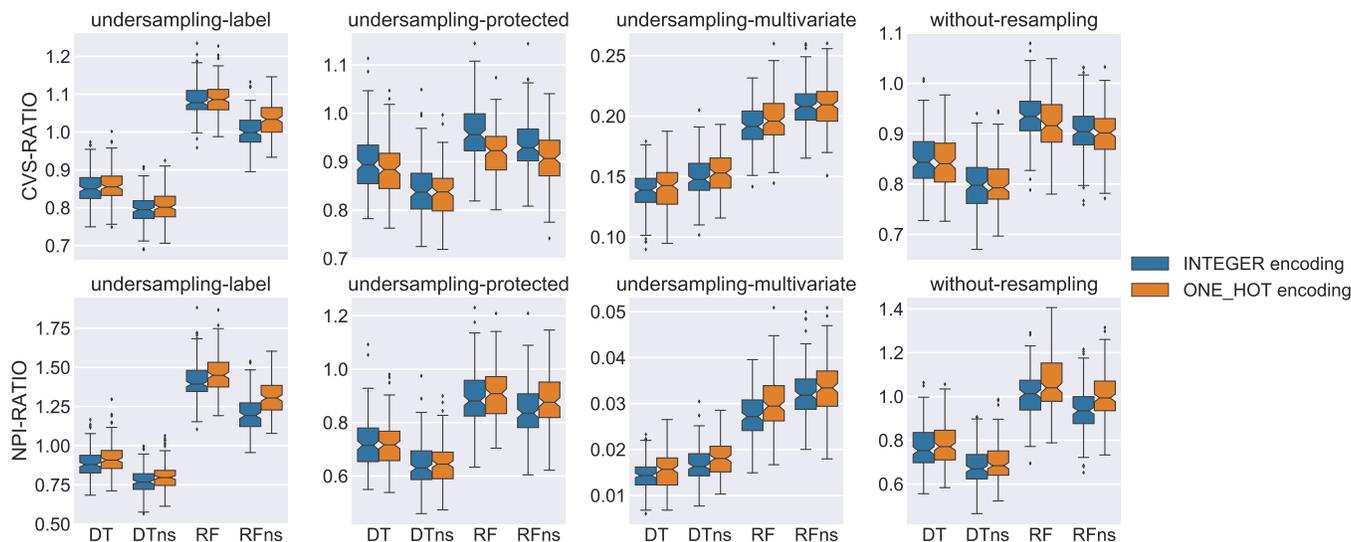


Fig. 5. Distributions of the CVS Ratio (top row) and the NPI Ratio (bottom row) for the Adult Income dataset.

DTns tend to reduce the unfairness of their predictions, while RF and RFns tend to increase it.

The NPI Ratio suggests that the unfairness in the predictions of models trained with `undersampling-protected` is smaller than the one in the data used to train them. Similarly to what was observed with the CVS Ratio, the unfairness found in the predictions of models trained with `undersampling-multivariate` is larger than the one in the training data. The unfairness in the predictions of models trained with `undersampling-protected` and `undersampling-multivariate` is approximately the same, since the NPI of the training data with `undersampling-protected` is, on average, between 0.0485 and 0.0534.

When it comes to `undersampling-label` and `without-resampling`, DT and DTns tend to reduce the unfairness in the data used to train them. On the other hand, the combination of `undersampling-label` with RF and RFns tends to result in models whose predictions increase the unfairness in the training data, while the NPI of the predictions made by RF and RFns trained `without-resampling` tends to be closer to the NPI of the data.

Fig. 6 represents the distributions of the ratios when a `stratified-cv` was used with German Credit Data. Extreme outliers, mainly detected with the NPI Ratio, are not represented to allow for a better visualisation. Unless stated otherwise, similar results were obtained with a `normal-cv`.

Similar to Adult Income, applying `undersampling-multivariate` always results in models whose predictions are more unfair than the data used to train them.

When it comes to `without-resampling`, the results for the CVS Ratio suggest that the predictions of DT and models trained without the sensitive attribute (DTns and RFns) tend to be more fair than the training data. However, the unfairness of the predictions made by RF may be closer to or higher than that of the training data.

The unfairness in the predictions made by models trained with `undersampling-protected` tends to be lower than

that in the data used to train them. However, in the particular case of stratified folds, the unfairness of the predictions made by models trained with the sensitive attribute may be closer to that of the training data.

With `undersampling-label`, RF tend to make predictions with a higher unfairness than that of the training data, while those made by DTns tend to have a lower unfairness. This is also observed with the NPI Ratio. The behaviour of DT and RFns is identical, with their predictions tending to be as unfair as the training data. The encoding seems to have more impact with a `normal-cv` than with a `stratified-cv`: DT and RFns trained with integer encoded data may increase the unfairness in the data, while those trained with one-hot encoded data may reduce it.

The results for the NPI Ratio suggest that DT, DTns, and RFns trained with `undersampling-protected` or `without-resampling` are able to reduce the unfairness in the training data. However, RF trained `without-resampling` tend to increase the unfairness in the training data, while RF trained with `undersampling-protected` tend to reduce it. The only exception, not observed with a `normal-cv`, are RF trained with integer encoded data for which the unfairness in the predictions is similar to that of the data used to train them.

Regarding `undersampling-label`, DT trained with integer encoded data tend to make predictions more unfair than the training data, while if trained with one-hot encoding their predictions tend to be fairer. A similar behaviour is observed when combining RFns with a `normal-cv`. However, when combining RFns with a `stratified-cv`, the unfairness in the predictions tends to be lower than that in the integer encoded data used to train them, but closer to the unfairness in the one-hot encoded data.

The analysis of the distributions of the CVS Ratio and the NPI Ratio leads to some conclusions that are also supported by the results presented in Section IV-A. The unfairness of the predictions made by models trained with `undersampling-`

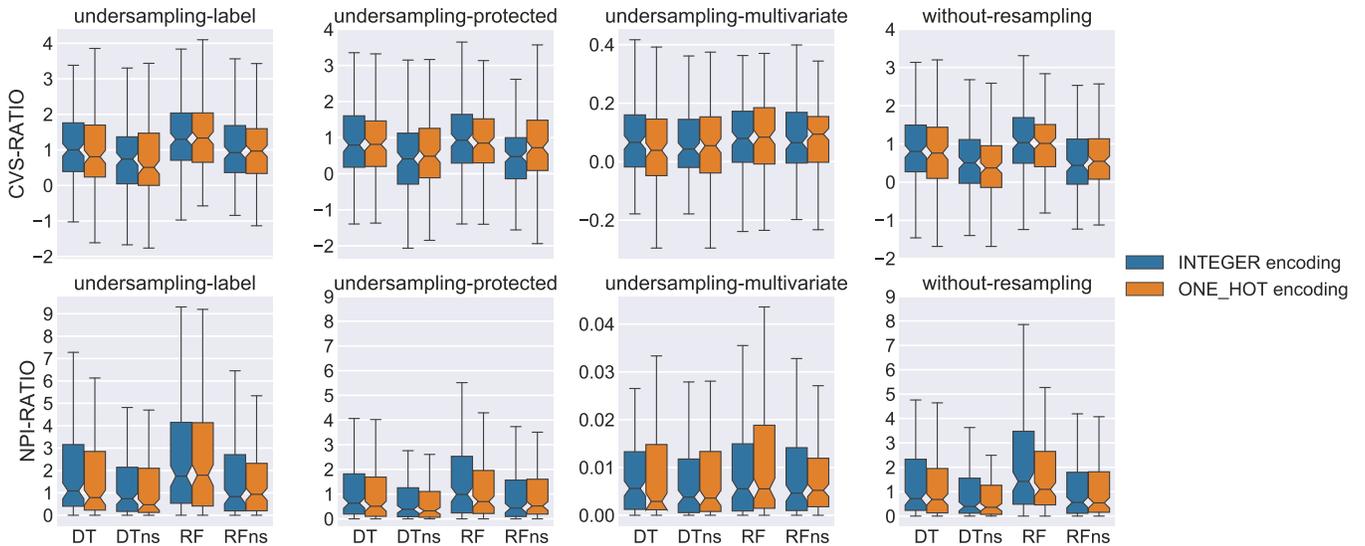


Fig. 6. Distributions of the CVS Ratio (top row) and the NPI Ratio (bottom row) for the German Credit Data dataset.

protected tends to be lower than the unfairness in the data used to train them. These results are similar to those obtained without-resampling, but some configurations may be worse with the latter. Even though undersampling-multivariate increases the unfairness in comparison to the training data, the unfairness in the predictions is similar to that of models trained with undersampling-protected. The worst configuration seems to be the combination of undersampling-label with Random Forests. The discrepancies between the CVS Ratio and the NPI Ratio are more accentuated with the smaller dataset.

V. DISCUSSION AND RECOMMENDATIONS

The **removal of the sensitive attribute**, together with the learning algorithm, **is one of the factors that influence the fairness of an ML model the most**. However, removing this attribute so that the learning algorithms do not have direct access to it does not always lead to models that make fairer predictions. This behaviour, usually exhibited when performing undersampling-multivariate, is somewhat counter-intuitive. When we perform undersampling-multivariate, the training data can be considered fair, at least according to the data level metrics. In addition, we are removing the main source of direct unfairness, i.e. the sensitive attribute. We believe this behaviour is caused by indirect prejudice, due to the presence of other attributes highly associated with the sensitive attribute. This indirect prejudice may become more apparent after the removal of the sensitive attribute, an hypothesis that requires further investigation by looking at the structure of the resulting trees.

The removal of the sensitive attribute seems to have low impact on predictive performance, with the models actually having better performance in some cases. Thus, in general, we can consider that this data preparation procedure does not penalize performance to gain in fairness.

For Adult Income, the choice of the better **encoding of the categorical attributes** depends on the fairness metric and

the learning algorithm, but most of the reported differences seem negligible. According to the NPI and DI, opting for an integer encoding seems to be the safest option for the majority of the configurations that were tested. When using CVS as the fairness metric, the results vary greatly, being difficult to find a more general pattern. Despite the occurrence of some exceptions, an integer encoding of the German Credit Data dataset tends to build models capable of making fairer predictions, independently of the fairness metric. There is no clear answer when it comes to the best encoding in terms of the models' performance.

When performing cross-validation, opting for using **stratification seems to have a minor impact on the models' fairness**. While no relevant differences were found for Adult Income, some minor changes were observed for German Credit Data. This may be due to the dataset's size and the under-representation of some classes. Furthermore, no stratification is made with respect to the sensitive attribute. Contrary to the sampling strategy, whose primary goal is to enhance the models' performance, stratification is meant to maintain the distribution of classes between folds and in comparison to the complete dataset, so as to get more accurate estimates.

For Adult Income, the better sampling strategy in terms of fairness is dependent on the metric used to perform the analysis. The results with CVS and the NPI suggest that undersampling-label is likely to lead to the worst results, while the results with DI indicate that the worst choice is undersampling-protected. For German Credit Data, all metrics suggest that undersampling-label tend to lead to the worst results. Even though some exceptions may occur, **undersampling-multivariate should be performed if one's aim is to build fairer models**. As already mentioned, the training data is completely fair after applying undersampling-multivariate, under the selected metrics. However, the learning algorithm is still able to explore the inherent unfairness of the original dataset, making predictions with some degree of unfairness. This behaviour

suggests the existence of indirect prejudice and shows the limitations of the fairness metrics applied at the data level.

The sampling method is actually one of the factors that affects the predictive performance of the models the most. However, for Adult Income, the best sampling method varies with performance metric. Taking the F1-score as a more suitable metric, `undersampling-label` is definitively the best strategy. As for German Credit Data, performing no random undersampling seems to be the best choice. The behaviour resulting from the application of `undersampling-multivariate` appears to be quite unstable, but its negative impact may be justified by the more drastic reduction in the number of instances used to train the models.

In general, **models based on Decision Trees produce fairer predictions than those based on Random Forests**, which means that model complexity may be a problem for fairness and needs to be further investigated. A more in-depth analysis of the resulting trees could allow for a better understanding of this behaviour, but we believe it may be due to the randomisation introduced by Random Forests during splitting. However, **using Random Forests instead of Decision Trees seems to be the decisive factor to achieve a better predictive performance**, which is expected. These observations highlight the trade-offs an organization may face when deploying a model into production.

An analysis based on the 80% rule, and the consequent decision on whether to consider a model to be fair, is highly dependent on the dataset. We can illustrate this by comparing the results on the two datasets used in our experiments: for Adult Income, no configuration allows for the creation of a fair model, while for German Credit Data most of the tested configurations originated fair models. We would also like to emphasize that when dealing with data imbalance, it is very unlikely to find a dataset with optimal NPI, since this metric is quite sensitive to small changes in class distributions. The sensitivity of this metric is also exacerbated by the presence of extreme outliers, mainly for the smaller dataset.

Based on our observations, we would suggest opting for Decision Trees and for following the standard procedure of removing the sensitive attribute to build fairer models. Even though performing random undersampling w.r.t. both the true labels and the sensitive attribute can lead to satisfactory fairness results, it may have a significant impact on the models' performance. **The best encoding of categorical attributes is data-dependent and different possibilities should be evaluated** instead of choosing an encoding *a priori*. In imbalanced scenarios, stratification is recommended, not only because of it being a good practice, but for its minimal impact on fairness. Combining `undersampling-label` with Random Forests should be avoided since other configurations are likely to offer a better trade-off between predictive performance and fairness.

Caution must be taken with the choice of performance metric, especially when dealing with imbalanced datasets, as could be observed with Adult Income. When fairness concerns are also being considered, one should analyse not only the class imbalance with respect to the true labels, as typically done, but also the disproportion between privileged

and unprivileged groups. One should also bear in mind that a false negative (for instance, some person being incorrectly classified as bad credit risk) is sometimes more costly than a false positive.

Although being aware of the drawbacks of fairness metrics like statistical parity, as these have been widely discussed in the literature [10], [23], we wanted to perform a not so common analysis of fairness that allowed us to compare the unfairness found in the predictions made by an ML model to that found in the data used to train that model. Nevertheless, our experiments have shown the brittleness of these metrics, as even those which were expected to show similar behaviours, such as CVS and DI, sometimes presented contradictory results [16]. This gives further confirmation that there is still room for improvement and progress when it comes to defining new fairness metrics. Specially in imbalanced scenarios, we believe that adopting more recently proposed fairness metrics based on group-conditioned performance [16] might be a small but crucial step towards achieving this goal. There is also room for improvement when it comes to the definition of individual fairness metrics, which seeks to treat similar individuals in a similar way [27]. The challenge here is in finding a suitable measure of the similarity between individuals [27].

VI. CONCLUSION AND FUTURE WORK

Our goal was to assess the potential impact that data preparation techniques commonly used in a machine learning pipeline may have on the fairness of a software system. Rather than focusing on fairness-aware methods, we first tried to understand how standard procedures influence fairness, as these are often used to train models without taking fairness concerns into account. Our findings suggest that the removal of the sensitive attribute and the learning algorithm are the factors that impact the fairness of a system the most. Despite potentially improving the overall performance of a model in imbalanced contexts, random undersampling must be performed with caution, since it may negatively impact the system's fairness. As future work, we plan to analyse the structure of the resulting trees, without neglecting the characteristics of the data used to train the models. We also want to extend the analysis to a wider range of fairness metrics, bearing the possible data imbalance in mind.

ACKNOWLEDGEMENT

This work has been partially supported by the project **ATMOSPHERE** (atmosphere-eubrazil.eu), funded by the Brazilian Ministry of Science, Technology and Innovation (Project 51119 - MCTI/RNP 4th Coordinated Call) and by the European Commission under the Cooperation Programme, Horizon 2020 grant agreement no 777154. It is also partially supported by the project **METRICS**, funded by the Portuguese Foundation for Science and Technology (FCT) – agreement no POCI-01-0145-FEDER-032504.

REFERENCES

- [1] European Parliament, “General Data Protection Regulation,” 2016. [Online]. Available: <http://data.europa.eu/eli/reg/2016/679/oj/eng>
- [2] N. Antunes, L. Balby, F. Figueiredo, N. Lourenco, W. Meira, and W. Santos, “Fairness and Transparency of Machine Learning for Trustworthy Cloud Services,” in *48th Annual IEEE/IFIP International Conference on Dependable Systems and Networks Workshops (DSN-W)*. Luxembourg: IEEE, 2018, pp. 188–193.
- [3] ACM, “Joint Statement on Algorithmic Transparency and Accountability by USACM and EUACM,” May 2017. [Online]. Available: https://www.acm.org/binaries/content/assets/public-policy/2017_joint_statement_algorithms.pdf
- [4] S. Marsland, *Machine Learning: An Algorithmic Perspective, Second Edition*, 2nd ed. Chapman & Hall/CRC, 2014.
- [5] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag, 2006.
- [6] T. Hastie, R. Tibshirani, and J. H. Friedman, *The elements of statistical learning: data mining, inference, and prediction, 2nd Edition*, ser. Springer series in statistics. Springer, 2009. [Online]. Available: <http://www.worldcat.org/oclc/300478243>
- [7] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An Introduction to Statistical Learning: With Applications in R*. Springer Publishing Company, Incorporated, 2014.
- [8] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, “Scikit-learn: Machine learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [9] S. Shalev-Shwartz and S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*. New York, NY, USA: Cambridge University Press, 2014.
- [10] R. Binns, “Fairness in machine learning: Lessons from political philosophy,” in *Conference on Fairness, Accountability and Transparency, FAT 2018, 23-24 February 2018, New York, NY, USA, 2018*, pp. 149–159. [Online]. Available: <http://proceedings.mlr.press/v81/binns18a.html>
- [11] M. B. Zafar, I. Valera, M. Gomez-Rodriguez, and K. P. Gummadi, “Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment,” in *Proceedings of the 26th International Conference on World Wide Web, WWW 2017, Perth, Australia, April 3-7, 2017*, 2017, pp. 1171–1180. [Online]. Available: <https://doi.org/10.1145/3038912.3052660>
- [12] D. Xu, S. Yuan, L. Zhang, and X. Wu, “FairGAN: Fairness-aware Generative Adversarial Networks,” *CoRR*, vol. abs/1805.11202, 2018. [Online]. Available: <http://arxiv.org/abs/1805.11202>
- [13] T. Calders and S. Verwer, “Three naive bayes approaches for discrimination-free classification,” *Data Mining and Knowledge Discovery*, vol. 21, no. 2, pp. 277–292, Sep 2010. [Online]. Available: <https://doi.org/10.1007/s10618-010-0190-x>
- [14] A. Romei and S. Ruggieri, “A multidisciplinary survey on discrimination analysis,” *Knowledge Eng. Review*, vol. 29, no. 5, pp. 582–638, 2014. [Online]. Available: <https://doi.org/10.1017/S0269888913000039>
- [15] M. Feldman, S. A. Friedler, J. Moeller, C. Scheidegger, and S. Venkatasubramanian, “Certifying and removing disparate impact,” in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia, August 10-13, 2015*, 2015, pp. 259–268. [Online]. Available: <https://doi.org/10.1145/2783258.2783311>
- [16] S. A. Friedler, C. Scheidegger, S. Venkatasubramanian, S. Choudhary, E. P. Hamilton, and D. Roth, “A comparative study of fairness-enhancing interventions in machine learning,” in *Proceedings of the Conference on Fairness, Accountability, and Transparency*, ser. FAT* ’19. New York, NY, USA: ACM, 2019, pp. 329–338. [Online]. Available: <http://doi.acm.org/10.1145/3287560.3287589>
- [17] M. B. Zafar, I. Valera, M. Gomez-Rodriguez, and K. P. Gummadi, “Fairness constraints: Mechanisms for fair classification,” in *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, AISTATS 2017, 20-22 April 2017, Fort Lauderdale, FL, USA, 2017*, pp. 962–970. [Online]. Available: <http://proceedings.mlr.press/v54/zafar17a.html>
- [18] “Uniform Guidelines on Employment Selection Procedures, 29 C.F.R. § 1607.4.” [Online]. Available: <https://www.law.cornell.edu/cfr/text/29/1607.4>
- [19] T. Kamishima, S. Akaho, H. Asoh, and J. Sakuma, “Fairness-aware classifier with prejudice remover regularizer,” in *Machine Learning and Knowledge Discovery in Databases*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 35–50.
- [20] F. Kamiran and T. Calders, “Data preprocessing techniques for classification without discrimination,” *Knowl. Inf. Syst.*, vol. 33, no. 1, pp. 1–33, Oct. 2012. [Online]. Available: <https://doi.org/10.1007/s10115-011-0463-8>
- [21] F. Kamiran, T. Calders, and M. Pechenizkiy, “Discrimination aware decision tree learning,” in *Proceedings of the 2010 IEEE International Conference on Data Mining*, ser. ICDM ’10. Washington, DC, USA: IEEE Computer Society, 2010, pp. 869–874. [Online]. Available: <http://dx.doi.org/10.1109/ICDM.2010.50>
- [22] E. Raff, J. Sylvester, and S. Mills, “Fair forests: Regularized tree induction to minimize model bias,” in *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, ser. AIES ’18. New York, NY, USA: ACM, 2018, pp. 243–250. [Online]. Available: <http://doi.acm.org/10.1145/3278721.3278742>
- [23] M. Hardt, E. Price, and N. Srebro, “Equality of opportunity in supervised learning,” in *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain, 2016*, pp. 3315–3323. [Online]. Available: <http://papers.nips.cc/paper/6374-equality-of-opportunity-in-supervised-learning>
- [24] T. Speicher, H. Heidari, N. Grgic-Hlaca, K. P. Gummadi, A. Singla, A. Weller, and M. B. Zafar, “A unified approach to quantifying algorithmic unfairness: Measuring individual & group unfairness via inequality indices,” in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, ser. KDD ’18. New York, NY, USA: ACM, 2018, pp. 2239–2248. [Online]. Available: <http://doi.acm.org/10.1145/3219819.3220046>
- [25] D. Dua and C. Graff, “UCI machine learning repository,” 2019. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [26] F. Kamiran and T. Calders, “Classifying without discriminating,” in *2009 2nd International Conference on Computer, Control and Communication*, Feb 2009, pp. 1–6.
- [27] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. S. Zemel, “Fairness through awareness,” in *Innovations in Theoretical Computer Science 2012, Cambridge, MA, USA, January 8-10, 2012*, 2012, pp. 214–226. [Online]. Available: <https://doi.org/10.1145/2090236.2090255>