



UNIVERSITY OF LEEDS

This is a repository copy of *Explaining contributions of features towards unfairness in classifiers: A novel threshold-dependent Shapley value-based approach*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/217962/>

Version: Accepted Version

Article:

Pelegrina, G.D., Siraj, S. orcid.org/0000-0002-7962-9930, Duarte, L.T. et al. (1 more author) (2024) Explaining contributions of features towards unfairness in classifiers: A novel threshold-dependent Shapley value-based approach. *Engineering Applications of Artificial Intelligence*, 138 (Part B). 109427. ISSN 0952-1976

<https://doi.org/10.1016/j.engappai.2024.109427>

© 2024, Elsevier. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0/>. This is an author produced version of an article published in *Engineering Applications of Artificial Intelligence*. Uploaded in accordance with the publisher's self-archiving policy.

Reuse

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) licence. This licence only allows you to download this work and share it with others as long as you credit the authors, but you can't change the article in any way or use it commercially. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

Explaining contributions of features towards unfairness in classifiers: A novel threshold-dependent Shapley value-based approach

Guilherme Dean Pelegrina, Sajid Siraj, Leonardo Tomazeli Duarte, Michel Grabisch

^a*Mackenzie Presbyterian University (UPM), 930 Consolação St, São Paulo, 01302-907, São Paulo, Brazil*

^b*Centre for Decision Research, Leeds University Business School, Leeds, LS6 1AN, UK*

^c*School of Applied Sciences - University of Campinas (UNICAMP), 1300 Pedro Zaccaria St, Limeira, 13484-350, São Paulo, Brazil*

^d*Centre d'Économie de la Sorbonne - Université Paris I Panthéon-Sorbonne, 106-112 bd de l'Hôpital, Paris, 75013, France*

Abstract

A number of approaches has been proposed to investigate and mitigate unfairness in machine learning algorithms. However, as the definition and understanding of fairness may vary in different situations, the study of ethical disparities remains an open area of research. Besides the importance of analyzing ethical disparities, explainability in machine learning is also a relevant issue in Trustworthy Artificial Intelligence. Usually, both fairness and explainability analysis are based on a fixed decision threshold, which differentiates the positive cases from the negative ones according to the predicted probabilities. In this paper, we investigate how changes in this threshold can impact the fairness of predictions between protected and other groups and how features contribute towards such a measure. We propose a novel Shapley value-based approach as a tool to investigate how changes in the threshold values change the contribution of each feature towards unfairness. This gives

us an ability to evaluate how fairness measures vary for different threshold values and which features have the higher (or lower) impact on creating ethical disparities. We demonstrate this using three different case studies that are carefully chosen to highlight different unfairness scenarios and features contributions. We also applied our proposal as a feature selection strategy, which contributed to decrease unfair results substantially.

Keywords:

Interpretable Machine Learning, Shapley value, Fairness, Feature contribution

1. Introduction

It is often argued that decision making by people contain several psychological biases [1, 2] and that may be avoided if we rely on machine-based intelligence. Based on the level of human involvement, the authors in [3] explained and structured operations into different types. On one hand, they define the fully automated tasks as “out of the loop” meaning that humans are not involved, while on the other hand, the “on the loop” tasks were defined as those where decisions are solely depending on humans. Although it appears that excluding humans from decision process should remove the issues and biases related to humans, it is not that simple as machines are often trained with the data that contain biases in itself [4]. The issue of unfairness in Machine Learning (ML) algorithms has been widely debated in recent years and a number of approaches has been proposed to mitigate unfairness [5, 6]. One may cite, for instance, works that exploit the concepts of equal opportunity [7] and accuracy parity [8]. However, as the concept

16 and definition of fairness may vary in different context, this remains an open
17 area of research. Regardless of its definition, fairness is closely related to the
18 explainability of algorithms. This is because to investigate fairness, there
19 should be a way to explain the results and performance of machine learning
20 algorithms.

21 One can define the term explainability as the ability of an algorithm to
22 provide information that can help evaluate causalities, similarities, and/or
23 uncertainties which in turn helps decision makers towards understanding the
24 model outcomes [9, 10]. Miller [11] described this field as the intersection
25 of ML with the subject areas of human-computer interaction and social sci-
26 ences. Some of the explainability tools focus on improving the interpretability
27 of data, for example, by reducing the dimension of feature space [12]. Others
28 focus on the machine learning models, for example, by explaining the out-
29 comes generated by ML algorithms [13]. Among such approaches, one may
30 cite the well-known method proposed by [14], called SHAP, which has been
31 addressed in several recent works [15, 16, 17].

32 The term explainable AI, or XAI in short, can be used at two very dif-
33 ferent levels: the global and the local ones [13, 10]. The global level inter-
34 pretations are important to assess and audit the explainability of algorithms
35 overall, for example, which feature contributes more towards explaining the
36 outcomes [18]. The local level interpretations, however, are more suited to
37 explain results (or predictions) provided for local samples [14], for example,
38 why a loan has been refused to a certain applicant. While both levels of
39 explainability can be useful in different situations, a single instance of local
40 explainability might not be sufficient for evaluating fairness. The notion of

group fairness demands investigation of multiple cases in order to compare the algorithm’s performance for different groups (for example, male group versus female group). Although some of the explainability strategies are only applicable to specific algorithms [19], our focus remains on model-agnostic algorithms that can be applied to different machine learning techniques equally. In this context, the use of Shapley value-based approaches has gained more attention due to its algorithm-agnostic characteristic [14], versatility [20, 21] and other useful properties [22, 23].

The issue of fairness in machine learning has also gained attention in recent years [24, 25]. For example, this has been proposed as a constrained optimization problem where the objective is to minimize the mis-classification probability whilst imposing an upper bound on an unfairness measure [26]. Fair solutions have also been exploited by means of multi-objective approaches [27, 28]. In this case, both algorithm performance and fairness concerns are optimized simultaneously. Although mitigation of unfairness in ML algorithms has been investigated, there is a little done in the field of explaining the contribution of features towards unfair results generated by these algorithms. Usually explainability for ML classifiers is done considering a single (pre-assigned/fixed) value of threshold that is used to differentiate the positive cases from the negative ones. Changing this threshold might not only impact the performance of classifiers but can also affect fairness for a protected group against the other groups. [15] proposed the use of Shapley values to investigate fairness in machine learning models and to explain the trade-off between accuracy and fairness. However, they did not investigate thresholds used for classification and their impact on fairness, whether using equalized

odds [7] or any other method of quantifying fairness.

A recent work [29] involved generating various different decision thresholds to estimate the contributions of features towards the quality of predictions, measured in terms of Receiver-operating Characteristics (ROC) curve and the Area under the ROC curve (AUC). Although the explanation of these curves can help analysts in feature engineering, it may also help investigate the issues of fairness, for example, by explaining the contribution of features towards disparities in quality of predictions between the protected and the privileged groups. Therefore, we see a gap in the literature in providing a mechanism to evaluate features contributions toward unfair results along with thresholds¹. Instead of an analysis for a single (and predefined) threshold, it is of interest to verify how fairness and features contributions varies for different threshold values. Aiming at overcoming this gap, in this paper, we investigate the range of decision thresholds to identify different levels of fairness along with the performance of classifiers. Our proposal is based on the Shapley value, which indicates features contributions towards both performance (for protected and privileged groups) and fairness measures. From the experimental results, we attest that our proposal can be useful to (i) evaluate which features impact disparate results, (ii) investigate the presence of features acting as proxies and (iii) observe how features contribute differently towards different sensitive groups. Moreover, we also apply our proposal as

¹Indeed, most of the explainable AI techniques, such as LIME [30] or SHAP [14], deal with local interpretability, i.e., with the purpose of explaining the outcome of a sample of interest. Although SHAP can be extended to global interpretability, it provides the contributions toward the predicted classes, instead of group fairness measures along with thresholds.

87 a tool to assist feature selection. Indeed, by removing features with relevant
 88 contributions toward unfairness, we are able to mitigate ethical disparities.

89 The next section describes the background with related literature, fol-
 90 lowed by the proposed technique (Section 3). The numerical experiments are
 91 presented in Section 4, and then, Section 5 concludes the discussion.

92 2. Background

93 This section discusses the theoretical background used in our proposal.
 94 Firstly, we present the adopted notations and, thereafter, we discuss some
 95 performance and fairness metrics frequently used to evaluate machine learn-
 96 ing models. Finally, we define the Shapley value and how it has been used
 97 as a feature attribution method for ML explainability.

98 2.1. Setup and Notations

99 Assume a binary classification setting where $\mathbf{X} \in \mathbb{R}^{n \times m}$ represents the
 100 m -dimensional dataset with n samples and $\mathbf{y} \in \{0, 1\}^n$ is the associated
 101 vector of labels. Generally, class 1 (the “positive” class) indicates a benefit.
 102 Consider $\mathbf{d} \in [0, 1]^n$ as the vector of predicted probabilities which indicate
 103 the likelihood of belonging to class 1. Given an instance $x^{(i)}$ and a predefined
 104 threshold t , $x^{(i)}$ is classified as class 1 if $d^{(i)} \geq t$. The vector of all predicted
 105 labels is represented by $\hat{\mathbf{y}}$, i.e., for a given instance $x^{(i)}$, if $d^{(i)} \geq t$, $\hat{y}^{(i)} = 1$.

106 Suppose a classification task whose samples can be split into sensitive
 107 groups defined by $G \in \{a, b\}$ (e.g., blacks and whites or men and women).
 108 Without loss of generality, consider that a and b are the protected and the
 109 privileged groups, respectively. By splitting \mathbf{X} into these groups, one may
 110 define $\mathbf{X} = [\mathbf{X}_a; \mathbf{X}_b]$, where $\mathbf{X}_a \in \mathbb{R}^{n_a \times m}$ and $\mathbf{X}_b \in \mathbb{R}^{n_b \times m}$ are m -dimensional

111 datasets with n_a and n_b samples, respectively. The vector of labels, predicted
112 probabilities and predicted labels can also be split similarly.

113 2.2. Performance metrics

114 In order to evaluate the performance of a classifier, one generally considers
115 metrics used to construct the confusion matrix. They are defined as follows:

- 116 • True positive (TP): Number of instances correctly classified as class 1.
- 117 • True negative (TN): Number of instances correctly classified as class
118 0.
- 119 • False positive (FP): Number of instances wrongly classified as class 1.
- 120 • False negative (FN): Number of instances wrongly classified as class 0.

121 Based on these metrics, the classifier's performance is evaluated by means
122 of the following rates:

- 123 • Positive predictive value ($PPV = \frac{TP}{TP+FP}$): ratio between the number
124 of class 1 correctly classified as class 1 and the total number of instances
125 classified as class 1. This ratio is also called Precision.
- 126 • Negative predictive value ($NPV = \frac{TN}{TN+FN}$): ratio between the number
127 of class 0 correctly classified as class 0 and the total number of instances
128 classified as class 0.
- 129 • True positive rate ($TPR = \frac{TP}{TP+FN}$): percentage of class 1 correctly
130 classified as class 1. It is also called Sensitivity or Recall.

- False positive rate ($FPR = \frac{FP}{FP+TN}$): percentage of class 0 wrongly classified as class 1.

All the aforementioned metrics can also be defined by conditioning on the sensitive groups. For example, TP_a means the number of instances in group a correctly classified as class 1 and FPR_b is the percentage of class 0 wrongly classified as class 1 for the individuals belonging to group b .

2.3. Measures of fairness

Besides the performance metrics presented in the previous subsection, machine learning models have been evaluated in terms of disparate results against protected groups. We present in the sequel a brief description of some fairness metrics used in our analysis (see [31, 24] for further details and other metrics).

2.3.1. Statistical parity (or Demographic parity)

A classifier satisfies statistical parity (SP) [32] if both groups have the same probability of being classified as the positive class, i.e., $\frac{TP_G+FP_G}{n_G}$ must be the same regardless group G . Therefore, the SP strategy is associated with the minimization of the following cost function:

$$f_{SP} = \left| \frac{TP_a + FP_a}{n_a} - \frac{TP_b + FP_b}{n_b} \right|. \quad (1)$$

The idea supported by this definition is that equivalent positive outcomes (such as receiving a credit) should be similar for individuals, regardless if its group.

151 *2.3.2. Predictive equality*

152 A classifier satisfies predictive equality (*PE*) if both groups have equal
 153 *FPR*. In other words, $P(\hat{Y} = 1 \| Y = 0, G = 1) = P(\hat{Y} = 1 \| Y = 0, G = 0)$.
 154 One here attempts to minimize

$$f_{PE} = \left| \frac{FP_1}{FP_1 + TN_1} - \frac{FP_0}{FP_0 + TN_0} \right|. \quad (2)$$

155 In this metric, the aim is to assign the same amount of positive outcomes to
 156 individuals of different groups that actually belong to the negative class.

157 *2.3.3. Equalized odds*

158 A classifier satisfies equalized odds (*EO*) if both groups have equal *TPR*
 159 and *FPR*. In other words, $P(\hat{Y} = 1 \| Y = 1, G = 1) = P(\hat{Y} = 1 \| Y = 1, G =$
 160 $0)$ and $P(\hat{Y} = 1 \| Y = 0, G = 1) = P(\hat{Y} = 1 \| Y = 0, G = 0)$. The definition
 161 of equalized odds was proposed by [7] to remedy previously noted flaws with
 162 demographic parity [32]. The goal is to minimize

$$f_{EO} = \left| \frac{TP_1}{TP_1 + FN_1} - \frac{TP_0}{TP_0 + FN_0} \right| + \left| \frac{FP_1}{FP_1 + TN_1} - \frac{FP_0}{FP_0 + TN_0} \right|. \quad (3)$$

163 *2.4. Shapley values in machine learning interpretability*

164 The Shapley value [33] is a classical solution concept in game theory.
 165 Consider a scenario in which a set $M = \{1, 2, \dots, m\}$ of m players join a
 166 coalition in order to achieve a common goal. For example, energy storage
 167 systems owners could join coalitions in order to save individual costs [34].
 168 In such a scenario, the Shapley value ϕ_j associated with each player j will
 169 indicate how much he/she should receive when sharing the whole benefit

170 achieved by the coalition of all players. Mathematically, it is defined as
 171 follows:

$$\phi_j = \sum_{A \subseteq M \setminus \{j\}} \frac{(m - |A| - 1)! |A|!}{m!} [v(A \cup \{j\}) - v(A)], \quad (4)$$

172 where $|A|$ indicates the cardinality of subset A and $v(A)$ is the game payoff
 173 (or benefit) when only players in A join the coalition. For the coalition of all
 174 players, $v(M)$ indicates the total benefit.

175 Among the several properties satisfied by the Shapley value (see [35] for
 176 further details), one is of interest in machine learning interpretability: effi-
 177 ciency. This property states that the payoff of the grand coalition $v(M)$ can
 178 be decomposed into the sum of the individuals Shapley values. Mathemati-
 179 cally, efficiency means $\sum_{j=1}^m \phi_j = v(M) - v(\emptyset) = v(M)$ (in game theory, one
 180 generally assumes $v(\emptyset) = 0$). In machine learning, this property allows us to
 181 explain the contributions of features from a predefined baseline ($v(\emptyset)$) to an
 182 achieved goal ($v(M)$). Therefore, how to defined both baseline and goal is
 183 an important task when adopting the Shapley value for ML interpretability
 184 (see [29] for further details).

185 One associates both baseline and goal in machine learning to the game
 186 payoff $v(\cdot)$ in game theory. Clearly, the definition of both elements depends
 187 on what one would like to explain. For instance, [36] used the Shapley value to
 188 evaluate the contribution of features towards the coefficient of determination
 189 in linear regression models. In this case, the payoff $v(A)$ is the coefficient of
 190 determination when features in A are available in the linear regression model.
 191 When no feature is available, $v(\emptyset) = 0$, and when all feature is available,

192 the coefficient of determination $v(M)$ is maximal. In the well-know SHAP
 193 method [14], the authors adopted the Shapley value as a model-agnostic fea-
 194 ture attribution approach to explain local predictions. In this case, given
 195 a sample of interest x^* , one may interpret how much the associated char-
 196 acteristics (features values $x_1^*, x_2^*, \dots, x_m^*$) contributes towards the obtained
 197 prediction or classification. Then, $v(\emptyset)$ represents the expected prediction
 198 when all features values in x^* are unknown and $v(M)$ is the actual predic-
 199 tion assigned to x^* .

200 Most of the work in explainable machine learning uses the Shapley value
 201 to interpret features contributions toward performance metrics such as accu-
 202 racy. However, only few works attempt to explain such contributions towards
 203 unfair results. An example is the work conducted by Begley et al. [15], where
 204 the authors aggregate local Shapley values in order to globally explain both
 205 performance and fairness measures. Moreover, most of the explainable ap-
 206 proaches only focus on results achieved by assuming a single threshold (and,
 207 therefore, performance or fairness metrics in a single scenario). In this pa-
 208 per, we borrowed the idea proposed by Pelegina and Siraj [29] to explain the
 209 overall performance of classifiers based on ROC curve and on the area un-
 210 der the ROC curve. As for such analysis where the authors investigated the
 211 classifier performance along with thresholds, our proposal in this work is to
 212 evaluate the impact of each feature in biased results for different predefined
 213 thresholds. We detail our proposal in the next section.

214 3. Explaining unfair results through Shapley values

215 In this section, we discuss the use of the Shapley values in our proposal
216 to assign contributions of features towards both performance and fairness
217 measures.

218 3.1. Features contributions towards performance measures

219 Generally, fairness implies equal performance measures for the considered
220 groups of individuals. For instance, to ensure the fairness of a model in
221 terms of predictive equality, it is imperative that both groups of individuals
222 exhibit equal false positive rates. Nevertheless, disparate model performances
223 among diverse groups of individuals are a common occurrence in numerous
224 applications. In this context, a crucial aspect involves interpreting how each
225 feature contributes to the observed unfair result. Such an interpretation is
226 vital for the redesign of the machine learning model, enabling the mitigation
227 of inherent inequalities.

228 In Section 1, we highlighted some existing approaches that are used to
229 evaluate features contributions towards performance measures. Typically,
230 these approaches calculate the impact of features on model performance us-
231 ing a single predefined threshold. However, in this paper, as delineated in [29],
232 we extend our analysis to incorporate contributions across varying thresh-
233 olds. Assume $p_{t,G}(A)$ as the model performance for group G and predefined
234 threshold t when a set of features, expressed by the set A , is considered in
235 the training step. We define the payoff as follows:

$$v_{t,G}^{PERF}(A) = p_{t,G}(A) - p_{t,G}(\emptyset), \quad (5)$$

236 where $p_{t,G}(\emptyset)$ represents the performance of a random classifier (i.e., when
 237 no features is available in training step). Note that, in accordance with
 238 the definition of a game, $v_{t,G}^{PERF}(\emptyset) = p_{t,G}^{PERF}(\emptyset) - p_{t,G}^{PERF}(\emptyset) = 0$. Based on
 239 $v_{t,G}^{PERF}(A)$ for all $A \in \mathcal{P}(M)$, where $\mathcal{P}(M)$ is the power set of M , it is possible
 240 to calculate the Shapley value of feature j as

$$\phi_j^{PERF,t,G} = \sum_{A \subseteq M \setminus \{j\}} \frac{(m - |A| - 1)! |A|!}{m!} [v_{t,G}^{PERF}(A \cup \{j\}) - v_{t,G}^{PERF}(A)] . \quad (6)$$

241 Due to the efficiency property, $\sum_{j=1}^m \phi_j^{PERF,t,G} = v_{t,G}^{PERF}(M) - v_{t,G}^{PERF}(\emptyset) =$
 242 $v_{t,G}^{PERF}(M)$, i.e., the model performance for group G when all features are
 243 available ($v_t(M)$) can be decomposed by the sum of the marginal contri-
 244 butions of each feature. Note that, depending on the adopted performance
 245 measure, $\phi_j^{PERF,t,G}$ may be either positive or negative. For instance, if we
 246 are evaluating overall accuracy and by assuming that the inclusion of a fea-
 247 ture into the model training would not decrease its performance, we would
 248 expect $0 \leq \phi_j^{PERF,t,G} \leq v_t(M)$. In this scenario, we have clear bounds, as the
 249 game is non-decreasing and normalized by the accuracy based on all features
 250 (see [37] for more details about such a game, usually called capacity). How-
 251 ever, if we are looking at the area under the Precision-Recall curve (see [29]
 252 for further details), $\phi_j^{PERF,t,G}$ could be negative and, therefore, it is difficult
 253 to define bound conditions for such a value.

254 In a biased scenario, the performances $v_{t,G}^{PERF}(M)$ as well as the features
 255 contributions $\phi_j^{PERF,t,G}$ for $G = \{a, b\}$ may be different and, therefore, we
 256 may interpret which features are creating disparate results. Figure 1 illus-
 257 trates the process of interpreting the disparate results by comparing the

258 performance measures of different groups of individuals. In summary, the
 259 process involves calculating performance measures and features contributions
 260 for all coalitions of features $A \in \mathcal{P}(M)$, thresholds t_0, t_1, \dots, t_s and groups
 261 a and b . Let us recall that groups a and b as the protected and the privi-
 262 leged groups, respectively. As in this process one calculates $p_{t,a}^{PERF}(M)$ and
 263 $p_{t,b}^{PERF}(M)$, one can compare the performance measure for both groups and
 264 visualize disparities along with thresholds. One illustrates this result on the
 265 center right plot of Figure 1, which assumes the true positive rate as the per-
 266 formance measure². Moreover, in the bottom right plot in Figure 1, one also
 267 illustrates a comparison between features contributions towards the TPR for
 268 groups a and b along with thresholds. Note that, for a fixed threshold t , the
 269 sum of the Shapley values is equal to the difference between the actual TPR
 270 (by using all features) and the random classifier performance which, in this
 271 case, is given by $1 - t$. It is worth mentioning that the shaded area in both
 272 figures indicate the standard deviation from the averaged value by taking
 273 into account the k -fold cross-validation strategy.

274 So far, we discussed how to use the Shapley values in order to compare
 275 performance measures for different groups of individuals and interpret dis-
 276 parities in features contributions. However, unfair results in machine learning
 277 are frequently evaluated by means of fairness measures. Therefore, one may
 278 also interpret features contributions directly on fairness measures. We elab-
 279 orate this in the sequel.

²It is worth highlighting that the overall TPRs are close to the TPRs from the privileged group due to an imbalance in the dataset used to create this illustrating example. Indeed, there are more instances from the privileged group in comparison with the protected group.

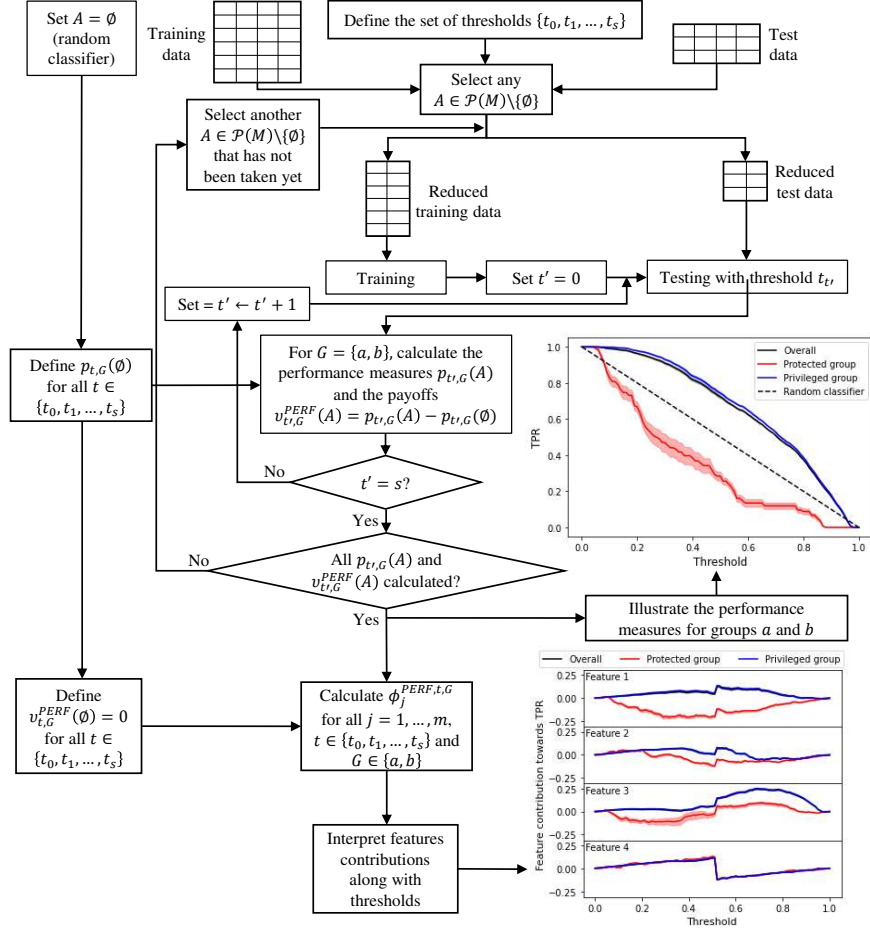


Figure 1: The proposed scheme to evaluate features contributions towards disparities in performance measures.

3.2. Features contributions towards fairness measures

Besides evaluating features contributions towards the model performance, one may also interpret their impacts on fairness measures. Assume $f_{t,p}(A)$ as the fairness measure associated with performance measures $p_{t,a}^{PERF}(A)$ and $p_{t,b}^{PERF}(A)$, for a predefined threshold t when only features in A are available

285 in training step. The payoff is defined by

$$v_t^{FAIR}(A) = f_{t,p}(A) - f_{t,p}(\emptyset), \quad (7)$$

286 where $f_{t,p}(\emptyset)$ represents the fairness measure of a random classifier³. One
 287 may also note that $v_t^{FAIR}(\emptyset) = f_{t,p}^{FAIR}(\emptyset) - f_{t,p}^{FAIR}(\emptyset) = 0$, i.e., in accordance
 288 with the definition of a game. Based on $v_t^{FAIR}(A)$ for all $A \in \mathcal{P}(M)$, we may
 289 calculate the Shapley value of feature j as

$$\phi_j^{FAIR,t} = \sum_{A \subseteq M \setminus \{j\}} \frac{(m - |A| - 1)! |A|!}{m!} [v_t^{FAIR}(A \cup \{j\}) - v_t^{FAIR}(A)]. \quad (8)$$

290 The interpretation in the considered fairness measure is given by the sum of
 291 the marginal contribution of each feature. Indeed, given the efficient property
 292 of Shapley values, $\sum_{j=1}^m \phi_j^{FAIR,t} = v_t^{FAIR}(M) - v_t^{FAIR}(\emptyset) = v_t^{FAIR}(M) =$
 293 $f_{t,p}(M) - f_{t,p}(\emptyset)$. In the fairest scenario, we normally expects $f_{t,p}(M) = f_{t,p}(\emptyset)$
 294 and therefore, $v_t^{FAIR}(M) = 0$. However, if $v_t^{FAIR}(M) > 0$, there is a disparity
 295 between groups a and b . The marginal contributions of features given by
 296 $\phi_j^{FAIR,t}$ will then highlight which features are creating disparate results and
 297 can be seen as a source of bias. As in the case of the performance measure,
 298 $\phi_j^{FAIR,t}$ can be either positive (feature j entail disparate results) or negative
 299 (feature j contributes to improve fairness).

300 We illustrate in Figure 2 the steps to interpret the contributions of fea-

³Although several fairness measure should be zero when considering a random classifier, we decided to keep $f_{t,p}(\emptyset)$ in our proposal in order to generalize the idea for any fairness measure.

301 tures on fairness measures. As in the Subsection 3.1⁴, the process here also
 302 involves calculations on all coalitions of features $A \in \mathcal{P}(M)$ and thresholds
 303 t_0, t_1, \dots, t_s . As a first interpretation, one may visualize the fairness measure
 304 along with the thresholds. This result is presented on the center right plot of
 305 Figure 2, which assumes the statistical parity as the fairness measure. The
 306 fairness measure along with the thresholds can be decomposed on individual
 307 contributions of features. The bottom right plots in Figure 2 illustrates the
 308 marginal contribution of each feature towards the fairness measure. Note
 309 that the sum of these marginal contributions is equal to the curve presented
 310 in the center right plot of Figure 2. Moreover, if there is an interest in analyz-
 311 ing the Shapley values for a single threshold t , one may use a waterfall plot
 312 as illustrated in the bottom left plot of Figure 2. As in Figure 1, the shaded
 313 area in both center and bottom right plots and the whiskers in the waterfall
 314 plot indicate the standard deviation from the averaged value by taking into
 315 account the k -fold cross-validation strategy.

316 4. Experiments

317 This section outlines the experiments conducted to validate the efficacy of
 318 our proposed method in interpreting feature contributions toward disparate
 319 outcomes in machine learning. While our experiments utilized the Random
 320 Forest classifier [38] with 5-fold cross-validation, it is essential to note that our
 321 approach is model-agnostic, as discussed in Section 3. Consequently, other

⁴Clearly, in a computational point of view and in order to avoid double calculating the performance measures, the schemes presented in Figures 1 and 2 could be merged into a single one. However, aiming at providing an easier visualization of both proposals, prefer to split both of them into two processes.

322 classifier can be employed. The explanations derived from our approach
323 reveal the contributions of features towards the unfair outcomes produced
324 by the trained classifier.

325 In order to evaluate our proposal, we examined three real datasets: COM-
326 PAS [39], Law School Admission Council (LSAC) [40] and Adult income [41].
327 For each dataset, we computed disparities related to a sensitive feature and
328 interpreted the individual contributions of each feature towards such dispar-
329 ities. For this purpose, we only consider the sensitive feature to split the
330 dataset into two groups when calculating the performance/fairness measures
331 (i.e., we do not use such a feature in training step). The remaining features
332 will, then, explain the achieved disparity for each threshold. We also con-
333 ducted a preprocessing step in all datasets in order to ensure that the two
334 classes are balanced. In this case, we used a re-sampling strategy that ran-
335 domly eliminates samples from the over-represented class until balancing the
336 dataset.

337 Subsequently, we delve into the results obtained for each dataset. All
338 codes and datasets are openly accessible at the following URL: [https://](https://github.com/shaprob/FairShap)
339 github.com/shaprob/FairShap.

340 4.1. COMPAS dataset

341 As a first experiment, we considered the COMPAS dataset [39], released
342 by ProPublica⁵ in 2016. In this dataset, one assigns recidivism risk scores
343 to defendants based on a set of numerical and categorical features describ-
344 ing them. We considered seven input features, namely *sex* (male or female),

⁵<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

345 *age_cat* (age intervals - less than 25, greater than 45 and between 25 and
 346 45 years old), *juv_fel_count* (number of juvenile felonies), *juv_misd_count*
 347 (number of juvenile misdemeanors), *juv_other_count* (number of other prior
 348 juvenile convictions), *priors_count* (number of prior crimes) and *c_charge_degree*
 349 (crime degree - misdemeanor or felony). As the COMPAS dataset also pro-
 350 vides the *race* of each defendant, we used this information as the sensible
 351 feature. For the purpose of evaluating unfair results between blacks and
 352 whites, we only considered African-American (protected group) or Caucasian
 353 (privileged group) defendants. After under-sampling the data by removing
 354 defendants from other races, we achieved 5048 samples.

355 It is known from the literature [39] that in this dataset one has disparities
 356 in false positive rates when comparing blacks and whites defendants. In
 357 other terms, the rate that blacks are wrongly classified as recidivists is higher
 358 than this rate for Caucasians. In order to investigate this unfair result, we
 359 provide in Figure 3 some interpretations with respect to the false positive
 360 rates and Predictive Equality, along with thresholds. One may clearly see
 361 from Figure 3a that, regardless the adopted threshold, blacks achieved higher
 362 FPRs in comparison with whites. An interesting result was achieved in terms
 363 of the feature contributions. As can be seen in Figure 3c, *age_cat* led to the
 364 highest contributions toward the unfair results. This finding can be explained
 365 by the data distribution with respect to the age categories, race and classes.
 366 For instance, among the defendants under 25 years old, 66% are considered
 367 recidivists. Moreover, within this age category, 70% are blacks. On the
 368 other hand, for those greater than 45 years old, only 24% are considered
 369 recidivists and 42% are blacks. Although race was removed from the training

step, the data distribution of age category and recidivism risk scores carries information from race and, therefore, the age category can be seen as a source of bias towards race.

The unfair scenario can also be interpreted directly from the fairness measure. As showed in Figures 3e and 3d, age category is the feature highly associated with disparate outcomes. It is interesting to remark here that there are a lot of dispersion around the averaged Shapley values for *priors_count* and *c_charge_degree* (see the shaded are in Figures 3b and 3d). Therefore, there are uncertainties in the Shapley values estimation for these two feature. However, there are less uncertainties for age category, with positive contributions regardless the adopted threshold. Moreover, note that the disparity decreases as the threshold value increases. For instance, for $t = 0.8$, one practically achieved fairness for predictive equality. However, for this threshold we possibly pay the price of decreasing the TPR. Another choice could be $t \approx 0.35$, which significantly reduces the disparity on FPRs (around 50%) in comparison with the default choice of $t = 0.5$.

4.2. LSAC dataset

In this second application of our proposal, we deal with interpretability in unfair results associated with the LSAC (Law School Admission Council) dataset [40] (see [42] for further details). The goal is to predict whether a student will pass the bar exam on the first try. For this purpose, the students are described by the following features: *decile1b* (decile based on the grades in the first year), *decile3* (decile based on the grades in the thrid year), *lsat* (score), *ugpa* (undergraduate GPA), *zfygpa* (first year law school GPA), *zgpa* (cumulative law school GPA), *fulltime* (full-time or part-time work) and

395 *fam_inc* (family income bracket). Besides these features, *male* (whether the
396 student is male or female) and *race* (white or non-white) is also provided. In
397 our analysis, both of them are assumed as sensible features. However, when
398 evaluating disparities with respect to race (resp., gender), we considered *male*
399 (resp., *race*) as an input feature in model training. This dataset is highly
400 unbalanced and, after under-sampling, we obtained 3672 samples (originally,
401 there were more than 45k samples).

402 Unfair results in the LSAC dataset are frequently associated with race.
403 We provide in Figures 4 and 5 the interpretations on the probability of success
404 in the bar exam and on the Statistical Parity, respectively. One may see
405 that, when evaluating fairness in race, there is a huge disparity between the
406 two groups (see Figure 4a). The probability of positive outcome (i.e., being
407 classified as success in the exam) for non-whites is much lower than for whites.
408 This disparate result can be explained by the features contributions presented
409 in Figure 4c, where most of them (specially *lsat*) contribute more to classify
410 whites as the positive class than non-whites. On the other hand, as can be
411 seen from Figure 4b, the probability of favorable outcomes is practically the
412 same regardless the gender and the adopted threshold. Moreover, in this
413 scenario, all features contribute equally to classify as successfully passing the
414 bar exam both males and females (see Figure 4d). Therefore, we could note
415 that features have different impacts when evaluating disparities for different
416 sensitive features. Moreover, the performance and fairness measures have
417 different shapes, which indicate that, in a threshold analysis, the choice of an
418 appropriate threshold to enhance fairness should be conducted individually.
419 However, in a scenario where two or more sensitive features are considered

420 simultaneously (what we generally refer as to intersectional fairness), a trade-
421 off analysis should be conducted.

422 The aforementioned findings can also be seen in Figure 5. The statisti-
423 cal parity presented in Figures 5a and 5b attests the existence of disparities
424 between races and the absence of unfair outcomes related to gender, respec-
425 tively. With respect to the contributions of features (see Figures 5c and 5d),
426 while there are contributions towards unfair results for race, they are almost
427 zero for gender. Finally Figure 5e presents the contributions of features
428 along with thresholds, where *lsat* appears as the most relevant information
429 that leads to the disparate results.

430 4.3. Adult income dataset

431 In this last scenario, we considered the Adult income dataset. The goal in
432 this dataset is to predict whether a person makes over 50K per year. As *gen-*
433 *der* is one of the available information, disparities can be noted when compar-
434 ing males and females. Therefore, we assumed gender as the sensible feature
435 in our analysis, with female and male being the protected and the privileged
436 groups, respectively. The remaining features used to train the model are⁶ *age*
437 (intervals - less than 25, greater than 60 and between 25 and 60 years old),
438 *workclass* (Private or Non-private), *educational-num* (numerical value associ-
439 ated with the education degree - the greater the better), *marital-status* (mar-
440 ried, never-married or other), *relationship* (Wife, Own-child, Husband, Not-
441 in-family, Other-relative or Unmarried), *race* (White, Asian-Pac-Islander,

⁶We preprocessed this dataset in order to group categories of some features. Moreover, we also removed categorical features with high number of categories. Most of these changes have already been discussed in [42].

442 Amer-Indian-Eskimo, Black or Other), *capital-gain*, *capital-loss*, *hours-per-*
443 *week* (hours of work per week) and *native-country* (US or non-US). The total
444 number of samples are 22416.

445 As in this dataset one observes less females being classified as making
446 over 50k per year (in comparison with males), we investigated such a dis-
447 parity by means of the Equalized Odds. Recall that in Equalized Odds one
448 considers true positive and false positive rates. Figures 6a and 6b present
449 a comparison between sensitive groups with respect to FPRs and TPRs, re-
450 spectively. Both figures attest that there are less females classified as making
451 over 50k per year than males, either if the classification is correct (the TPRs)
452 or incorrect (the FPRs). The interpretations in terms of equalized odds are
453 presented in Figure 6c. The unfair outcomes can be explained by the features
454 contributions provided in Figures 6e and 6d. It is interesting to note that
455 both *marital-status* and *relationship* are the features that contribute the most
456 toward the unfair result. Indeed, as some categories in *relationship* describe
457 marital status, such features are somehow redundant. This explains the sim-
458 ilar results in terms of contribution towards the equalized odds. Moreover,
459 *relationship* is a feature that can be seen as a proxy for gender. For instance,
460 the indication that a person is a Wife or a Husband is practically the same
461 that saying the this person is a woman or a man, respectively. Therefore,
462 even if gender is removed from training step, the use of proxy features such
463 as *relationship* is a source of bias that can lead to disparate results.

464 4.3.1. Conducting feature selection

465 Although the main goal of this paper is to propose an approach to evalu-
466 ate features contributions toward unfair results along with threshold values,

our method could be helpful in feature engineering. As we detect features with high impact on disparate results, this information can be useful to conduct feature selection aiming at improving fairness. Indeed, by removing features that contribute to increase unfairness, one expects to mitigate ethical disparities.

In this subsection, we attest our proposal as a feature selection strategy on the Adult income dataset. As achieved in the previous section, both *marital-status* and *relationship* have high impact towards unfair results. Therefore, aiming at improving fairness, we could remove both features from the analysis⁷. By adopting the remaining features into the classifier, the obtained FPRs, TPRs and Equalized Odds are presented in Figure 7. As can be seen in Figures 7a and 7b, the performance of both protected and privileged groups are much closer in comparison with the model with features *marital-status* and *relationship* (see, for instance, Figures 6a and 6b). Therefore, we can attest the reduction of disparate results when features with high impact on unfairness are removed from the dataset. This finding can also be verified in Figure 7c, where the equalized odds decreased along with thresholds.

5. Conclusions

In this paper, we propose to investigate performance and disparity across various decision threshold(s) and quantifying the contribution of different features towards these two objectives. We have demonstrated the usefulness of our proposed approach with the help of three different case studies in-

⁷For the scope of this paper, we only consider improving fairness. Therefore, we do not further evaluate the impact on the performance measure by removing both features.

489 volving real datasets. For example, in the recidivism dataset, we showed the
490 disparity decreased with the increasing value of threshold. In the adult in-
491 come dataset, our proposed approach identified that the use of proxy features
492 such as relationship is a source of bias that can lead to disparate results, even
493 if the gender attributes is excluded whilst training the model. In the LSAC
494 dataset, we noted that features may have different contributions towards
495 different sensitive features. Indeed, in this dataset, some features impacted
496 the statistical parity associated with race, however, they do not contribute
497 towards disparity associated with gender.

498 In summary, this provided us a tool to identify the trade-offs between the
499 quality of prediction and the disparities between protected and other groups.
500 Moreover, one may identify which features contribute the most for both mea-
501 sures. This highlights an imminent use of our proposal. As illustrated in the
502 Adult income dataset, the user may adopt our proposal to help conducting
503 feature engineering by selecting features that do not entail disparate out-
504 comes. As a result, one may improve fairness. However, it is important
505 to see this feature selection impact into the model performance. Therefore,
506 we see as a future perspective a trade-off analysis looking performance and
507 fairness when conducting feature selection based on features contributions.

508 As another future perspective, we believe that it will be useful to validate
509 the practical usefulness of our proposal through experimental studies, for
510 example, to assess its impact in real world problems by doing survey-based
511 studies and collecting feedback from stakeholders. Another important area
512 of work will be to develop an interactive software tool for stakeholders that
513 can help visualize the trade-off between performance and disparity in machine

514 learning classifier; and therefore, enabling them to choose a decision threshold
515 based on their preferences. Finally, we would like to extend our approach
516 to deal with multi-class classification problems. In this case, once we adapt
517 both performance and fairness measures to have a single measure for each
518 coalition of features, we are able to calculate the payoff are, therefore, the
519 Shapley value and feature contributions.

520 **Acknowledgments**

521 Work supported by São Paulo Research Foundation (FAPESP) under the
522 grants #2020/09838-0 (BIOS - Brazilian Institute of Data Science), #2020/10572-
523 5 and #2021/11086-0. L. T. Duarte would like to thank the National Council
524 for Scientific and Technological Development (CNPq, Brazil) for the financial
525 support.

526 **References**

- 527 [1] J. M. Sheffield, R. Smith, P. Suthaharan, P. Leptourgos, P. R. Corlett,
528 Relationships between cognitive biases, decision-making, and delusions,
529 Scientific Reports 13 (2023) 9485.
- 530 [2] B. De Martino, D. Kumaran, B. Seymour, R. J. Dolan, Frames, biases,
531 and rational decision-making in the human brain, science 313 (2006)
532 684–687.
- 533 [3] S. H. Ivanov, Automated decision-making, foresight 25 (2023) 4–19.
- 534 [4] S. Barocas, A. D. Selbst, Big data’s disparate impact, California Law

- 535 Review 104 (2016) 671. URL: [http://lawcat.berkeley.edu/record/](http://lawcat.berkeley.edu/record/1127463)
536 [1127463](http://lawcat.berkeley.edu/record/1127463).
- 537 [5] A. Petrović, M. Nikolić, M. Jovanović, M. Bijanić, B. Delibašić, Fair clas-
538 sification via Monte Carlo policy gradient method, Engineering Appli-
539 cations of Artificial Intelligence 104 (2021). doi:[10.1016/j.engappai.](https://doi.org/10.1016/j.engappai.2021.104398)
540 [2021.104398](https://doi.org/10.1016/j.engappai.2021.104398).
- 541 [6] P. A. Grabowicz, N. Perello, A. Mishra, Marrying Fairness and Ex-
542 plainability in Supervised Learning, in: Proceedings of the 2022 ACM
543 Conference on Fairness, Accountability, and Transparency, Association
544 for Computing Machinery, Seoul, Republic of Korea, 2022, pp. 1905–
545 1916. doi:[10.1145/3531146.3533236](https://doi.org/10.1145/3531146.3533236).
- 546 [7] M. Hardt, E. Price, N. Srebro, Equality of opportunity in supervised
547 learning, in: Advances in Neural Information Processing Systems 29
548 (NIPS 2016), Barcelona, Spain, 2016, pp. 3315–3323.
- 549 [8] T. Quan, F. Zhu, Q. Liu, F. Li, Learning fair representations for ac-
550 curacy parity, Engineering Applications of Artificial Intelligence 119
551 (2023). doi:[10.1016/j.engappai.2023.105819](https://doi.org/10.1016/j.engappai.2023.105819).
- 552 [9] G. Xu, T. D. Duong, Q. Li, S. Liu, X. Wang, Causality Learning: A New
553 Perspective for Interpretable Machine Learning, ArXiv ID: 2006.16789
554 (2020). URL: <http://arxiv.org/abs/2006.16789>.
- 555 [10] N. Burkart, M. F. Huber, A survey on the explainability of supervised
556 machine learning, Journal of Artificial Intelligence Research 70 (2021)
557 245–317. doi:[10.1613/JAIR.1.12228](https://doi.org/10.1613/JAIR.1.12228).

- 558 [11] T. Miller, Explanation in artificial intelligence: Insights from the social
559 sciences, *Artificial Intelligence* 267 (2019) 1–38. doi:[10.1016/j.artint.](https://doi.org/10.1016/j.artint.2018.07.007)
560 [2018.07.007](https://doi.org/10.1016/j.artint.2018.07.007).
- 561 [12] D. Charte, F. Charte, M. J. del Jesus, F. Herrera, An analysis on the
562 use of autoencoders for representation learning: Fundamentals, learning
563 task case studies, explainability and challenges, *Neurocomputing* 404
564 (2020) 93–107. doi:[10.1016/j.neucom.2020.04.057](https://doi.org/10.1016/j.neucom.2020.04.057).
- 565 [13] C. Molnar, *Interpretable machine learning*, 2021. URL: [https://](https://christophm.github.io/interpretable-ml-book/)
566 christophm.github.io/interpretable-ml-book/.
- 567 [14] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model pre-
568 dictions, *Advances in neural information processing systems* 30 (2017).
- 569 [15] T. Begley, T. Schwedes, C. Frye, I. Feige, Explainability for fair machine
570 learning, *arXiv preprint arXiv:2010.07389* (2020).
- 571 [16] S. F. Nimmy, O. K. Hussain, R. K. Chakraborty, F. K. Hussain,
572 M. Saberi, Interpreting the antecedents of a predicted output by captur-
573 ing the interdependencies among the system features and their evolution
574 over time, *Engineering Applications of Artificial Intelligence* 117 (2023)
575 105596. doi:[10.1016/j.engappai.2022.105596](https://doi.org/10.1016/j.engappai.2022.105596).
- 576 [17] W. Cai, A. B. Kordabad, S. Gros, Energy management in residential
577 microgrid using model predictive control-based reinforcement learning
578 and Shapley value, *Engineering Applications of Artificial Intelligence*
579 119 (2023) 105793. doi:[10.1016/j.engappai.2022.105793](https://doi.org/10.1016/j.engappai.2022.105793).

- 580 [18] B. Rozemberczki, L. Watson, P. Bayer, H. T. Yang, O. Kiss, S. Nilsson,
581 R. Sarkar, The Shapley Value in Machine Learning, in: IJCAI Interna-
582 tional Joint Conference on Artificial Intelligence, Vienna, Austria, 2022,
583 pp. 5572–5579. doi:[10.24963/ijcai.2022/778](https://doi.org/10.24963/ijcai.2022/778).
- 584 [19] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair,
585 R. Katz, J. Himmelfarb, N. Bansal, S.-I. Lee, From local explanations
586 to global understanding with explainable AI for trees, *Nature Machine*
587 *Intelligence* 2 (2020) 56–67. doi:[10.1038/s42256-019-0138-9](https://doi.org/10.1038/s42256-019-0138-9).
- 588 [20] E. Albini, J. Long, D. Dervovic, D. Magazzeni, Counterfactual Shapley
589 Additive Explanations, in: *Proceedings of the 2022 ACM Conference on*
590 *Fairness, Accountability, and Transparency*, Association for Computing
591 Machinery, Seoul, Republic of Korea, 2022, pp. 1054–1070. doi:[10.1145/](https://doi.org/10.1145/3531146.3533168)
592 [3531146.3533168](https://doi.org/10.1145/3531146.3533168).
- 593 [21] D. Watson, Rational Shapley Values, in: *Proceedings of the 2022 ACM*
594 *Conference on Fairness, Accountability, and Transparency*, Seoul, Re-
595 public of Korea, 2022, pp. 1083–1094. doi:[10.1145/3531146.3533170](https://doi.org/10.1145/3531146.3533170).
- 596 [22] K. Aas, M. Jullum, A. Løland, Explaining individual predictions when
597 features are dependent: More accurate approximations to Shapley val-
598 ues, *Artificial Intelligence* 298 (2021) 103502. doi:[10.1016/j.artint.](https://doi.org/10.1016/j.artint.2021.103502)
599 [2021.103502](https://doi.org/10.1016/j.artint.2021.103502).
- 600 [23] G. D. Pelegina, L. T. Duarte, M. Grabisch, A k-additive Choquet
601 integral-based approach to approximate the SHAP values for local inter-

- 602 pretability in machine learning, *Artificial Intelligence* 325 (2023) 104014.
603 doi:[10.1016/j.artint.2023.104014](https://doi.org/10.1016/j.artint.2023.104014).
- 604 [24] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, A. Galstyan, A
605 survey on bias and fairness in machine learning, *arXiv preprint*
606 arXiv:1908.09635 (2019). URL: <http://arxiv.org/abs/1908.09635>.
- 607 [25] M. Kearns, A. Roth, *The ethical algorithm: The science of socially*
608 *aware algorithm design*, Oxford University Press, New York, USA, 2019.
609 doi:[10.1017/CB09781107415324.004](https://doi.org/10.1017/CB09781107415324.004).
- 610 [26] Y. Wang, V. A. Nguyen, G. A. Hanasusanto, Wasserstein robust classi-
611 fication with fairness constraints, 2021. [arXiv:2103.06828](https://arxiv.org/abs/2103.06828).
- 612 [27] N. Martinez, M. Bertran, G. Sapiro, Minimax pareto fairness: A multi
613 objective perspective, in: *Proceedings of the 37th International Confer-*
614 *ence on Machine Learning (ICML 2020)*, 2020, pp. 6711–6720.
- 615 [28] S. Liu, L. N. Vicente, Accuracy and fairness trade-offs in machine learn-
616 ing: a stochastic multi-objective approach, *Computational Management*
617 *Science* 19 (2022) 513–537. doi:[10.1007/s10287-022-00425-z](https://doi.org/10.1007/s10287-022-00425-z).
- 618 [29] G. D. Pelegina, S. Siraj, Shapley value-based approaches to explain
619 the quality of predictions by classifiers, *IEEE Transactions on Artificial*
620 *Intelligence* (2024) 1–15. doi:[10.1109/TAI.2024.3365082](https://doi.org/10.1109/TAI.2024.3365082).
- 621 [30] M. T. Ribeiro, S. Singh, C. Guestrin, "Why should I trust you?" Ex-
622 plaining the predictions of any classifier, in: *Proceedings of the 22nd*
623 *ACM SIGKDD international conference on knowledge discovery and*
624 *data mining*, 2016, pp. 1135–1144.

- 625 [31] S. Verma, J. Rubin, Fairness definitions explained, in: 2018 ACM/IEEE
626 International Workshop on Software Fairness, IEEE, 2018, pp. 1–7.
- 627 [32] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, R. Zemel, Fairness
628 through awareness, in: 3rd Innovations in Theoretical Computer Sci-
629 ence Conference (ITCS 2012), Cambridge, MA, USA, 2012, pp. 214–226.
630 doi:[10.1145/2090236.2090255](https://doi.org/10.1145/2090236.2090255). [arXiv:1104.3913](https://arxiv.org/abs/1104.3913).
- 631 [33] L. S. Shapley, A value for n-person games, in: W. Kuhn, A. W. Tucker
632 (Eds.), Annals of mathematics studies: Vol. 28. Contributions to the
633 theory of games, Vol. II, Princeton University Press, Princeton, 1953,
634 pp. 307–317.
- 635 [34] L. Han, T. Morstyn, M. McCulloch, Incentivizing Prosumer Coalitions
636 With Energy Management Using Cooperative Game Theory, IEEE
637 Transactions on Power Systems 34 (2019) 303–313.
- 638 [35] H. P. Young, Monotonic solutions of cooperative games, International
639 Journal of Game Theory 14 (1985) 65–72.
- 640 [36] S. Lipovetsky, M. Conklin, Analysis of regression in game theory ap-
641 proach, Applied Stochastic Models in Business and Industry 17 (2001)
642 319–330.
- 643 [37] M. Grabisch, P. Miranda, Exact bounds of the Möbius inverse of mono-
644 tone set functions, Discrete Applied Mathematics 186 (2015) 7–12.
645 doi:[10.1016/j.dam.2015.01.016](https://doi.org/10.1016/j.dam.2015.01.016).
- 646 [38] L. Breiman, Random forests, Machine Learning 45 (2001) 5–32.

- 647 [39] J. Angwin, J. Larson, S. Mattu, L. Kirchner, Machine bias - propublica,
648 Propublica (2016). URL: [https://www.propublica.org/article/
649 machine-bias-risk-assessments-in-criminal-sentencing](https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing).
- 650 [40] L. F. Wightman, LSAC national longitudinal bar passage study, 1998.
- 651 [41] B. Becker, R. Kohavi, Adult, UCI Machine Learning Repository, 1996.
652 doi:doi.org/10.24432/C5XW20.
- 653 [42] T. Le Quy, A. Roy, V. Iosifidis, W. Zhang, E. Ntoutsi, A survey on
654 datasets for fairness-aware machine learning, Wiley Interdisciplinary
655 Reviews: Data Mining and Knowledge Discovery 12 (2022) 1–59.

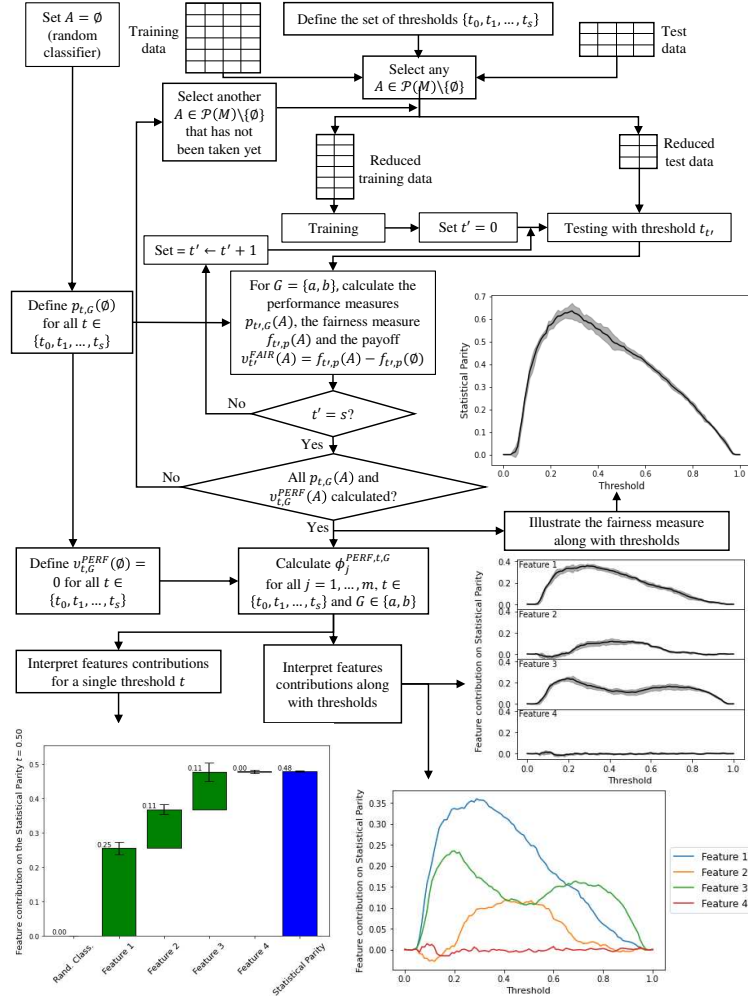
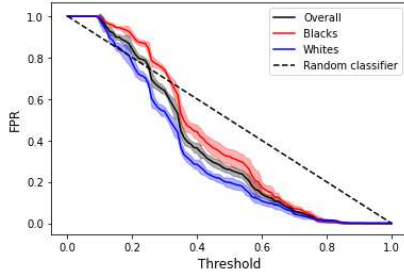
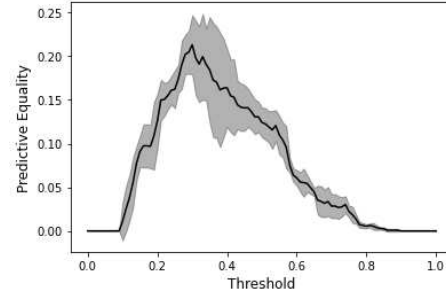


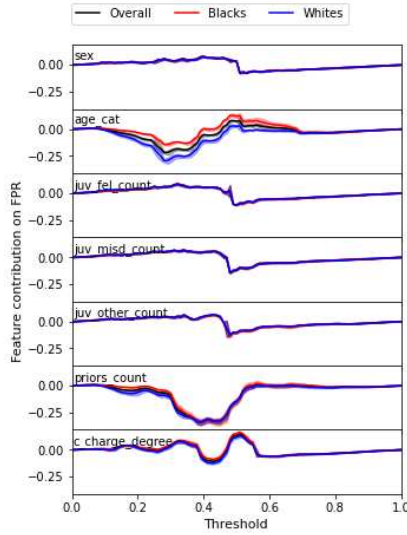
Figure 2: The proposed scheme to evaluate features contributions towards fairness measures.



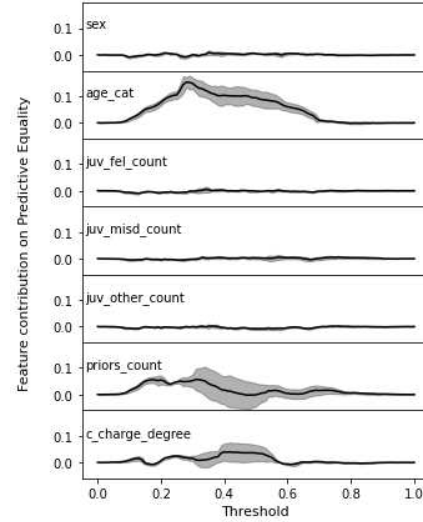
(a) FPRs along with thresholds.



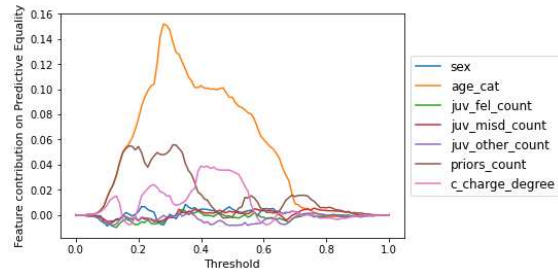
(b) Predictive Equality along with thresholds.



(c) Contributions of features towards FPRs (split).

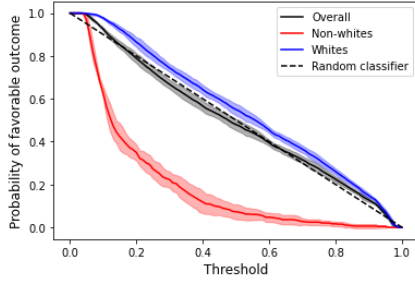


(d) Contributions of features towards Predictive Equality (split).

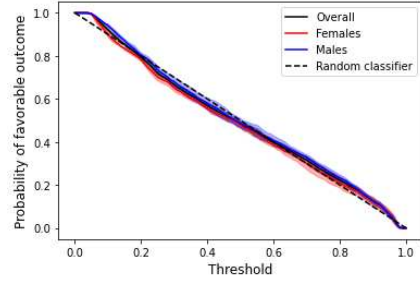


(e) Contributions of features towards Predictive Equality.

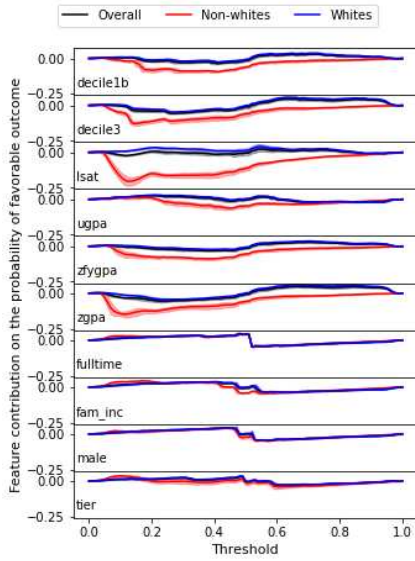
Figure 3: Interpreting disparate results towards FPRs and Predictive Equality - COMPAS dataset.



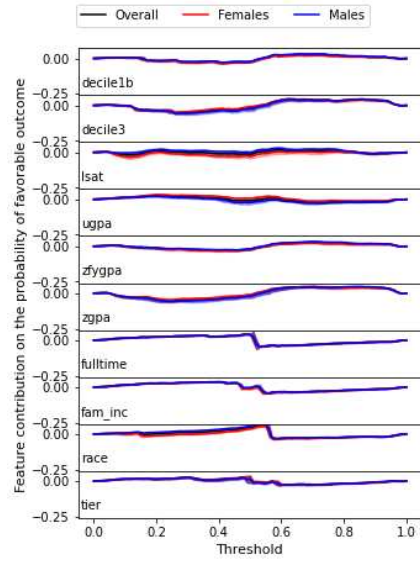
(a) Disparities between whites and non-whites.



(b) Disparities between males and females.

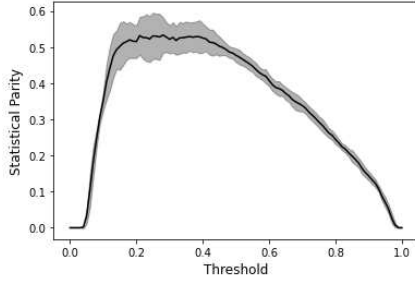


(c) Contributions of features (split) - Whites and non-whites.

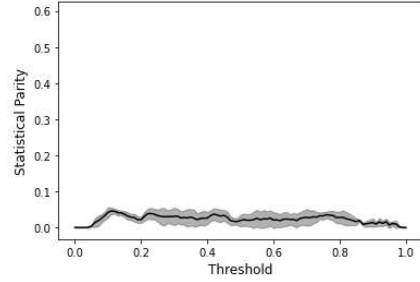


(d) Contributions of features (split) - Males and females.

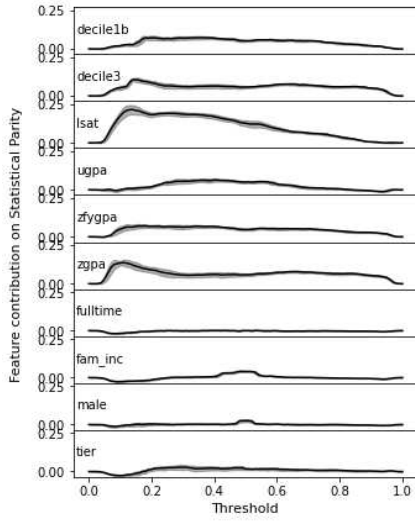
Figure 4: Interpreting disparate results towards the probability of success in the bar exam - LSAC dataset.



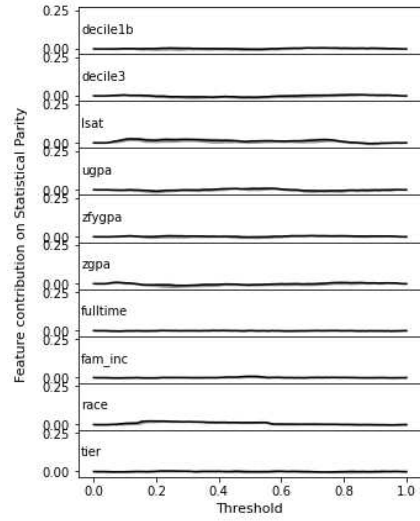
(a) Disparities between whites and non-whites.



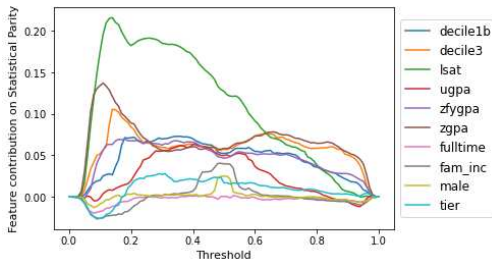
(b) Disparities between males and females.



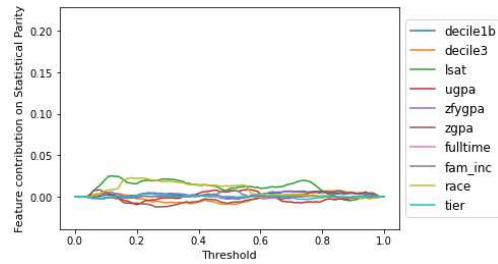
(c) Contributions of features (split) - Whites and non-whites.



(d) Contributions of features (split) - Males and females.

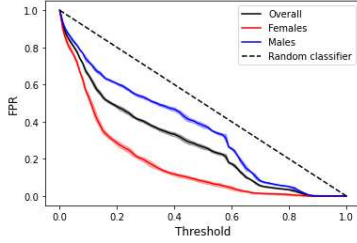


(e) Contributions of features - Whites and non-whites.

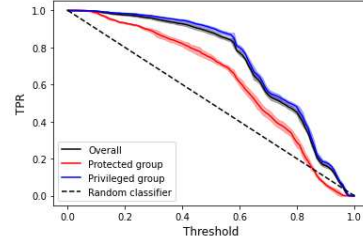


(f) Contributions of features - Males and females.

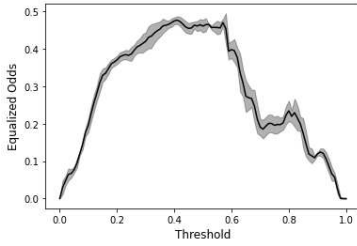
Figure 5: Interpreting disparate results towards the Statistical Parity - LSAC dataset.



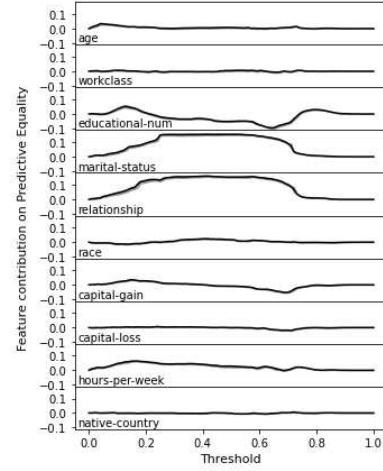
(a) FPRs along with thresholds.



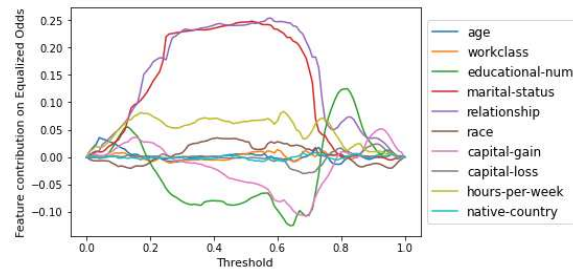
(b) TPRs along with thresholds.



(c) Equalized Odds along with thresholds.

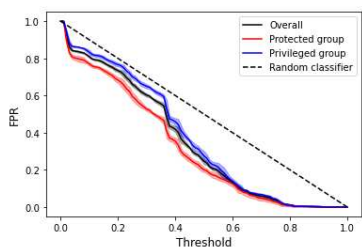


(d) Contributions of features towards Equalized Odds (split).

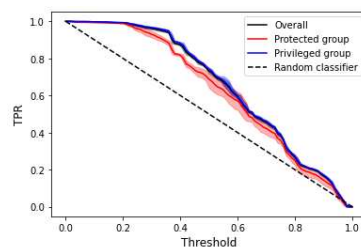


(e) Contributions of features towards Equalized Odds.

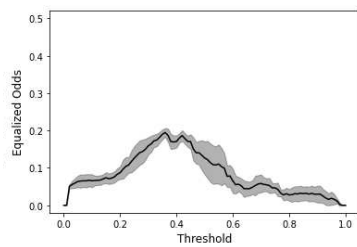
Figure 6: Interpreting disparate results towards FPRs, TRPs and Equalized Odds - Adult income dataset.



(a) FPRs along with thresholds.



(b) TPRs along with thresholds.



(c) Equalized Odds along with thresholds.

Figure 7: textcoloredVisualizing FPRs, TRPs and Equalized Odds after feature selection - Adult income dataset.