



UvA-DARE (Digital Academic Repository)

PEBAM: A Profile-Based Evaluation Method for Bias Assessment on Mixed Datasets

Wilms, M.; Sileno, G.; Haned, H.

DOI

[10.1007/978-3-031-15791-2_17](https://doi.org/10.1007/978-3-031-15791-2_17)

Publication date

2022

Document Version

Final published version

Published in

KI 2022: Advances in Artificial Intelligence

License

Article 25fa Dutch Copyright Act (<https://www.openaccess.nl/en/in-the-netherlands/you-share-we-take-care>)

[Link to publication](#)

Citation for published version (APA):

Wilms, M., Sileno, G., & Haned, H. (2022). PEBAM: A Profile-Based Evaluation Method for Bias Assessment on Mixed Datasets. In R. Bergmann, L. Malburg, S. C. Rodermund, & I. J. Timm (Eds.), *KI 2022: Advances in Artificial Intelligence: 45th German Conference on AI, Trier, Germany, September 19–23, 2022 : proceedings* (pp. 209-223). (Lecture Notes in Computer Science; Vol. 13404), (Lecture Notes in Artificial Intelligence). Springer. https://doi.org/10.1007/978-3-031-15791-2_17

General rights

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

Disclaimer/Complaints regulations

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

UvA-DARE is a service provided by the library of the University of Amsterdam (<https://dare.uva.nl>)



PEBAM: A Profile-Based Evaluation Method for Bias Assessment on Mixed Datasets

Mieke Wilms, Giovanni Sileno^(✉), and Hinda Haned

University of Amsterdam, Amsterdam, The Netherlands
mieke.wilms@student.uva.nl g.sileno@uva.nl

Abstract. Bias evaluation methods focus either on individual bias or on group bias, where groups are defined based on protected attributes such as gender or ethnicity. More generally, however, descriptively relevant combinations of feature values in the data space (*profiles*) may serve also as anchors for biased decisions. This paper introduces therefore a semi-hierarchical clustering method for profile extraction from mixed datasets. It elaborates on how profiles can be used to reveal *historical*, *representational*, *aggregation* and *evaluation biases* in algorithmic decision-making models, taking as example the German credit data set. Our experiments show that the proposed profile-based evaluation method for bias assessment on mixed datasets (PEBAM) can reveal forms of bias towards profiles expressed by the dataset that are undetected when using individual- or group-bias metrics alone.

Keywords: Algorithmic fairness · Bias prevention · Bias evaluation · Clustering · Domain analysis

1 Introduction

The wider introduction of machine learning algorithms in decision-making processes feeds an ongoing debate over algorithmic decision-making (ADM). To prevent or correct ADM from taking biased decisions several fairness-aware machine learning algorithms have been proposed [8]. However, these algorithms are not always accessible to practitioners due to their ‘black-box’ nature [2]; they are highly dependent on the data preprocessing phase [8]; and, at more fundamental level, fairness and bias can be given technical meanings but cannot be captured by one single definition [14, 16]. Contemporary bias evaluation methods used for fairness analysis generally focus either on individual bias or on group bias, where groups are defined based on protected attributes such as gender or ethnicity, but this is not without drawbacks. For instance, analysing the German credit dataset—a real world dataset¹ collecting features of loan applicants and a credit risk label *good* or *bad* assigned to them—the group of young individuals

¹ Available at: <https://www.kaggle.com/uciml/german-credit>.

(age below 25) obtains more often a false negative label than the group above 25, hence young individuals are discriminated when applying for a loan [11]. The simplest solution would be to take this sensitive attribute out of consideration, however there are multiple attributes that correlate with the “age” attribute (e.g. “own house” [12]). From a more general standpoint, one may ask whether there exist relevant descriptive combinations of feature values in the data space, that we will call here *profiles*, which may act as anchors for (assessing the presence of) biased decisions. As an additional source of complexity, we need also to take into account that data is commonly presented in form of mixed datasets (i.e. including both categorical and numeric features). Discretization of numeric dimensions, or embedding of categorical dimensions, add further complexity and potentially undesired effects. Given this context, we address the following research questions: *How can profiles be defined? How can we extract profiles from mixed datasets? How can profiles be used to assess biases? How does a profile-based assessment compare with existing individual- or group-based methods?* The goal of this paper is to develop and test a Profile-based Evaluation method for Bias Assessment of algorithmic decision-making on Mixed datasets (PEBAM). Our contribution is twofold: (i) we present an effective and computationally efficient method for profile extraction on mixed datasets based on clustering; (ii) we show how profiles can be utilized to evaluate various forms of biases—most of them associated to trained ADM models. The paper is structured as follows. Section 2 provide a brief overview of relevant concepts. Section 3 presents the proposed methodology. Section 4 elaborates on the experiments and results on the German credit dataset. A note on future work ends the paper.

2 Theoretical Background

Types of Bias. Several types of bias have been identified in the literature (see e.g. the 23 types in [14]), but for the scope of this research we will focus in particular on biases that can arise during a ML-product lifecycle (see e.g. [16]). For instance, during the *data generation* process, we may have: *historical bias*, produced by the world as it is, and occurring even if data is perfectly measured and sampled; *representation bias*, occurring when the training data for the ML model under-represents parts of the population the algorithm will be used on. During the *model building* and *implementation* phases, we may have: *aggregation bias*, arising when a general model is used for all groups, while in reality different groups have a different mapping from input features to labels (e.g. some ethnic groups can have different indicators for a disease than others); *evaluation bias*, occurring when the data on which the model is evaluated is a misrepresentation of the target population. These four types of bias do not cover all possible sources of bias, but they will be used as relevant examples about how to set up a profile-based evaluation.

Bias Evaluation Methods for Algorithmic Fairness. In their extensive literature review, Mehrabi et al. [14] give an overview of the most widely used definitions of fairness within machine learning, providing 3 definitions focused on individual fairness, 6 on group fairness, in which groups are defined by protected attribute classes (e.g. sex, ethnicity, etc.), 1 on subgroup fairness. In this work we will build upon two (group-fairness) measures. The first is *equal opportunity* [10], a criterion for fairness in binary algorithms. Reading the outcome $y = 1$ as the “advantaged” outcome, and A as the protected class attribute, we have:

Definition 1. *Equal opportunity* *A binary predictor \hat{y} satisfies equal opportunity w.r.t. attribute A and ground truth y iff: $Pr\{\hat{y} = 1 \mid A = 0, y = 1\} = Pr\{\hat{y} = 1 \mid A = 1, y = 1\}$.*

The second is *contextual demographic (dis)parity* (CDD)—based on *conditional (non-)discrimination* by [12]—a measure found to be the most compatible with the decisions of the European Court of Justice on cases of discrimination [18].

Definition 2. *Conditional Demographic Disparity* *Let R be a given set of attributes, A_r be the proportion of people belonging to a protected class in the advantaged group and with attribute $r \in R$, and let D_r be the proportion of people of protected class in the disadvantaged group with attribute r . A decision-making process exhibits conditional demographic disparity iff: $\forall r \in R : D_r > A_r$.*

The conditions r in R should be *explanatory* [12], i.e. they should hypothetically explain the outcome even in the absence of discrimination against the protected class (e.g., different salaries between men and women might be due to different working hours). Under this view, R is derived from domain expert knowledge.

Clustering Algorithm for Mixed Data. In ADM one very often has to deal with mixed datasets, i.e. datasets that consist of both categorical and numerical features. Various solutions have been proposed in the literature to the known difficulty to capture distributions on mixed datasets [15]; the present work will rely in particular on *k-medoids clustering* [5]. The main benefit of *k-medoids clustering* over *k-means* is that it is more robust to noise and outliers; we also do not have to come up with a measure to compute the mean for categorical features. On the other hand, the *k-medoids clustering* problem is NP-hard to solve exactly. For this reason, in our work we will make use of the heuristic Partitioning Around Medoids (PAM) algorithm [4].

3 Methodology

PEBAM (*Profile-based Evaluation for Bias Assessment for Mixed datasets*) is a method consisting of three main steps: (1) a profile selection—based on the iteration of clustering controlled by a measure of variability—to extract profiles representative of the target domain from an input dataset; (2) profiles are evaluated in terms of stability over repetitions of extractions; (3) a given ADM classification model is evaluated for bias against those profiles.

3.1 Profile-Selection Based on Clustering

Informally, profiles can be seen as relevant descriptive elements that, as a group, act as a “summary” of the data space. Because individuals sharing to an adequate extent similar attributes should intuitively be assigned to the same profile, clustering can be deemed compatible with a profile selection task. We consider then three different clustering algorithms to implement profile selection: (i) *simple clustering*; (ii) a form of *hierarchical clustering* based on the iteration of the first; and (iii) a novel *semi-hierarchical clustering* method based on adding static and dynamic constraints to control the second. The first two algorithms will be used as baselines to evaluate the third one, and will be described succinctly.

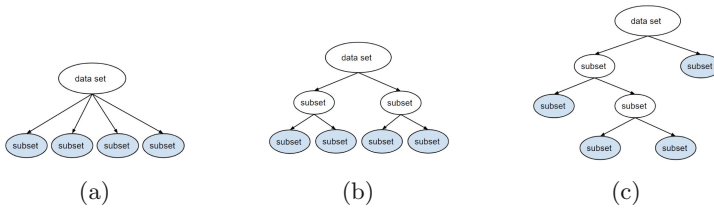


Fig. 1. Schematic overview of the simple clustering algorithm with $k = 4$ (a), the hierarchical clustering algorithm with $k = 2$ and $l = 2$ (b) and a possible outcome of the semi-hierarchical clustering algorithm (c).

The **simple clustering** algorithm consists in the application of the chosen clustering algorithm (in our case, k -medoids, Sect. 2), and results in a flat partition of the sample space (see e.g. Fig. 1a). The challenge is to determine the right number of clusters k . **Hierarchical clustering** consists in the nested iteration of the previous algorithm set with k clusters over l layers (see e.g. Figure 1b). The benefit of this method is that the resulting tree-based structure gives further insight on the basis of which feature clusters are created. The downside is that we have to tune an extra hyperparameter besides k : the number of layers l .

The **semi-hierarchical clustering** algorithm is an automatically controlled version of the second algorithm. It is based on performing iteratively two operations. First, we apply a clustering method with $k = 2$, i.e. at every step we divide input subsets into two new subsets (clusters). Then, we test each new subset to decide whether we should continue clustering by looking at the *variability* of the features it expresses, or at its cardinality. For variability of a feature f w.r.t. a dataset D we intend a measure between 0 and 1 that indicates to what extent the feature f is varying within D . A value close to 0 indicates that f is rather stable over D , and a value close to 1 indicates that f mostly varies in D (and so it is not relevant for describing it). The intuition behind using this measure is that stable (i.e. not varying) features help to discriminate a profile over another one. For instance, coffees are characterized by a black (or dark brown) colour, so the colour feature is very stable to support discriminating coffee from other

drinks; in contrast, colour is not a relevant feature to discriminate books from other objects. In general, any feature can then be: (1) an *irrelevant feature*, when the variability of a feature exceeds an upper bound c_u , is deemed not to be a characteristic property of the cluster; (2) a *relevant feature*: when the variability of a feature is smaller than a lower bound c_l , it means that the feature is strongly characteristic of that cluster. When all features expressed by a subset satisfy either case (1) or case (2), or the subset has less than a fixed amount of elements n_{stop} , there is no need to continue further clustering the input subset, and therefore it is selected as a *profile*. The resulting structure is then a binary tree, possibly unbalanced (see Fig. 1c), whose leaves are the selected profiles. The benefit of semi-hierarchical clustering over simple and hierarchical clustering is that we do not have to decide the numbers of clusters and layers in advance, but requires setting the variability thresholds values c_u and c_l , as well as the threshold cluster cardinality n_{stop} .

In quantitative terms, when *numerical features* are not significant, we expect that their distribution should approximate a uniform distribution. Let us assume we have a numerical non-constant feature X ; we can normalize X between 0 and 1 via $(X - X_{\min}) / (X_{\max} - X_{\min})$, and then compute the sample standard deviation s from the normalized samples. Theoretically, for a random variable $U \sim \text{Uniform}(0, 1)$, we have $\mu = 1/2$, and $\text{Var}(U) = \mathbb{E}[(U - \mu)^2] = \int_0^1 (x - 1/2)^2 dx = 1/12$. Thus, the standard deviation of a random variable uniformly distributed is $\sqrt{1/12} \approx 0.29$. Therefore, if the sample standard deviation s approximates 0.29, we can assume that the feature X is uniformly distributed across the given cluster and is therefore not a unique property of species within the cluster. On the other hand, when the standard deviation is close to zero, this means that most sample points are close to the mean. This indicates that feature X is very discriminating for that specific cluster. To obtain a measure of variability, we need to compute the standardized standard deviation $s_s = s/0.29$.

For *categorical features*, we consider the variability measure proposed in [1]. Let X a n -dimensional categorical variable consisting of q categories labelled as $1, 2, \dots, q$. The relative frequency of each category $i = 1, \dots, q$ is given by $f_i = n_i/n$, where n_i is the number of samples that belongs to category i and $n = \sum_{i=1}^q n_i$. Let $\vec{f} = (f_1, f_2, \dots, f_q)$ be the vector with all the relative frequencies. We define the variability of X as: $v_q = 1 - \|\vec{f}\|_q$. Allaj [1] shows that the variability is bounded by 0 and $1 - 1/q$, where an outcome close to $1 - 1/q$ associates to high variability. We can also compute the standardized variability: $v_{q,s} = \frac{v_q}{1 - 1/\sqrt{q}}$, such that the variability lies between 0 and 1 for all number of categories q where again a variability close to 1 implies a high variability and hence a non-characteristic feature. A variability close to 0 indicates that the feature is highly characteristic for that sample set.

3.2 Evaluation of Clustering Methods for Profile Selection

Given a clustering method, we need to evaluate whether it is working properly with respect to the profile selection task. Unfortunately, since this task is an unsupervised problem, we do not have access to the ground truth, but we can still focus on a related problem: *Is our method stable?* The **stability** of a clustering method towards initialization of the clusters can be tested by running the algorithm multiple times with different initial cluster settings (in our case, different datapoints selected as the initial medoids), and check whether we end up with the same sets of clusters (hereby called *cluster combinations*). When the stability analysis results in two or more different cluster combinations, we want to know how similar these combinations are, and for doing this, we will introduce a measure of inter-clustering similarity.

Inter-clustering similarity is a similarity score that tells to what extent different outcomes of clustering are similar, by comparing how many elements clusters belonging to the two clustering outputs have in common with each other. This score can be computed by comparing the distribution of the elements over the clusters of two combinations. We present the process through an example:

Example: Let us consider a dataset with 20 data points. Suppose the clustering algorithm returns the following two different cluster combinations C^1 and C^2 : where $C^1 = (c_1^1, c_2^1, c_3^1) = ((1, 4, 6, 7, 13, 17, 18, 20), (3, 5, 11, 12, 15, 16), (2, 8, 9, 10, 14, 19))$, and $C^2 = (c_1^2, c_2^2, c_3^2) = ((1, 4, 6, 7, 13, 17, 18, 20), (3, 11, 12, 14, 15, 16, 19), (2, 5, 8, 9, 10))$. For every cluster c_i^1 in C_1 we compute its overlap with each cluster c_j^2 in C_2 . For instance, for the first cluster of C^1 , $c_1^1 = (1, 4, 6, 7, 13, 17, 18, 20)$ the max overlap is 1, since c_2^1 of C^2 is exactly the same. For the second cluster c_2^1 we have:

$$\max\left(\frac{|c_2^1 \cap c_1^2|}{|c_2^1|}, \frac{|c_2^1 \cap c_2^2|}{|c_2^1|}, \frac{|c_2^1 \cap c_3^2|}{|c_2^1|}\right) = \max(0/6, 5/6, 1/6) = 0.83$$

Applying the same calculation on c_3^1 returns 0.67. The similarity score of cluster-combination C^1 with respect to cluster-combination C^2 is given by the mean over all the three maximum overlap values, which in this case is 0.83.

3.3 Profile-Based Evaluation of a Given Classifier

By means of profiles, we can assess the *historical*, *representational*, *aggregation* and *evaluation biases* (Sect. 2) of a certain ADM classification model.

Historical bias arises on the dataset used for training. We measure it using conditional demographic disparity (Def. 2); however, the resulting value may be wrong if an attribute r is not a relevant characteristic for grouping individuals. Therefore, we consider a *profile-conditioned demographic disparity*, differing from [12, 18] in as much each profile label c_i acts as attribute r . In this way we capture behavioural attitudes (w.r.t. assigning or not an advantageous label) of the labeller-oracle towards elements belonging to that profile.

For **Representational bias**, by clustering the dataset around profiles, we may get insights if there are parts of the population which are overrepresented and parts which are underrepresented. A domain expert can compare the resulting distribution of profiles with the expected distribution for the population on which the decision-making algorithm is going to be used. When these two distributions differ much from each other, one can *bootstrap sampling* from profiles to get a more correct distribution in the training dataset. If no domain expertise is available, one can consider using this method during deployment. By collecting the data where the algorithm is applied on, once the dataset is sufficiently large (e.g. about the size of the training dataset), one can divide the samples of the collected dataset over the different known profiles based on distance towards the medoids associated to profiles. If the distribution of the collected dataset over the profiles relevantly differs from the distribution of the training dataset, we can conclude that there is representation bias. Alternatively, one can repeat the profile selection on the collected data, and evaluate how much they differ from the ones identified in the training dataset.

In order to evaluate **Aggregation bias**, we need a good metric to evaluate the performance of the trained model under assessment. Our goal is to evaluate the model against all profiles, i.e. to test whether the model works equally well on individuals with different profiles (i.e. individuals from different clusters). We start from the definition of equal opportunity (Def. 1), but we reformulate it in a way that equal opportunity is computed with respect to profiles instead of the protected class:

Definition 3 (Equal opportunity w.r.t. a single profile). *We say a binary predictor \hat{y} satisfies equal opportunity with respect to a profile C equal to i and outcome/ground truth y iff: $Pr\{\hat{y} = 1|C = i, y = 1\} = Pr\{\hat{y} = 1|C \neq i, y = 1\}$*

Definition 4 (General equal opportunity). *A binary predictor \hat{y} satisfies general equal opportunity iff it satisfies equal opportunity with respect to all profiles $C \in \{1, \dots, k\}$ and ground truth y . In formula: $Pr\{\hat{y} = 1|C = 1, y = 1\} = \dots = Pr\{\hat{y} = 1|C = k, y = 1\}$*

In some cases, it might be that getting a wrong prediction for a ground truth or positive outcome label occurs more often than with a negative outcome label, and may be more valuable (e.g. in some medical disease treatment); looking at distinct values of y may give insights on the overall functioning of the model.

Evaluation bias arises when the model is evaluated on a population which differs in distribution from the data that was used for training the model. It has been shown [8] that fairness-preserving algorithms tend to be sensitive to fluctuations in dataset composition, i.e. the performance of the algorithms is affected by the train-test split. To ensure that we do not have evaluation bias, we run a Monte Carlo simulation of the decision-making algorithm. This means that we make M different train-test splits of the dataset. For each train-test split, we train the decision-making algorithm on the train set and use the test set for evaluation. For the model-evaluation, we use general equal opportunity (Def. 4). This gives us insights on which profiles are more sensitive towards train-test splitting (and thus to evaluation bias).

4 Experiments and Results

We conducted two experiments to evaluate PEBAM.² In the first, we considered a small artificial dataset to test the processing pipeline in a context in which we knew the ground truth. In the second, we focused on the German credit dataset, used in multiple researches on fairness [7, 8, 11, 13]. This dataset contains 8 attributes (both numerical and categorical) of 1000 individuals, including an association of each individual to a credit risk score (good or bad). For the sake of brevity, we will limit our focus here on the German credit dataset. All Python-code that we run for obtaining the results is publicly available.³

Table 1. Stability analysis of the cluster combinations for all three clustering algorithms applied on the German credit data set with varying parameter settings. *For the semi-hierarchical clustering the number of clusters is not fixed, we reported the number of clusters for the cluster combination that occurs most.

Algorithm	l	k	# Clusters	# Cluster comb	Freq. most occurring comb	Running time
Simple	1	20	20	15	43	3.956 ± 0.376 s
Simple	1	30	30	50	9	8.170 ± 0.777 s
Simple	1	40	40	17	26	12.068 ± 0.875 s
Hierarchical	2	4	16	31	24	1.119 ± 0.078 s
Hierarchical	2	5	25	21	28	1.152 ± 0.077 s
Hierarchical	2	6	36	18	32	1.298 ± 0.113 s
Hierarchical	3	3	27	54	12	1.144 ± 0.047 s
Hierarchical	4	2	16	18	29	1.450 ± 0.083 s
Hierarchical	5	2	32	39	20	1.614 ± 0.071 s
Semi-hierarchical	–	–	34*	47	23	2.264 ± 0.112 s

4.1 Evaluation of Clustering for Profile Selection

We evaluate the three clustering algorithm for profile selection specified in Sect. 3.1 following the method described in Sect. 3.2. For all three clustering algorithms, we used a k -medoids clustering algorithm with the Gower distance [3, 9]. The simple algorithm and the hierarchical clustering require the tuning of hyperparameters as k (number of clusters), and l (number of layers), and therefore have a fixed number of final clusters (in this context seen as profiles). Since we do not know the correct number of profiles in advance, we tried several hyperparameters. The semi-hierarchical clustering algorithm needs instead three other parameters: c_u , c_l (upper and lower bounds of variability), and n_{stop} (the threshold for cluster cardinality). For our experiments, we chose to set $c_u = 0.9$,

² Experimental setup: Intel Core i7-10510u, 16 GB RAM, Windows-10 64-bit.

³ <https://github.com/mcwilms/PEBAM>.

$c_l = 0.1$, and n_{stop} to 5% of the size of the dataset (i.e. 50 for the German credit dataset).

As a first step, we test if clustering methods are *stable* enough. Table 1 gives a summary of the stability analysis performed for all three clustering methods on the German credit dataset, reporting (when relevant) the parameters l and c , the number of clusters (e.g. k^l), the number of different outcomes after running the clustering algorithm 100 times with different random initializations, how often the most occurring cluster combinations occurs in these 100 runs, and the mean running time. Amongst other things, Table 1 shows that, when running the stability analysis with semi-hierarchical clustering, the German credit dataset produces 47 different cluster combinations. However, several combinations occur only once, and only one of the combination (number 0, the first one) occurs significantly more often than the other combinations, see Fig. 2a.

As a second step, we compute the *inter-clustering similarity* to test if the different cluster combinations are adequately similar. Figure 2b shows the inter-clustering similarity of the different cluster combinations we obtain on the German credit dataset via the semi-hierarchical clustering algorithm, showing only the cluster combinations that occur more than once. A dark blue tile means that two clusters are very similar (max 100%), and a white tile means that they have 50% or less of the clusters in common.

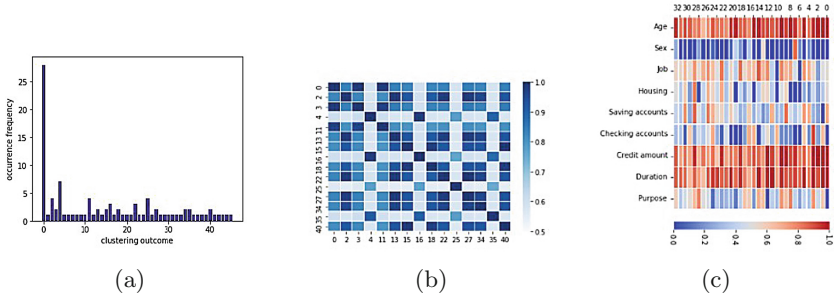


Fig. 2. Stability and variability analysis of the semi-hierarchical clustering algorithm applied on the German credit dataset: (a) frequency of each clustering outcome (or cluster combination) obtained over 100 runs; (b) inter-clustering similarity for cluster combinations that occur more than once; (c) variability of the different features for each profile in the most recurrent cluster combination.

As a confirmation that the algorithms end up on profiles which are descriptive attractors, we compute the feature *variability* of each cluster (supposedly a profile) within the most occurring clustering outcome. The variability plot of Fig. 2c shows for each profile (columns) the variability of the features (rows), where dark blue indicates a low variability and dark red a high variability. In tendency, qualitative features becomes stable, whereas numerical features show

still a certain variability at the level of profiles. For each profile, however, the majority of features becomes stable.

4.2 Profile-Based Evaluation of Bias

We now apply the bias evaluation methods described in Sect. 3.3 with the profiles obtained by applying the semi-hierarchical clustering algorithm on classifiers trained on the German credit data set via three commonly used machine learning algorithms: *logistic regression classifier*, *XGboost classifier*, and *support vector machine (SVM) classifier* (e.g. [2]).⁴ Following the standard practice of removing protected attributes (gender, ethnicity, etc.) as input features during training, we do not use the feature “Sex” provided by the German credit dataset for training the classifiers.

For **Representational bias**, Fig. 3 gives an overview of the presence of the identified profiles within the German credit dataset. We see that not all profiles are equally frequent; this is not necessarily an error, as long as this profile distribution is a good representation of the data on which ADM will be applied in practice. Expert knowledge or actual data collection can be used to test this assumption.

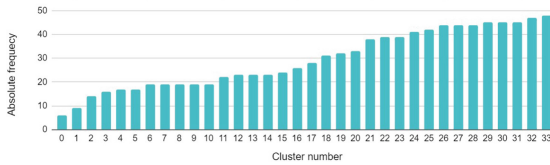


Fig. 3. Absolute frequencies of profiles obtained after performing the semi-hierarchical clustering algorithm on the German credit dataset.

To assess **Historical bias**, we first test for (general) *demographic disparity* with respect to protected attributes. A disadvantaged group with attribute x (DG_x) is the group with risk-label ‘bad’ and attribute x , whereas an advantaged group with attribute x (AG_x) is the group with risk-label ‘good’ and attribute x . Denoting with A the proportion of people from the full dataset belonging to the protected class (female) in the advantaged group over all people (female and male) in the advantaged group, and D for the proportion of people belonging to the protected class in the disadvantaged group over all people, we find:

$$D = \frac{\#DG_f}{\#DG_{f+m}} = 0.36 > 0.29 = \frac{\#AG_f}{\#AG_{f+m}} = A$$

and hence we conclude that, at aggregate level, the German credit data exhibit demographic disparity. We now will do the same computation for each profile

⁴ Note however that the same approach will apply with any choice of profile selection method or of ML method used to train the classifier.

subset of the German credit dataset, to test whether there is *demographic disparity for profiles*. We write A_c for the proportion of people belonging to the protected class (female) in the advantaged group of cluster c over all people in the advantaged group of cluster c , and D_c for the proportion of people belonging to the protected class in the disadvantaged group of cluster c over all disadvantaged people in cluster c . Table 2 shows that there are 3 profiles (7, 20 and 28) that show demographic disparity. The fact that profiles show demographic disparity indicates that it might be possible that for some profiles other (not-protected) attributes correlate with the protected attribute, and so the protected attribute can indirectly be used in the training-process of the model.

In the computation of the (dis)-advantage fraction of Table 2 we still looked at the protected group *female*, however, we can also compute the measures A_c^* (D_c^*) as the fraction of (dis)advantaged individuals in a profile c over the total individuals within that profile (without distinguishing the protected class in it):

$$A_c^* = \frac{\#AG_c}{\#AG_c + \#DG_c} \quad D_c^* = \frac{\#DG_c}{\#AG_c + \#DG_c}$$

By doing so, we get an indication of how informative a profile is for belonging to the (dis)advantaged group. Table 3 shows the fraction of advantaged and disadvantaged individuals for each profile. Note that there are profiles for which the majority of the samples is clearly advantaged (e.g. 0, 1, 2, ...), a few have some tendency towards disadvantaged outcomes (e.g. 3, 15), but in comparison could be put together with other profiles that have no clear majority (e.g. 9, 10, ...). Plausibly, for profiles exhibiting a mixed distribution of the risk label, there may be factors outside the given dataset that determine the label. Since the ADM models also do not have access to these external features, it may be relevant to evaluate performance on these profiles to evaluate this hypothesis.

Table 2. Fractions of disadvantaged (D) and advantaged (A) individuals with protected attribute *female* in each profile c .

c	D_c	A_c	c	D_c	A_c	c	D_c	A_c	c	D_c	A_c	c	D_c	A_c
0	0.00	0.00	7	0.36	0.25	14	0.00	0.19	21	0.00	0.03	28	0.96	0.90
1	0.00	0.00	8	0.00	0.00	15	1.00	1.00	22	0.00	0.00	29	0.00	0.00
2	0.00	0.00	9	0.00	0.00	16	1.00	1.00	23	0.00	0.00	30	1.00	1.00
3	0.00	0.00	10	1.00	1.00	17	1.00	0.86	24	0.00	0.00	31	1.00	1.00
4	0.00	0.07	11	0.00	0.00	18	0.00	0.00	25	0.00	0.03	32	0.00	0.00
5	1.00	1.00	12	0.00	0.00	19	0.06	0.06	26	1.00	1.00	33	0.00	0.09
6	0.00	0.08	13	0.00	0.00	20	0.24	0.00	27	0.00	0.00			

Table 3. Fractions of disadvantaged (D^*) and advantaged (A^*) individuals in each profile c .

c	D^*	A^*	c	D^*	A^*	c	D^*	A^*	c	D^*	A^*	c	D^*	A^*
0	0.00	1.00	7	0.58	0.42	14	0.09	0.91	21	0.05	0.95	28	0.52	0.48
1	0.11	0.89	8	0.58	0.42	15	0.62	0.38	22	0.36	0.64	29	0.40	0.60
2	0.14	0.86	9	0.47	0.53	16	0.38	0.62	23	0.15	0.85	30	0.09	0.91
3	0.62	0.38	10	0.53	0.47	17	0.25	0.75	24	0.37	0.63	31	0.22	0.78
4	0.12	0.88	11	0.18	0.82	18	0.19	0.81	25	0.20	0.71	32	0.13	0.87
5	0.12	0.88	12	0.35	0.65	19	0.50	0.50	26	0.45	0.55	33	0.06	0.94
6	0.32	0.68	13	0.52	0.48	20	0.52	0.48	27	0.14	0.86			

For the **Aggregation bias** we look at the blue dots in Fig. 4a, which indicate the mean performances of the algorithm over training the algorithm 100 times on different train-test splits. Looking at performance over each profile gives us a visual way to see to what extent general equal opportunity (Def. 4) is satisfied; we consider the average to provide a more robust indication. We see that the XGboost classifier performs the best of the three algorithms with respect to predicting the labels correctly, however we also observe some difference in performance depending on profile. In contrast, the SVM classifier has very low probabilities of getting an unjustified disadvantage label (Fig. 4b), while the probability of getting a correct label is not very high.

For the **Evaluation bias**, we look at the performance ranges of the different classification methods (visualized in terms of standard deviations). We see that the SVM classifier is the least sensitive towards the train-test split. The logistic regression classifier is already slightly more sensitive, however the XGboost classifier is by far the most sensitive towards the train-test split. All three algorithms are equally sensitive towards small profiles as much as larger profiles.

5 Conclusion

The paper introduced PEBAM: a new method for evaluating biases in ADM models trained on mixed datasets, focusing in particular on profiles extracted through a novel (semi-hierarchical) clustering method. Although we have proven the feasibility of the overall pipeline, several aspects need further consolidation, as for instance testing other measures of variability (e.g. to be compared with entropy-based forms of clustering, e.g. [6]), similarity scores, and distance measures. Yet, the method was already able to find biases that were not revealed by most used bias evaluation methods, since they would not test for biased decisions against groups of individuals that are regrouped by non-protected attributed values only. For instance, profile 7, exhibiting demographic disparity against women as historical bias, refers to applicants with little saving/checking accounts and renting their house, who are asking credit for cars (see Appendix for details).

Why, *ceteris paribus* (all other things being the same), men are preferred to women for access to credit for buying cars, if not in presence of a prejudice?

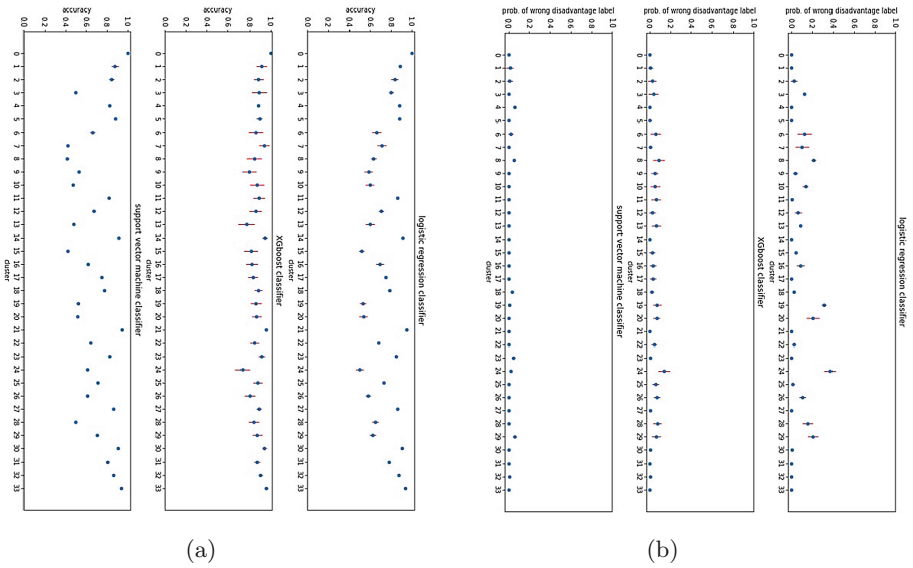


Fig. 4. (a) Probability of getting a correct prediction label. (b) Probability of getting a disadvantage (‘bad’) label when the true label is the advantage (‘good’), for *logistic regression classifier*, *XGboost classifier*, and *SVM classifier* on the different profiles.

At a technical level, although the proposed semi-hierarchical clustering algorithms has shown a shorter running time than the baseline on the German credit dataset, the PAM algorithm does not scale well to larger datasets. Tiwari et al. propose BanditPAM as alternative for PAM [17], a method that reduces the complexity of each PAM iteration from $\mathcal{O}(n^2)$ to $\mathcal{O}(n \log n)$. When using PEBAM on large datasets one might consider using BanditPAM over PAM. This will be investigated in future work.

Acknowledgments. Giovanni Sileno was partly funded by the Dutch Research Council (NWO) for the HUMAINER AI project (KIVI.2019.006).

A Profiles on the German Credit Dataset

The following table reports the profiles selected on the German credit dataset by applying the semi-hierarchical clustering proposed in the paper, as described by their medoids:

Profile	Age	Sex	Job	Housing	Saving accounts	Checking account	Credit amount	Duration	Purpose	Sample
0	26	Male	2	Rent	Moderate	Unknown	3577	9	Car	859
1	37	Male	2	Own	Unknown	Unknown	7409	36	Business	868
2	39	Male	3	Own	Little	Unknown	6458	18	Car	106
3	26	Male	2	Own	Little	Little	4370	42	Radio/TV	639
4	31	Male	2	Own	Quite rich	Unknown	3430	24	Radio/TV	19
5	38	Female	2	Own	Unknown	Unknown	1240	12	Radio/TV	135
6	43	Male	1	Own	Little	Little	1344	12	Car	929
7	36	Male	2	Rent	Little	Little	2799	9	Car	586
8	39	Male	2	Own	Little	Little	2522	30	Radio/TV	239
9	31	Male	2	Own	Little	Moderate	1935	24	Business	169
10	33	Female	2	Own	Little	Little	1131	18	Furniture/equipment	166
11	26	Male	1	Own	Little	Moderate	625	12	Radio/TV	220
12	23	Male	2	Own	Unknown	Moderate	1444	15	Radio/TV	632
13	42	Male	2	Own	Little	Little	4153	18	Furniture/equipment	899
14	29	Male	2	Own	Unknown	Unknown	3556	15	Car	962
15	37	Female	2	Own	Little	Moderate	3612	18	Furniture/equipment	537
16	27	Female	2	Own	Little	Little	2389	18	Radio/TV	866
17	26	Female	2	Rent	Little	Unknown	1388	9	Furniture/equipment	582
18	29	Male	2	Own	Little	Unknown	2743	28	Radio/TV	426
19	53	Male	2	Free	Little	Little	4870	24	Car	4
20	36	Male	2	Own	Little	Little	1721	15	Car	461
21	38	Male	2	Own	Little	Unknown	804	12	Radio/TV	997
22	29	Male	2	Own	Little	Moderate	1103	12	Radio/TV	696
23	43	Male	2	Own	Unknown	Unknown	2197	24	Car	406
24	27	Male	2	Own	Little	Little	3552	24	Furniture/equipment	558
25	30	Male	2	Own	Little	Moderate	1056	18	Car	580
26	24	Female	2	Own	Little	Moderate	2150	30	Car	252
27	34	Male	2	Own	Little	Unknown	2759	12	Furniture/equipment	452
28	24	Female	2	Rent	Little	Little	2124	18	Furniture/equipment	761
29	34	Male	2	Own	Little	Moderate	5800	36	Car	893
30	34	Female	2	Own	Little	Unknown	1493	12	Radio/TV	638
31	30	Female	2	Own	Little	Unknown	1055	18	Car	161
32	35	Male	2	Own	Little	Unknown	2346	24	Car	654
33	35	Male	2	Own	Unknown	Unknown	1979	15	Radio/TV	625

References

1. Allaj, E.: Two simple measures of variability for categorical data. *J. Appl. Stat.* **45**(8), 1497–1516 (2018)
2. Belle, V., Papantonis, I.: Principles and practice of explainable machine learning. *Front. Big Data* **4** (2021)
3. Ben Ali, B., Massmoudi, Y.: K-means clustering based on Gower similarity coefficient: a comparative study. In: 2013 5th International Conference on Modeling, Simulation and Applied Optimization (ICMSAO), pp. 1–5. IEEE (2013)
4. Budiaji, W., Leisch, F.: Simple k-medoids partitioning algorithm for mixed variable data. *Algorithms* **12**(9), 177 (2019)
5. Caruso, G., Gattone, S., Fortuna, F., Di Battista, T.: Cluster analysis for mixed data: an application to credit risk evaluation. *Soc.-Econ. Plan. Sci.* **73**, 100850 (2021)
6. Cheng, C.H., Fu, A.W., Zhang, Y.: Entropy-based subspace clustering for mining numerical data. In: Proceedings of the 5th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining, KDD 2009, pp. 84–93. Association for Computing Machinery, New York (1999)
7. Feldman, M., Friedler, S., Moeller, J., Scheidegger, C., Venkatasubramanian, S.: Certifying and removing disparate impact. In: Proceedings of the 21th ACM

- SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2015, pp. 259–268. ACM (2015)
8. Friedler, S., Scheidegger, C., Venkatasubramanian, S., Choudhary, S., Hamilton, E., Roth, D.: A comparative study of fairness-enhancing interventions in machine learning. In: Proceedings of the Conference on fairness, accountability, and transparency, FAT 2019, pp. 329–338. ACM (2019)
 9. Gower, J.C.: A general coefficient of similarity and some of its properties. *Biometrics* **27**(4), 857–871 (1971)
 10. Hardt, M., Price, E., Price, E., Srebro, N.: Equality of opportunity in supervised learning. In: Advances in Neural Information Processing Systems (NIPS) (2016)
 11. Kamiran, F., Calders, T.: Classifying without discriminating. In: Proceedings of 2nd IEEE International Conference on Computer, Control and Communication (2009)
 12. Kamiran, F., Žliobaitė, I., Calders, T.: Quantifying explainable discrimination and removing illegal discrimination in automated decision making. *Knowl. Inf. Syst.* **35**(3), 613–644 (2013)
 13. Kamishima, T., Akaho, S., Sakuma, J.: Fairness-aware learning through regularization approach. In: 2011 IEEE 11th International Conference on Data Mining Workshops, pp. 643–650. IEEE (2011)
 14. Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., Galstyan, A.: A survey on bias and fairness in machine learning. *ACM Comput. Surv.* **54**(6), 1–35 (2021)
 15. Pleis, J.: Mixtures of Discrete and Continuous Variables: Considerations for Dimension Reduction. Ph.D. thesis, University of Pittsburgh (2018)
 16. Suresh, H., Gutttag, J.: A framework for understanding sources of harm throughout the machine learning life cycle. In: Equity and Access in Algorithms, Mechanisms, and Optimization, EAAMO 2021 (2021)
 17. Tiwari, M., Zhang, M.J., Mayclin, J., Thrun, S., Piech, C., Shomorony, I.: Banditpam: almost linear time k-medoids clustering via multi-armed bandits. In: Advances in Neural Information Processing Systems (NIPS) (2020)
 18. Wachter, S., Mittelstadt, B., Russell, C.: Why fairness cannot be automated: bridging the gap between EU non-discrimination law and AI. *Comput. Law Secur. Rev.* **41**, 105567 (2021)