

Data Preprocessing Techniques for Classification with Discrimination

Faisal Kamran - Toon Calders / 2021

- Problem: Training data exhibit unlawful discrimination
- Task:
 - Optimize accuracy
 - Classifier doesn't have discrimination on test data
- Method: Preprocess the data to remove discrimination
- Application: A binary sensitive attribute + 2-class classification problem



- Discrimination: Unfair or unequal treatment of individuals of a certain group based solely on their affiliation to that group

• Discrimination-aware classification

- Input: A labeled dataset with one or more sensitive attributes
- Output: A classifier that doesn't correlate with the sensitive attribute
- Quality measurement: Accuracy and discrimination

• Discrimination metric:

- Sensitive attribute S
- Target attribute X
- b: Deprived community
w: Favored community
- +: Desirable outcome
-: Non-desirable outcome

→ Discrimination of a classifier (evaluation on the test set!)

$$\text{disc}_{S=b} = P(C(X) = + | X(S) = w) - P(C(X) = + | X(S) = b)$$

↳ The deprived and the favored communities should have the same proportion of desirable outcome for the classifier to be discrimination free ($\text{disc}_{S=b} = 0$)

⚠️ Vocab: In the paper, a classifier is a classifier that only outputs the classes $\rightarrow +$
a ranker is a classifier that outputs probabilities $\rightarrow -$

Methods for preprocessing the data (non-detailed version)

- **Suppression:** Removes and the attributes correlated to it → Establish a baseline
- **Managing the data:** Change the labels of some data points to remove discrimination
A classifier/ranker is used to choose the data points whose labels we're going to change
- **Reweighing:** Weights are assigned to each data point → used for training a classifier
- **Sampling:** To be used when reweighing's not possible.
Resample the combinations of $S = \{w, b\}$ and $\text{Class} = \{+, -\}$
Till making the dataset discrimination free

$$\begin{cases} \cdot b+ : \text{upsample} \\ \cdot b- : \text{downsample} \end{cases} \quad \begin{cases} \cdot w+ : \text{downsample} \\ \cdot w- : \text{upsample} \end{cases}$$

2 proposed sampling methods → Uniform sampling (US)
↳ Preferred sampling (PS)

Discrimination in a dataset

$$\hookrightarrow \text{disc}_{S=b} = P(X(\text{class})=+) \times (S=w) - P(X(\text{class})=+) \times (S=b) \quad \} \text{ can be negative as well !!}$$

The prob of having a + outcome The prob of having a positive outcome
for someone in the favored group for someone in the deprived group

→ For discrimination of a classifier, the \neq is :

→ Instead of $X(\text{class}) \rightarrow C(x)$, C = the classifier
→ The evaluation is done on the test set

→ Discrimination-aware classifier ↳ The accuracy of the classifier C is high
 ↳ The discrimination of the classifier is less

↳ 3 assumptions made in the article:

→ A1: The goal is to learn a classifier ↳ highest accuracy
→ A2: Don't use the sensitive attribute at prediction time

→ A3: The total ratio of + predictions of C should be similar to that of D

• Some theory

$\rightarrow C \text{ dominates } C' \text{ when } \begin{cases} \text{Accuracy}(C) > \text{Accuracy}(C') \\ \text{Disc}(C) \leq \text{Disc}(C') \end{cases}$

$\rightarrow C \text{ strictly dominates } C' \text{ when one of these 2 inequalities is strict}$

$\rightarrow \text{DA-optimality}$: A classifier C in a set of classifiers \mathcal{C} is DA-optimal if there is no other classifier in \mathcal{C} that strictly dominates C

In other words, there should be no classifier with both a better accuracy and a lower discrimination than C (where one of them is strict)

• Some notation

$\mathcal{C}_{\text{all}} = \text{set of all classifiers}$ Same ratio of + cases
 $\mathcal{C}_{\text{opt}}^* = \text{set of all classifiers with } P(C(x)=+ | x \in D) = P(x/\text{class})=+ | x \in D$

• Accuracy-Discrimination trade-off:

A perfect classifier $C \rightarrow C^{\text{perf}}(x) = x(\text{class}), x \in D$

Theorem 1 A classifier C is DA-optimal in \mathcal{C}_{all} iff

$$\text{acc}(C^{\text{perf}}) - \text{acc}(C) = \frac{\min(d_b, d_w)}{d} (\text{disc}(C^{\text{perf}}) - \text{disc}(C))$$

A classifier C is DA-optimal in $\mathcal{C}_{\text{all}}^*$ iff

$$\text{acc}(C^{\text{perf}}) - \text{acc}(C) = 2 \frac{d_b}{d} \frac{d_w}{d} (\text{disc}(C^{\text{perf}}) - \text{disc}(C))$$

\Rightarrow Discrimination-accuracy
Tradeoff for DA-optimal
classifiers is linear

- Imperfect classifiers \Rightarrow The relationship remains linear
- Using a ranker allows for a relationship that is sub-linear

- Solutions: Data preprocessing techniques

1. Massaging → changing the labels to get a discrimination of 0
 ↳ is a 2-step method

- Step 1: Define M , the # of data points whose labels have to be changed for each group ($b \& w$)

$$M = \frac{\text{Disc}(D) \times |D_b| \times |D_w|}{D}$$

▷ M depends only on D

▷ The distribution of labels doesn't change as the same # of points change labels + ↔ -

- Step 2: Learn a ranker to choose which points to modify in each group ($b \& w$)

We define two sets of points that are considered for label change

$$\begin{aligned} \hookrightarrow \text{promotion candidates (pr)} &= \{x \in D \mid x(\text{class}) = -8 \times S = b\} \\ \hookrightarrow \text{demotion candidates (dem)} &= \{x \in D \mid x(\text{class}) = +8 \times S = w\} \end{aligned}$$

- We build a ranker (R)

↳ We order pr in ascending order

↳ We order dem in descending order

↳ The points at the top of each class are close to the borderline,

we take M points of each set and we change their labels

Having the def in the ranker, D is discrimination free ^_~

2. Reweighting → Weights are attached to the dataset to make it disc-free

- Idea: In a discrimination-free dataset, the sensitive attribute and the label should be independent*

expected $P_{\text{exp}}(S=b, X(\text{Class})=+) = P(S=b) \cdot P(X(\text{Class})=+)$

If the dataset is disc-free, the observed prob should = P_{exp}

We define weights for all cases: (+, b), (+, w), (-, b), (-, w)

$$W(x) = \frac{P_{\text{exp}}(S=x(s), \text{Class}=x(\text{Class}))}{\text{Prob}(S=x(s), \text{Class}=x(\text{Class}))}$$

⚠ Reweighting only works with models that can use weights in the training

⚠ The authors define discrimination as the absence of correlation between the sensitive attribute and the label. Nonetheless, when talking about reweighting, they use independence!!!

3. Sampling → Instead of reweighting, over and under sample datapoints as not all models can use weights

This method is a 2-step method

Step 1: Calculate the weight for each case and multiply the # of occurrences of each case by its corresponding weight to get the new # of samples

$$W(s, c) = \frac{|\{x \in D | x(s)=s\} \times |\{x \in D | x(\text{Class})=c\}|}{|D| \times |\{x | s=x, x(\text{Class})=c\}|}$$

• Step 2: Sample till getting the expected # of samples, using either uniforming sampling (US) or preferential Sampling (PS)

Uniform Sampling (US): Sample data points with replacement from each class randomly until getting the expected # of samples

Preferential Sampling (PS):

- A rank is used to rank data points among \neq classes $(+, b), (+, w), (-, b), (-, w)$

PS starts from the original training dataset and iteratively duplicates (for the groups DP and FN) and removes objects (for the groups DN and FP) in the following way:

- Decreasing the size of a group is always done by removing the data objects closest to the boundary; i.e., the top elements in the ranked list.
- Increasing the sample size is done by duplication of the data object closest to the boundary. When an object has been duplicated, it is moved, together with its duplicate, to the bottom of the ranking. We repeat this procedure until the desired number of objects is obtained.



Conclusions:

- Removing the sensitive attribute from the data is not enough
- Massaging the data and Preferential Sampling usually work best
- The preprocessing works best when using tree-based classifiers down the line (decision trees, KNN with a small K)