

---

# On the reproducibility of "Exacerbating Algorithmic Bias through Fairness Attacks"

---

Anonymous Author(s)

Affiliation

Address

email

## Reproducibility Summary

1

### 2 **Scope of Reproducibility**

3 The paper presents two novel kinds of adversarial attacks against fairness: the IAF attack and the anchoring attacks.  
4 Our goal is to reproduce the five main claims of the paper. The first claim states that using the novel IAF attack we can  
5 directly control the trade-off between the test error and fairness bias metrics when attacking. Claims two to five suggest  
6 a superior performance of the novel IAF and anchoring attacks over the two baseline models. We also extend the work  
7 of the authors by implementing a different stopping method, which changes the effectiveness of some attacks.

### 8 **Methodology**

9 To reproduce the results, we use the open-source implementation provided by the authors as the main resource, although  
10 many modifications were necessary. Additionally, we implement the two baseline attacks which we compare to the  
11 novel proposed attacks. Since the assumed classifier model is a support vector machine, it is not computationally  
12 expensive to train. Therefore, we used a modern local machine and performed all of the attacks on the CPU.

### 13 **Results**

14 Due to many missing implementation details, it is not possible to reproduce the original results using the paper alone.  
15 However, in a specific setting motivated by the authors' code (more details in section 3), we managed to obtain results  
16 that support 3 out of 5 claims. Even though the IAF and anchoring attacks outperform the baselines in certain scenarios,  
17 our findings suggest that the superiority of the proposed attacks is not as strong as presented in the original paper.

### 18 **What was easy**

19 The novel attacks proposed in the paper are presented intuitively, so even with the lack of background in topics such as  
20 fairness, we managed to easily grasp the core ideas of the paper.

### 21 **What was difficult**

22 The reproduction of the results requires much more details than presented in the paper. Thus, we were forced to make  
23 many educated guesses regarding classifier details, defense mechanisms, and many hyperparameters. The authors  
24 also provide an open-source implementation of the code, but the code uses outdated dependencies and has many  
25 implementation faults, which made it hard to use as given.

### 26 **Communication with original authors**

27 Contact was made with the authors on two occasions. First, we asked for some clarifications regarding the provided  
28 environment. They promptly replied with lengthy answers, which allowed us to correctly run their code. Then, we  
29 requested additional details concerning the pre-processing of the datasets. The authors pointed at some of their previous  
30 projects, where we could find further information on the processing pipeline.

## 31 1 Introduction

32 Machine Learning models have shown impressive performance in countless domains in the last decade. However, it has  
33 been demonstrated that an adversary can input carefully-crafted perturbations to subvert the predictions of these models.  
34 The area of Adversarial Machine Learning has emerged to study vulnerabilities of machine learning approaches in  
35 adversarial settings and to develop techniques that make them robust against malicious attacks.

36 Most of the research has focused on studying malign interventions that degrade the accuracy of a system: imagine, for  
37 example, the consequences of inducing wrong predictions in an autonomous driving system. Only recently, fairness  
38 has become a rising concern for the performance of machine learning models, especially for sensitive fields such as  
39 criminal justice and loan decisions. Along these lines, “Exacerbating Algorithmic Bias through Fairness Attacks” [1]  
40 proposes two families of poisoning attacks that inject malicious points into the models’ training sets and intentionally  
41 target the fairness of a classification model.

42 The first, the *influence* attack, extends the optimization-based technique introduced by Koh et al. [2] by incorporating  
43 in the loss function a constraint for fair classification. An attacker can hence harm both accuracy and fairness  
44 simultaneously, with a trade-off regularized via a parameter  $\lambda$ . The second type of attack, the *anchoring* attack, affects  
45 solely fairness and aims to place poisoned data points to bias the decision boundary without modifying the attacker loss.  
46 Depending on whether the target point is chosen at random, anchoring attacks are classified as *random* or *non-random*.

## 47 2 Scope of reproducibility

48 This report investigates the reproducibility of the original paper by Mehrabi et al. and aims to verify its main claims.  
49 Since these heavily rely on the datasets and metrics used by the authors, the reader is invited to consult Sections 3.2 and  
50 3.3 – respectively – for a refresh of such concepts. Then, the main claims can be summarized as follows:

### 51 – *Influence Attack on Fairness (IAF)*:

- 52 • *Claim 1*: Increasing the parameter  $\lambda$  results in stronger attacks against fairness. Contrarily, for lower values  
53 the model acts similarly to the original influence attack [2] targeted towards accuracy;
- 54 • *Claim 2*: The proposed IAF outperforms the attack of Koh et al. [2] in affecting both fairness metrics (SPD  
55 and EOD), on all three datasets;
- 56 • *Claim 3*: The proposed IAF also outperforms the attack based on the loss function proposed by Solans et al.  
57 [3] in affecting SPD and EOD, on all tested datasets.

### 58 – *Anchoring Attack*:

- 59 • *Claim 4*: Both random and non-random anchoring attacks (RAA and NRAA, respectively) outperform Koh  
60 et al. [2] in degrading the SPD and EOD of the classification model, on all three datasets;
- 61 • *Claim 5*: On the German and Drug Consumption datasets, RAA and NRAA have a greater impact on fairness  
62 metrics (SPD and EOD) compared to the attack based on Solans et al. [3]. However, the latter outperforms the  
63 proposed anchoring attack in affecting fairness when classification is performed on the COMPAS dataset.

## 64 3 Methodology

65 The authors provided an open-source implementation of their code on GitHub [4]. Unfortunately, the repository has  
66 several issues: dependencies are not sufficiently specified, and simply running the code in the given environment results  
67 in conflicts. Furthermore, the code does not provide an option to run baseline methods used in the paper, nor does it  
68 include the essential hyperparameter  $\lambda$ , which is used in the experiments. The majority of the code is based on Koh  
69 et al. [2]’s public implementation [5], and a code coverage analysis revealed that more than 50% is not used for running  
70 experiments related to this paper<sup>1</sup>. Moreover, the repository comes with pre-processed datasets and while this may  
71 sound advantageous, there is no mention of the processing procedure in the paper nor on GitHub. Finally, the code is  
72 generally complex and hard to understand due to insufficient comments and documentation.

73 Therefore, we used the codebase provided by the authors and customized it for our purposes. First, to aid maintainability  
74 and scalability, as well as to ensure future reproducibility of the original experiments, the code was modernized and  
75 made compatible with the latest version of *every* dependency. This involved major changes to migrate from TensorFlow

---

<sup>1</sup>The coverage.py tool [6] was used to measure code coverage, and the study was performed considering all possible attacks-datasets combinations.

76 1.12.0 to 2.6.2 and to update CVXPY from version 0.4.11 to 1.1.18<sup>2</sup>. Secondly, datasets were downloaded from the  
77 original sources [7, 8] and processed from scratch. The procedure is thoroughly reported in Section 3.2. Furthermore,  
78 the code was trimmed down to the essential, and the user was given the option to choose any of the available models and  
79 the corresponding parameters. Lastly, we added comprehensive documentation to make the code more interpretable.

### 80 3.1 Model descriptions

81 It appears that the authors of the original paper do not specify the model that they use for the given classification  
82 task. From the implementation details given in Koh et al. [2], as well as from [1]'s codebase, we assume the use  
83 of a Support Vector Machine (SVM) trained with a smooth hinge loss and L2 regularization (refer to further  
84 details). Additionally, the optimization algorithm is not indicated; we assumed it to be Newton's Conjugate Gradient  
85 (Newton-CG) method, as suggested by the codebase. Such a method is used for both the minimization of the parameters  
86 on the training set and the update step of the poisoned points (for attacks utilizing an adversarial loss). The gradient  
87 is computed using the full datasets, i.e., without using mini-batches. Although hardly recognizable, this follows the  
88 implementation of the original paper: from our interpretation of the code, it seems that the authors define a variable  
89 containing the size of the mini-batch size and the necessary functionality, but then never use it.

90 Our base algorithmic setup for the IAF, RAA, and NRAA attacks is described in the Methods section of the original  
91 paper. However, the authors omitted important details that we consequently had to assume based on more or less  
92 concrete evidence. First, an advantaged and disadvantaged group for the sensitive attribute (i.e., gender, as per the  
93 original work) has to be specified for all attacks. Since the rationale behind this choice does not seem to be included  
94 in the paper, we infer from the codebase that the authors did it automatically and deduced it from the datasets. More  
95 specifically, we assume that the advantaged group is chosen as the group with the highest ratio of data points with  
96 positive label  $\hat{y} = 1$ , regardless of the actual class label it corresponds to. This method is simple yet fallacious: for  
97 instance, it means that the group taking on the label "likely to perform a crime soon" more often (in the context of the  
98 COMPAS dataset) is considered "advantaged" in terms of the algorithm.

99 Secondly, for the computation of the feasible set using an anomaly detector we assume that the intersection of the  
100 Slab defense and the L2 defense was originally employed, as described in Koh et al. [2]. For reprojecting poisoned data  
101 points into the feasible set, we again use the approach of [2], which incorporates LP rounding for discrete variables.

102 Moreover, we implement two baselines. The three proposed attacks are compared against the original accuracy-targeting  
103 attack proposed by Koh et al. [2], and another attack that uses a loss function proposed by Solans [3], which  
104 targets fairness<sup>3</sup>. Lastly, the model-specific changes/improvements are presented below:

105 IAF. As mentioned before, we modified the code to include the hyperparameter  $\alpha$  which controls the trade-off between  
106 the accuracy and the fairness loss in the adversarial loss.

107 Koh attack. We were not able to find a way of running this baseline attack using the given codebase. We have decided  
108 to implement it from scratch, treating it as the limiting case of the IAF attack when  $\alpha = 0$  (meaning no fairness loss in  
109 the adversarial loss function). Consequently, it is not exactly as presented in the original Koh attack sampling, the  
110 initial poisoned points are not drawn from advantaged and disadvantaged groups, contrary to the IAF attack. However,  
111 we argue that equalizing the sampling method provides a stronger comparison between the two methods, as we alleviate  
112 the issue of the missing inductive bias from the original Koh inference attack.

113 Solans attack. This attack serves as the second baseline. We could not find it in the codebase, thus we implemented it  
114 by replacing the adversarial loss in the IAF attack with a weighted sum loss, as presented in [3]. Implementing this  
115 change posed a bigger issue than expected, due to the infeasibility of the TensorFlow-based implementation. Thus,  
116 major revisions were required.

### 117 3.2 Datasets

118 The authors provide compressed copies of the three real-world datasets used for their experiments – the German  
119 Credit Dataset [7], the COMPAS Dataset [8] and the Drug Consumption Dataset [7]. However, these are already  
120 pre-processed, and the processing procedure is not reported nor documented in the code. This constitutes an important  
121 reproducibility barrier, because raw datasets are not directly usable with the given codebase.

<sup>2</sup>In our repository we provide a YAML configuration file to quickly set up the required environment.

<sup>3</sup>For simplicity, we will refer to the inference attack presented in [2] as the Koh attack, and we will also refer to the attack presented in [3] as the Solans attack.

<sup>4</sup>The German Credit Dataset and the Drug Consumption Dataset can be downloaded from the UCI machine learning repository [7], while the COMPAS can be found in the corresponding GitHub repository [8].

122 In this section, we present our pre-processing pipeline, which was mainly determined by reverse engineering of the  
123 given les. Like the authors, we provide a set of les containing already-processed data to run our implementation,  
124 but we also include the scripts used to pre-process each dataset. Custom\_data\_preprocessing directory. Lastly,  
125 to run the attacks, we assume that the advantaged and disadvantaged groups are males and females respectively. We  
126 accordingly map them to 0 and 1 to create a group\_label binary array.

127 In the rest of this section, we outline our dataset-specific details of the pre-processing pipeline and the assumptions that  
128 were made for the sake of reproducibility of the original results.

129 German Credit Dataset. The dataset contains the credit profile of 1000 individuals with 20 attributes associated with  
130 each person. In our experiments, we use all of them, as in the original paper. The attributes are both numerical and categorical, and  
131 we assumed the original authors used one-hot representations to encode the latter. The assumption was based on an  
132 extensive study of the provided datasets, with particular attention to their shapes. We then autonomously standardize  
133 the data, as it is common practice in Machine Learning, and split the data into an 80-20 train and test split, as indicated  
134 in the original paper.

135 COMPAS Dataset. ProPublica's COMPAS dataset [1] contains information about 7214 defendants from Broward  
136 County. We use the features specified in Table 1 of the original paper. In this case, based on the provided dataset, we concluded that  
137 the authors must have used numerical label encoding to represent the categorical attributes. Finally, we standardize the  
138 data and split it into an 80-20 train and test split.

139 Drug Consumption Dataset. The dataset contains information about the drug consumption of 1885 individuals [2].  
140 We use the attributes indicated in Table 1 of the original paper. The pre-processing procedure is as follows: first, we  
141 binarize the categorical data linked to cocaine consumption into users and non-users. Intuitively, non-users should be  
142 mapped to 0 (and 1 in the opposite case), but an inspection of the provided data suggests that the authors reversed  
143 the mapping. We decided to adhere to their choice for the sake of reproducibility. Moreover, we suspect that the dataset  
144 was shuffled before splitting it into training and test sets. By doing so, we obtain similar results in the experiments.  
145 Finally, we standardize the data. The original processing of this dataset was particularly difficult to replicate, because  
146 contrary to what was reported in the paper, the authors did not follow an exact 80-20 train and test split. Rather, the two  
147 contained 1500 and 385 data points respectively.

148 To conclude, it is noteworthy that even the pre-processed datasets provided by the authors are not immediately usable:  
149 the position (specified as index) of the sensitive feature (i.e., gender) is different for each dataset and is only given for  
150 the German dataset in the running instructions. To account for this unnecessary confusion, our custom pre-processing  
151 procedure includes the moving of the gender column to the index, which is taken as default by the main function. In  
152 this way, we simplify the running instructions and make them coherent across datasets. Still, the user is given the ability  
153 to pass the sensitive feature index as an argument, to facilitate future experiments on different and untested data.

### 154 3.3 Metrics

155 The attacks are evaluated in terms of accuracy and fairness. Along with classification (test) error, the original paper  
156 uses two important metrics to evaluate the attack in terms of fairness: Statistical Parity Difference and Equality of  
157 Opportunity Difference.

158 Statistical Parity Difference. Statistical Parity Difference (SPD) was first introduced by Dworkin et al. [3] and is used  
159 to capture the predictive outcome differences between different (advantaged and disadvantaged) demographic groups.  
160 The mathematical formulation is reported in Equation 1.

$$161 \text{ SPD} = \mathbb{P}(\hat{Y} = +1 | \mathbf{x} \in \mathcal{D}_a) - \mathbb{P}(\hat{Y} = +1 | \mathbf{x} \in \mathcal{D}_d) \quad (1)$$

161 where  $\mathcal{D}_a$  denotes the advantageous group and  $\mathcal{D}_d$  denotes the disadvantageous group.

162 Equality of Opportunity Difference. Equality of Opportunity Difference (EOD) (Hardt et al. [4]) captures differences  
163 in the true positive rate between different (advantaged and disadvantaged) demographic groups. It is defined as shown  
164 in Equation 2.

$$165 \text{ EOD} = \mathbb{P}(\hat{Y} = +1 | \mathbf{x} \in \mathcal{D}_a; Y = +1) - \mathbb{P}(\hat{Y} = +1 | \mathbf{x} \in \mathcal{D}_d; Y = +1) \quad (2)$$

---

<sup>5</sup>The main author followed a similar pre-processing procedure in another project that is publicly available on their GitHub [5].

### 165 3.4 Experimental setup and hyperparameters

166 All experiments shown in this paper can easily be reproduced using our code, which is publicly available on [GitHub](#)<sup>6</sup>.  
167 There we also provide technical details on how to run experiments and test different attacks in various settings. In this  
168 section, however, we list some additional details necessary to replicate the exact setup.

- 169 • The original code constrains the maximum iterations of an attack to 10000 and uses early stopping to interrupt  
170 training if the accuracy on the test set does not decrease for a specific number of iterations, which is hardcoded to  
171 be 2. We follow this strategy but adapt it for our experiments. First, we implement early stopping on both accuracy  
172 and fairness, meaning that the user can also choose to stop training in the absence of changes in fairness. We utilize  
173 average fairness  $(SPD + EOD) = 2$  as the stopping criterion<sup>7</sup> since the two metrics have similar behavior and equal  
174 range  $[0, 1]$ . Then, we set the early stopping patience as a controllable hyperparameter.
- 175 • It is unclear from the paper how the best-performing model was selected by the authors. The code suggests the usage  
176 of the model after the last attack iteration and training of the model parameters. Instead, we decided to save the  
177 best-performing model on the test set according to the chosen stopping metric (average fairness or accuracy), to better  
178 reflect the actual best performance. By selecting the best model based on fairness, we hope to choose more relevant  
179 states of the poisoned data affecting the fairness metrics. We compare the results in Section 4.
- 180 • The computation of the feasible set and the reprojection of poisoned points onto it is handled as a convex optimization  
181 problem (see [\[1\]](#)). Since we upgraded CVXPY to its newest version, we can let the library select the most appropriate  
182 solver for the given problem, instead of specifying one (the authors of [\[1\]](#) seem to have used [SCS](#) however).
- 183 • Following the original implementation, we utilize the `train_ncg` optimizer of the `scipy` library [\[13\]](#) for the Newton-  
184 CG optimization. We comply with the choices of the authors and set the convergence threshold of the `fmin` optimizer  
185 to  $10^{-8}$ , and the maximum number of iterations to 100. We follow the implementation details specified in [\[1\]](#)  
186 computing the inverse Hessian-vector.
- 187 • During training, the temperature of the smooth hinge loss is chosen `1004` as found hardcoded in the original  
188 implementation. The value for the weight decay is set to 0.09 for all datasets (apart from the code of the authors, this  
189 assumption is also backed up by the main experiments of Koh [\[2\]](#)). The step size utilized in the IAF algorithm  
190 (and thus also in the Koh and Solans attack) is set to 1 for all experiments, as found in the codebase.
- 191 • The threshold of the anomaly detector (see [\[2\]](#)) is controlled by a hyperparameter named `percentile`, which  
192 specifies the percentage of the data left after applying the anomaly detector. We first experimented with a value of 95  
193 as suggested by Koh et al. [\[2\]](#) but, as this seemed to lead to some training failings, we settled on 90 (the default value  
194 given in the codebase).
- 195 • The number of injected poisoned points is proportional to the number of clean data points,  $|D_p| \propto |D_c|$   
196 (where  $D_c$  and  $D_p$  are the set of clean and poisoned data points respectively). The authors control such quantity  
197 by using the proportionality factor as a changeable parameter. Accordingly, we do the same and also make  
198 it a controllable parameter.
- 199 • After careful inspection and testing of the authors' code, the EOD metric calculation was found to be faulty and was  
200 consequently re-implemented. Our adaptation is based on the paper that originally proposed [\[3\]](#) and inspired by  
201 the implementation found in the `IAF360` library [\[14\]](#).
- 202 • Finally, the distance to original points in anchoring attacks was set to 0 for all experiments, as in the original paper.
- 203 • The random seed in all experiments was set to 1

### 204 3.5 Computational requirements

205 To give a complete overview of our experimental setup, we collect the average runtimes per iteration for different  
206 datasets and types of attacks. These are presented in Table 1. All models have been trained on a local machine with an  
207 AMD Ryzen 5 5600x CPU (6 cores, Base clock 3.7 GHz). Since the datasets are small, there is no need for more than  
208 4Gb of RAM. In this sense, training should be virtually possible on any entry-level PC.

## 209 4 Results

### 210 4.1 Results reproducing original paper

211 As stated in Section 2, our main claims were identified in the original paper. In our specific setting, we were able to  
212 reproduce three of these, as summarized in Table 2. In this section we elaborate on our reproduction results: first, in

<sup>6</sup>[https://anonymous.4open.science/r/MLRC2021\\_fairness\\_attack/](https://anonymous.4open.science/r/MLRC2021_fairness_attack/)

<sup>7</sup>In the rest of the paper, we might refer to it simply as fairness stopping metric.



| Attack | German dataset[s] | COMPAS dataset[s] | Drug dataset[s] | Claim   | Reproducible? |
|--------|-------------------|-------------------|-----------------|---------|---------------|
| IAF    | 0:870             | 0265              | 0312            | Claim 1 | Yes           |
| NRAA   | 1:123             | 10623             | 3678            | Claim 2 | Yes           |
| RAA    | 0:934             | 0306              | 0324            | Claim 3 | No            |
| Koh    | 0:474             | 0267              | 0201            | Claim 4 | Yes           |
| Solans | 0:862             | 0332              | 0262            | Claim 5 | No            |

Table 1: Average runtime per iteration for different attack types and dataset. All values are stated in units of seconds.

Table 2: Summary of the claims investigation under our specific setup.

Figure 1: Influence of  $\alpha$  on the different metrics for different attacks on the German dataset, using accuracy as the stopping criteria during training.

213 section 4.1.1 we show the effect of the hyperparameter  $\alpha$  on various metrics (Claim 1). In section 4.1.2 we compare the  
 214 newly proposed attacks and the baselines (Claims 2-5).

#### 215 4.1.1 Effect of $\alpha$ on the different metrics

216 To verify Claim 1, we conducted the same experiment as the authors. We run an IAF attack for each dataset using  
 217 different  $\alpha$  values and increasing  $\alpha$  to recreate Figure 3 of the original paper (see Appendix B.3, Fig. 8). However,  
 218 compared to the original experiment we test a larger range of  $\alpha$  values (from 0.0 to 2.0) to gain better insights into its  
 219 effects. As depicted in Figure 1, increasing  $\alpha$  does result in stronger attacks against fairness. Here we use the German  
 220 dataset and accuracy as the stopping metric, but similar trends were observed on the other datasets and using fairness  
 221 for early stopping. The plots are included in Appendix B.1 for the sake of completeness. Therefore in this specific  
 222 setup, we were able to reproduce the claim.

#### 223 4.1.2 Comparison between the proposed attacks and the baselines

224 To investigate Claims 2-5 we design an experiment that is heavily inspired by the work of the authors. We perform  
 225 each attack on each dataset,  $\alpha = 1$  and gradually increasing  $\alpha$  (from 0.0 to 1.0, with steps of 0.1), and repeat this  
 226 procedure for each stopping metric. The results essentially replicate Figure 2 of the original paper (as seen in Appendix  
 227 B.3, Fig. 7) and are collected in Figures 5 and 6 of Appendix B.2. However, to facilitate a comparative study between  
 228 the proposed attacks and the baselines, we average the metrics over  $\alpha$  values and report the results in Table 3. In this  
 229 way, we can base our observations on quantifiable measures instead of solely using visual inspection.  
 230 Assuming that the authors used accuracy as the early stopping criteria, the corresponding values in the table reveal that –  
 231 in this specific setting:

- 232 • Claim 2 is reproducible. On average, IAF has a much stronger influence on SPD and EOD compared to Koh's  
 233 attack, on all three datasets.
- 234 • Claim 3 is not reproducible, because Solan's attack outperformed IAF in affecting the EOD on the Compas dataset.
- 235 • Claim 4 is reproducible. NRAA and RAA were found to degrade the fairness metrics (SPD and EOD) more than  
 236 Koh's attack, on all three datasets.
- 237 • Claim 5 is not reproducible. Solans' attack had a greater impact on the SPD than NRAA on the German and greater  
 238 impact than NRAA on both SPD and EOD on the Compas dataset. It also has a greater impact on the EOD than the  
 239 RAA attack on the Compas dataset.

| Attack | German Dataset                         |           |           | Compas Dataset                         |           |           | Drug Dataset                           |           |           |
|--------|--|-----------|-----------|--|-----------|-----------|--|-----------|-----------|
|        | Test error                             | SPD       | EOD       | Test Error                             | SPD       | EOD       | Test error                             | SPD       | EOD       |
|        | (Stopping metric: Fairness / Accuracy) |           |           | (Stopping metric: Fairness / Accuracy) |           |           | (Stopping metric: Fairness / Accuracy) |           |           |
| IAF    | 0:40=0:47                              | 0:84=0:68 | 0:88=0:74 | 0:46=0:47                              | 0:83=0:75 | 0:87=0:77 | 0:43=0:45                              | 0:89=0:75 | 0:90=0:76 |
| NRAA   | 0:26=0:26                              | 0:26=0:25 | 0:36=0:33 | 0:41=0:42                              | 0:59=0:59 | 0:64=0:64 | 0:39=0:39                              | 0:53=0:53 | 0:53=0:53 |
| RAA    | 0:27=0:28                              | 0:24=0:17 | 0:36=0:19 | 0:47=0:47                              | 0:84=0:73 | 0:87=0:75 | 0:42=0:44                              | 0:66=0:55 | 0:68=0:57 |
| Koh    | 0:27=0:61                              | 0:17=0:08 | 0:13=0:12 | 0:45=0:53                              | 0:81=0:46 | 0:85=0:48 | 0:40=0:56                              | 0:56=0:26 | 0:56=0:29 |
| Solans | 0:40=0:48                              | 0:65=0:44 | 0:49=0:16 | 0:44=0:45                              | 0:76=0:73 | 0:83=0:78 | 0:40=0:56                              | 0:53=0:28 | 0:55=0:32 |

Table 3: Average metrics overvalues, obtained for each measure-attack combination and each dataset. We report one pair of values in each entry, corresponding to the two stopping criteria (average fairness and accuracy), and highlight the greatest one.

| Value                                       | German (Solans) | Drug (IAF) |
|---|-----------------|------------|
| Min. test accuracy                          | 0.465           | 0.506      |
| Avg. fairness at the point of min. accuracy | 0.229           | 0.822      |
| Actual max. average fairness                | 0.619           | 1.000      |

Table 4: Minimum accuracy, the value of the average fairness at the point of minimum accuracy, and maximum achievable average fairness of the plots of Figure 2.

Figure 2: Difference between the two stopping metrics (accuracy and average fairness) for the Solans attack on the German dataset (left), IAF attack on the Drug dataset (right).

## 4.2 Results beyond the original paper: using fairness as the early stopping metric

While the original codebase seems to use accuracy as the early stopping metric (and hence for selecting and saving the best model), we investigate the change in the results if fairness is used instead. The main motivation behind such an experiment lies in the assumption that interrupting training based on the fairness measures supposedly yields more relevant states of the poisoned data, effectively resulting in more efficient attacks against fairness. Since the SPD and EOD have similar behavior and equal range, we employ average fairness  $(SPD + EOD) = 2$  for the task at hand.

Figure 2 depicts the test accuracy and the average fairness over epochs for two different dataset-attack combinations. An analysis of the curves confirms that the maximum achievable average fairness is much greater than the same measure at the point of minimal accuracy (see Table 4). The same phenomenon is observed for the other dataset-attack combinations, as reported in Table 3: fairness undergoes a stronger degradation if average fairness is used to interrupt the training process and save the best model. This is reflected in the corresponding values of the fairness measures, which appear much higher compared to when accuracy is used.

## 5 Discussion

Our reproduction reveals that although the proposed methods represent valid novel attacks against the fairness of a model, they are not always superior to other methods in the literature. IAF showed important performance in terms of SPD and EOD degradation, but anchoring attacks were outperformed by the baseline models on multiple occasions. This result conflicts with the findings of the main paper (see Appendix B.3, Fig. 7) where the baselines are generally inferior to the proposed attacks. We had to make several assumptions to solve issues and inconsistencies between the original paper and corresponding implementation (many of which have already been mentioned throughout the report, but we systematically collect them in Appendix A). These assumptions are, by definition, uncertain and might have been the cause of the discrepant results. To better understand the source of discrepancy, we initially planned to perform an ablation study, which would have also unveiled more information regarding the model's behavior. This was ultimately not possible, given the time constraints and the contingencies encountered in the reproduction process.

In the remainder of this section, we elaborate on the main claims and our ability to reproduce them. We then present some personal reflections on the overall execution of the work and conclude with a summary and look into future works.

265 5.1 Discussion of the results

266 The first claim was found to be reproducible under our experimental setup, as we expected. The parameter  
267 specially designed to control the trade-off between accuracy and fairness, hence a rejection of the claim would have  
268 implied a major flaw in the core idea of the paper. The other claims focused on the comparison with the two baselines  
269 and, while the results presented in Section 4.1.2 are explicative enough, some remarks are still noteworthy.

270 In general, better statistics of the results would give us a clearer insight into the relative performance of the models.  
271 However, only four weeks were allocated for this project and we were unable to re-run the experiments with multiple  
272 seeds. For example, the Solans attack outperformed the IAF attack in terms of EOD metric on the Compas dataset  
273 (when using accuracy as the stopping method) and led to the non-reproducibility of 3. Yet, this difference is  
274 relatively small and a measure of uncertainty could potentially reverse our decision.

275 Furthermore, it was shown that the final fairness metrics can highly vary depending on the chosen stopping method.  
276 This is especially prominent for Claim 4, which was accepted under the assumption that accuracy was used for stopping  
277 and saving the best model. In reality, Koh attack outperforms NRAA on both Compas and Drug datasets in the terms of  
278 SPD/EOD metrics, if fairness is used instead. Since the validity of the claim depends on the stopping metric of choice,  
279 we argue that the claim is much weaker than originally proposed. Similarly, compare the IAF and the Koh attack in  
280 terms of fairness measures, using accuracy as the stopping criteria. On the Drug dataset, IAF's SPD/EOD metrics are  
281 respectively 2:89 / 2:62 higher than Koh's. This gap tightens if fairness is used: IAF's SPD/EOD metrics become  
282 1:022 / 1:024 higher. Although these numbers indicate the same result, we find the claim to be weaker than proposed,  
283 as the superior performance of the IAF attack is diminished by the use of a different stopping metric.

284 Finally it is important to notice the different behavior of the test accuracy and the average fairness (Fig. 2) used as  
285 stopping criteria. While the latter has a relatively high variance, the former is pretty constant, meaning that using  
286 fairness as the stopping metric does not result in significant variations in the model's accuracy. Contrarily, as empirically  
287 proved by our experiments, it can be highly beneficial for the fairness measures.

288 5.2 Reflection: What was easy? What was difficult?

289 The new methods presented in the paper were described both intuitively and formally, with a clear mathematical  
290 structure. The authors also provided figures to aid the intuition on how new attacks can affect decision boundaries,  
291 which allowed us to easily understand the core novel ideas presented in the publication.

292 However, it was not trivial to re-implement the proposed methods, because many details required for the implementation  
293 do not appear in the paper. The provided open-source implementation was ultimately hard to follow due to its convoluted  
294 organization, lack of documentation, poorly named functions/variables, and abundance of unused code. Even setting  
295 up a working environment using the authors-given dependencies took longer than one would expect, prompting us  
296 to get help from the authors. Eventually, the hope to aid future experiments motivated the decision to make the code  
297 compatible with up-to-date dependencies. This was one of the biggest struggles because the codebase heavily relies on  
298 packages that underwent major updates (TensorFlow and CVXPY).

299 The authors also provided pre-processed datasets. We spent a considerable amount of time trying to replicate their exact  
300 pipeline through reverse-engineering of the given files. Additionally, after recognizing some imperfections in the code  
301 and inconsistencies with the paper, we verified all of the existing implementation details to make sure that no further  
302 errors were made. This was a daunting task, given the complete lack of documentation and intuitive variable use.

303 5.3 Communication with original authors

304 To reiterate, we have initially contacted the main author to aid us with the dependency issues, who helped us with  
305 setting up a working environment. We then had additional contacts regarding the dataset pre-processing procedure. The  
306 author provided us with some indications on the pipeline and pointed at some useful resources. Eventually, we decided  
307 to gain a better understanding of the datasets through reverse-engineering.

308 5.4 Conclusion

309 In this paper, we have presented a reproducibility study of "Exacerbating Algorithmic Bias through Fairness Attacks",  
310 whereon we can draw some conclusions. Due to all the mentioned issues and inconsistencies (collected in Appendix A),  
311 we find it not possible to reproduce the original results from sole use of the paper, and difficult even in possession of  
312 the provided codebase. Yet, we managed to obtain similar findings that supported three out of the five main claims of  
313 the publication, albeit using partial re-implementations and numerous assumptions. Ascertaining the validity of such  
314 assumptions is therefore important for future works. Moreover, further studies could extend the classifier to work with  
315 multiple demographic groups and investigate the results using different fairness metrics.



## References

- 316
- 317 [1] Ninareh Mehrabi, Muhammad Naveed, Fred Morstatter, and Aram Galstyan. Exacerbating algorithmic bias  
318 through fairness attacks, 2020.
- 319 [2] Pang Wei Koh, Jacob Steinhardt, and Percy Liang. Stronger Data Poisoning Attacks Break Data Sanitization  
320 Defenses. arXiv:1811.00741, November 2018.
- 321 [3] David Solans, Battista Biggio, and Carlos Castillo. Poisoning attacks on algorithmic fairness, 2020.
- 322 [4] Ninareh Mehrabi, Muhammad Naveed, Fred Morstatter, and Aram Galstyan. Exacerbating algorithmic bias  
323 through fairness attacks. <https://github.com/Ninarehm/attack>, 2020. Accessed on: 23-01-2022.
- 324 [5] Pang Wei Koh, Jacob Steinhardt, and Percy Liang. Stronger Data Poisoning Attacks Break Data Sanitization  
325 Defenses. <https://github.com/kohpangwei/data-poisoning-journal-release>, 2018. Accessed on:  
326 23-01-2022.
- 327 [6] Ned Batchelder and other 103 contributors. Coverage. [coverage.readthedocs.io](https://coverage.readthedocs.io), 2021. Accessed on:  
328 23-01-2022.
- 329 [7] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.  
330 Last accessed on: 23-01-2022.
- 331 [8] Jeff Larson, Marjorie Roswell, and Vaggelis Atlidakis. Compas analysis. [https://github.com/propublica/  
332 compas-analysis/](https://github.com/propublica/compas-analysis/), 2016. Last accessed on: 23-01-2022.
- 333 [9] E. Fehrman, A. K. Muhammad, E. M. Mirkes, V. Egan, and A. N. Gorban. The ve factor model of personality  
334 and evaluation of drug consumption risk, 2017.
- 335 [10] Ninareh Mehrabi. Fairness- statistical equity: A fairness classification objective. [https://github.com/  
336 Ninarehm/Fairness](https://github.com/Ninarehm/Fairness), 2020. Accessed on: 23-01-2022.
- 337 [11] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. Zemel. Fairness through awareness. pages 214–226, 2012.  
338 URL <https://arxiv.org/abs/1104.3913>.
- 339 [12] Moritz Hardt, Eric Price, and Nathan Srebro. Equality of opportunity in supervised learning. [arXiv:  
340 abs/1610.02413](https://arxiv.org/abs/1610.02413), 2016. URL <http://arxiv.org/abs/1610.02413>.
- 341 [13] SciPy v1.7.1 Manual. scipy.optimize.fmin\_ncg. URL [https://docs.scipy.org/doc/scipy/reference/  
342 generated/scipy.optimize.fmin\\_ncg.html](https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.fmin_ncg.html). Accessed on: 31-01-2022.
- 343 [14] Rachel K. E. Bellamy, Kuntal Dey, Michael Hind, Samuel C. Hoffman, Stephanie Houde, Kalapriya Kannan,  
344 Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilovic, Seema Nagar, Karthikeyan Natesan  
345 Ramamurthy, John Richards, Diptikalyan Saha, Prasanna Sattigeri, Moninder Singh, Kush R. Varshney, and  
346 Yunfeng Zhang. AI Fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted  
347 algorithmic bias, October 2018. URL <https://arxiv.org/abs/1810.01943>.

348 **A Table of issues**

| Issue  | Our contribution   |
|--|--|
| Running the code in the given environment results in conflicts   | Code modernized and made compatible with the latest version of every dependency  |
| The code is generally complex and hard to understand due to insufficient comments and documentation as well as leftover code   | Trimmed down the code to the essential, included option to choose any of the available models and the corresponding parameters. Added comprehensive documentation to make the code more interpretable  |
| It appears that the pre-processing pipeline of the given datasets is not specified   | Made the scripts we used to pre-process each dataset available as well as a detailed description   |
| It appears that the position (i.e. index) of the sensitive feature for the COMPAS and Drug Consumption datasets is not indicated, posing a challenge to reproduce the author's results   | Moved the sensitive feature (i.e. gender) of every dataset to the 0th index, which is taken as default by the main function  |
| The advantaged and disadvantaged groups for the sensitive attribute (gender) has not been specified for any attack   | Assumed from the codebase that the authors did this automatically and inferred it from the dataset (the advantaged group is chosen as the group with a higher ratio of datapoints with the positive label ( $y=1$ ), regardless of the actual class label it corresponds to) to be specified for all attacks |
| 349 The code does not provide an option to run baseline methods used in the paper, nor does it include the hyperparameter  | Included option to run baseline methods (Koh attack, Solans attack) and to include in IAF attack   |
| The code implements a deterministic point sampling in the anchoring attacks (RAA, NRAA) due to the same seed being reset in every attack iteration. Thus the sampling yields the same point every iteration not properly applying the randomness | Fixed the issue so that randomness takes effect  |
| The code makes use of a faulty EOD metric calculation  | Re-implemented the EOD metric calculation to fix the issue   |
| The paper specifies the feasible set computation to be done on the union of the clean dataset and the initial poisoned points. The original code however does this on the clean data only when using the running commands given by the authors   | Implemented the feasible set as specified in the paper   |
| It appears that the model used for the given classification task is not specified  | Assumed they used a Support Vector Machine (SVM) trained with a smooth hinge loss and L2 regularization  |
| The optimization algorithm is not indicated  | Assumed it to be Newton's Conjugate Gradient (Newton-CG) method, as suggested by the codebase  |
| It is unclear how the best performing model was selected   | Saved best performing model on the test set according to the chosen stopping metric  |

350 **B Additional figures**

351 Here we collect additional figures that support the results discussed above.

352 **B.1 Effect of  $\alpha$  on the different metrics**

353 Figure 3 shows the influence of  $\alpha$  on the different metrics when accuracy is used as the stopping criteria. The experiment is repeated using average fairness as the stopping metric, and the results are collected in Figure 4. These results support Claim 1 of Section 2, effectively proving it.

356 B.2 Comparative study between the proposed attacks and the baselines

357 We report the results of the experiment designed to support Claims 2-5 of Section 2. We perform each attack (IAF,  
358 NRAA, RAA, Koh, Solans) on each dataset (German, Compas, Drug),  $\alpha$  and gradually increasing from 0.0  
359 to 1.0, with steps of 0.1. We repeat this procedure for each stopping metric (average fairness and accuracy). The results  
360 are respectively collected in Figures 4 and 5.

361 B.3 Figures of the original paper

362 For the sake of self-containedness of this reproducibility study, we report the two main figures of the original paper.  
363 Figures 7 and 8 correspond – respectively – to Figures 2 and 3 of "Exacerbating Algorithmic Bias through Fairness  
364 Attacks".

Figure 3: Results obtained for different attacks with regards to accuracy (test error) and fairness (SPD and EOD) measures on German Credit, COMPAS, and Drug Consumption databases with different  $\alpha$  and with accuracy as the stopping method.

Figure 4: Results obtained for different attacks with regards to accuracy (test error) and fairness (SPD and EOD) measures on German Credit, COMPAS, and Drug Consumption databases with different levels and with average fairness as the stopping method.

Figure 5: Accuracy (test error) and fairness (SPD and EOD) measures obtained after the IAF attack the on German Credit, COMPAS, and Drug Consumption databases for different increasing values, with accuracy as the stopping method.



Figure 6: Accuracy (test error) and fairness (SPD and EOD) measures obtained after the IAF attack the on German Credit, COMPAS, and Drug Consumption databases for different increasing values, with average fairness as the stopping method.

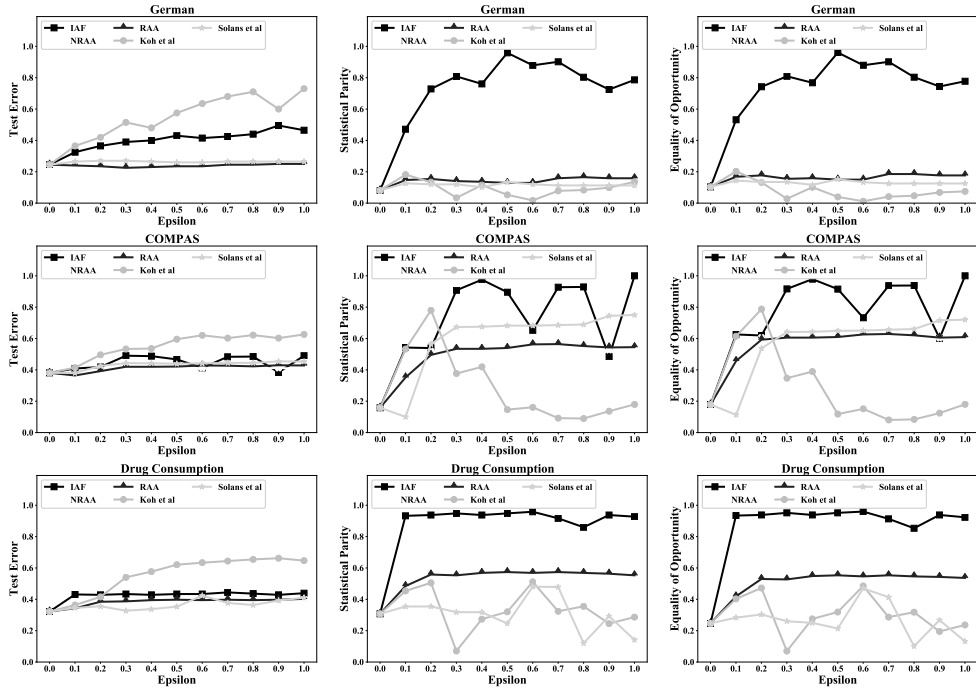


Figure 7: Results of the original paper obtained for different attacks with regards to different fairness (SPD and EOD) and accuracy (test error) measures on three different datasets (German Credit, COMPAS, and Drug Consumption) with different  $\epsilon$  values. Retrieved from [1].

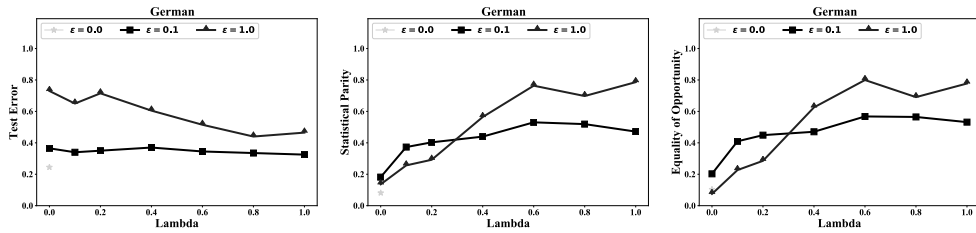


Figure 8: Results of the original paper obtained for different  $\epsilon$  values for the IAF attack with regards to different fairness (SPD and EOD) and accuracy (test error) measures on three different datasets (German Credit, COMPAS, and Drug Consumption) with different  $\lambda$ . Retrieved from [1].