

---

# Reproducing "Fair Selective Classification via Sufficiency"

---

Anonymous Author(s)

Affiliation

Address

email

## Reproducibility Summary

1

### 2 **Scope of Reproducibility**

3 In this reproducibility study we focus on the paper "Fair Selective Classification via Sufficiency". Our experiments  
4 focus on the following claims: 1. Sufficiency is able to mitigate disparities in precision across the entire coverage scale  
5 and in margin distributions, and will not increase these disparities compared to a baseline selective classification model  
6 in any case. 2. Using sufficiency may decrease overall accuracy in some cases, but still mitigates the disparity between  
7 groups when looking at individual classification scores. 3. The sufficiency-regularised classifier exhibits better fairness  
8 performance on traditional fairness datasets.

### 9 **Methodology**

10 As the authors have not made their code publicly available, all code was written from scratch, based on the instructions  
11 and pseudocode given in the original paper. Our code reconstruction contains code for training both the sufficiency  
12 model and a baseline model performing standard selective classification.

### 13 **Results**

14 We were not able to fully reproduce the results of the original paper in this setting. The numbers (accuracies, precisions  
15 and margin distributions) obtained in our experiments differ significantly from those reported in the original paper.  
16 Though differences between the baseline model and the sufficiency model are not as significant as in the original paper,  
17 our results do support the main claims about sufficiency being able to increase the worst-group precision and thus  
18 causing disparities between groups to decrease.

### 19 **What was easy**

20 The authors made the importance of implementing fair selective classification with sufficiency very clear. Moreover, the  
21 authors provided an in-depth mathematical background to sufficiency and selective classification, making their reasoning  
22 explicit. Finally, the authors presented their results in such a manner that allowed for straightforward comparison once  
23 we had trained the model.

### 24 **What was difficult**

25 Many technical details and model parameters were not specified in the original paper, and as no code was provided  
26 by the authors, these initially had to be determined by experimentation. Furthermore, some of the figures in the paper  
27 caused confusion about the exact implementation of the model.

### 28 **Communication with original authors**

29 As soon as we noticed we needed clarification on the hyperparameters, datasets and models, we contacted the authors  
30 via email. Initially we did not receive a reply, and eventually the authors were only able to answer some of our questions  
31 on the Tuesday before the deadline. While we re-implemented our model based on the newly supplied information,  
32 time was too short to fix the new issues that became apparent with the new model.

## 33 1 Introduction

34 Fair classification problems emerge when one wishes to ensure that underprivileged groups sharing some sensitive  
35 attribute, such as race or gender, are not disadvantaged against any other group with the same sensitive attribute (Lee  
36 et al., 2021). A variant of the fair classification problem is selective classification, where a model is allowed to abstain  
37 from making a decision. This is usually implemented via confidence thresholding. When the confidence threshold is  
38 higher, one should expect to see better performance on the remaining samples, as the system is only making decisions  
39 when it is very confident with regards to some confidence measure (Jones et al., 2020). However, it has been shown that  
40 while decreasing the coverage can increase overall performance, it can additionally magnify disparities between groups  
41 (Jones et al., 2020).

42 The paper that is central to this reproducibility study by Lee et al. (2021) proposes a method for enforcing fairness  
43 during selective classification, consisting of a sufficiency criterion and a regulariser based on mutual information. The  
44 authors claim that the method ensures that a classifier is fair, even if it abstains from classifying on a large number of  
45 samples. They demonstrate their method on four datasets, each consisting of a different type of data.

## 46 2 Scope of reproducibility

47 In this reproducibility study we focus on several claims. The first claim is that sufficiency is able to mitigate disparities  
48 in precision across the entire coverage scale and in margin distributions, and will not increase these disparities compared  
49 to a baseline selective classification model in any case. The second claim is that using sufficiency may decrease overall  
50 accuracy in some cases, but still mitigates the disparity between groups when looking at individual classification scores.  
51 Finally, the authors claim that sufficiency-regularised classifier exhibits better fairness performance on traditional  
52 fairness datasets.

53 Our study consists of two components:

- 54 • Code reconstruction: Since the author’s code is not publicly available, all code was written from scratch in  
55 Python 3, using the instructions and pseudocode as described in the paper. Models, code and datasets are  
56 described in Section 3. Our code can be found on GitHub<sup>1</sup>.
- 57 • Replication: The main part of our study is focused on reproducing the results in Lee et al. (2021), and to  
58 validate their observations and conclusions. Our replication results are presented in Section 4.

## 59 3 Methodology

60 First, an overview of the general sufficiency model is given, after which we discuss how the model was adapted to each  
61 of the datasets. This is followed by a discussion on how we evaluated our implementation, and finally we discuss the  
62 computational requirements.

### 63 3.1 Model descriptions

64 As mentioned before, the original paper uses a selective classification model to which the sufficiency criterion has been  
65 applied during training. The sufficiency criterion ensures that the predictive accuracy is the same for each group at  
66 each confidence level, that precision increases for each group when using selective classification and helps prevent  
67 disparities between groups when decreasing coverage.

68 For a binary target  $Y$  and sensitive attribute  $D$ , the sufficiency criterion imposes a conditional independence between  $Y$   
69 and  $D$  conditioned on the learned features  $\Phi$ , thus requiring:

$$P(Y = 1|\Phi(x), D = a) = P(Y = 1|\Phi(x), D = b), \forall a, b \in D.$$

70 An overview of the general sufficiency model is given in Figure 1. When training the model, depending on which data  
71 set is used, the data is first passed through either or both a pre-trained deep neural network and a two-layer neural

---

<sup>1</sup><https://github.com/MLRC2022FSCS/FSCS>, accessed 04-02-22

72 network. The first layer serves as a feature extractor and the second one serves as a classifier. From these features,  
 73 in addition to training a joint classifier, a group-specific classifier is trained for each  $d \in D$ . For each data point, a  
 74 group-specific loss and a group-agnostic loss are computed. To obtain the first, the datapoint is assigned to the correct  
 75 group-specific classifier, that is the one corresponding to the input’s sensitive attribute  $D = d$ , while for the second the  
 76 input is assigned to either of the classifiers based on the marginal distribution  $P_D$ . A combination of these losses is then  
 77 used as a sufficiency regulariser:

$$L_R \frac{1}{n} \sum_{i=1}^n (\log q(y_i | \Phi(x_i); \theta_{d_i}) - \log q(y_i | \Phi(x_i); \theta_{d_i}^{\sim}))$$

78 The overall loss function then becomes:

$$\min \frac{1}{n} \sum_{i=1}^n (L(T(\Phi(x_i)), y_i) + \lambda \log q(y_i | \Phi(x_i); \theta_{d_i}) - \lambda \log q(y_i | \Phi(x_i); \theta_{d_i}^{\sim}))$$

79 and is used to update the feature extractor and joint classifier.

80 By minimising the difference between the group-specific and group-agnostic loss,  $\Phi(x)$  will be such that the group-  
 specific models trained on it will decrease their individual biases and converge towards the same model. In binary

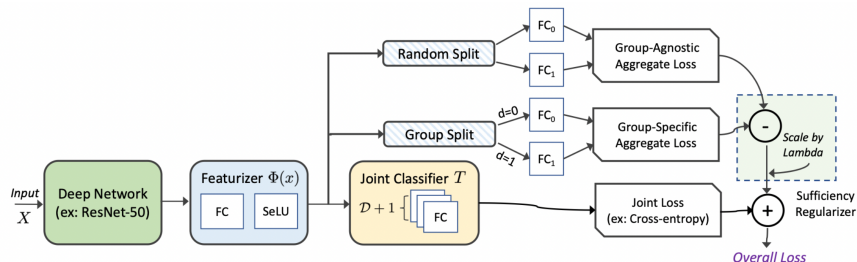


Figure 1: Overview of the sufficiency model. Obtained from the original paper by (Lee et al., 2021).

81 selective classification, an input  $X$  is classified as belonging to a certain class when the confidence exceeds some  
 82 threshold. The softmax response  $s(x)$  is monotonically mapped to the confidence score  $\kappa(x)$  with the following  
 83 formula, which maps  $[0.5, 1]$  to  $[0, \infty]$  and provides much higher resolution on the values close to 1 (Lee et al., 2021):  
 84

$$\kappa(x) = \frac{1}{2} \log \left( \frac{s(x)}{1 - s(x)} \right)$$

85 When  $\hat{y} = y$ , the margin  $M(x)$  is  $\kappa(x)$  and  $-\kappa(x)$  otherwise. Given a threshold  $\tau$ , the classifier makes a correct  
 86 prediction when  $M(x) \geq \tau$  and an incorrect prediction when  $M(x) \leq -\tau$ .

### 87 3.2 Code reconstruction

88 Following the paper, our code was implemented in PyTorch<sup>2</sup>. This was achieved by creating the featuriser for each  
 89 dataset, a joint classifier and two fully connected layers: one for the privileged and one for the unprivileged protected  
 90 group. No activation layers were added to the joint classifier and group-layers, since cross-entropy loss requires logits  
 91 as input. However, for evaluation of the model a softmax layer was applied to the predictions of the joint classifier as  
 92 this was required for selective classification.

93 Three separate Adam optimisers were used: one for the featuriser, one for the joint classifier and one for both layers in  
 94 the group classifiers. The loss regulariser  $\lambda$  was set to 0.7 for all datasets as was stated in the paper. The learning rate  
 95 was not provided in the paper, but later clarified by the authors to be 0.001 for each of the three featurisers. Moreover,  
 96 there was no mention of what range was used for the confidence threshold, which determines the coverage. As such,  
 97 testing starts with a threshold  $\tau$  of 0 (i.e. with a coverage of 100%), and increases with some threshold step size until  
 98 we reach a coverage of under 0.19. The cut-off point for coverage at 0.19 was chosen somewhat arbitrarily, as it lies a  
 99 little past 0.20, which seems to be roughly the point beyond which neither the accuracy, nor the precision change much  
 100 at all. This is in line with both our observations, as well as the data presented by Lee et al. (2021).

<sup>2</sup><https://pytorch.org/docs/stable/index.html>, accessed 04-02-22

### 101 3.3 Dataset-specific models

102 We ran the experiments on three of the four binary classification datasets used in the paper. For each of the datasets, we  
103 used the predetermined train/test splits if available.

#### 104 3.3.1 Adult dataset

105 The Adult<sup>3</sup> dataset (Kohavi et al., 1996) consists of 48,842 entries containing tabular census data, such as age, sex  
106 and education. The first step in preprocessing the dataset was removing all entries with missing values. The data was  
107 then split into 29,092 training examples and 15,060 test examples. Categorical variables within the data were one-hot  
108 encoded, and the continuous variables were normalised to be between 0 and 1, the latter of which was done to remove  
109 the outliers that could incorrectly skew the gradient learning of the parameters. Following the original paper, we only  
110 kept the first 50 samples for women with a high income, that is  $D = 0$  and  $Y = 1$ , to stimulate disparities between  
111 groups. The resulting data,  $X$ , was used to predict the target label  $Y$ , which in this case is an individual’s income.  
112 Classification is binary: an income of over 50k is viewed as high income and assigned label 1, and every other income  
113 was assigned label 0. Sex was designated as the sensitive attribute  $D$ .

114 For this dataset, we followed the original paper and used a two-layer neural network with a hidden layer consisting of  
115 80 nodes. The first layer is a feature extractor using a Scaled Exponential Linear Unit (SELU) activation function, and  
116 the second layer serves as the joint classifier. The network is trained for 20 epochs.

#### 117 3.3.2 CelebA dataset

118 The CelebA<sup>4</sup> dataset (Liu et al., 2015) consists of 202,599 images of 10,177 different celebrities, along with a list of  
119 attributes depicted in the images. It was not clear from the original paper whether the aligned dataset or the original  
120 one was used. Moreover, it was only specified that the images were resized to 224x224, but it was not explained how  
121 this was done and whether there were any other preprocessing steps, such as normalisation. After communication  
122 with the authors it became clear that the aligned dataset was used, and that the images were to be normalised with 0.5  
123 mean and 0.5 standard deviation. Because the Pytorch dataloader loaded in 38 extra columns that were unnecessary in  
124 our research, we manually resized the images to 224x224 with a Pytorch transformation. In order to be able to use  
125 the cross-entropy function in a later stage, all -1 values of the binary ‘blond’ and ‘male’ variables were mapped to 0.  
126 The resulting images were used as data  $X$ , the hair colour (blond or not) was used as the target variable  $Y$  and the sex  
127 variable was used as the sensitive attribute  $D$ .

128 To obtain features from the images, we trained a ResNet-50 model (He et al., 2016), initialised with the pre-trained  
129 ImageNet weights, for 10 epochs. The features were then extracted from the second to last layer. The last layer was  
130 removed and replaced with a layer consisting of two output nodes to form the classifier.

#### 131 3.3.3 Civil Comments dataset

132 The Civil Comments dataset<sup>5</sup> (Borkan et al., 2019) is a text-based dataset consisting of 1,999,514 online comments on  
133 various news articles, along with metadata about the commenter and a label indicating whether the comment displays  
134 toxicity or not. The Kaggle repository does not provide a test set with labels nor a validation set. This meant that  
135 exclusively datapoints from the training set were used in our study. Following Lee et al. (2021), we let  $X$  be the  
136 comment text,  $Y$  be the binary toxicity label, and  $D$  be whether Christianity is mentioned. The dataset does not include  
137 mention-of-Christianity values for each data point and therefore all data points without one were dropped. A total  
138 of 235,087 comments remained. These were subsequently split into a training, validation and test set using ratios of  
139 0.8, 0.1, 0.1 respectively. Additionally, the targets  $Y$  and mentions of Christianity  $D$  were converted to binary values,  
140 where values above or equal to 0.5 were mapped to 1 and values below 0.5 were mapped to 0. Lastly, the comments  $X$   
141 were tokenised using the BERT tokeniser<sup>6</sup>, with max length set to 512, truncation set to true, and padding set to the max  
142 length.

---

<sup>3</sup><https://archive.ics.uci.edu/ml/datasets/adult>, accessed 04-02-22

<sup>4</sup><http://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>, accessed 04-02-22

<sup>5</sup><https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification/data>, accessed 04-02-22

<sup>6</sup>[https://huggingface.co/docs/transformers/model\\_doc/bert#transformers.BertTokenizer](https://huggingface.co/docs/transformers/model_doc/bert#transformers.BertTokenizer), accessed 04-02-

143 To obtain features from the texts, the tokenised data was passed through a BERT model (Devlin et al., 2018) using the  
144 pretrained parameters. The exact BERT model was not specified in the original paper. Due to time constraints, the  
145 Tiny BERT model from Hugging Face<sup>7</sup> (Turc et al., 2019; Bhargava et al., 2021) was used, which had previously been  
146 adapted to Pytorch. Similarly to the Adult dataset, we then applied a two-layer neural network to the BERT output with  
147 80 nodes in the hidden layer. The first layer was treated as a feature extractor and the second layer as the classifier.  
148 Following the original paper, we trained the model for 20 epochs.

### 149 3.4 Evaluation

150 To make sure our implementation is correct, we also implement a standard classification baseline where we only  
151 optimise the cross-entropy loss function. This can be observed in the lower part of Figure 1. Moreover, we plot the  
152 margin distributions of our sufficiency implementation and compare them to that of the original paper.

153 To measure the effectiveness of our selective classification implementation, we follow the evaluation method of the  
154 authors and plot the accuracy-coverage and precision-coverage curves, and then compute the area under the curves to  
155 summarise the performance across coverage values.

### 156 3.5 Computational requirements

157 The experiments were run using a Nvidia RTX 3090 with 24 GB VRAM at 1785 MHz. The batch sizes were not  
158 provided in the original paper, and so they were chosen based on memory constraints. As the Adult dataset consists of  
159 relatively little data, the batch size was set to 32 in order to perform enough gradient steps to fit the parameters. This  
160 resulted in a total training runtime of about 10 seconds for the baseline and 5 minutes for the sufficiency implementation  
161 across all 20 epochs. For the CelebA dataset, the largest batch size that fit in memory was 96, which resulted in a total  
162 training runtime of 1 hour and 30 minutes for the baseline and 3 hours and 30 minutes for sufficiency for 10 epochs.  
163 For the Civil Comments model, a batch size of 48 was used, resulting in a total of 30 minutes of training time for the  
164 baseline model and 1 hour and 38 minutes for the sufficiency implementation when running for 20 epochs.

## 165 4 Results

### 166 4.1 Overall accuracy-coverage graphs

167 Figure 2 displays the overall accuracy plotted against the coverage for different datasets and for both the baseline and  
168 the sufficiency-regularised model. From the Adult dataset graph, we can infer that accuracies are the same for both  
169 models across all coverages. For the CelebA dataset, the sufficiency model increases the accuracy for most of the  
170 coverage scale, only converging with the baseline at a coverage of around 0.25. In the Civil Comments dataset graph,  
171 the baseline model outperforms the sufficiency model across the entire coverage scale.

172 In the original paper, the authors claim that sufficiency may decrease accuracy in some cases. Specifically, they show  
173 that the baseline model outperforms the sufficiency model on overall accuracy for the Adult dataset. Our results do  
174 support this specific result on the CelebA dataset, though for the Adult dataset the baseline and sufficiency models  
175 perform equally. For the Civil dataset, however, it is the case that sufficiency decreases accuracy, which thus confirms  
176 the general claim that sufficiency does not necessarily improve accuracy.

### 177 4.2 Group-specific precision-coverage curves

178 Figure 3 shows the group-specific precisions across the entire coverage scale. When comparing the baseline model  
179 to the sufficiency model, Figure 3a shows that, from a coverage of below about 0.7, sufficiency leads to a smaller  
180 gap between the female and male precisions on the Adult dataset. The worst-group precision, i.e. the male precision,  
181 improves most.

182 For the CelebA dataset in Figure 3b, we observe both groups' precisions increasing when using sufficiency. The  
183 precisions increase equally however, causing the gap between the genders to remain the same.

184 As was to be expected from Figure 2c, both precisions decrease when using sufficiency on the Civil Comments dataset.  
185 However, the gap between the two groups very slightly decreases for coverages between 0.7 and 1.0.

<sup>7</sup><https://huggingface.co/prajjwal1/bert-tiny>, accessed 04-02-22

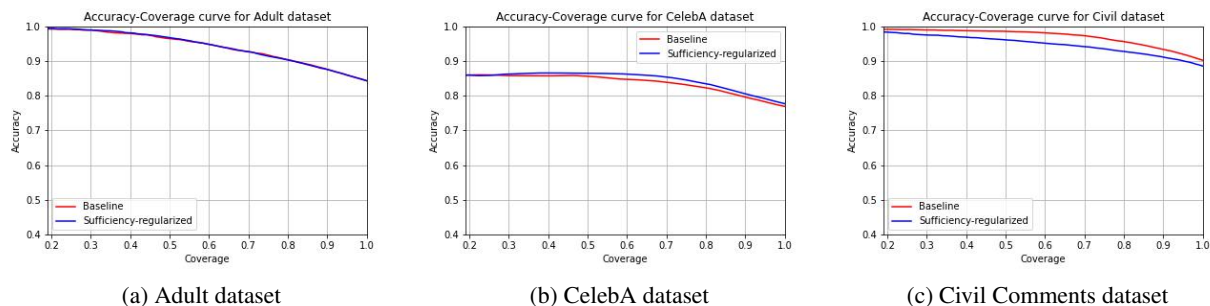


Figure 2: Overall accuracy-coverage graphs.

186 These findings are mostly in line with the findings in the original paper: while for the Adult dataset and the Civil  
 187 Comments dataset the gaps between the two group do decrease when including sufficiency, and while for the CelebA  
 188 dataset this is not the case, sufficiency does not increase the gap but does significantly improve accuracy. These results  
 189 neither confirm nor deny the authors’ claim that the sufficiency criterion introduces a method for mitigating the disparity  
 190 in precision, though we do note that the differences in precision in our results are much less significant than those as  
 191 reported in the original paper.

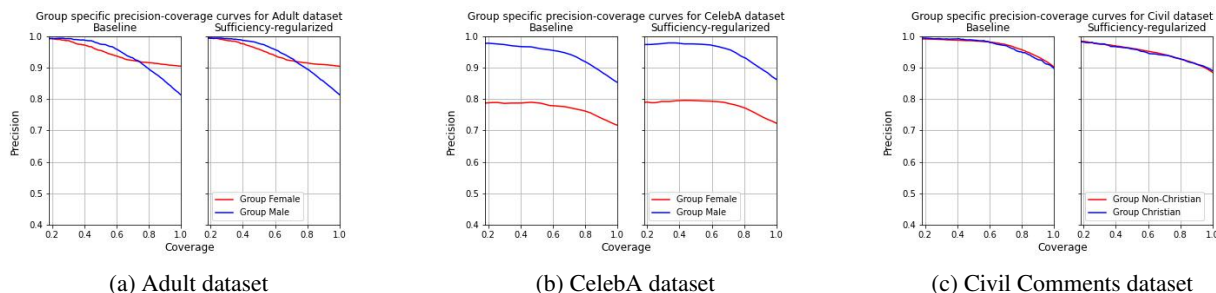


Figure 3: Group-specific precision-coverage graphs.

### 192 4.3 Margin distributions

193 In Figure 4, the margin distributions for both groups are displayed for each of the datasets. For the Adult dataset, the  
 194 margins do not appear to be affected much by sufficiency.

195 Conversely, in the CelebA dataset, both margin distributions become more positively centred, causing the distributions  
 196 to be more similar. Especially the male group margin shifts more towards the positive side, obtaining a smaller range in  
 197 the negative region and a wider peak in the positive region. The number of samples in the female group with a negative  
 198 margin has decreased. We also observe an increase in the number of outliers in the positive region.

199 Finally, the Civil Comments dataset shows the Non-Christian group’s margin becoming more normally centred around  
 200 a margin of around 1, and also shows the two groups’ distributions becoming more aligned.

201 Our results show sufficiency mitigating and in any case not worsening disparities between the two groups, with the  
 202 Adult dataset distributions staying the same and the other two datasets confirming that sufficiency causes a slight  
 203 reduction of the gap between the margin distributions of different groups. This thus confirms the claim that sufficiency  
 204 helps mitigate disparities in margin distributions, however, again, the differences between the models’ distributions are  
 205 not as clear as in the original paper.

### 206 4.4 Numerical evaluations

207 In Table 1, the areas under the accuracy curves and the areas between the precision curves are presented for each of the  
 208 datasets. For the Adult dataset, the area under the accuracy curve virtually remains the same when using sufficiency, in  
 209 line with Figure 2a. The area between the precision curves slightly increases. While this seems to contradict Figure  
 210 4a, note that we only observed a decrease in the precision gap for coverages of below 0.7, and the numbers in Table 1

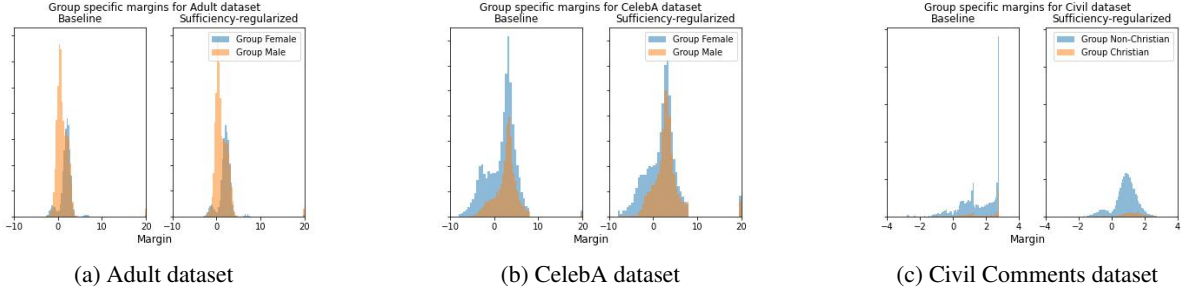


Figure 4: Margin distribution graphs.

211 concern the entire coverage score.

212 For the CelebA dataset, the area under the accuracy curve increases, resulting in an increase in overall accuracy as  
 213 previously observed in Figure 2b. However, as already suggested by Figure 4b, the area between the precision curves  
 214 stays virtually the same when using the sufficiency method.

215 Once again confirming the results observed in Figure 3, when using sufficiency, the area under the Civil Comments  
 216 dataset’s accuracy curve decreases. The area between the precision curves effectively stays the same.

217 As mentioned before, in the original paper sufficiency causes the Adult dataset accuracy to diminish, while in our results  
 218 both models achieve the same performance. In contrast, while for the Civil Comments dataset the original paper’s  
 219 accuracy increases, our results show a decrease in accuracy. The CelebA results both exhibit an increase in accuracy,  
 220 though this increase is more prominent in the original paper.

221 The area between the precision curves significantly decreases for the Adult dataset in the original paper, which is not the  
 222 case for our results. The same holds for the precision curves of the the CelebA and Civil Comments dataset: although  
 223 less so than for the Adult dataset, the original paper’s result show that disparities are decreased when using sufficiency,  
 224 while our results do not show any significant change.

225 The Civil Comments results show that accuracy can reduce when using sufficiency, but disparities will not increase.  
 226 Furthermore, although the claim about disparity in precision mitigating is not directly confirmed by our results as the  
 227 areas stay the same, it does show that sufficiency will not (significantly) worsen disparities in any case.

Dataset	Method	Area under accuracy curve	Area between precision curves
Adult	Baseline	0.931	0.220
	<b>Reproduced baseline</b>	<b>0.941</b>	<b>0.004</b>
	Sufficiency	0.887	0.021
	<b>Reproduced sufficiency</b>	<b>0.942</b>	<b>0.005</b>
CelebA	Baseline	0.852	0.094
	<b>Reproduced baseline</b>	<b>0.855</b>	<b>0.141</b>
	Sufficiency	0.975	0.013
	<b>Reproduced sufficiency</b>	<b>0.863</b>	<b>0.142</b>
Civil Comments	Baseline	0.888	0.026
	<b>Reproduced baseline</b>	<b>0.973</b>	<b>0.0012</b>
	Sufficiency	0.943	0.010
	<b>Reproduced sufficiency</b>	<b>0.954</b>	<b>0.0010</b>

Table 1: Numerical comparison between original paper and reproduction.

## 228 5 Discussion

229 To summarise, the numbers (accuracies, precisions, margin distributions etc.) obtained in our experiments differ  
 230 significantly from those reported in the original paper. However, although differences between the baseline model  
 231 and the sufficiency model are not as significant as in the original paper, our results do support the main claims about  
 232 sufficiency being able to increase the worst-group precision and thus causing disparities between groups to decrease.

233 It is worth mentioning that the Figures 4b and 4c show the largest increase in margin alignment, and these are also  
234 the datasets that either improve in overall accuracy, or decrease in disparities between groups. Moreover, the authors  
235 claimed that the sufficiency-regularised classifier exhibited better fairness performance on traditional fairness datasets.  
236 Though we were not able to reproduce their results in this study, we do believe we can validate this claim, as sufficiency  
237 is either able to decrease the disparities in precision between groups (Figures 3a and 3c), or increase the precision for  
238 both groups in an equal manner as we traverse the coverage scale, meaning that no group is penalised for the sake of  
239 improving the other group’s precision.

240 The fact that we were not able to precisely reproduce the results from the original paper is likely due to the fact that not  
241 all technical details required to fully replicate the original paper were provided by the authors in the paper. Specifically  
242 the learning rate, selective classification threshold and optimiser algorithms had to be decided upon ourselves. While a  
243 well-informed guess of what parameters to use was made possible due to experimentation, it could well be possible that  
244 the authors’ implementation differs on these fronts, and that this caused our results to differ from the ones in the paper.  
245 We also did not have time to run all the experiments we would have liked to, such as the fourth CheXpert<sup>8</sup> dataset or  
246 experiments beyond replication, such as applying the sufficiency method to a new dataset. This was due to the fact that  
247 we spent a large amount of time trying to improve our original results, because we wanted to make sure these were  
248 stable before attempting to generalise further.

## 249 **5.1 Reproducibility of the paper**

### 250 **5.1.1 What was easy**

251 The authors provided a strong and logically structured theoretical background that made the importance of implementing  
252 fair selective classification with sufficiency very clear. Moreover, the authors provided an in-depth mathematical  
253 background to sufficiency and selective classification, making their reasoning explicit. Finally, the authors provided  
254 clear explanations of the evaluation method and presented their results in such a manner that allowed for straightforward  
255 comparison once we had trained the model.

### 256 **5.1.2 What was difficult**

257 As mentioned previously, many crucial technical details (e.g. pretrained models and hyperparameters) required to  
258 replicate the original paper were not provided by the authors. Furthermore, we found the overview of the model shown  
259 in Figure 1 (Figure 2 in the original paper) difficult to interpret. This caused the implementation of the model to take  
260 more time than we had anticipated. The first issue was the use of "ex" in the deep network and joint loss depictions,  
261 which is generally short for "excluding". In section 4.1 of the original paper, it appears that the ResNet-50 model is  
262 modified in place, leading to the features being extracted and classified within ResNet-50 itself. This would indeed  
263 indicate 'ex' meaning 'excluding', as there is no separate featuriser in this case. However, this interpretation means that  
264 cross-entropy is excluded from the joint loss, though it is explicitly mentioned in section 4.1. This would indicate "ex"  
265 is short for 'exemplum', which is a contradiction. Moreover, the image does not make immediately clear that the fully  
266 connected layers  $FC_0$  and  $FC_1$  are the same for both the group-specific and the group-agnostic classifier. There was  
267 also no mention of the loss functions or activation functions used for the fully connected layers in the group-specific  
268 classifiers. Finally, it was not explicitly mentioned whether the featurisers were the same for the Adult and Civil  
269 datasets.

## 270 **5.2 Communication with original authors**

271 As soon as we noticed we were missing crucial information about the hyperparameters and the CelebA dataset and  
272 we needed some clarifications on the workings of the model, we contacted the authors via email. Initially we did not  
273 receive a reply, and so we sent a follow-up email. We received an answer from the authors that they needed more time  
274 to verify the information we asked for and were currently working towards a deadline themselves and we eventually  
275 received an email on 01-02-2022. In this email, the authors were only able to answer some of our questions. While  
276 we re-implemented our model based on the newly supplied information, time was too short to fix the new issues that  
277 became apparent with the new model.

---

<sup>8</sup><https://stanfordmlgroup.github.io/competitions/chexpert/>, accessed 04-02-22



278 **References**

- 279 Bhargava, P., Drozd, A., and Rogers, A. (2021). Generalization in nli: Ways (not) to go beyond simple heuristics.
- 280 Borkan, D., Dixon, L., Sorensen, J., Thain, N., and Vasserman, L. (2019). Nuanced metrics for measuring unintended  
281 bias with real data for text classification. In *Companion proceedings of the 2019 world wide web conference*, pages  
282 491–500.
- 283 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for  
284 language understanding. *arXiv preprint arXiv:1810.04805*.
- 285 He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *2016 IEEE Conference  
286 on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778.
- 287 Jones, E., Sagawa, S., Koh, P. W., Kumar, A., and Liang, P. (2020). Selective classification can magnify disparities  
288 across groups. *arXiv preprint arXiv:2010.14134*.
- 289 Kohavi, R. et al. (1996). Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. In *Kdd*, volume 96,  
290 pages 202–207.
- 291 Lee, J. K., Bu, Y., Rajan, D., Sattigeri, P., Panda, R., Das, S., and Wornell, G. W. (2021). Fair selective classification  
292 via sufficiency. In Meila, M. and Zhang, T., editors, *Proceedings of the 38th International Conference on Machine  
293 Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 6076–6086. PMLR.
- 294 Liu, Z., Luo, P., Wang, X., and Tang, X. (2015). Deep learning face attributes in the wild. In *Proceedings of International  
295 Conference on Computer Vision (ICCV)*.
- 296 Turc, I., Chang, M., Lee, K., and Toutanova, K. (2019). Well-read students learn better: The impact of student  
297 initialization on knowledge distillation. *CoRR*, abs/1908.08962.