

---

# [Re] Fair Selective Classification Via Sufficiency

---

Anonymous Author(s)

Affiliation

Address

email

## Reproducibility Summary

1

### 2 Scope of Reproducibility

3 Bu, Lee et al. (2021) introduced a method for enforcing fairness in selective classification, deriving a novel upper bound  
4 for the conditional mutual information from the sufficiency criterion. We attempt to verify the second claim that: "[this  
5 novel upper bound] can be used as a regularizer to enforce the sufficiency criteria, [and] then **show that it works to**  
6 **mitigate the disparities on real-world datasets.**"[4]

### 7 Methodology

8 To verify the author's claim, we implemented the model and regularizer described in the original paper. We wrote  
9 the code from scratch, since there was no code available. We train both a baseline and regularized model on three of  
10 the four datasets used by the authors: Adult, CelebA, and CheXpert. The Civil Comments dataset which the authors  
11 also used was computationally too expensive. We trained using the Adam optimizer and a constant learning rate. The  
12 training of the regularized models for Adult, Chexpert and CelebA takes under a minute, under 2 hours and under 4  
13 hours respectively on a single NVIDIA Titan RTX.

### 14 Results

15 We found that we could not reproduce the original paper's results. While the area between precision curves decreases  
16 somewhat for the CelebA and CheXpert datasets, it increases for the Adult experiment. Also, the analysis of our margin  
17 distributions between the baseline and regularized models do not seem to indicate an increase in overlap between  
18 groups. Due to these results, we cannot conclude the effectiveness of the regularizer in reducing the disparity between  
19 two groups demonstrated by the authors of *Fair Selective Classification via Sufficiency*, using our implementation.

### 20 What was easy

21 Implementing the regularizer and group specific models in PyTorch was relatively straight forward, using the well  
22 documented loss functions and Algorithm 1 from the original paper [4].

### 23 What was difficult

24 We found that the results for all datasets were sensitive to preprocessing and hyperparameter tuning. Since the authors  
25 specified very little in this regard, experimenting with the dataset specific preprocessing steps, and the hyperparameter  
26 tuning for three datasets took us a considerable amount of time.

### 27 Communication with original authors

28 The authors of the original paper were emailed with multiple questions about preprocessing, training and the baseline.  
29 Unfortunately, because they also had an important deadline, they responded three days before the deadline giving us  
30 little time to make changes.

# 31 1 Introduction

32 Machine learning algorithms are being used to solve more and more diverse problems, and are fulfilling tasks in  
33 increasingly difficult situations. One way to improve the performance of classification models is to use selective  
34 classification [4]. This means that models are allowed to abstain when their prediction confidence is low. However,  
35 abstaining does come at the cost of coverage (the ratio of samples for which a decision is made). Previous work has  
36 shown that classifying selectively does not always affect all distinguishable groups within the data evenly, for instance  
37 in the CelebA<sup>1</sup> and CivilComments<sup>2</sup> datasets [3]. Implementing a model that is able to classify selectively in a fair way,  
38 not discounting certain groups within the data therefore is an active challenge for the artificial intelligence research  
39 community. The authors of the *Fair Selective Classification Via Sufficiency* paper<sup>3</sup> propose a method to improve  
40 fairness in selective classification accuracy between groups by using the sufficiency criterion.

41 The contribution of the authors of this paper is twofold. Firstly, they prove a novel upper bound for conditional mutual  
42 information. Secondly, they use this result to introduce a regularization technique that forces a model to be more fair to  
43 all protected groups when classifying selectively. These protected groups can be selected based on sensitive attributes  
44 (e.g. race, gender). They report improved overall group specific performance relative to a baseline method where they  
45 only optimize the cross-entropy loss. Furthermore, they improve relative to the group DRO method which has been  
46 shown to mitigate the disparity in recall rates between groups in selective classification [5]. In this study we verify  
47 the second claim that "[... their regularizer] works to mitigate the disparities on real-world datasets" by building a  
48 sufficiency regularized classifier that is more fair to underrepresented groups in selective classification.

## 49 2 Scope of reproducibility

50 In this paper we aim to reproduce the second claim from the original paper, which states: sufficiency can be used to  
51 train fairer selective classifiers which ensure that precision always increases for all groups as coverage is decreased. The  
52 authors support their claim by evaluating on the positive predictive parity, also called precision, by looking at the area  
53 under the curve for the accuracy for two groups within the Adult<sup>4</sup>, CelebA, CheXpert<sup>5</sup> and CivilComments datasets.  
54 Since we found the CivilComments dataset and corresponding model to be too computationally expensive, we aim to  
55 reproduce the results on the first three datasets. The original authors did not publish any of their code. The scope of this  
56 reproducibility report is thus to write all necessary code and train and evaluate both the baselines and the regularized  
57 models for the Adult, CelebA and CheXpert datasets.

## 58 3 Methodology

59 This section discusses the methodology and experimental setup used to reproduce the paper *Fair Selective Classification*  
60 *via Sufficiency*. Firstly, the model is discussed, after which we go over the dataset specifics, the evaluation metrics, and  
61 the computational requirements.

### 62 3.1 Model description

63 The architecture comprises of three distinct components: the featurizer, classifier and regularizer. Using these  
64 components, the modelling objective can be described as finding model parameters  $\theta_T$  and  $\theta_\Phi$  for the classifier and  
65 featurizer such that the following equation is satisfied:

$$\min_{\theta_T, \theta_\Phi} \frac{1}{n} \sum_{i=1}^n \left( L(T(\Phi(x_i)), y_i) + R(\Phi(x_i), y_i, \theta_d, \theta_{\tilde{d}}) \right) \quad (1)$$

66 Where  $\Phi(x)$ ,  $T(\Phi(x))$  and  $R(\Phi(x), y, \theta_d, \theta_{\tilde{d}})$  represent the featurizer, classifier and regularizer respectively. These  
67 individual components will be discussed in the following sections. An overview can be found in Figure 1. The loss  
68 function  $L$  is not specified in the original paper, therefore we assume the cross-entropy loss is used, due to its popularity

<sup>1</sup><https://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>

<sup>2</sup><https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification/>

<sup>3</sup><https://proceedings.mlr.press/v139/lee21b.html>

<sup>4</sup><https://archive.ics.uci.edu/ml/datasets/Adult>

<sup>5</sup><https://stanfordmlgroup.github.io/competitions/chexpert/>

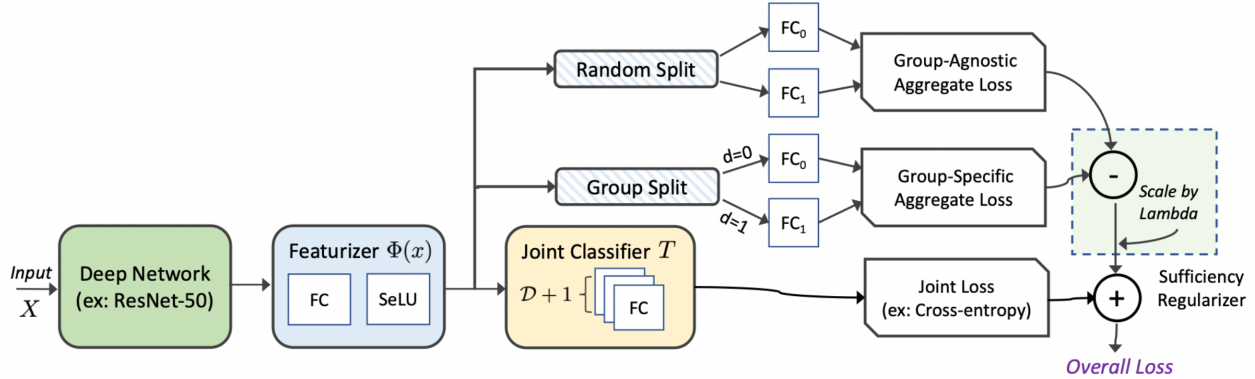


Figure 1: The model architecture from the original paper.

69 This entire Fair Selective Classifier ensures an improvement in precision for all groups as coverage decreases by  
 70 applying the sufficiency criteria to the learned features. This sufficiency is enforced by regularizing the model through a  
 71 novel upper bound of the conditional mutual information (CMI). In this section we describe the classifier architecture  
 72 including the CMI regularization.

### 73 3.1.1 Featurizer

74 The featurizer  $\Phi(x)$  for each dataset is trained to be predictive about  $Y$  while allowing for the classification to be  
 75 calibrated by group based on the sensitive features. Under that constraint for all groups  $d \in \mathcal{D}$  we have that classification  
 76 performance for specific groups is never sacrificed to increase overall performance. This part of the architecture is  
 77 dataset specific.

### 78 3.1.2 Classifier

79 The structure of the *joint classifier*  $T(\Phi(x))$  is vaguely specified in the proposed architecture, specifying a number  
 80  $D + 1$  and displaying a parallel stack of fully connected parallel layers in Figure 2 of the original paper. We chose  
 81 to implement it as a single linear layer taking as input the extracted features from the featurizer. The output size is  
 82 dependent on the classification task, with a single node for binary classification and  $n$  nodes for multi-categorical  
 83 classification. The output of this part of the model is used for inference and therefore also final evaluation.

### 84 3.1.3 Regularizer

85 In order to impose the *sufficiency condition*, a regularization term  $R(\Phi(x), y, \theta_d, \theta_{\tilde{d}})$  is added to the joint classification  
 86 loss. This regularization loss (Equation 2) is calculated according to the following equation.

$$87 \quad R(\mathbf{x}, y_i, \theta_{d_i}, \theta_{\tilde{d}_i}) = \lambda \log q(y_i | \Phi(x_i); \theta_{d_i}) - \lambda \log q(y_i | \Phi(x_i); \theta_{\tilde{d}_i}) \quad (2)$$

88 The architecture achieves this in the following way. Firstly, the sample features are split in two ways: based on their  
 89 true sensitive attribute values  $d$  (the group split) and based on a randomized sensitive attribute  $\tilde{d}$  (the random split),  
 90 where  $\tilde{d}$  is sampled randomly from empirical distribution  $\hat{P}_d$ . Both these splits are subsequently activated by the  
 91 same attribute-specific linear layers  $q(y | \Phi(x); \theta_d)$  to perform group-specific and group-agnostic classification. After  
 92 classification, the loss term  $L_d$  is calculated for each split using the cross-entropy loss, scaled by  $\lambda$ , and subtracted to  
 capture the difference (Equation 2).

## 93 3.2 Datasets

94 This section goes over each dataset, and gives a description of the data, the selected sensitive attribute, and the  
 95 featurization. For the experiments we used the same datasets as the authors. To handle the the different types of datasets,  
 96 four different featurization architectures are used. For a more detailed description of the datasets see Appendix A.

97 The Adult dataset<sup>6</sup> records census data. Each data point has both continuous attributes and categorical attributes. The  
98 goal for this dataset is to predict whether an individual makes over 50k per year. The selected sensitive attribute for this  
99 dataset is sex. To preprocess the data we normalized the continuous attributes to have zero mean and unit variance. We  
100 introduced a bias in the data the same way as in the original paper by removing all but the first 50 rows for which the  
101 protected group  $D = 0$  and the target  $Y = 1$  (what this means is unclear from the paper; we used "Female" and ">50K"  
102 respectively). The Adult Dataset comprises of tabular data, and therefore the featurization component is a single linear  
103 layer. The layer has 80 output nodes and is followed by the SeLU activation function.

104 The CelebA dataset<sup>7</sup> contains RGB images of celebrities. Each image is provided with annotations about the appearance  
105 of the celebrities. The task at hand is to predict whether a celebrity has blond hair and the selected sensitive attribute  
106 is sex. We followed the original paper and resized the images to 224 by 224. The ResNet-50 model has its weights  
107 pre-trained on ImageNet and the classification layer is dropped to output the 2048 features of the penultimate layer.

108 The CheXpert dataset<sup>8</sup> consists of chest x-ray images and corresponding attributes annotated by experts. The task  
109 for this experiment is to predict the presence of pleural effusion (a lung disease), and the sensitive attribute is the  
110 presence of a support device. To preprocess CheXpert we removed all data points where either the sensitive attribute  
111 or the target was labeled uncertain. Again following the original paper, the x-ray images were all resized to 224 by  
112 224 pixels and stacked to simulate the red, green and blue color channels that Densenet121 expects. For the CheXpert  
113 image dataset a DenseNet-121 model (pre-trained on image net) is used as featurizer. The classification layer is dropped  
114 and an AvgPool2d layer is added to output the penultimate layer's feature vector of size 1024, in accordance with the  
115 DenseNet121 architecture [2].

116 The Civil comments dataset<sup>9</sup> contains comments on news articles. For this dataset the task is to predict whether a  
117 comment was toxic or not. The selected sensitive attribute was the commenter being Christian. To preprocess the data,  
118 all comments with unknown religious background were removed. The comments were tokenized using the BERT  
119 tokenizer and the tokenized comments were padded to be able to fit the dimensions of the BERT model as it takes  
120 tensors with a size of 512 tokens. The final classification layer of the BERT-model is replaced with an extra linear layer  
121 to output a feature vector of size 80. The specific version of BERT was not specified, it can be safe to assume that it was  
122 the uncased BERT-Base. The model can be downloaded from the PyTorch-Transformers library. We chose a model  
123 pre-trained for sequence classification.

124 During experimentation we concluded that the computational requirements were too high for our time constraints and  
125 setup. We therefore decided to train models for the Adult, CelebA and CheXpert dataset, but we provide code for  
126 CivilComments as well.

### 127 3.3 Hyperparameters

128 The original paper specifies  $\lambda = 0.7$  and the dataset specific number of epochs. The learning rates and optimizer were  
129 not specified. Therefore, we decided to grid search the learning rate(s). They use three different learning rates in their  
130 notation:  $\eta_d, \eta_f, \eta$  for the group specific models, the featurizer and the joint classifier respectively. However, trying  
131 multiple different combinations of these learning rates would require a lot of computational resources, so we decided to  
132 only perform a grid search over a general learning rate  $\eta$ , which sets the same value for  $\eta_d, \eta_f$  and  $\eta$ .

133 In the process of finetuning the model we found that depending on the dataset learning rates between 0.001 and 0.0001  
134 led to stable and generalizing models. Shortly before the submission deadline for the Machine Learning Reproducibility  
135 Challenge 2021 we were informed by the authors that they used the Adam optimizer, and used a learning rate of 0.001.  
136 In Table 1 we specify the entire set of hyperparameters for every per dataset.

### 137 3.4 Experimental setup

138 In this section we go over the experimental setup. We first explain the training and the evaluation process. After that,  
139 we also go over the hardware and software used for this study.

---

<sup>6</sup><https://archive.ics.uci.edu/ml/datasets/Adult>

<sup>7</sup><https://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>

<sup>8</sup><https://stanfordmlgroup.github.io/competitions/chexpert/>

<sup>9</sup><https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification/>

Parameters	Adult	CelebA	CheXpert
$\lambda$	{0, 0.7}	{0, 0.7}	{0, 0.7}
$\eta$	0.001	0.001	0.001
nr. epochs	20	10	10
batch size	32	128	64

Table 1: The set of hyperparameters for every dataset.

### 140 3.4.1 Model Training

141 During model training, we alternate between two backpropagation steps following the original paper. Firstly we fit the  
 142 group-specific models for each batch and secondly we update the feature extractor and joint classifier.

### 143 3.4.2 Evaluation metrics

144 This section goes over the evaluation metrics of the original paper by first discussing the confidence score and the  
 145 margin. These concepts are then used to derive the final evaluation metrics: the area under the accuracy curve and the  
 146 area between precision curves. For more details, see Section 2 of the original paper.

147 The classifier has the possibility to abstain from the decision based on a confidence score  $\kappa(x)$  and a threshold  $\tau$ . The  
 148 used confidence score is defined as the monotonic mapping of the softmax response  $s(x)$

$$\kappa(x) = \frac{1}{2} \log \left( \frac{s(x)}{1 - s(x)} \right) \quad (3)$$

149 which maps  $[0.5, 1]$  to  $[0, \infty]$  to provide a high resolution on the values close to 1. Since we can map a softmax  
 150 response to the interval  $[0.5, 1]$  for both targets ( $s(x) = s(x)$  for  $s(x) \geq 0.5$ , and  $s(x) = 1 - s(x)$  for  $s(x) < 0.5$ ), it  
 151 is possible to use this function.

152 The confidence score is used to define the margin  $M$ , such that is defined as  $M(x) = \kappa(x)$  if  $\hat{y}(x) = y$  and as  
 153  $M(x) = -\kappa(x)$  otherwise. If we then use  $\tau$  as our threshold for abstaining, the selective classifier makes correct  
 154 predictions when  $M(x) \geq \tau$  and incorrect predictions when  $M(x) \leq -\tau$ .

155 The model is evaluated with different values of  $\tau$ . The selective accuracy is computed for the different coverages,  
 156 caused by the different values of  $\tau$ . The selective precision is computed similarly, conditioning on  $\hat{Y} = 1$ . To measure  
 157 the effectiveness of the selective classifier at different coverage levels, the area under this curve is computed. The  
 158 difference in precision across groups sometimes reveals some disparities that are not revealed by only considering the  
 159 difference in accuracy. Therefore the precision-coverage curves are also plotted per group. The difference between  
 160 these curves is computed to encapsulate the difference across different coverages.

### 161 3.4.3 Computational requirements & Code

162 All training is done on GPU nodes of a cluster which contains multiple kinds of Nvidia GPU’s. For our training we  
 163 used the Titan RTX nodes. These GPU’s have 24GB of GDDR6 memory. The training of the regularized model took  
 164 under 2 minutes for Adult, under 2 hours for Chexpert and under 5 hours on the CelebA dataset. To be able to store all  
 165 the datasets around 20GB of storage is necessary.

166 All code used for the data preprocessing, model training and model evaluation can be found in our Github repository<sup>10</sup>.

## 167 4 Results

### 168 4.1 Area under the Curve & Area between Precision Curves

169 Table 2 lists the area under the accuracy curve (*auc*) and the area between precision curves (*abc*) for the baseline and  
 170 regularized model for every dataset. The change in *auc* shows that the introduction of the regularization does not harm  
 171 (or even improves) the overall performance of the models for all datasets. This is in line with the results of the original  
 172 paper.

<sup>10</sup><https://anonymous.4open.science/r/FSCS-4F57/README.md>

Dataset	Method	Area under accuracy curve	Area between precision curves
Adult	Baseline	$0.93 \pm 0.0002$	$0.056 \pm 0.0008$
	Regularized	$0.93 \pm 0.006$	$0.065 \pm 0.009$
CelebA	Baseline	$0.93 \pm 0.040$	$0.010 \pm 0.0002$
	Regularized	$0.99 \pm 0.0001$	$0.006 \pm 0.0001$
CheXpert	Baseline	$0.83 \pm 0.013$	$0.075 \pm 0.036$
	Regularized	$0.84 \pm 0.02$	$0.070 \pm 0.0009$

Table 2: Area under accuracy curve results for all datasets

Our results for the *abc*, however, diverge. The original paper reports a decrease for the *abc* by a factor of 10, 9 and 2 going from the baseline to the regularized model, for Adult, CelebA and CheXpert respectively. The *abc* of our regularized model for Adult is higher than its baseline, and for CelebA and CheXpert the improvement is a factor of 2 and 1.1 over the baseline. Another point of interest is the high standard deviation for the CheXpert baseline which makes the improvement of the regularized model questionable.

## 4.2 Margins and Precision-Coverage Plots

The dataset specific margins of both the baseline and the sufficiency regularized model can be found in Figure 2 and the precision-coverage plots can be found in Figure 3. Our plots show different characteristics compared to the original paper, and the difference between the unregularized and regularized models is small for all datasets.

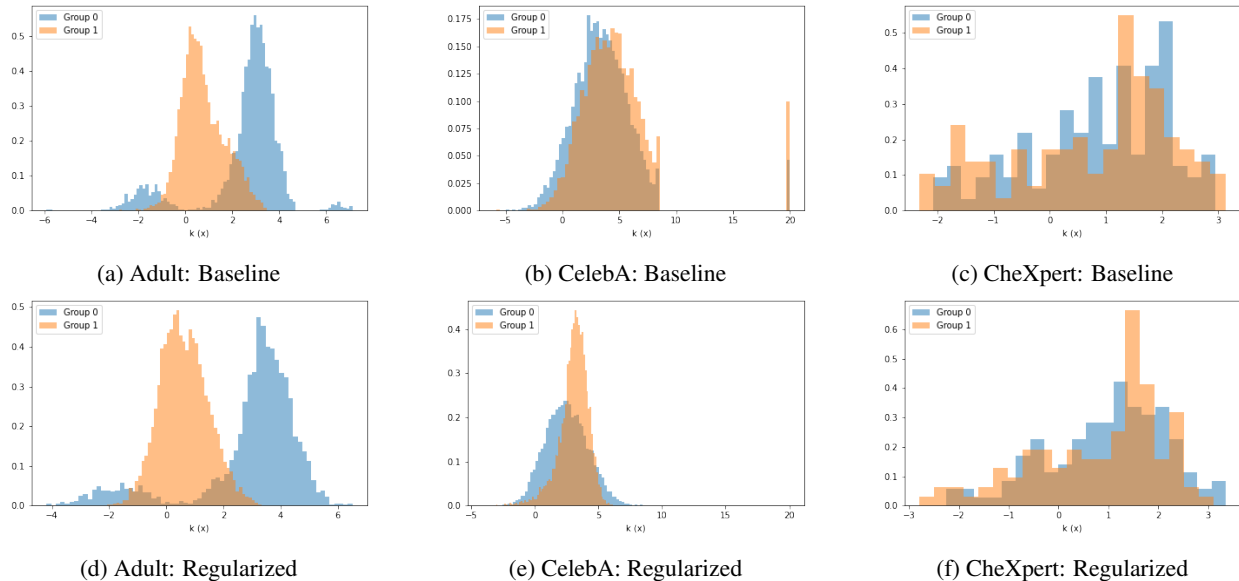


Figure 2: Margin distributions for the datasets for the baseline and regularized methods. For CelebA, all confidence scores were capped to 20.

## 5 Discussion

Table 2 shows the method comparison between the baseline and regularized models for each dataset. While *abc* values decrease somewhat for the CelebA and CheXpert datasets, it increases for the Adult experiment. The margins and precision-coverage plots of the original paper show a clear improvement for the worst case group going from the baseline to the regularized version. Their regularized models show more overlap between the distributions of the two groups, and the precision-coverage curves are closer to each other for the regularized models than for the baseline. The

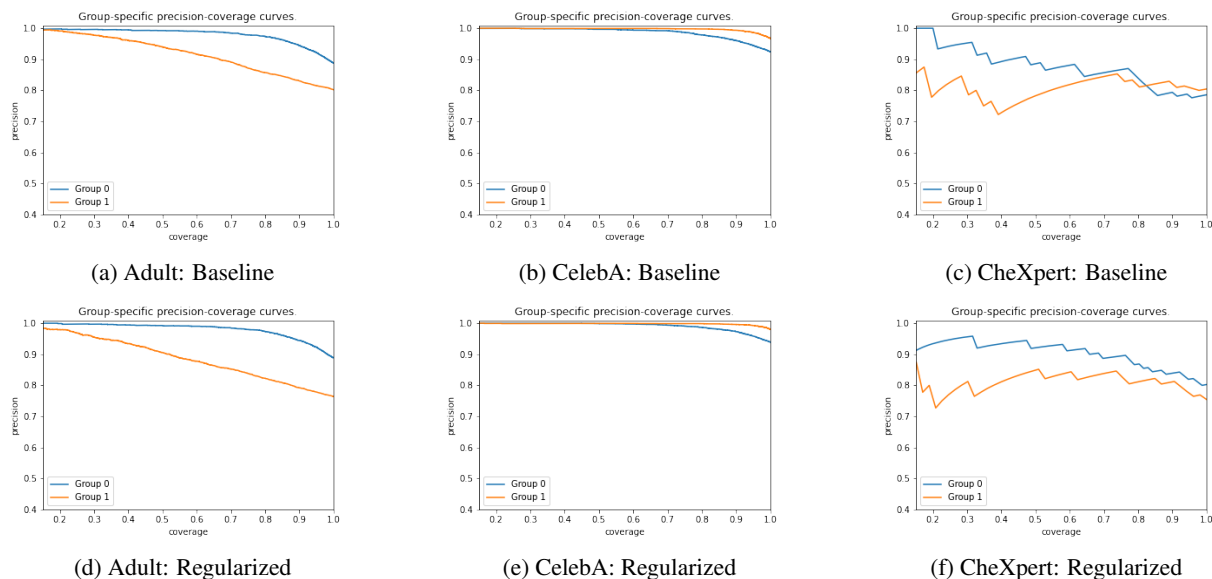


Figure 3: Group-specific precision-coverage curves for the baseline and regularized methods.

188 analysis of our margin distributions between the baseline and regularized models (Figure 2) do not seem to indicate  
 189 an increase in overlap between groups. This means that these margin and precision-coverage comparisons for our  
 190 experiments do not give a conclusive result on the effect of regularizing the classifier.

191 Due to these contrary and/or inconclusive results on the impact on the regularization on selective classification across  
 192 datasets, we cannot conclude its effectiveness in reducing the disparity between two groups demonstrated by the authors  
 193 of *Fair Selective Classification via Sufficiency* using our implementation.

### 194 5.1 What was easy

195 Implementing the regularizer and group specific models in PyTorch was relatively straight forward, using the well  
 196 documented loss functions and Algorithm 1 from the original paper. It was also relatively easy to configure the  
 197 featurizers for each dataset as they were clearly described as well. Implementing the architecture was straight forward  
 198 using the clear figure from the original paper.

### 199 5.2 What was difficult

200 We found that the results for all datasets were sensitive to preprocessing and hyperparameter tuning. Since the authors  
 201 specified very little in this regard, experimenting with the dataset specific preprocessing steps, and the hyperparameter  
 202 tuning for three datasets took us a considerable amount of time.

### 203 5.3 Communication with original authors

204 We sent an email to the original authors. We asked them about preprocessing the CheXpert and CelebA dataset because  
 205 we didn't know exactly how they preprocessed the images in this dataset and how they used some of the attributes. We  
 206 also wondered how Figure 2 in the original paper should be interpreted.

207 Furthermore we had questions about training the models. We did not know what learning rates and optimizers they used  
 208 for training. Knowing this would have saved us from the time-consuming task of running a grid search. Shortly before  
 209 our submission deadline, we received an email from the authors with answers to our questions. This gave us some  
 210 insights into which hyperparameters and optimizer they used. This was just in time to rerun all models for multiple  
 211 seeds.

212 **Acknowledgement**

213 We would like to thank [name to be added] for his guidance during the replication process.

214 **References**

- 215 [1] R. K. E. Bellamy, K. Dey, M. Hind, S. C. Hoffman, S. Houde, K. Kannan, P. Lohia, J. Martino, S. Mehta,  
216 A. Mojsilovic, S. Nagar, K. N. Ramamurthy, J. Richards, D. Saha, P. Sattigeri, M. Singh, K. R. Varshney, and  
217 Y. Zhang. Ai fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic  
218 bias, 2018.
- 219 [2] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger. Densely connected convolutional networks, 2018.
- 220 [3] E. Jones, S. Sagawa, P. W. Koh, A. Kumar, and P. Liang. Selective classification can magnify disparities across  
221 groups, 2021.
- 222 [4] J. K. Lee, Y. Bu, D. Rajan, P. Sattigeri, R. Panda, S. Das, and G. W. Wornell. Fair selective classification via  
223 sufficiency. In M. Meila and T. Zhang, editors, *Proceedings of the 38th International Conference on Machine  
224 Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 6076–6086. PMLR, 18–24 Jul 2021.
- 225 [5] S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang. Distributionally robust neural networks for group shifts: On  
226 the importance of regularization for worst-case generalization, 2020.



## 227 A Dataset Analysis

### 228 A.1 Descriptive analysis datasets

#### 229 A.1.1 Civil Dataset

230 The Civil Comments dataset consists of comments on news article's which were collected on the Civil Comments  
231 platform. This is a platform in which people are able to post a comment after they verify two other comments. This  
232 makes the users self moderate the comments on the platform. All comments are annotated by multiple users and the  
233 toxicity score is an average of binary toxicity classifications. In figure 4 we plotted the amount of comments for each  
234 length. As you can see there is a clear spike of comments at a length of 1000 words. This can be explained due to a  
235 maximum amount of words per comment. Because the BERT model takes 512 words at maximum we truncated every  
comment which was longer than this.

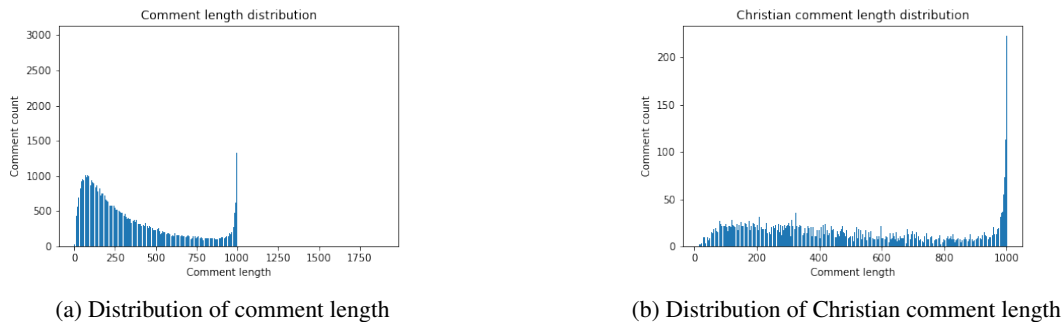


Figure 4: Two comment length distributions of the civil comments dataset

236

237 The Civil comments dataset was available through a challenge on the Kaggle dataset platform<sup>11</sup>. It contains 1971916  
238 comments from the Civil comments platform. This was a platform created to make comments more civilized by letting  
239 users that want to comment first rate other people's comments. The dataset was split into a training set and a test set,  
240 with 90 percent of the data being the training set and 10 percent the test set. Each comment row contained information  
241 about other peoples like reactions, multiple forms of toxicity and information about the person who commented. For this  
242 dataset the task is to predict whether a comment was toxic or not. The selected sensitive attribute was the commenter  
243 being Christian which is reported as a binary Christian or not Christian (1 or 0, respectively).

244 The following steps were taken to preprocess the data. In 78 percent of the comments it was not known whether the  
245 commenter was Christian or not. Therefore these comments were removed from the dataset and it left us with 370646  
246 comments in the train set and 5811 comments in the test set. The comments also need to be tokenized to be able to use  
247 it for the BERT model and we used the BERT tokenizer to do this. The tokenized comments also needed to be padded  
248 to be able to fit the dimensions of the BERT model as it takes tensors with a size of 512 tokens.

#### 249 A.1.2 CelebA Dataset

250 The CelebA dataset was obtained through the Large-scale CelebFaces Attributes Dataset website<sup>12</sup> but it was stored on  
251 Google Drive. The dataset contains 202599 RGB images of 10177 celebrities. The dataset was split into a training  
252 set of 162770 images, a validation set of 19868 images and a test set of 19961 images. There were three types of  
253 annotations available: landmark annotations, attributes annotations and identities annotations. We only used the attribute  
254 annotations, these contain annotations about their appearance. The task at hand is to predict whether a celebrity has  
255 blond hair. The selected sensitive attribute (by the original paper) is sex just like in the Adult Dataset. It is reported as a  
256 binary for the attribute "Male".

<sup>11</sup><https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification/>

<sup>12</sup><https://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>

257 The only necessary preprocessing step was to resize the images. All images in the dataset have the dimension 178 x 218  
258 x 3. But ResNet50 takes dimension 224 x 224 x 3 as input. Therefore each image had to be resized to this dimension.

259 To get some insights into the data we decided to calculate the average images for some specific subgroups as you  
260 can see in figure 5. In subplot 5b you can clearly see a female average for all celebrities with blond faces. This was  
261 confirmed when we calculated the percentage of blond celebrities and the percentage of blond male celebrities. From  
262 the complete dataset 14 percent was blond and of this subgroup 5 percent was male. This clearly indicates a minority  
group in the dataset.

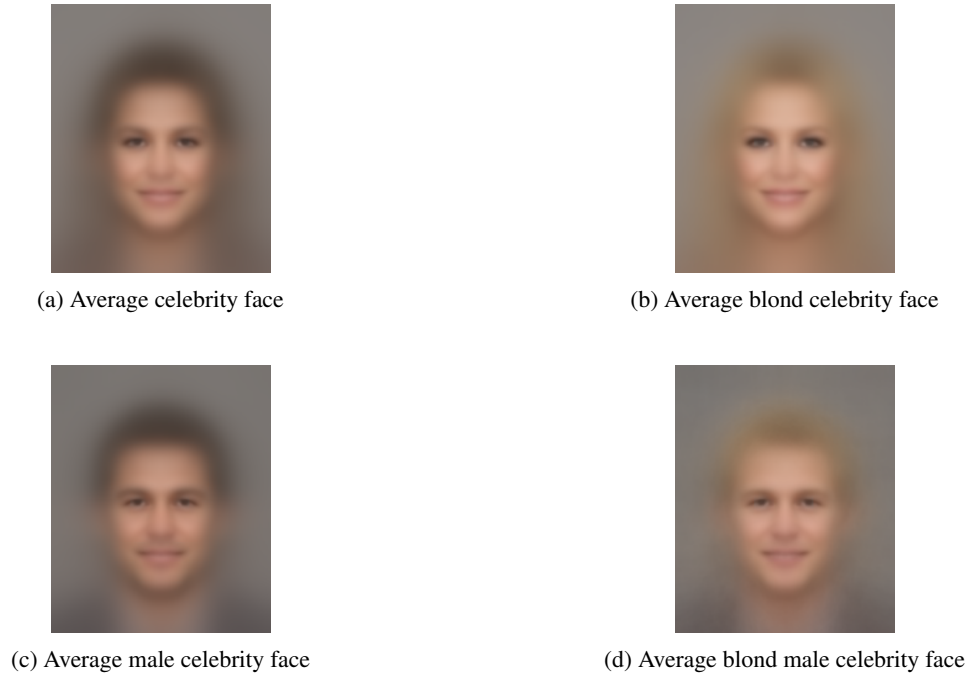


Figure 5: Average faces of different subgroups of the dataset

263

### 264 A.1.3 Adult Dataset

265 The Adult dataset was obtained through the AI Fairness 360 toolkit [1]. The dataset records census data of 45222  
266 individuals (after removing incomplete rows), split in a training set of 29100 and a test set of 15060 data points. Each  
267 data point has both continuous attributes such as age, capital gain and capital loss, and categorical attributes such as sex,  
268 marital status and native country. The task corresponding to this dataset is to predict whether each individual earns  
269 more than 50K per year. The selected sensitive attribute (by the original paper) for this dataset is sex which is reported  
270 as a binary: "Male" and "Female".

271 We took the following steps to preprocess the data. Firstly, we normalized the continuous attributes to have zero mean  
272 and unit variance. Secondly, we introduced a bias in the data the same way as in the original paper by removing all but  
273 the first 50 rows for which the protected group  $D = 0$  (our experiments showed that this was probably "Female"), and  
274 the target is ">50k". This removes 1062 samples from the training set.

### 275 A.1.4 Chexpert Dataset

276 We obtained the small CheXpert dataset from the website of the Stanford ML Group<sup>13</sup>. This dataset consists of chest  
277 x-ray images and corresponding attributes annotated by experts. These attributes indicate the presence of diseases  
278 (pathology), which are all possible classification targets, and the presence of a support device. The task for this  
279 experiment is to predict the presence of pleural effusion, and the chosen sensitive attribute is the presence of a support  
280 device.

<sup>13</sup><https://stanfordmlgroup.github.io/competitions/chexpert/>

281 The preprocessing of CheXpert was relatively straight forward. We removed all data points where either the sensitive  
282 attribute or the target was labeled uncertain. The x-ray images were all resized to 224 by 224 pixels and stacked to  
283 simulate the red, green and blue color channels that Densnet121 expects.