
[RE] An Implementation of Fair Robust Learning

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1

2 **Scope of Reproducibility**

3 This work attempts to reproduce the results of the 2021 ICML paper To be Robust or to be Fair: Towards Fairness
4 in Adversarial Training. I first reproduce classwise accuracy and robustness discrepancies resulting from adversarial
5 training, and then implement the authors' proposed Fair Robust Learning (FRL) algorithms for correcting this bias.

6 **Methodology**

7 In the spirit of education and public accessibility, this work attempts to replicate the results of the paper from first
8 principles using Google Colab resources. To account for the limitations imposed by Colab, a much smaller model and
9 dataset are used. All results can be replicated in approximately 10 GPU hours, within the usual timeout window of an
10 active Colab session. Serialization is also built into the example notebooks in the case of crashes to prevent too much
11 loss, and serialized models are also included in the repository to allow others to explore the results without having to
12 run hours of code.

13 **Results**

14 This work finds that (1) adversarial training does in fact lead to classwise performance discrepancies not only in standard
15 error (accuracy) but also in attack robustness, (2) these discrepancies exacerbate existing biases in the model, (3)
16 upweighting the standard and robust errors of poorly performing classes during training decreased this discrepancy for
17 both both the standard error and robustness and (4) increasing the attack margin for poorly performing classes during
18 training also decreased these discrepancies, at the cost of some performance. (1) (2) and (3) match the conclusions of
19 the original paper, while (4) deviated in that it was unsuccessful in helping increasing the robustness the most poorly
20 performing classes. Because the model and datasets used were totally different from the original paper's, it is hard to
21 quantify the exact similarity of our results. Conceptually however, I find very similar conclusions.

22 **What was easy**

23 It was easy to identify the unfairness resulting from existing adversarial training methods and implement the authors'
24 FRL (reweight) and FRL (remargin) approaches for combating this bias. The algorithm and training approaches are well
25 outlined in the original paper, and are relatively accessible even for those with little experience in adversarial training.

26 **What was difficult**

27 Because of the resource limitations imposed, I was unable to successfully implement the suggested training process
28 using the authors' specific model and dataset. Also, even with a smaller model and dataset it was difficult to thoroughly
29 tune the hyperparameters of the model and algorithm.

30 **Communication with original authors**

31 I did not have contact with the authors during the process of this reproduction. I reached out for feedback once I had a
32 draft of the report, but did not hear back.

33 1 Introduction

34 The advent of adversarial examples (1)(2) has motivated the need for procedures which decrease the sensitivity to noise
35 of learned models (which I will call adversarial robustness or simply robustness.) Once such method is adversarial
36 training (3)(4), in which adversarial examples are generated during the training process and are mixed in with "clean"
37 examples to create mixed training batches of both manipulated and unmanipulated images. Learning on these batches
38 has been shown to improve the robustness of models to adversarial attacks, often at a slight cost to standard performance
39 (accuracy.)

40 To be Robust or to be Fair: Towards Fairness in Adversarial Training identifies that adversarial training creates unfairness
41 in the resulting robust model. While the overall robustness of the model improves, some classes in the resulting model
42 are more robust to adversarial attacks than others. Not only are the robustness benefits unfairly distributed, so too are
43 the standard performance losses; the classes which are less robust at the end of the procedure tend to be the ones which
44 suffer more in terms of standard performance. Moreover, these classes tend to be the ones which were harder to learn
45 before adversarial training. As Xu et al. describe it: "adversarial training tends to make the hard classes even harder to
46 be classified or robustly classified."

47 Motivated by this unfairness, Xu et al. conduct a theoretical analysis of the problem to explain this empirically observed
48 phenomenon. They then draw on (5) to describe robust error in terms of the sum of standard errors (i.e. the probability
49 that a class will be incorrectly classified without manipulation) and boundary errors (i.e. the probability that there exists
50 some ϵ -ball attack which can change a classifier's decision on a given class.) Using this description, they reformulate
51 the learning problem into a series of cost-sensitive classification problems that can be penalized for violating fairness
52 constraints. With this reformulation, they present two FRL algorithms for making adversarial training more fair: one
53 which upweights the error of classes which violate the fairness constraints during training, and one which increases the
54 attack radius for classes which violate fairness constraints during training.

55 2 Scope of reproducibility

56 The focus of this reproduction will be attempting to demonstrate the following:

- 57 • Claim 1, which is supported by Experiment 1 in Figure 1, is that adversarial training creates unfair outcomes
58 in terms of both robustness and standard error.
- 59 • Claim 2, which is also supported by Experiment 1 in Figure 1, is that this unfairness exacerbates existing
60 biases in model performance.
- 61 • Claim 3, which is supported by Experiment 2 in Figure 2, is that upweighting the error of classes which
62 violates fairness constraints (using the authors' FRL: reweight algorithm) can improve the both the standard
63 errors for the most poorly performing classes, and to a lesser degree their robustness.
- 64 • Claim 4, which is explored by Experiment 3 in Figure 3, is that increasing the margin of attack for classes
65 which violates fairness constraints (using the authors' FRL: remargin algorithm) can also improve the fairness
66 of the model— perhaps more effectively than reweighting.

67 3 Methodology

68 As an educational exercise, I aimed to re-implement the authors' training approaches from their descriptions in the
69 paper. Because of the limitation imposed on the resources, however, I opted to use a simpler model and dataset in my
70 experiments.

71 3.1 Model descriptions

72 The paper used the PreAct-ResNet18 and WRN28 architectures for their experimentation; I opted for the LeNet-5
73 architecture in the interest of efficiency. Though it is a much simpler model than the paper's originals, it provided
74 enough complexity to conduct my experiments.

75 3.2 Datasets

76 The paper used the CIFAR10 and SVHN datasets for their experimentation; I used the Fashion-MNIST dataset. The
77 train set is comprised of 60,000 examples, the test set 10,000. Both have a uniform label distribution across all 10
78 classes. The original train and test sets are used Experiment 1, while Experiments 2 and 3 split the train set into an
79 80/20 train/validation set for the FRL process. The only preprocessing done was to resize the images from 28x28 to
80 32x32. The data is freely available here.

81 3.3 Hyperparameters

82 The fairness tolerance hyperparameter was selected based on the recommendations in the paper (5%), as was the
83 baseline ϵ (8/255 for the PGD attack.) For Experiment 1 I used a learning rate of 1e-3 for regular training and adversarial
84 training, as the paper recommended. Due to resource constraints I had to limit the number of epochs I trained for to 15,
85 and from convergence behavior I decayed the learning rate more often than the original paper (every 4 rounds by a
86 factor of 3, as opposed to every 40 rounds by a factor of 10.) For the simpler model and dataset, this worked well.

87 For Experiments 2 and 3 I used a baseline learning rate of 1e-4, which I selected based on unstable behavior at a rate of
88 1e-3. I suspect this is due to differences in the model and dataset used, as well as the way I implemented the reweighting
89 and remargining systems.

90 I utilized the results of the fairness evaluation (ϕ values) in the training process by applying a Softmax function to
91 creating cross-entropy loss weightings, and as such the α values were different than the original paper's. I tried a variety
92 of α values in the space of (1, 2, 5, 10,) and a variety of ratios of natural- α s to boundary- α s. The best results came from
93 a ratio of 5:1 natural:boundary error weighting, which decreased the worst-case standard error by 25%, and the worst
94 case robust error by 11%.

95 3.4 Experimental setup and code

96 For Experiment 1, I defined the LeNet-5 architecture and trained a classifier on the Fashion-MNIST dataset for 15
97 epochs at a learning rate of 1e-3. I then adversarially trained a new LeNet-5 model using a PDG attack for the same
98 number of epochs at the same learning rate, with a 50/50 mixture of clean and manipulated images. I then compared the
99 classwise standard accuracy (i.e. ability to predict a "clean" image correctly) and robust accuracy (i.e. ability to predict
100 a image correctly despite manipulation) of the natural model and adversarially trained model. The results are recorded
101 in Figure 1.

102 For Experiments 2, I retrained the unfair adversarially-trained model under the FRL (reweight) paradigm. During this
103 procedure, I recorded the overall and classwise standard and boundary errors of the model during each batch, and based
104 on these errors I re-calculated loss weights for each class. The loss function used was the sum of the standard loss
105 and the loss for adversarially manipulated images with respect to the predictions on their unmanipulated counterparts
106 (corresponding to standard error and boundary error, respectively.) Classes were penalized based on violations of
107 fairness constraints, i.e. how greatly they differed from the average standard and boundary errors for all classes. I ran
108 10 rounds of retraining, and then compared the original unfair adversarially trained model with its retrained counterpart,
109 comparing classwise standard and robust accuracy. These results can be found in Figure 2.

110 Experiment 3 was much the same as Experiment 2, the only difference being that instead of simply upweighting
111 the loss of classes which violated fairness constraints, the radius of a class' attack during training was increased or
112 decreased based on the size of their violation. Again, I ran 10 rounds of retraining, and then compared the original
113 unfair adversarially trained model with its retrained counterpart, comparing classwise standard and robust accuracy.
114 These results can be found in Figure 3.

115 All the code for these experiments, as well as example notebooks that walk through the procedure, can be found here.

116 3.5 Computational requirements

117 As mentioned, I used Google Colab for all of the experimentation. As such, it is difficult to describe the exact hardware
118 that was used, or to even be confident of the consistency of the hardware throughout this process. I did use GPU
119 resources, though I cannot speak to any specific type.

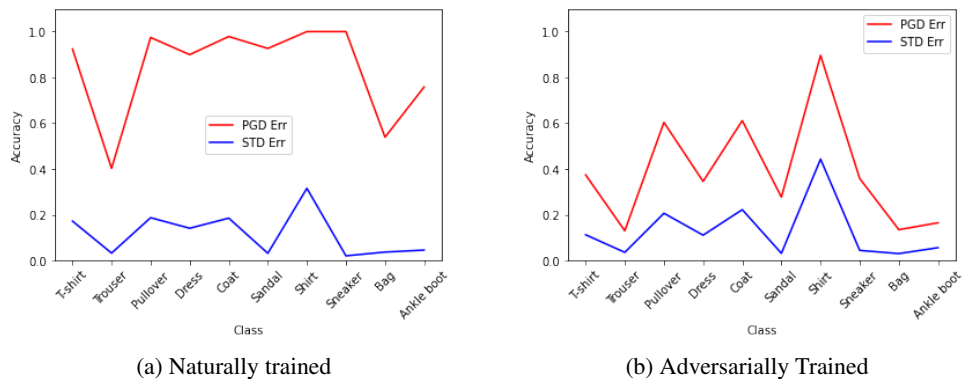


Figure 1: Adversarial training produces unfair outcomes across classes, and worsens existing performance discrepancies

120 Experiment 1 can be run in approximately 15 minutes of GPU time. Experiment 2 can be run in approximately 5 hours
 121 of GPU time (for all alpha-combinations) and Experiment 3 can be run in approximately 3 hours.

122 All three notebooks can sometimes be run in parallel, but not always. Colab can be a bit unpredictable.

123 4 Results

124 In my experiments, I found that:

- 125 • Adversarial training does in fact lead to classwise discrepancies in standard error and adversarial robustness,
 126 that the least robust classes in the resulting model are the ones the model originally had a hard time learning,
 127 and that the penalties to standard performance brought on by adversarial training exacerbate existing biases in
 128 model performance.
- 129 • Reweighting the natural and boundary errors to penalize classes violating fairness constraints during adversarial
 130 retraining can improve the fairness of the model with respect to standard error, and to a lesser degree robust
 131 error.
- 132 • Remargining the attack radius for classes violating fairness constraints during adversarial retraining can also
 133 improve the fairness (i.e. lower the variance across classes) of the model’s robustness (at a cost to robust
 134 performance) as well as improve the standard error.

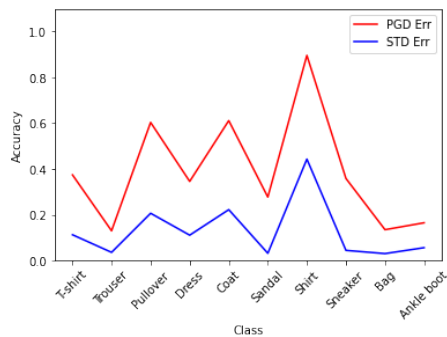
135 Most of these results agree with the paper’s conclusions, although the results in Experiment 3 differ in that I was not
 136 able to improve the robustness of the model with remargining as well as I could with reweighting. The original paper
 137 showed the opposite: that reweighting was unable to improve robustness for the most poorly performing classes. One
 138 experiment I did not conduct was to try both reweighting and remargining together, which the authors suggest might be
 139 fruitful. I leave that as a further exercise.

140 4.1 Results reproducing original paper

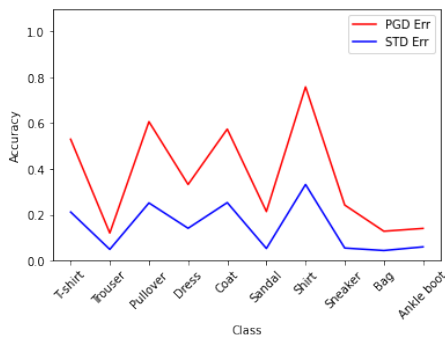
141 4.1.1 Result 1

142 The result of Experiment 1 (shown in Figure 1) relates to claims 1 and 2 in Section 2. I found that in the naturally
 143 trained model, the standard error is quite low and the adversarial error (PGD error) is quite high. The adversarial error
 144 is not quite as uniform as in the original paper, I suspect because of the simplicity of the dataset and model I used.

145 Still it is observable that after adversarial training, the model’s adversarial error is much lower across the board, but
 146 not in a fair way. Certain classes are much more robust to attack than others, and in particular the classes which had
 147 poorer initial standard performance are the ones with worse adversarial robustness. Moreover, we can see that there are
 148 penalties to standard performance incurred as a result of adversarial training, and the classes which suffer the most are
 149 the ones the natural model already had a hard time learning.



(a) "Vanilla" adversarially-trained model



(b) After the FRL (Reweight) procedure

Figure 2: FRL Reweight is able to mitigate standard performance losses, while also increasing the robustness of the most difficult class

150 Indeed, as Xu et al. put it, "adversarial training tends to make the hard classes even harder to be classified or robustly
 151 classified." This is exactly what I found, even with a totally different model and dataset.

152 **4.1.2 Result 2**

153 The result of Experiment 2 (shown in Figure 2) relates to claim 3 in Section 2. Here we can see the result of my best
 154 attempt at reweighting the loss of classes during adversarial retraining based on their violation of fairness constraints.
 155 As per the paper's FRL retraining algorithm, I began with an adversarially trained model and iteratively tried to retrain
 156 it, adjusting the loss of each class as I went depending on whether it violated fairness, and to what degree. As such, I
 157 compared the "vanilla" adversarially trained model with the resulting model after retraining.

158 I observed that for the hardest class to classify, there is a 25% reduction in standard error (bringing it nearly in line with
 159 the naturally trained model) and an 11% reduction in robust error. This is not totally free; we can observe, for example,
 160 that the standard and robust error for some of the easier classes suffers as a result. Still, the resulting model is fairer
 161 than it originally was.

162 These results seem relatively in-line with the original paper's, though again because of the different model and dataset
 163 selected it is hard to quantify the exact similarity. The overall conclusion is much the same though: reweighting is
 164 hugely successful in decreasing the classwise standard error discrepancies brought on by adversarial training, and to a
 165 lesser degree in decreasing classwise robustness discrepancies.

166 **4.1.3 Result 3**

167 The result of Experiment 3 (shown in Figure 3) relates to claim 4 in Section 2. This is the result of my best attempt at
 168 remargining during the retraining procedure. I observed a slight improvement in the worst-case standard error, but little
 169 to no improvement in the worst case robustness, and indeed a general degradation in robustness across most classes.

170 These results were not in line with the paper's, which found FRL (Remargin) to be more effective than FRL (Reweight.)
 171 This may be due to differences in our datasets, or artifacts of my implementation. It should be noted that because of
 172 the greater expense of this procedure, it was harder to thoroughly explore its hyperparameters, and this is still an
 173 interesting area of exploration for me.

174 **5 Discussion**

175 I believe that overall my results are quite in line with the original paper's. I found that adversarial training does produce
 176 unfair results, both in the improvements to robustness the model receives as well as the degradation of standard error
 177 it experiences. I also found that these unequal costs penalize classes that are harder for the model to learn, making
 178 it worse at what classifying what it already had trouble with. Finally, I found that the FRL (reweight) approach was

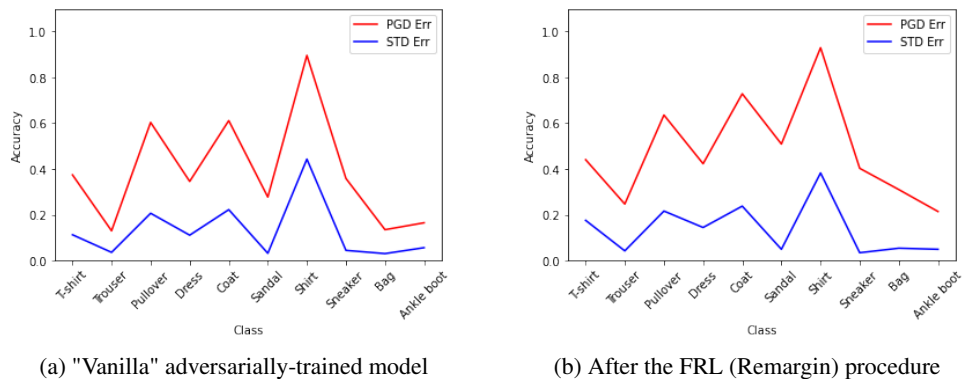


Figure 3: FRL Remargin was able to slightly improve the standard performance of the hardest class, but decreased the overall robustness of the model

179 able to mitigate most of the degradation in standard performance for the hardest to learn classes, and to a lesser degree
 180 improve the robustness for that class as well as well as the overall robustness.

181 One weak point of my implementation was in the FRL (Remargin) procedure. I was unable to successfully improve the
 182 model’s robustness via remargining, though I am not confident that I thoroughly explored the space. It was the most
 183 costly procedure I ran, and it ran into its fair share of Colab timeouts, making hyperparameter tuning tricky.

184 One last experiment I did not have time for was a combination of reweighting and remargining, which Xu et al. suggest
 185 is the most effective means increasing adversarial fairness. This is because I wanted positive results in remargining
 186 before attempting to combine the two approaches, which I was unfortunately unable to achieve. This is still an open
 187 question to pursue.

188 5.1 What was easy

189 One of the paper’s easiest claims to verify was that adversarial training creates the unfair outcomes described above.
 190 Even with little experience in adversarial training, we found that with only a bit of effort I could observe this phenomenon
 191 myself.

192 It was also fairly easy to implement Xu et al’s FRL algorithms; the remargining and reweighting procedures are very
 193 clearly explained in the paper and were straightforward to put into code. One aspect of the paper not discussed in
 194 this report is their theoretical analysis, which was also very clear and helped motivate and explain the FRL problem
 195 formulation.

196 5.2 What was difficult

197 As mentioned above, the part I had the most difficulty with was the remargining procedure. It took much longer
 198 than anticipated, and its expense made automated hyperparameter searches difficult. Because I was unsuccessful in
 199 improving the model’s robustness with remargining, I was also hesitant to implement a combined FRL (Reweight) and
 200 FRL (Remargin) approach, which the authors suggest might be the most effective result. As mentioned, this is an area
 201 in which I am still actively exploring. Hopefully in the future I can replicate their success there too.

202 5.3 Communication with original authors

203 As mentioned in my summary, I did not have contact with the authors throughout this process. It was only upon drafting
 204 my report that I learned it was encouraged to contact the original authors; in the future, I think it would be a great idea
 205 to communicate with them sooner. I reached out with a preprint of the report for any feedback or suggestions, but did
 206 not hear back.

207 **References**

- 208 [1] Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. arXiv preprint
209 arXiv:1412.6572, 2014.
- 210 [2] Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties
211 of neural networks. arXiv preprint arXiv:1312.6199, 2013.
- 212 [3] Kurakin, A., Goodfellow, I., and Bengio, S. Adversarial machine learning at scale. arXiv preprint arXiv:1611.01236,
213 2016.
- 214 [4] Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to
215 adversarial attacks. arXiv preprint arXiv:1706.06083, 2017.
- 216 [5] Zhang, H., Yu, Y., Jiao, J., Xing, E. P., Ghaoui, L. E., and Jordan, M. I. Theoretically principled trade-off between
217 robustness and accuracy. arXiv preprint arXiv:1901.08573, 2019b.