Fairness Through Robustness: Investigating Robustness Disparity in Deep Learning

Vedant Nanda*, Samuel Dooley*, Sahil Singla, Soheil Feizi, John P. Dickerson

9 March 2021

FAccT 2021



ROBUSTNESS BIAS

Most notions of fairness are defined by system outputs

Noise can cause a system to treat groups unfairly



Are some types of points systematically more robust than others?



The closer points are more likely to be misclassified

DATA PRELIMINARIES – PARTITION

Partition by Predictor Variable



Dataset



Partition by Sensitive Attribute



INTUITION



INTUITION



INTUITION



Predictor Partition	Robustness	Sensitive Partition	Robustness
Red	2/3	Circles	2/3
Blue	2/3	X's	2/3

ROBUSTNESS BIAS

Depends on the Data, Model, and Partition *and also on tau



METRICS

Given a dataset *D*, $\tau > 0$, and partition $P \subset D$. -

 $\widehat{I_P}(\tau) =$

1.0

0.8

0.6

0.2

0.0

0

2

10 0.4

- Let $d_{\theta}(x)$ be the minimal distance from point x to a decision boundary.
- Calculate -

Proportion of that are more than tau from decision boundary

Proportion $d_{\theta} > \tau$

4

τ

6

Relative average distance to the decision boundary





RESEARCH QUESTIONS

- Does robustness bias exist in the wild? -
- How much does it depend on the dataset? the *partition*? the **model**? -



BUT WE CAN'T COMPUTE $d_{\theta}(x)$

- All our models (except MLP) do not permit direct computation
- We can approximate $d_{\theta}(x)$ with upper or lower bounds



DeepFool and CarliniWagner



Randomized Smoothing

RESULTS

RECALL: Red means partition is <u>more</u> vulnerable Blue means partition is <u>less</u> vulnerable



- 2 - 1 - 0.5

0

-0.5

--1 --2

RESULTS

RECALL: Red means partition is more vulnerable Blue means partition is less vulnerable



RESULTS

Our results suggest that:

- Robustness Bias appears to be present in the wild.
- Some classes appear to uniformly be more/less robust across all models.
- Some partitions of datasets appear to exhibit robustness bias more than others.
- DeepFool and CarliniWagner agree on the sign of the bias (with \$p<0.001\$ in the Pearson's Chi-squared).
- There is some evidence of agreement between the upper and lower bounds.

CAN WE MITIGATE ROBUSTNESS BIAS?

Attempted the obvious regularization fix.



1.0

0.8

0.6

0.4

0.2

0.0

0.0

(Fu

FUTURE DIRECTIONS

- Use adversarial training to balance mitigation with overall robustness hit
- What about on deployed models?
 - Quantifying harms ethical concerns for such measurements makes it hard
- How do other methods to mitigate fairness affect robustness bias?
- How does robustness bias change for robust models?
- Benchmarking current models for other kinds of perturbation based biases
 - Synthetic corruptions eg: ImageNet-C
 - Natural corruptions eg: ImageNetV2

Check out our paper!

https://arxiv.org/abs/2006.12621 Get in touch: {vedant,sdooley1}@cs.umd.edu



