



EMORY

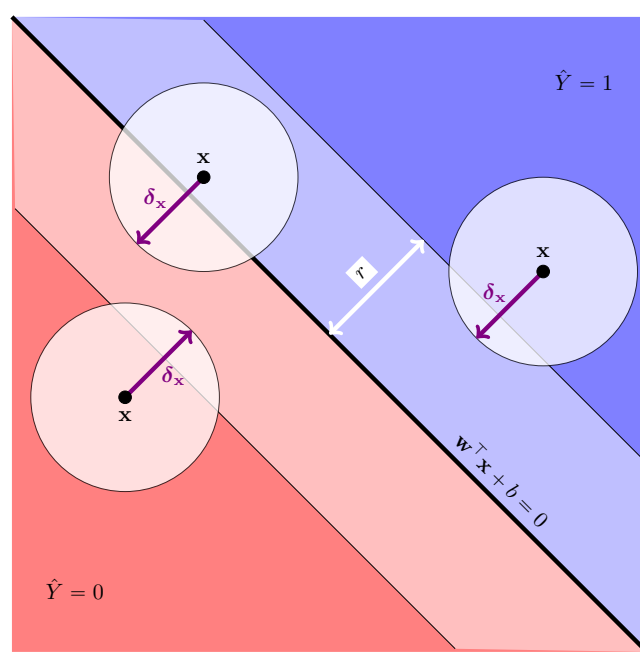
FAST & FAIR: EFFICIENT SECOND-ORDER ROBUST OPTIMIZATION FOR FAIRNESS IN ML

Allen Minch¹, Hung Anh Vu², Annie Warren³, advised by Dr. Elizabeth Newman⁴¹ Brandeis University, ² University of Maryland, ³ University of Minnesota, ⁴ Emory University

Abstract.

This project explores adversarial training techniques to develop fairer Deep Neural Networks (DNNs) to mitigate the inherent bias they are known to exhibit. DNNs are susceptible to inheriting bias with respect to sensitive attributes such as race and gender, which can lead to life-altering outcomes (e.g., demographic bias in facial recognition software used to arrest a suspect). We propose a robust optimization problem to improve fairness in DNNs, and leveraging second-order information, we are able to efficiently find a solution.

Adversarial training



Our research investigates whether implementing robust optimization - so that two nearby points are more likely to be classified similarly - would improve fairness. Robust optimization takes into consideration at a radius r around a data point (see fig. 1).

Optimization problem:

$$\min_{\theta} \frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{T}} \left[\max_{\|\delta_{\mathbf{x}}\| \leq r} L(f_{\theta}(\mathbf{x} + \delta_{\mathbf{x}}), \mathbf{y}) \right] + R(\theta)$$

Figure 1: Robust optimization

Solving the inner optimization problem

Projected Gradient Descent: We use a first-order method with iterates

$$\delta_{\mathbf{x}}^{(k+1)} = P \left[\delta_{\mathbf{x}}^{(k)} + \alpha^{(k)} \cdot \nabla_{\mathbf{x}} L(f_{\theta}(\mathbf{x} + \delta_{\mathbf{x}}), \mathbf{y}) \right]$$

where P is a projection operator such that the constraint, $\|\delta_{\mathbf{x}}\| \leq r$, is satisfied.

Trust Region Subproblem: We use a second-order Taylor approximation of our loss function and solve

$$\delta_{\mathbf{x}}(\lambda) = -(\nabla_{\mathbf{x}}^2 L(f_{\theta}(\mathbf{x}), \mathbf{y}) + \lambda \mathbf{I})^{-1} \nabla_{\mathbf{x}} L(f_{\theta}(\mathbf{x}), \mathbf{y})$$

where $\lambda \geq 0$ is chosen such that the constraint is satisfied.

Error in trust region subproblem solution

For binary classification, we use a logistic regression loss function with a sigmoid activation function $\sigma(z) = \frac{1}{1+e^{-z}}$, so we have the following inner optimization problem:

$$\max_{\|\delta_{\mathbf{x}}\| \leq r} \left[-y \ln(\sigma(\mathbf{w}^{\top}(\mathbf{x} + \delta_{\mathbf{x}}) + b)) - (1-y) \ln(1 - \sigma(\mathbf{w}^{\top}(\mathbf{x} + \delta_{\mathbf{x}}) + b)) \right]$$

When we use a second-order approximation (RHS) instead of the true loss function (LHS), we end up solving two slightly different problems with differing terms

$$\sigma(\mathbf{w}^{\top}(\mathbf{x} + \delta_{\mathbf{x}}) + b) \neq \sigma(\mathbf{w}^{\top} \mathbf{x} + b) + \sigma'(\mathbf{w}^{\top} \mathbf{x} + b) \mathbf{w}^{\top} \delta_{\mathbf{x}}$$

Taylor expanding the LHS, we can obtain the following:

$$\sigma(\mathbf{w}^{\top}(\mathbf{x} + \delta_{\mathbf{x}}) + b) \approx \sigma(\mathbf{w}^{\top} \mathbf{x} + b) + \sigma'(\mathbf{w}^{\top} \mathbf{x} + b) \mathbf{w}^{\top} \delta_{\mathbf{x}} + \frac{1}{2} \delta_{\mathbf{x}}^{\top} \mathbf{w} \sigma''(\mathbf{w}^{\top} \mathbf{x} + b) \mathbf{w}^{\top} \delta_{\mathbf{x}} + \dots$$

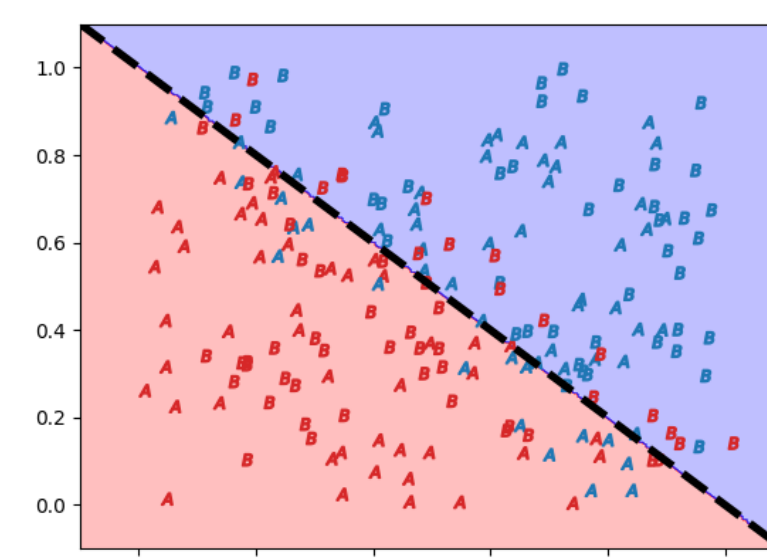
The error of solving the trust region subproblem comes from truncating the higher-order terms in the Taylor expansion. The error from truncating the quadratic term depends on the magnitude $|\sigma''(z)| \leq 0.1$ and the $\delta_{\mathbf{x}}$ for which we solved.

Fairness metrics

Let \hat{Y} be the classifier's prediction, Y be the true class, and s be a binary sensitive attribute. How do we measure the **fairness** of a binary classifier with respect to s ?

$$\begin{aligned} \text{Independence:} \quad & \mathbb{P}(\hat{Y} = 1 | s = 0) = \mathbb{P}(\hat{Y} = 1 | s = 1) \\ \text{Separation:} \quad & \mathbb{P}(\hat{Y} = 1 | Y = 1, s = 0) = \mathbb{P}(\hat{Y} = 1 | Y = 1, s = 1) \\ \text{Sufficiency:} \quad & \mathbb{P}(Y = 1 | \hat{Y} = 1, s = 0) = \mathbb{P}(Y = 1 | \hat{Y} = 1, s = 1) \end{aligned}$$

Hiring (synthetic data)

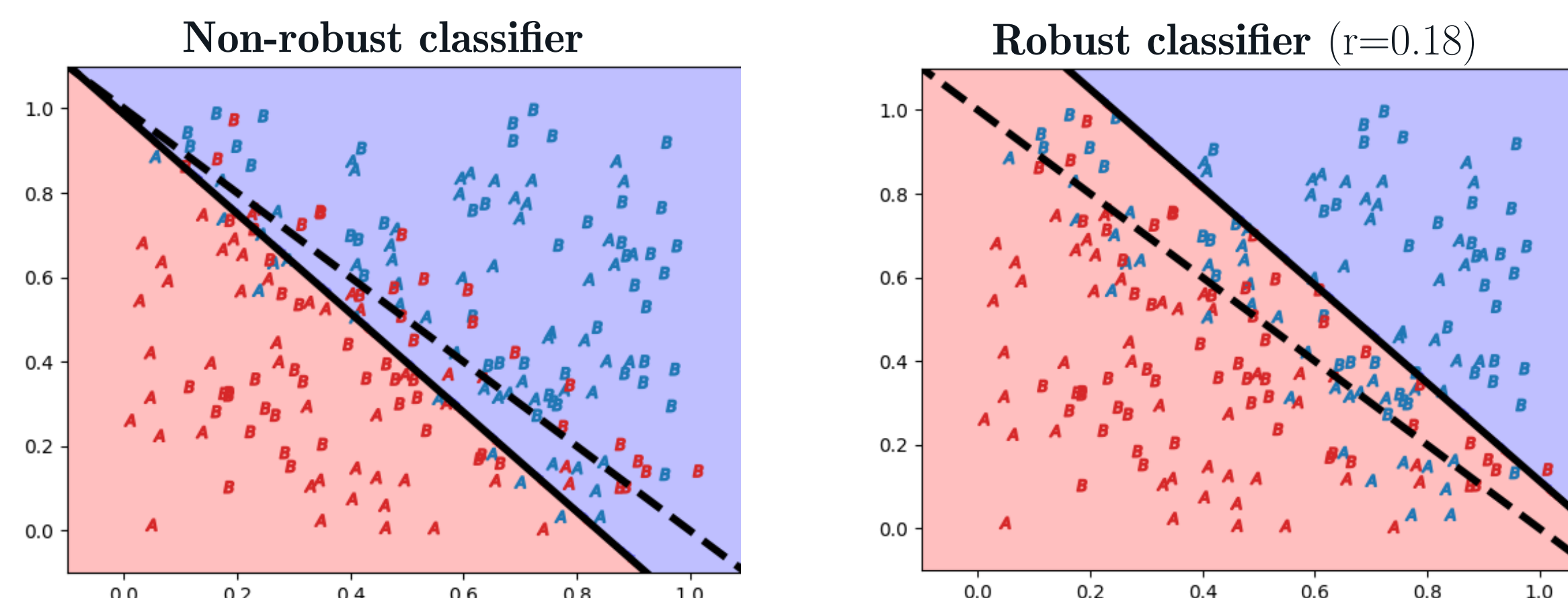


- Y : should be **hired** or **not hired**
- \hat{Y} : predicted to be **hired** or **not hired**
- s : sensitive attribute A or B
- **unfairness:** B s shifted up and right
As shifted down and left

All red individuals in the blue region who would be hired in error are B s, while all blue individuals in the red region who would be incorrectly not hired are A s.

Figure 2: Unfair setup of data

Hiring results

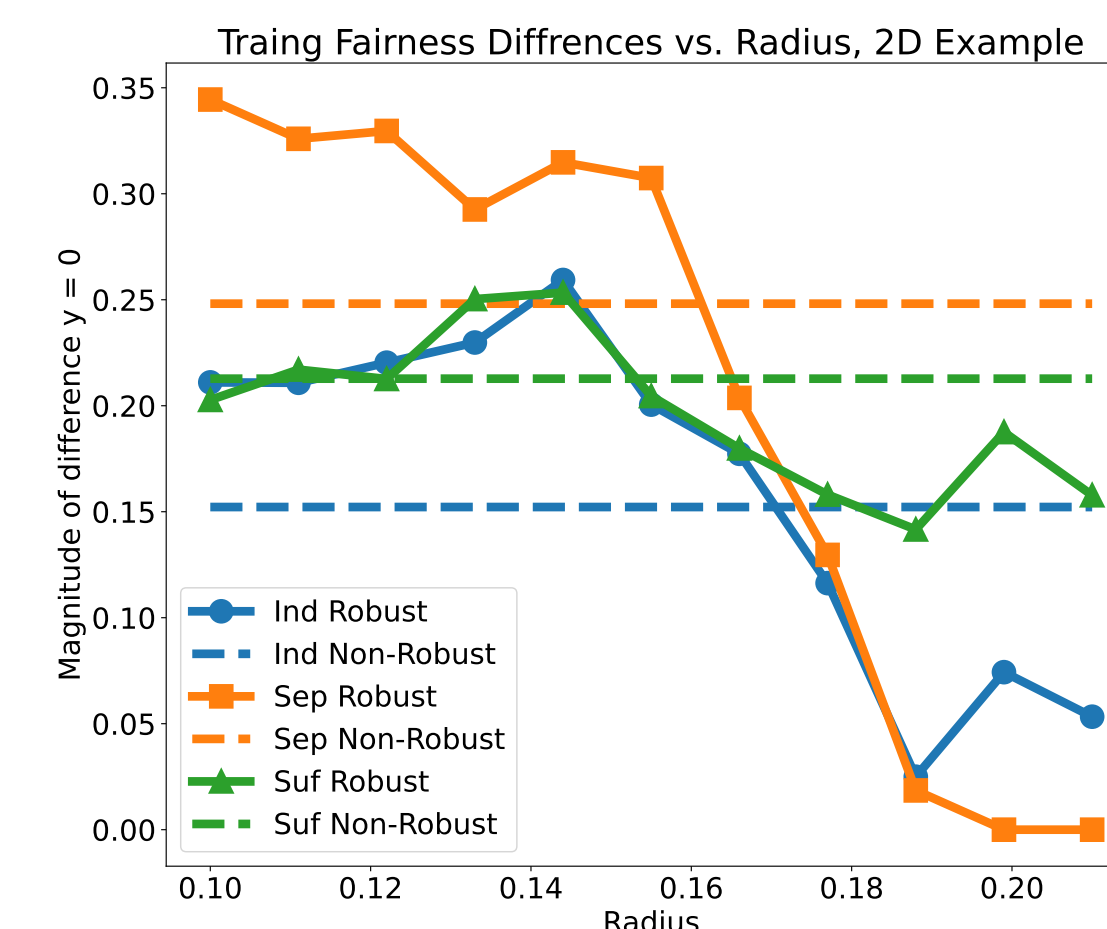
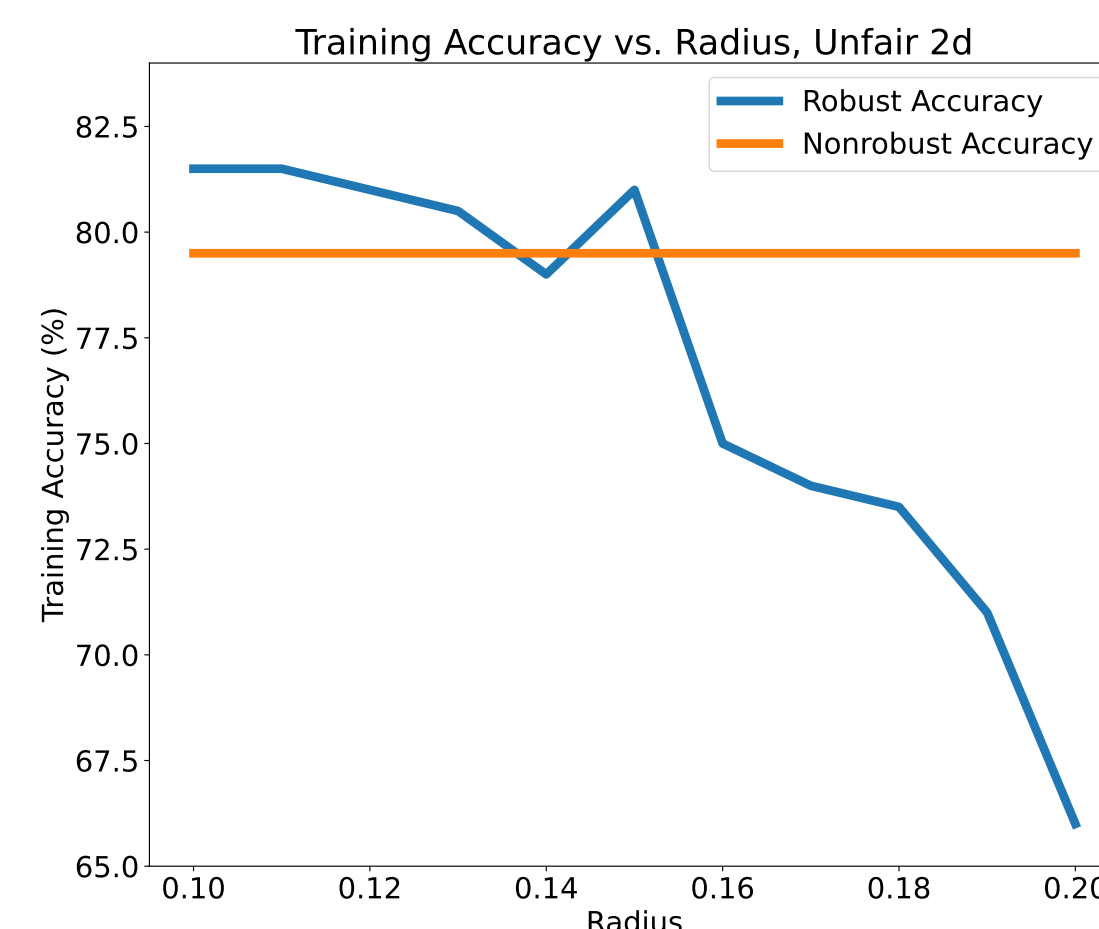


Diff:	Y=0 S1-S0	Y=1 S1-S0
Ind.	0.152	0.152
Sep.	0.248	0.179
Suff.	0.213	0.190

Training Accuracy: 79.5%
Test Accuracy: 78.0%

Diff:	Y=0 S1-S0	Y=1 S1-S0
Ind.	0.025	0.025
Sep.	0.019	0.127
Suff.	0.142	0.038

Training Accuracy: 73.5%
Test Accuracy: 73.0%



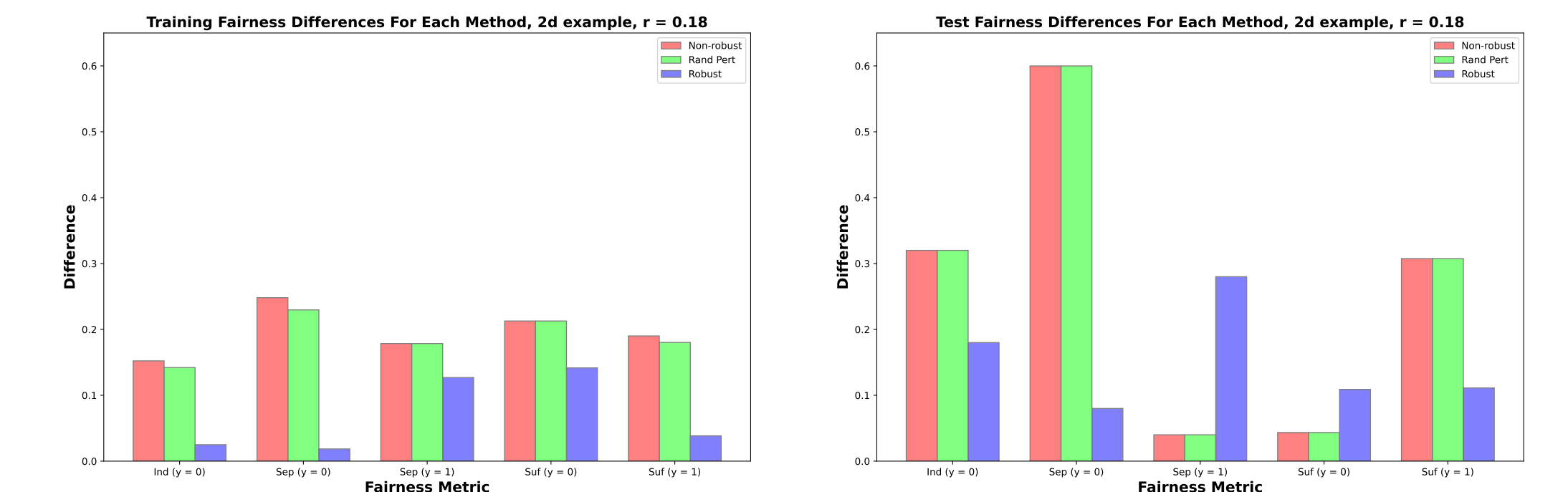
Which method is faster?

Avg. Epoch Time $\frac{PGD}{Trust}$

Dataset:	Min	Max
Synthetic	1.377	3.130
LSAT	2.852	9.639
Adult	8.497	31.407

Unlike in PGD, where we may have to do many gradient computations, with the second-order method, we at worst solve a system of linear equations and use a bisection method. The trust region method has proven to be the more efficient method.

Which method is fairer?



Conclusions

- Utilizing the trust region subproblem method **significantly improves efficiency**: computing second-order information using **hessQuik** outperforms first-order PGD across all perturbation radii on three different data sets.
- **Robustness can improve fairness, but potentially at the cost of accuracy.** Fairness improves as the perturbation radius increases but accuracy decreases in both training and testing data as the radius increases.
- If using robust training with a certain radius improves fairness, it appears to improve fairness by larger margins compared to random perturbation; **solving the optimization problem well is worthwhile.**

Acknowledgements

- Thank you to Dr. Elizabeth Newman, our mentor, for her guidance and support.
- This work is supported in part by the US NSF award DMS-2051019.

References

- [1] Martin Arjovsky et al. *Invariant Risk Minimization*. 2020.
- [2] Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning: Limitations and Opportunities*. <http://www.fairmlbook.org>. fairmlbook.org, 2019.
- [3] Elizabeth Newman and Lars Ruthotto. "hessQuik": Fast Hessian computation of composite functions". In: *Journal of Open Source Software* 7.72 (2022), p. 4171.
- [4] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. 2e. New York, NY, USA: Springer, 2006.
- [5] Tai Le Quy et al. "A survey on datasets for fairness-aware machine learning". In: *WIREs Data Mining and Knowledge Discovery* 12.3 (2022).