

Fast & Fair: Efficient Second-Order Robust Optimization for Fairness in Machine Learning*

Allen Joseph Minch[†], Hung Anh Dinh Vu[‡], and Anne Marie Warren[§]

Project Advisor: Dr. Elizabeth Newman[¶]

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, which we demonstrate can improve fairness in several datasets, both synthetic and real-world, using an affine linear model. Leveraging second order information, we are able to find a solution to our optimization problem more efficiently than with a purely first order method.

Key words. machine learning, fairness, robust optimization, adversarial training, optimization

MSC codes. 65F10, 65F22, 65K05, 90C47

1. Introduction. Machine learning has become an integral part of data analysis with its powerful ability to reveal underlying patterns and structures in data. Deep Neural Networks (DNNs) in particular are the gold standard classifying complex data; however, there is a tendency for DNNs to inherit bias from the datasets on which they train. Bias in this sense is not statistical bias, but the ways in which individual advantages or disadvantages manifest in data. This can be especially problematic in areas where machine learning is used to make life-altering decisions such as criminal justice [4] and corporate hiring [7].

The ever-expanding use of machine learning poses a significant ethical question when models are known to perpetuate societal biases [7]. While an in-depth discussion of these ethical concerns is beyond the scope of this work, they motivate our efforts to improve fairness within the models themselves. The unfortunate truth is that the harmful biases we see in our models and in our data come from deep-rooted societal structures that are at present beyond the abilities of machine learning to correct. However, we feel that in the face of these larger issues it is our duty within our means to work towards fairer outcomes.

One way to potentially achieve fairer outcomes is to use adversarial training to introduce robustness to a model. Robust optimization aims to make the model less susceptible to small variations in data, known as adversarial attacks, but in doing so decreases model accuracy. Recently, the Fair-Robust-Learning framework was proposed to reduce this unfairness problem in adversarial training [13]. The authors demonstrated that a combination of fairness and adversarial regularization yielded fairer models on benchmark image classification datasets.

*Submitted to the editors January 31, 2024.

Funding: This work is supported in part by the US NSF award DMS-2051019.

[†]Department of Mathematics, Brandeis University, Waltham, MA (allenminch@brandeis.edu)

[‡]Department of Mathematics, University of Maryland, College Park, MD (hvu1@terpmail.umd.edu)

[§]University of Minnesota, Minneapolis, MN (warre659@umn.edu, <https://anniewarren.github.io>)

[¶]Department of Mathematics, Emory University, Atlanta, GA (elizabeth.newman@emory.edu).

35 Our research shares the goal of addressing fairness issues in DNNs through the use of
 36 adversarial training techniques, but focuses on an additive bias rather than out-of-distribution
 37 bias or other forms. We define fairness on different metrics (independence, separation, and
 38 sufficiency vs. average and worst-class boundary, robust, and standard errors) to measure
 39 additive bias with respect to sensitive ‘hidden’ attributes. Without aiming to cater to specific
 40 types of data, we explore the effects of adversarial training on this definition of fairness. A
 41 simultaneous focus is to improve the efficiency of solving robust optimization problems. To this
 42 end, we use second-order information to accelerate training, a concept that was not addressed
 43 in previous work [13].

44 We implement a second-order method, termed the “trust region subproblem” (TRS), de-
 45 signed explicitly to address inner optimization challenges encountered when introducing robust
 46 training. Our experiments, spanning both synthetic and real-world datasets, demonstrate the
 47 capabilities of robust optimization in enhancing fairness. We employ three distinct optimiz-
 48 ers for these tests, allowing us to compare their performance. Notably, the integration of
 49 `hessQuik` [8], has proven instrumental in efficiently deriving exact Hessians. This approach
 50 surpasses the projected gradient descent (PGD) method in terms of time efficiency while pro-
 51 ducing the same solution. For transparency and further community engagement, we’ve made
 52 our Python implementation, including all our experiments, available on our GitHub repository
 53 at Fast-N-Fair (<https://github.com/elizabethnewman/fast-n-fair>).

54 The paper is organized as follows: [Section 2](#) introduces DNNs and the necessary notation,
 55 robust optimization, and our choice of fairness metrics. [Section 3](#) describes our proposed sec-
 56 ond order method, our implementation, and is followed by an analysis of the error produced
 57 by approximations used in our methods ([Subsection 3.2](#)). In [Subsection 3.3](#), we introduce
 58 alternate methods of solving the robust optimization problem that are implemented as a com-
 59 parison to our proposed approach. [Section 4](#) first describes the setup of a synthetic dataset
 60 along with the preliminary fairness results, and then extends the discussion to several real
 61 world datasets. We also examine the relative computational efficiency of our different meth-
 62 ods of performing robust optimization. Lastly, [Section 5](#) concludes the paper and discusses
 63 potential future work.

64 **2. Background.** First we must discuss DNNs and some necessary notation, robust opti-
 65 mization, and our choice of fairness metrics.

66 **2.1. Notation.** Deep neural networks (DNNs) can be represented by a parameterized
 67 mapping $f_{\theta} : \mathcal{X} \times \Theta \rightarrow \mathcal{Y}$ from input-target pairs $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$, where $\mathcal{D} \subseteq \mathcal{X} \times \mathcal{Y}$ is the
 68 data space, $\mathcal{X} \subseteq \mathbb{R}^{n_{\text{in}}}$ is the input space, and $\mathcal{Y} \subseteq \mathbb{R}^{n_{\text{out}}}$ is the target space, and $\Theta \subset \mathbb{R}^{n_{\theta}}$
 69 is the parameter space. Our goal is to learn the weights $\theta \in \Theta$ such that $f_{\theta}(\mathbf{x}) \approx \mathbf{y}$ for all
 70 input-target pairs. Typically, learning the weights is posed as the optimization problem

$$71 \quad (2.1) \quad \min_{\theta} \frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{T}} L(f_{\theta}(\mathbf{x}), \mathbf{y}) + R(\theta)$$

72 where $\mathcal{T} \subset \mathcal{D}$ is the training set and $R : \Theta \rightarrow \mathbb{R}$ is a regularization term to enforce desirable
 73 properties on the weights.

74 For many problems with a well-chosen optimizer, we can solve (2.1) well. However, this
 75 can lead to problems such as overfitting, where the model fits the training data well but does

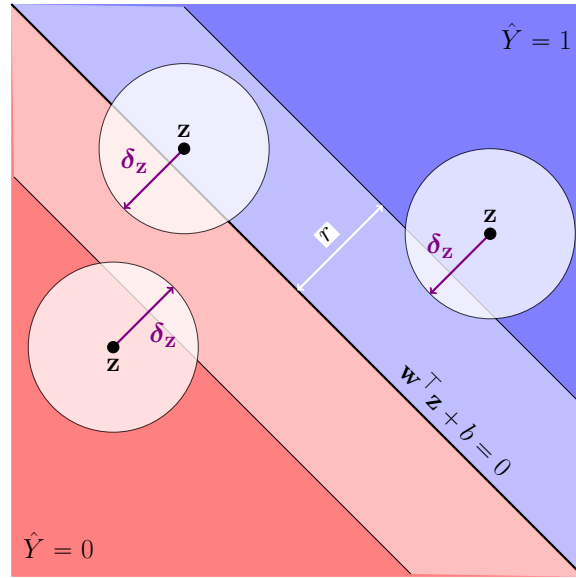


Figure 1: Robust optimization, visualized in the case of a linear classifier (black line) in two dimensions $\mathbf{w}^\top \mathbf{z} + b = 0$. The black data points $\mathbf{z} \equiv f_\theta(\mathbf{x})$ are the network outputs for various data inputs. The white circles indicate output features within a radius of r of the network outputs. The direction of perturbation $\delta_{\mathbf{z}}$ that maximizes the inner optimization problem is normal to the linear classifier defined by \mathbf{w} . Any network outputs in the white channel, r away from the linear classifier, change the predicted class. Robust optimization encourages network outputs to live outside of the white channel to avoid ambiguous class predictions.

76 not generalize to unseen data, or a lack of robustness, where small changes to the data result
77 in significantly different results (e.g., incorrect classifications).

78 **2.2. Robust Optimization.** Adversarial training promotes robustness in DNNs by intro-
79 ducing a perturbation $\delta_{\mathbf{x}}$ for each input \mathbf{x} and solving the minimax problem

$$80 \quad (2.2a) \quad \min_{\theta} \frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{T}} L(f_\theta(\mathbf{x} + \delta_{\mathbf{x}}(\theta)), \mathbf{y}) + R(\theta)$$

$$81 \quad (2.2b) \quad \text{s. t. } \delta_{\mathbf{x}}(\theta) \in \operatorname{argmax}_{\|\delta_{\mathbf{x}}\|_2 \leq r} L(f_\theta(\mathbf{x} + \delta_{\mathbf{x}}), \mathbf{y}) \quad \text{for each } (\mathbf{x}, \mathbf{y}) \in \mathcal{T}$$

82 We perturb the inputs \mathbf{x} by $\delta_{\mathbf{x}}$ and maximize the Euclidean norm of the perturbation $\|\delta_{\mathbf{x}}\|_2$
83 (inner optimization problem (2.2b)) while optimally fitting the data (outer minimization prob-
84 lem (2.2a)). We build neighborhoods of radius r around our training points where we can
85 rely on our model classifying anything within the neighborhood similarly. See Figure 1 for a
86 visualization.

87 The complexity of our new minimax problem is a large consideration for the applicability
88 of our results to large scale real-world situations. To address this, we use second order in-
89 formation to solve the inner optimization problem efficiently in terms of computational time.

90 Solving this problem well means satisfying first order optimality conditions. Following [2], we
 91 first negate the loss to produce an equivalent minimization problem and set up a Lagrangian.

$$92 \quad (2.3) \quad \mathcal{L}(\boldsymbol{\delta}_x, \lambda) = -L(f_{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\delta}_x), \mathbf{y}) + \lambda(\frac{1}{2}\|\boldsymbol{\delta}_x\|_2^2 - \frac{1}{2}r^2)$$

93 Here λ is a Lagrangian multiplier and we use an equivalent version of the constraint $\frac{1}{2}\|\boldsymbol{\delta}_x\|_2^2 \leq$
 94 $\frac{1}{2}r^2$ that we can differentiate more easily. The perturbation $\boldsymbol{\delta}_x$ that maximizes the inner
 95 optimization problem of (2.2) necessarily satisfies the Karush-Kuhn-Tucker (KKT) conditions
 96 below [9].

$$97 \quad (2.4a) \quad \nabla_{\boldsymbol{\delta}_x} \mathcal{L}(\boldsymbol{\delta}_x, \lambda) = -\nabla_{\boldsymbol{\delta}_x} L(f_{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\delta}_x), \mathbf{y}) + \lambda \boldsymbol{\delta}_x = \mathbf{0} \quad (\text{stationarity})$$

$$98 \quad (2.4b) \quad \|\boldsymbol{\delta}_x\|_2 \leq r \quad (\text{primal feasibility})$$

$$99 \quad (2.4c) \quad \lambda \geq 0 \quad (\text{dual feasibility})$$

$$100 \quad (2.4d) \quad \lambda(\|\boldsymbol{\delta}_x\|_2 - r) = 0 \quad (\text{complementary slackness})$$

101 Satisfying the KKT conditions ensures that gradients of the outer optimization problem are
 102 accurate; in particular, for each training sample, we have

$$103 \quad (2.5) \quad \nabla_{\boldsymbol{\theta}} L(f_{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\delta}_x(\boldsymbol{\theta})), \mathbf{y}) = [\nabla_{\boldsymbol{\theta}'} f_{\boldsymbol{\theta}'}(\mathbf{x} + \boldsymbol{\delta}_x(\boldsymbol{\theta})) \nabla_{\boldsymbol{\delta}_x} L(f_{\boldsymbol{\theta}'}(\mathbf{x} + \boldsymbol{\delta}_x(\boldsymbol{\theta})), \mathbf{y})]_{\boldsymbol{\theta}'=\boldsymbol{\theta}} \\ 104 \quad \quad \quad + [\nabla_{\boldsymbol{\theta}'} \boldsymbol{\delta}_x(\boldsymbol{\theta}') \nabla_{\boldsymbol{\delta}_x} L(f_{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\delta}_x(\boldsymbol{\theta}')), \mathbf{y})]_{\boldsymbol{\theta}'=\boldsymbol{\theta}}$$

105 The first term in (2.5) is the traditional gradient that we want to preserve. The second
 106 term comes from considering the perturbation as a function of the network weights, $\boldsymbol{\delta}_x(\boldsymbol{\theta})$.
 107 From the stationarity condition (2.4a), we get that $\nabla_{\boldsymbol{\delta}_x} L(f_{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\delta}_x), \mathbf{y})$ is parallel to $\boldsymbol{\delta}_x$ if
 108 the perturbation is a maximizer. If the constraint is inactive ($\lambda = 0$), then $\nabla_{\boldsymbol{\delta}_x} L(f_{\boldsymbol{\theta}}(\mathbf{x} +$
 109 $\boldsymbol{\delta}_x), \mathbf{y}) = \mathbf{0}$ and the second term is zero. If the constraint is active ($\lambda > 0$), then from
 110 primal feasibility (2.4b) we know that the perturbation must satisfy the constraint even when
 111 undergoing changes incurred from $[\nabla_{\boldsymbol{\theta}'} \boldsymbol{\delta}_x(\boldsymbol{\theta}')]_{\boldsymbol{\theta}'=\boldsymbol{\theta}}$. With a sufficiently small perturbation of
 112 the weights $\boldsymbol{\theta}$, the change in perturbation will follow the boundary of the constraint, nearly
 113 orthogonal to the direction of the gradient $\nabla_{\boldsymbol{\delta}_x} L(f_{\boldsymbol{\theta}}(\mathbf{x} + \boldsymbol{\delta}_x), \mathbf{y})$. This again makes the second
 114 term zero. Thus, if we solve the inner optimization problem well and thereby satisfy the KKT
 115 conditions, we can ignore the contribution of the second term.

116 **2.3. Fairness.** We use three different fairness metrics defined in [1] in our experiments. All
 117 of these fairness metrics pertain to fairness of a classifier with respect to a sensitive attribute,
 118 in terms of true labels against the classifier's predictions. In all of our experiments, the
 119 sensitive attribute s , true label Y , and classifier prediction \hat{Y} are all binary. For convenience,
 120 in defining the fairness metrics, we treat Y as a random variable representing an object's true
 121 label and \hat{Y} as a random variable representing its prediction.

122 **2.3.1. Independence.** For a classifier to satisfy independence its prediction \hat{Y} must be un-
 123 correlated with the sensitive attribute s . This requires an equal rate of positive classifications
 124 across all sensitive groups.

$$125 \quad (2.6) \quad P(\hat{Y} = 1 | s = 0) = P(\hat{Y} = 1 | s = 1) = P(\hat{Y} = 1)$$

126 For instance, if the classifier was being used to recommend hiring decisions (so $\hat{Y} = 1$ means a
 127 candidate should be hired, and $\hat{Y} = 0$ means a candidate should not), satisfying independence
 128 would mean that if the classifier hires 20% of applicants in class $s = 1$, then it also hires 20%
 129 of applicants in class $s = 0$.

130 **2.3.2. Separation.** Separation is similar to independence; for separation to be satisfied \hat{Y}
 131 must be conditionally independent of s given the value of Y .

$$132 \quad (2.7) \quad P(\hat{Y} = 1|Y = 1, s = 0) = P(\hat{Y} = 1|Y = 1, s = 1)$$

$$133 \quad P(\hat{Y} = 1|Y = 0, s = 0) = P(\hat{Y} = 1|Y = 0, s = 1)$$

134 Separation enforces equality of true and false positive rates. If again \hat{Y} determines hiring rec-
 135 ommendations, then Y might indicate an individual's true qualifications. Separation requires
 136 that individuals with similar qualifications have an equal chance of being hired, regardless of
 137 sensitive attribute.

138 **2.3.3. Sufficiency.** Sufficiency enforces the conditional independence of Y and s given \hat{Y} .

$$139 \quad (2.8) \quad P(Y = 1|\hat{Y} = 1, s = 0) = P(Y = 1|\hat{Y} = 1, s = 1)$$

$$140 \quad P(Y = 1|\hat{Y} = 0, s = 0) = P(Y = 1|\hat{Y} = 0, s = 1)$$

141 Sufficiency requires that the rates of individuals with the same predicted label also having the
 142 same true label is equal across different sensitive groups. If sufficiency is satisfied, then an
 143 individual from one group who is hired by the classifier is as likely to be truly qualified as a
 144 hired individual from another group.

145 **3. Our Approach.** Next we introduce our proposed second order method, and discuss its
 146 implementation. Our approach relies on approximation, so an analysis of the error produced by
 147 this approximation follows in [Subsection 3.2](#). Then in [Subsection 3.3](#), we introduce alternate
 148 methods of solving the robust optimization problem and their implementations to test against
 149 our proposed approach.

150 **3.1. Trust Region Subproblem (TRS).** Our main algorithm ([Algorithm 3.1](#)) solves an
 151 approximation of the inner optimization problem (2.2b) using second order information. For
 152 each training sample $(\mathbf{x}, \mathbf{y}) \in \mathcal{T}$, we fix θ and expand the loss function using a quadratic
 153 Taylor series approximation about \mathbf{x} in the direction of $\delta_{\mathbf{x}}$.

$$154 \quad (3.1) \quad \min_{\|\delta_{\mathbf{x}}\|_2 \leq r} -L(f_{\theta}(\mathbf{x}), \mathbf{y}) - (\nabla_{\mathbf{x}}L(f_{\theta}(\mathbf{x}), \mathbf{y}))^T \delta_{\mathbf{x}} - \frac{1}{2} \delta_{\mathbf{x}}^T \nabla_{\mathbf{x}}^2 L(f_{\theta}(\mathbf{x}), \mathbf{y}) \delta_{\mathbf{x}}$$

155 To fit our constraint, we construct a Lagrangian term by squaring our initial constraint and
 156 scaling the Lagrange multiplier by one-half. This gives us a function that depends on $\delta_{\mathbf{x}}$ and
 157 λ .

$$158 \quad (3.2) \quad \tilde{\mathcal{L}}(\delta_{\mathbf{x}}, \lambda) = -L(f_{\theta}(\mathbf{x}), \mathbf{y}) - (\nabla_{\mathbf{x}}L(f_{\theta}(\mathbf{x}), \mathbf{y}))^T \delta_{\mathbf{x}} - \frac{1}{2} \delta_{\mathbf{x}}^T \nabla_{\mathbf{x}}^2 L(f_{\theta}(\mathbf{x}), \mathbf{y}) \delta_{\mathbf{x}} + \frac{\lambda}{2} (\|\delta_{\mathbf{x}}\|^2 - r^2)$$

159 We approximate the optimal $\delta_{\mathbf{x}}^*$ to the inner optimization problem as the optimal $\delta_{\mathbf{x}}$ solution
 160 to the quadratic problem (3.1). The KKT conditions are the same as (2.4) except for the

Algorithm 3.1 Trust Region Subproblem

Require: network $f_\theta : \mathcal{X} \times \Theta \rightarrow \mathcal{Y}$, loss function $L : \mathbb{R}^{n_{\text{out}}} \times \mathcal{Y} \rightarrow \mathbb{R}$, batch $\mathcal{T}_i \subset \mathcal{T}$, trust region radius r

Ensure: Candidate perturbation per sample $\mathbf{S} \in \mathbb{R}^{n_{\text{in}} \times |\mathcal{T}_i|}$

- 1: Initialize \mathbf{S} as empty array
- 2: **for** each sample $(\mathbf{x}, \mathbf{y}) \in \mathcal{T}_i$ **do**
- 3: Evaluate loss and derivatives, $L(f_\theta(\mathbf{x}), \mathbf{y})$, $\nabla_{\mathbf{x}} L(f_\theta(\mathbf{x}), \mathbf{y})$, and $\nabla_{\mathbf{x}}^2 L(f_\theta(\mathbf{x}), \mathbf{y})$
- 4: Define perturbation as a function of Lagrange multiplier $\delta_{\mathbf{x}}(\lambda)$ ▷ Equation (3.4)
- 5: Compute unconstrained perturbation $\mathbf{s}_{\mathbf{x}} = \delta_{\mathbf{x}}(0)$
- 6: **if** $\|\mathbf{s}_{\mathbf{x}}\|_2 > r$ **then**
- 7: Set $\lambda_{\text{low}} = 0$, compute λ_{high} ▷ Subsection 3.1.1
- 8: Solve for λ^* using bisection method on the function $g(\lambda)$ ▷ Equation (3.5)
- 9: Choose optimal search direction $\mathbf{s}_{\mathbf{x}} = \delta_{\mathbf{x}}(\lambda^*)$
- 10: **end if**
- 11: Concatenate \mathbf{S} and $\mathbf{s}_{\mathbf{x}}$
- 12: **end for**

161 stationarity condition.

$$162 \quad (3.3) \quad -\nabla_{\mathbf{x}} L(f_\theta(\mathbf{x}), \mathbf{y}) - \nabla_{\mathbf{x}}^2 L(f_\theta(\mathbf{x}), \mathbf{y}) \delta_{\mathbf{x}} + \lambda \delta_{\mathbf{x}} = \mathbf{0} \quad (\text{stationarity})$$

163 This gives us an explicit relation of $\delta_{\mathbf{x}}$ to λ .

$$164 \quad (3.4) \quad \delta_{\mathbf{x}}(\lambda) = -(\nabla_{\mathbf{x}}^2 L(f_\theta(\mathbf{x}), \mathbf{y}) - \lambda I)^{-1} \nabla_{\mathbf{x}} L(f_\theta(\mathbf{x}), \mathbf{y})$$

165 There are two cases to (3.4). If $\lambda = 0$, then the optimal $\delta_{\mathbf{x}}$ for (3.1) can be found by solving a
 166 system of linear equations involving the gradient and the Hessian. Alternatively, if $\lambda \neq 0$, then
 167 complementary slackness enforces $\|\delta_{\mathbf{x}}\|_2 = r$, so we need to find λ such that $\|\delta_{\mathbf{x}}(\lambda)\|_2 = r$.

168 We note that Algorithm 3.1 and our derivation uses a “per-sample” option. This means
 169 the trust region constraint is applied independently for each sample in the input dataset; i.e.,
 170 for each data point, we solve a separate trust region optimization problem and perturb. This
 171 is beneficial when different data points require different level of adjustment. An alternative
 172 approach would be to use a “global” option; i.e., a single constraint is applied to the entire
 173 batch of data. The “per-sample” approach is beneficial because we can optimize the adversarial
 174 perturbation for each samples (i.e., increase robustness) and offers the potential to use a
 175 different trust region radius per sample (we have not programmed this adaptability into our
 176 code yet). The trade off is that the “per-sample” method is run sequentially over the batch
 177 samples, which can be slow. There is a potential for parallelization; however, this is non-trivial
 178 to code, particularly when ensuring the gradients track properly for automatic differentiation.
 179 We consider as a future improvement of the repository. During our experiments, we had a
 180 choice to use either but we mainly used “per-sample” which is why it was included.

181 **3.1.1. The Bisection Method Bracket.** To find a value for λ such that $\|\delta_{\mathbf{x}}(\lambda)\|_2 = r$, we
 182 build a univariate function $g(\lambda) := \|\delta_{\mathbf{x}}(\lambda)\|_2 - r$ and find a root of this function. Applying

183 some linear algebra to (3.4), one can show that

$$\begin{aligned}
184 \quad (3.5a) \quad & g(\lambda) = \|\delta_{\mathbf{x}}(\lambda)\|_2 - r \\
185 \quad (3.5b) \quad & = \left\| -(\nabla_{\mathbf{x}}^2 L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) - \lambda I)^{-1} \nabla_{\mathbf{x}} L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \right\|_2 - r \\
186 \quad (3.5c) \quad & = \left\| (QDQ^{\top} - \lambda I)^{-1} \nabla_{\mathbf{x}} L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \right\|_2 - r \\
187 \quad (3.5d) \quad & = \left\| Q(D - \lambda I)^{-1} Q^{\top} \nabla_{\mathbf{x}} L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \right\|_2 - r \\
188 \quad (3.5e) \quad & = \left\| (D - \lambda I)^{-1} Q^{\top} \nabla_{\mathbf{x}} L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \right\|_2 - r
\end{aligned}$$

189 where $\nabla_{\mathbf{x}}^2 L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) = QDQ^{\top}$ is the eigendecomposition of the Hessian. Because the Hessian
190 is symmetric, by the Spectral Theorem, we know Q is orthogonal and D is diagonal and real-
191 valued.

192 We use the bisection method [5] to find a root of $g(\lambda) = 0$. We do recognize that it may not
193 seem appropriate to use a bisection method here because g is not continuous at eigenvalues of
194 the Hessian of L . However, given the complexity of the function g and of robust optimization
195 in general, we wanted to use a straight-forward approach to find a root of g . In practice, we did
196 not seem to run into any numerical issues in our code using the bisection method. Moreover,
197 we obtained the same final fairness and accuracy results using a first order projected gradient
198 descent (PGD) method, outlined later, as using a bisection method with our new second-
199 order optimization approach. Considering these things, we considered the bisection method
200 to be adequate for the purposes of this work, and we leave it for future work to improve the
201 root-finding strategy of our optimization method.

202 In order to use the bisection method, we need to establish a bracket $[\lambda_{\text{low}}, \lambda_{\text{high}}]$ such that
203 g has different signs at the endpoints; that is, $g(\lambda_{\text{low}})g(\lambda_{\text{high}}) < 0$. If $\lambda = 0$, then constraint is
204 satisfied. Thus, $g(0) \geq 0$ and, in practice, positive, so $\lambda_{\text{low}} = 0$ is a good candidate. To find
205 the upper bound, we first bound the norm

$$\begin{aligned}
206 \quad (3.6) \quad & \left\| (D - \lambda I)^{-1} Q^{\top} \nabla_{\mathbf{x}} L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \right\|_2 \leq \left\| (D - \lambda I)^{-1} \right\|_2 \left\| \nabla_{\mathbf{x}} L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \right\|_2 \\
207 \quad (3.7) \quad & = \sqrt{\sum_{i=1}^{n_{\text{in}}} \frac{1}{(d_i - \lambda)^2}} \left\| \nabla_{\mathbf{x}} L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \right\|_2
\end{aligned}$$

208 Following from [9], as $\lambda \rightarrow d_{\text{max}}^+$, the upper bound approaches $+\infty$ and as $\lambda \rightarrow \infty$, the upper
209 bound approaches 0. This guarantees that there is some $\lambda \in (d_{\text{max}}, \infty)$ such that the upper
210 bound is less than r . If we let

$$211 \quad (3.8) \quad \lambda_{\text{high}} = |d_{\text{max}}| + \frac{\sqrt{n} \left\| \nabla_{\mathbf{x}} L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \right\|_2}{r}$$

212 then $(d_i - \lambda_{\text{high}})^2 \geq \frac{n \left\| \nabla_{\mathbf{x}} L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \right\|_2^2}{r^2}$ for $i = 1, \dots, n_{\text{in}}$. Substituting into the upper bound,
213 we get

$$214 \quad (3.9) \quad \sqrt{\sum_{i=1}^{n_{\text{in}}} \frac{1}{(d_i - \lambda)^2}} \left\| \nabla_{\mathbf{x}} L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \right\|_2 \leq \frac{\sqrt{nr}}{\sqrt{n} \left\| \nabla_{\mathbf{x}} L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \right\|_2} \left\| \nabla_{\mathbf{x}} L(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) \right\|_2 = r.$$

215 Thus, we have found a bracket for the bisection method.

216 **3.2. Algorithm Analysis.** When solving the inner optimization problem using a second
 217 order approximation, we would like to know how well this approximation actually solves
 218 this problem. For specific classes of models, loss functions, and activation functions, we can
 219 confine the error explicitly to depend on high orders of the perturbation $\delta_{\mathbf{x}}$ and loss function
 220 derivatives.

221 **3.2.1. Affine model.** An affine model $f_{\mathbf{w},b}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ with weight vector \mathbf{w} and a
 222 scalar bias b combined with a logistic regression loss function L and a sigmoid activation
 223 function σ is convex with respect to inputs. The loss is given explicitly as

$$224 \quad (3.10) \quad L(f_{\mathbf{w},b}(\mathbf{x}), y) = -y \ln[\sigma(\mathbf{w}^\top \mathbf{x} + b)] - (1 - y) \ln[1 - \sigma(\mathbf{w}^\top \mathbf{x} + b)]$$

225 with $\sigma(z) = (1 + e^{-z})^{-1}$. Introducing the perturbation $\delta_{\mathbf{x}}$ results in a specific case of the inner
 226 optimization problem in (2.2b). To solve this optimization problem without approximation,
 227 we introduce as before a Lagrange multiplier λ with the constraint $\|\delta_{\mathbf{x}}\|^2 - r^2 \leq 0$.

$$228 \quad (3.11) \quad \mathcal{L}_{\text{aff}}(\delta_{\mathbf{x}}, \lambda) = -y \ln(\sigma(\mathbf{w}^\top (\mathbf{x} + \delta_{\mathbf{x}}) + b)) - (1 - y) \ln(1 - \sigma(\mathbf{w}^\top (\mathbf{x} + \delta_{\mathbf{x}}) + b)) \\ + \frac{1}{2} \lambda (\|\delta_{\mathbf{x}}\|_2^2 - r^2)$$

229 The first order optimality conditions for (3.11) tell us that at the optimal $\delta_{\mathbf{x}}$,

$$230 \quad (3.12) \quad \nabla_{\delta_{\mathbf{x}}} \mathcal{L}_{\text{aff}}(\delta_{\mathbf{x}}, \lambda) = (-y + \sigma(\mathbf{w}^\top (\mathbf{x} + \delta_{\mathbf{x}}) + b)) \mathbf{w} + \lambda \delta_{\mathbf{x}} = \mathbf{0}.$$

231 To compare the exact solution from equation (3.12) to the approximation made when solving
 232 using the trust region method of section 3.1, we find the second order approximation of the
 233 loss function $L(f_{\mathbf{w},b}(\mathbf{x}), y)$ by a Taylor expansion in \mathbf{x} in the direction of $\delta_{\mathbf{x}}$.

$$234 \quad (3.13) \quad L(f_{\mathbf{w},b}(\mathbf{x} + \delta_{\mathbf{x}}), y) \approx L(f_{\mathbf{w},b}(\mathbf{x}), y) + \nabla_{\mathbf{x}}^\top L(f_{\mathbf{w},b}(\mathbf{x}), y) \delta_{\mathbf{x}} + \frac{1}{2} \delta_{\mathbf{x}}^\top \nabla_{\mathbf{x}}^2 L(f_{\mathbf{w},b}(\mathbf{x}), y) \delta_{\mathbf{x}}$$

235 where the gradient and Hessian are

$$236 \quad (3.14) \quad \nabla_{\mathbf{x}} L(f_{\mathbf{w},b}(\mathbf{x}), y) = (-y + \sigma(\mathbf{w}^\top \mathbf{x} + b)) \mathbf{w} \\ \nabla_{\mathbf{x}}^2 L(f_{\mathbf{w},b}(\mathbf{x}), y) = \mathbf{w} \sigma'(\mathbf{w}^\top \mathbf{x} + b) \mathbf{w}^\top.$$

237 The associated Lagrangian is

$$238 \quad (3.15) \quad \tilde{\mathcal{L}}_{\text{aff}}(\delta_{\mathbf{x}}, \lambda) = L(f_{\mathbf{w},b}(\mathbf{x}), y) + \nabla_{\mathbf{x}}^\top L(f_{\mathbf{w},b}(\mathbf{x}), y) \delta_{\mathbf{x}} + \frac{1}{2} \delta_{\mathbf{x}}^\top \nabla_{\mathbf{x}}^2 L(f_{\mathbf{w},b}(\mathbf{x}), y) \delta_{\mathbf{x}} + \frac{1}{2} \lambda (\|\delta_{\mathbf{x}}\|^2 - r^2).$$

239 As before, take the gradient with respect to $\delta_{\mathbf{x}}$ and set it equal to zero to solve using
 240 first-order optimization conditions.

$$241 \quad (3.16) \quad \nabla_{\delta_{\mathbf{x}}} \tilde{\mathcal{L}}(\delta_{\mathbf{x}}, \lambda) = (-y + \sigma(\mathbf{w}^\top \mathbf{x} + b)) \mathbf{w} + \mathbf{w} \sigma'(\mathbf{w}^\top \mathbf{x} + b) \mathbf{w}^\top \delta_{\mathbf{x}} + \lambda \delta_{\mathbf{x}} = \mathbf{0}$$

242 Comparing (3.12) (LHS) and (3.16) (RHS), the discrepancy between the exact solution and
 243 the approximation is restricted to

$$244 \quad (3.17) \quad \sigma(\mathbf{w}^\top (\mathbf{x} + \delta_{\mathbf{x}}) + b) \neq \sigma(\mathbf{w}^\top \mathbf{x} + b) + \sigma'(\mathbf{w}^\top \mathbf{x} + b) \mathbf{w}^\top \delta_{\mathbf{x}}.$$

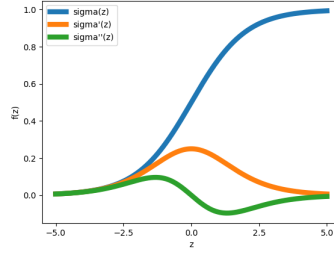


Figure 2: Flattening behavior of sigmoidal function, σ , derivatives.

245 Now, applying a Taylor expansion centered at \mathbf{x} in the ball $\mathbf{x} + \boldsymbol{\delta}_x$ to the [LHS](#), we obtain:

$$246 \quad (3.18) \quad \sigma(\mathbf{w}^\top(\mathbf{x} + \boldsymbol{\delta}_x) + b) \approx \sigma(\mathbf{w}^\top \mathbf{x} + b) + \sigma'(\mathbf{w}^\top \mathbf{x} + b) \mathbf{w}^\top \boldsymbol{\delta}_x + \frac{1}{2} \boldsymbol{\delta}_x^\top \mathbf{w} \sigma''(\mathbf{w}^\top \mathbf{x} + b) \mathbf{w}^\top \boldsymbol{\delta}_x$$

247 We have recovered the [RHS](#) of (3.17), so the error is due to the truncation of the second-order
 248 and higher terms of (3.18). In this case with a sigmoidal loss function, that means this error
 249 depends on the magnitude of $|\sigma''(z)|$ and the $\boldsymbol{\delta}_x$ for which we solved.

250 For any sigmoidal function, their structure gives first and second order derivatives of the
 251 classes shown in [Figure 2](#). For our sigmoid function defined as $\sigma(z) = (1 + e^{-z})^{-1}$ with
 252 $|\sigma''(z)| \leq 0.1$, and in general for any choice of sigmoidal function this flattening of higher-
 253 order derivatives will be observed. By nature $\boldsymbol{\delta}_x$ is bounded by the data since it defines the
 254 perturbation from a given point, and a perturbation larger than the size of the data space
 255 in any given dimension would be meaningless. For data in the form of continuous numerical
 256 values normalized to be between 0 and 1, as in our case, each component of $\boldsymbol{\delta}_x$ must be
 257 less than 1. In practice, $\boldsymbol{\delta}_x$ tends to be much smaller than that. This bound means higher
 258 order terms are generally quite small, and for this combination of loss function, activation
 259 function, model, and radii on the order of 10^{-1} , the approximation error is on the order of
 260 $\|\boldsymbol{\delta}_x\|^2 |\sigma''(z)| \approx 10^{-3}$.

261 **3.3. Other Methods.** We compare our trust region subproblem (TRS) algorithm to ran-
 262 dom perturbation to examine whether or not it is important to solve the inner optimization
 263 problem well. We also use random perturbation, along with a projected gradient descent
 264 (PGD) method, to verify that our second order TRS approach has greater computational
 265 efficiency than only using lower order information.

266 **3.3.1. Random Perturbation.** For each data point, we sample the perturbation $\boldsymbol{\delta}_x$ ran-
 267 domly from a multivariate standard normal distribution and rescale to the length of the trust
 268 region radius. This method acts as a control in our experiments to show the advantages of
 269 solving for an optimal perturbation.

270 **3.3.2. Projected Gradient Descent (PGD).** In order to test that our second order TRS
 271 approach is computationally faster than using purely first order information, we also imple-
 272 mented a version of gradient descent for our inner optimization problem. Since our inner
 273 problem has the constraint $\|\boldsymbol{\delta}_x\|_2 \leq r$, we cannot use vanilla gradient descent, and instead

274 use projected gradient descent (PGD) [3, 10]. PGD operates similarly to standard gradient
 275 descent, but once it has found its optimal step it projects the step onto the constrained set
 276 before returning it. Mathematically, this looks like:

$$277 \quad (3.19) \quad \boldsymbol{\delta}_x^{(k+1)} = P \left[\boldsymbol{\delta}_x^{(k)} + \alpha^{(k)} \cdot \nabla_x L(f_\theta(\mathbf{x} + \boldsymbol{\delta}_x), \mathbf{y}) \right]$$

278 where $\boldsymbol{\delta}_x^{(k)}$ is the k th iterate, $\alpha^{(k)}$ is the step size at the k th iteration, and P is the projection
 279 operator $P(\mathbf{x}) = \operatorname{argmin}_{\mathbf{z}: \|\mathbf{z}\|_2 \leq r} \|\mathbf{x} - \mathbf{z}\|_2^2$. This projection operator turns out to be very
 280 simple, returning the inputted point if the point already satisfies the constraint, or scaling the
 281 point inward to the boundary of the constraint if it is outside it. In particular,

$$282 \quad (3.20) \quad P(x) = \min \left\{ 1, \frac{r}{\|\mathbf{x}\|_2} \right\} \mathbf{x}.$$

283 Numerical experiments on how PGD performs with adversarial training [13, 6] have shown
 284 that PGD is a reliable method when it comes to solving robust optimization problems.

285 **4. Numerical Results.** Now we present results of our numerical experiments pertaining to
 286 fairness, accuracy, and computational time. Subsections 4.1, 4.2, and 4.3 discuss fairness and
 287 accuracy results on three datasets. For each dataset, we compute the fairness and accuracy re-
 288 sults for nonrobust training, for robust training with various radii, and random perturbations.
 289 We measure the absolute difference across sensitive attributes for analysis of each fairness
 290 metric, and the closer this difference is to zero, the fairer the classifier is with respect to that
 291 metric. Subsection 4.4 presents our results on relative computational time across our various
 292 methods for solving the inner optimization problem.

293 **4.1. Synthetic Data (Unfair2D).** The primary dataset we used for carrying out numerical
 294 experiments was a synthetic dataset. Individuals belonging to two different groups, labeled
 295 with respect to a sensitive attribute A or B , are being hired on the basis of two numeric scores
 296 x_1 and x_2 . Individuals have a binary label that is either “should be hired” or “shouldn’t be
 297 hired,” and we train a linear classifier to decide whether or not to hire an individual. The data
 298 is initially fair (Figure 3a), and we introduce unfair bias into the data by artificially raising
 299 the scores of all B s while lowering the scores of all A s (Figure 3b). In the real world, this
 300 could be a manifestation of structural unfairness in which B s are more likely to belong to a
 301 wealthy socioeconomic class, and thus can afford training that boosts their scores, whereas
 302 A s do not have this opportunity. In fact, A s may not only lack the advantage of B s, but also
 303 have an active disadvantage, such as an increased likelihood of needing to work longer hours,
 304 impeding time for study and test prep, lowering their scores.

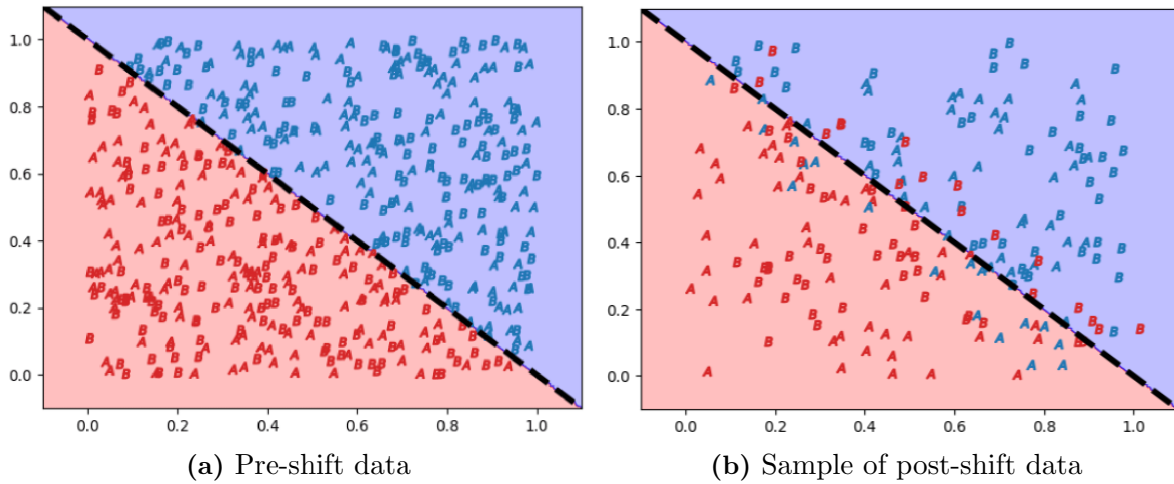


Figure 3: Points are colored based on their original location in the blue region ($Y = 1$) or red region ($Y = 0$). Post shift, note the unfair presence of red B s in the blue region and blue A s in the red region.

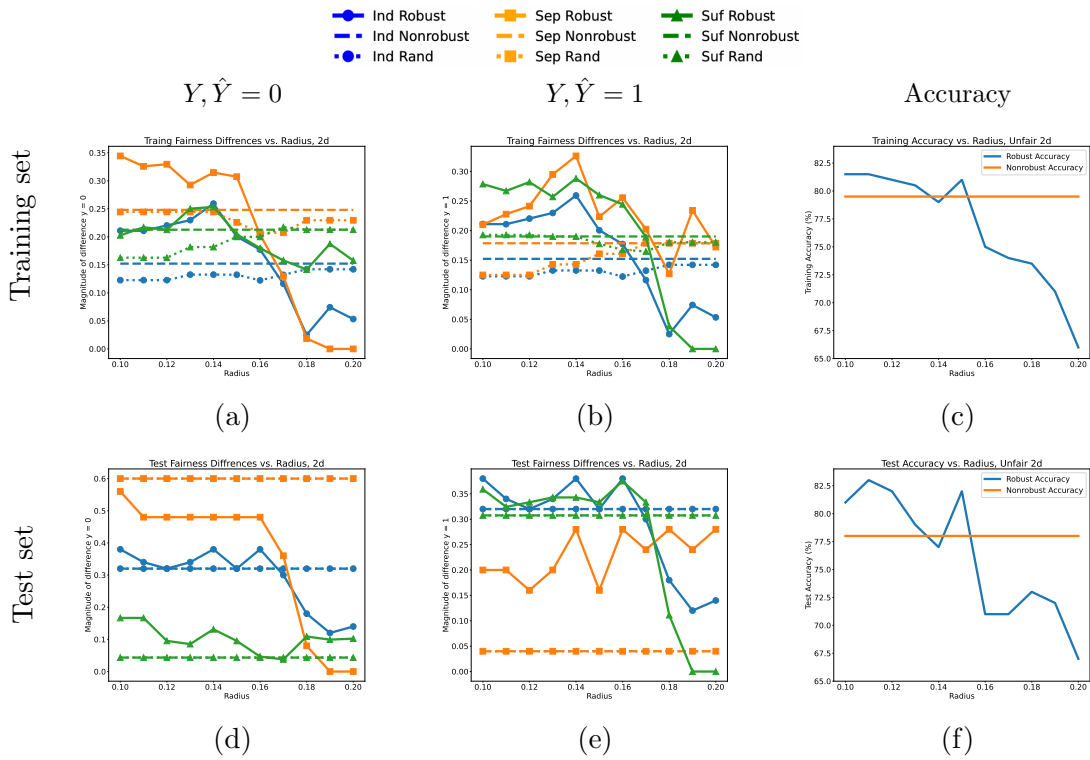
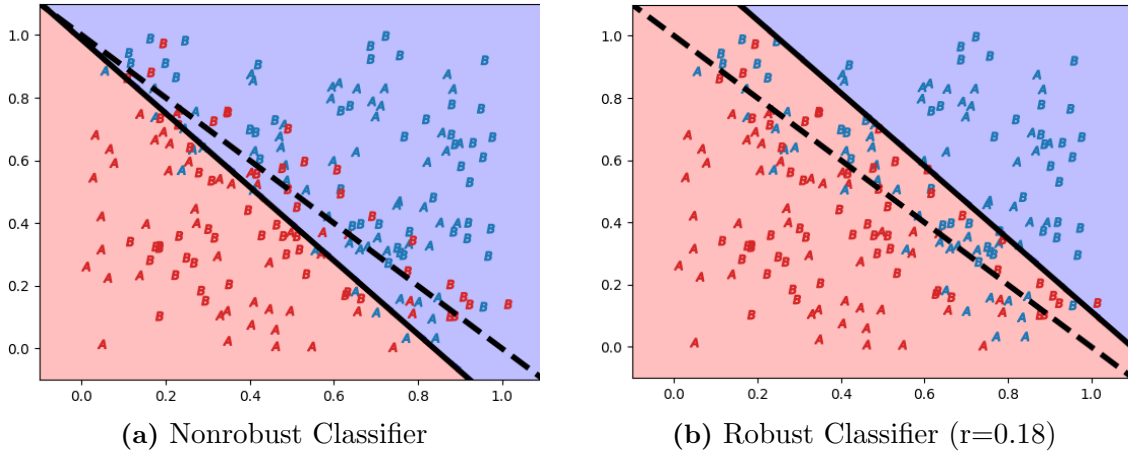


Figure 4: Synthetic fairness and accuracy results. For fairness, values closer to zero are desirable.



Diff:	$Y = 0, S1 - S0 $	$Y = 1, S1 - S0 $	Diff:	$Y = 0, S1 - S0 $	$Y = 1, S1 - S0 $
Ind.	0.152	0.152	Ind.	0.025	0.025
Sep.	0.248	0.179	Sep.	0.019	0.127
Suff.	0.213	0.190	Suff.	0.142	0.038

Training Accuracy: 79.5%
Test Accuracy: 78.0%

Training Accuracy: 73.5%
Test Accuracy: 73.0%

Figure 5: Comparative Analysis of Non-Robust and Robust Classifiers

305 We compare nonrobust training on this synthetic dataset with robust training (using the
306 TRS method) and random perturbation for 11 different perturbation radii, with the radius
307 increasing in 0.01 increments from 0.1 up to 0.2. These experiments were run with a total of
308 10 training epochs and a learning rate of 0.01 in the outer optimization problem. Four plots of
309 fairness metric differences versus radius are shown in Figure 4, as well as plots of the training
310 and testing accuracy versus radius. For many radii in the lower end of the plotted range,
311 many of the fairness differences are worse in the training data with robust training than with
312 nonrobust. However, some fairness improvement can be seen with robust training.

313 In all cases, at least two of the three fairness metrics show a downward trend for robust
314 training. This suggests that while robust training may worsen fairness for a very small radius,
315 fairness improvement is possible at more appropriate radii. At a radius of 0.18, all of the
316 robust fairness differences are better than the corresponding nonrobust ones in the training
317 data, although it does lead to a decrease in test accuracy from around 78% to 73% (Figure 5).
318 The nonrobust classifier is visualized in Figure 5a versus the robust classifier in Figure 5b.
319 In Figure 5b, robust optimization is improving fairness by raising the bar, giving a positive
320 classification to only the most qualified individuals. It eliminates nearly all of the false positive
321 *B*s that exist with the original dashed classifier and greatly increases the quantity of blue *B*s
322 in the red region, equalizing false negative rates. Increasing the radius even further to 0.2
323 with 20 epochs of training, the robust classifier eventually classifies nothing positively. Our
324 robust classifier is not helping the disadvantaged *A*s in the process of improving fairness – it

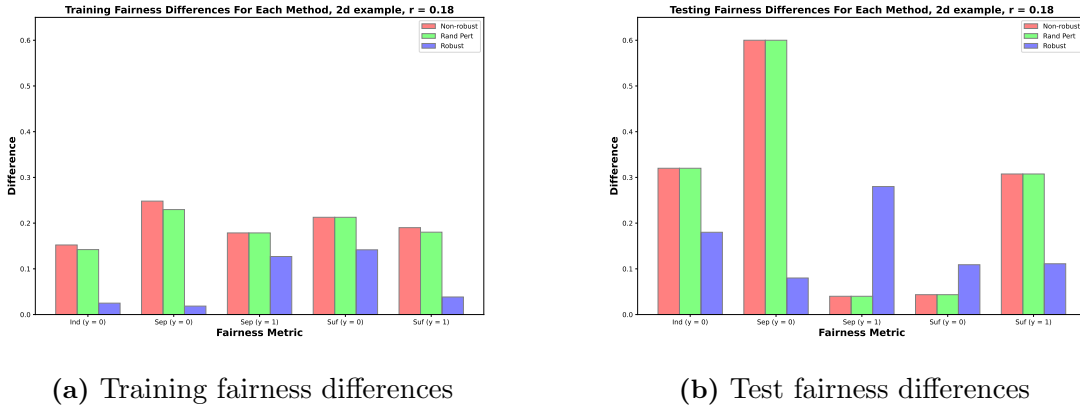


Figure 6: Fairness differences ($r=0.18$). Left bar is nonrobust, middle is random, right is robust.

325 is only hurting the advantaged B_s . While this does not provide an indication of how robust
 326 training would work on every dataset, it does illustrate that robust optimization may improve
 327 fairness in an unexpected or unintended way.

328 There is an advantage to solving the inner optimization problem well instead of just
 329 using random perturbations. In Figure 4, the fairness differences for random perturbations
 330 are either the same or stay close to the nonrobust differences. This stands in contrast to
 331 robust training, where fairness differences are initially high and then get significantly lower,
 332 surpassing nonrobust differences. Figure 6 also exhibits this advantage.

333 **4.2. Adult Dataset.** We also extended our numerical experiments to real-world datasets.
 334 The Adult dataset [11] consists of demographic data about individuals that are used to classify
 335 whether their annual income is more than \$50,000. Note that the dataset predominantly
 336 consists of white males in the age range of 25-60. It contains 48,842 instances and each
 337 instance is described using 15 attributes. We want to look at the 5 continuous numerical
 338 attributes (age, education-num, capital-gain, capital-loss, hours-per-week) for analysis. The
 339 income (salary) data is converted into binary form (1 for $\leq 50k$, 0 for $> 50K$), and the
 340 protected attribute can be sex or race.

341 Unlike our synthetic data, the adult example yielded mixed results in terms of fairness
 342 improvement. For the training data, Figure 7a and Figure 7b show that only three out of
 343 the six differences were measured to be better with robust training. There was a similar
 344 result for the test data as seen in Figure 7c and Figure 7d. Despite only having a 50%
 345 improvement rate, the majority of the fairness metrics exhibit a downward trend, and when
 346 there is an improvement robust optimization outperforms random perturbation significantly.
 347 The expected accuracy-robustness trade-off is present (Figure 7e), with both the training
 348 and test robust accuracy decreasing with increasing radii. However, unlike in the synthetic
 349 dataset, the decay appears to be linear and does not spike at certain radii, and does not yield

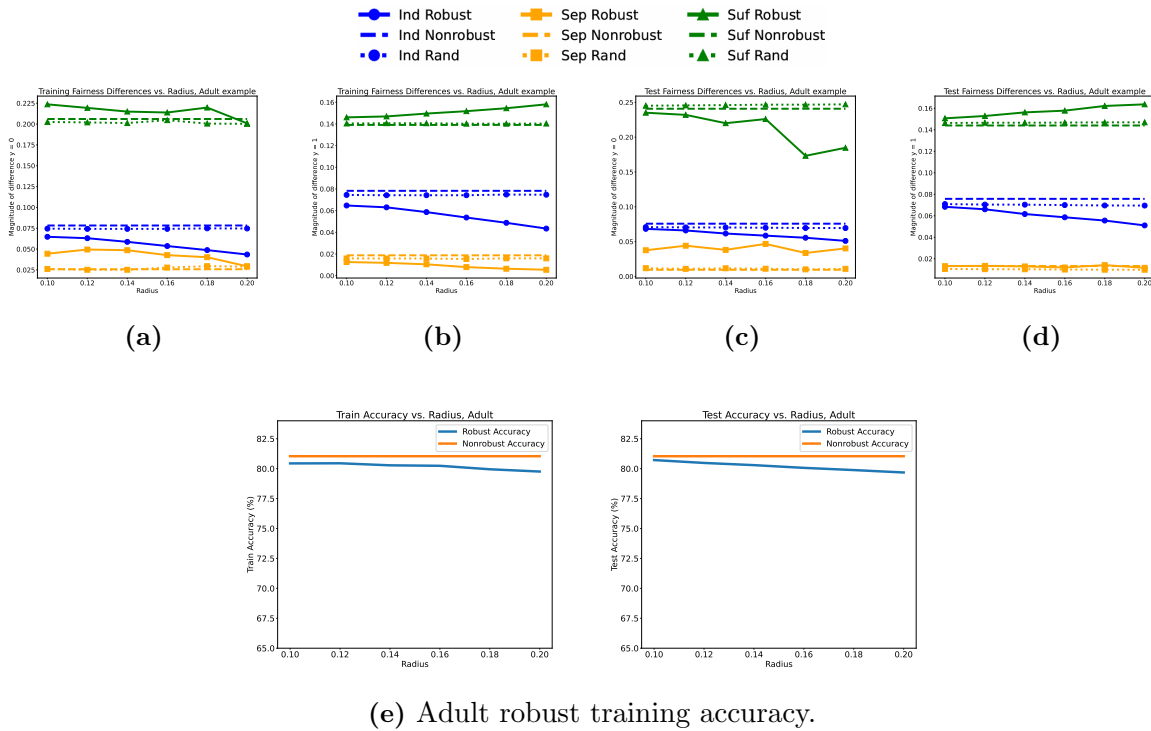


Figure 7: Fairness ((a)-(d)) and accuracy (e) trends in the adult dataset for nonrobust and robust training.

350 better accuracy for smaller radii. Overall, there is still reasonable case for improving fairness
 351 metrics at the expense of classifier accuracy.

352 **4.3. LSAT Data.** Another extension to a real-world dataset comes from the Law School
 353 Admissions Council (LSAC) [12]. This dataset was collected to explore the reasons behind
 354 low bar passage rates among racial and ethnic minorities. We train our classifier to predict
 355 whether or not a student will pass the bar, based on their Law School Admission Test (LSAT)
 356 score and undergraduate GPA. We are using GPA and LSAT scores because they are the
 357 strongest predictors for passing the bar examination. Our primary interest lies in examining
 358 five key features of the dataset: the bar exam pass/fail prediction made by a DNN, the gender
 359 of the student, their LSAT score, the true bar exam pass/fail value for the student, and their
 360 race. For the purpose of our experiment, the race feature is made binary to indicate a student
 361 as either white or non-white, which is used as the sensitive attribute.

362 Unlike the two previous examples, there is not a lot of fluctuation with LSAT robust
 363 results (Figure 8). Surprisingly, robust optimization seems not to deviate from nonrobust
 364 training. As we have seen, using a random perturbation yields fairness results that are very
 365 similar to the results of non-robust training. For this dataset, robust optimization also per-
 366 forms very similarly to just picking a random perturbation. This is especially the case with

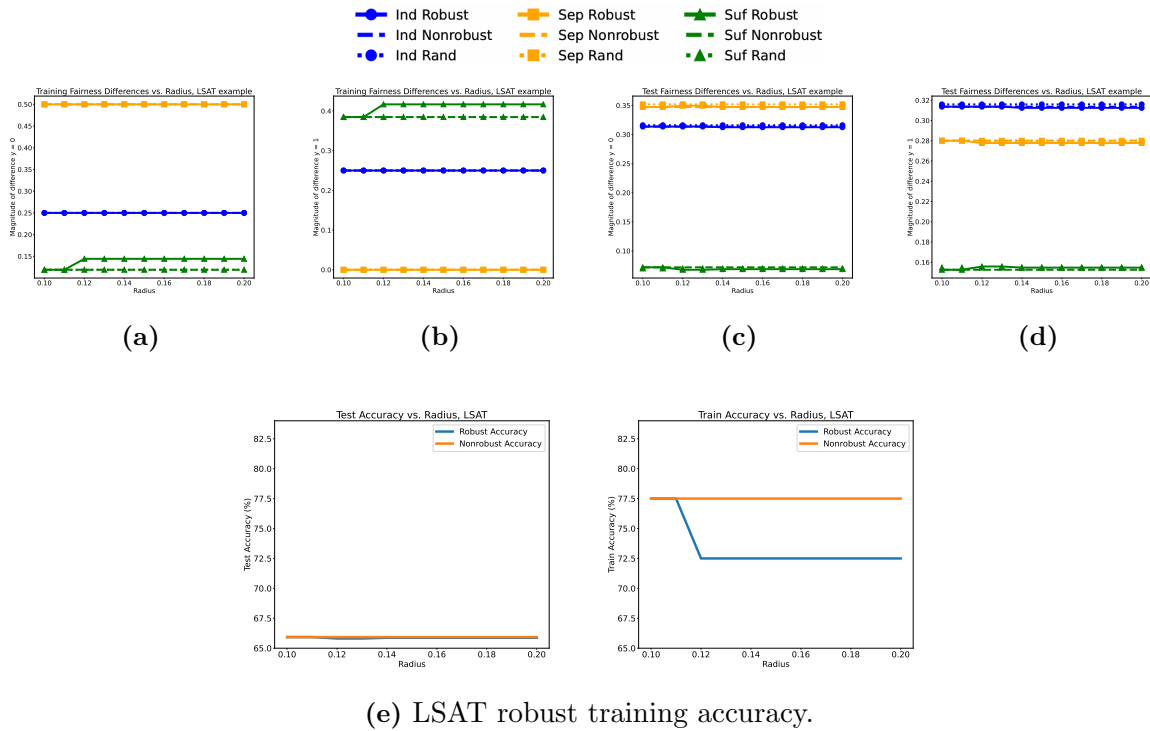


Figure 8: Fairness ((a)-(d)) and accuracy (e) trends in the LSAT dataset for nonrobust and robust training.

367 the independence and separation metrics. The only two notable deviations come from the
 368 sufficiency metric in [Figure 8a](#) and [Figure 8b](#) of the training dataset. This might be attribut-
 369 able to our definition of the sensitive attribute as white vs nonwhite, which creates a very small
 370 dataset due to the dominance of white individuals in the original dataset. For this specific
 371 example, robust optimization does not improve fairness, and in fact performs worse since it
 372 loses accuracy when yielding the same fairness results.

373 **4.4. Efficiency Comparison.** We gathered time data to see if the TRS method converged
 374 faster than PGD on our datasets. On each of the three datasets above - synthetic, Adult,
 375 and LSAT - we computed for each radius the average epoch time elapsed for TRS, PGD, and
 376 random perturbation. To compare the speed of the TRS method and PGD, we examine the
 377 ratio of the average PGD epoch time to the average TRS epoch time, looking at the extreme
 378 values of this ratio to get a range of how much faster the trust region subproblem was than
 379 PGD across all radii. The results are shown in [Table 1](#).

380 Random perturbation is the fastest adversarial training method in all three datasets. This
 381 is expected, as it does not actually solve the optimization problem; its only computation task
 382 is generating a random vector and rescaling it. It is also noteworthy that using the TRS
 383 method consistently is computationally faster than using PGD. Over all radii shown, training

Table 1: Average Epoch Times. For each dataset, the first three rows show the average epoch times for each of the three robust optimization approaches, where a lower value indicates faster computational performance. The fourth row shows the ratio of PGD to TRS time, where a ratio greater than 1 indicates that the TRS approach was faster than the PGD approach for the given radius. The gray values highlight the minimum ratio of PGD to TRS time over all radii, while the yellow values highlight the maximum ratio.

	Radii	.10	.12	.14	.16	.18	.20
Synthetic	PGD	2.680	3.597	4.228	4.486	5.061	5.080
	TRS	1.945	1.875	1.806	1.891	1.742	1.752
	RND	0.0387	0.0366	0.0387	0.0389	0.0387	0.0375
	PGD/TRS	1.377	1.919	2.340	2.373	2.904	2.900
Adult	PGD	512.970	593.173	697.399	1146.856	1689.761	1854.932
	TRS	60.372	61.557	57.723	58.707	57.492	59.061
	RND	0.0947	0.101	0.0972	0.0919	0.0974	0.0939
	PGD/TRS	8.497	9.636	12.082	19.535	29.391	31.407
LSAT	PGD	1.129	1.559	2.844	3.575	3.347	2.752
	TRS	0.396	0.396	0.421	0.371	0.414	0.385
	RND	0.0107	0.0123	0.0154	0.0122	0.0162	0.0104
	PGD/TRS	2.852	3.936	6.762	9.639	8.094	7.150

384 with TRS is between 1.4 and 2.9 times faster than PGD in the synthetic dataset, between
385 8.5 and 31.4 times faster in Adult, and between 2.9 and 9.6 times faster in LSAT. The very
386 short average epoch times for the LSAT dataset are due to the significantly smaller scale of
387 the input data. All of the smallest factors of time improvement of TRS relative to PGD
388 (highlighted in gray) are greater than 1 suggesting that the trust region subproblem has a
389 consistent advantage over PGD in computational speed.

390 Looking at the PGD/TRS ratios, the factor of improvement that TRS has in computa-
391 tional time over PGD appears to be higher in the real-world datasets than in the synthetic
392 dataset. The real-world datasets, and especially Adult, are trained on larger amounts of data,
393 so the advantage of TRS over PGD seems to scale with the size of the dataset. This advantage
394 of the trust region subproblem also improves with larger perturbation radii. In particular,
395 the minimum factor of improvement (gray) always occurs with the smallest radius, and the
396 maximum factor of improvement (yellow) always occurs with one of the three largest radii.

397 **5. Conclusion.** In our affine linear model setup, we were able to see improvement in
398 fairness by using robust optimization. In the synthetic dataset, whenever there was an im-
399 provement, the gain was a significant reduction in fairness difference magnitudes (which are
400 ideally zero). In our numerical experiments extending to real-world datasets, we have shown
401 that robust training performs similarly to non-robust training even in the worst-case scenario
402 (LSAT dataset). Across all three datasets, the accuracy of robust optimization decreased
403 as the radius increased, the majority of the fairness metrics displayed a downward trend as
404 the perturbation radius increased, and when fairness improved with robust training, precise

405 solutions to the inner optimization problem outperformed randomly selected solutions. Fur-
406 thermore, we were able to quantify the fact that, with the help of `hessQuik`, using second-order
407 information is much faster for solving our class of optimization problem.

408 We acknowledge that while we were able to achieve positive results with our experiments in
409 both synthetic and real-world datasets, there are a few mathematical limitations to our results
410 that prevent generalization to higher-dimensional applications. We used a neural network in
411 our training with only one hidden layer, our experiments were conducted using a linear and
412 binary classifier, and our sensitive attribute was binary. This motivates future exploration
413 of extending our approach to deeper neural networks, multinomial classification, and other
414 fairness metrics relevant to those cases. It may help to improve fairness even further to
415 introduce a regularization term to our approach to penalize violations of our fairness metrics,
416 which is another avenue for further work. There limitations of our implementation. For
417 PGD, we used an arbitrary step size instead of varying the step size as training proceeds.
418 Additionally, we did not solve our inner optimization problems in parallel. Parallelizing the
419 computations for our inner optimization problem could provide a significant reduction in
420 overall computation time.

421 Despite these limitations, this work demonstrates initial promise for the ability of robust
422 training to bring about fairness improvement in machine learning models, and motivates
423 further research on similar methodologies.

424 **Acknowledgments.** This work was supported by NSF award DMS-2051019 and was com-
425 pleted during the “Computational Mathematics for Data Science” REU/RET program in the
426 summer of 2023. We would like to thank Dr. Elizabeth Newman, our mentor, and the rest of
427 the faculty who participated in the program for their feedback and support.

428

REFERENCES

- 429 [1] S. BAROCAS, M. HARDT, AND A. NARAYANAN, *Fairness and Machine Learning: Limitations and Op-*
430 *portunities*, fairmlbook.org, 2019. <http://www.fairmlbook.org>.
- 431 [2] A. BECK, *Introduction to Nonlinear Optimization*, Society for Industrial and Applied Mathemat-
432 ics, Philadelphia, PA, 2014, <https://doi.org/10.1137/1.9781611973655>, <https://epubs.siam.org/doi/abs/10.1137/1.9781611973655>, <https://arxiv.org/abs/https://epubs.siam.org/doi/pdf/10.1137/1.9781611973655>.
- 433 [3] A. BECK, *First-order methods in optimization*, SIAM, 2017.
- 434 [4] N. FURL, P. PHILLIPS, AND A. J. O’TOOLE, *Face recognition algorithms and the other-race ef-*
435 *fect: computational mechanisms for a developmental contact hypothesis*, *Cognitive Science*, 26
436 (2002), pp. 797–815, [https://doi.org/https://doi.org/10.1016/S0364-0213\(02\)00084-8](https://doi.org/https://doi.org/10.1016/S0364-0213(02)00084-8), <https://www.sciencedirect.com/science/article/pii/S0364021302000848>.
- 437 [5] A. KAW, *Numerical methods with applications (kaw)*, University of South Florida, 2011, [https://math.libretexts.org/Under_Construction/Numerical_Methods_with_Applications_\(Kaw\)](https://math.libretexts.org/Under_Construction/Numerical_Methods_with_Applications_(Kaw)).
- 438 [6] A. MADRY, A. MAKELOV, L. SCHMIDT, D. TSIPRAS, AND A. VLADU, *Towards deep learning models*
439 *resistant to adversarial attacks*, in *International Conference on Learning Representations*, 2018, <https://openreview.net/forum?id=rJzIBfZAb>.
- 440 [7] N. MEHRABI, F. MORSTATTER, N. SAXENA, K. LERMAN, AND A. GALSTYAN, *A survey on bias and*
441 *fairness in machine learning*, *CoRR*, abs/1908.09635 (2019), <http://arxiv.org/abs/1908.09635>, <https://arxiv.org/abs/1908.09635>.
- 442 [8] E. NEWMAN AND L. RUTHOTTO, ‘hessquik’: *Fast hessian computation of composite functions*, *Journal*
443 *of Open Source Software*, 7 (2022), p. 4171, <https://doi.org/10.21105/joss.04171>, <https://doi.org/10.21105/joss.04171>.

- 450 21105/joss.04171.
- 451 [9] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer, New York, NY, USA, 2e ed., 2006.
- 452 [10] N. PARIKH, S. BOYD, ET AL., *Proximal algorithms*, Foundations and trends® in Optimization, 1 (2014),
- 453 pp. 127–239.
- 454 [11] T. L. QUY, A. ROY, V. IOSIFIDIS, W. ZHANG, AND E. NTOUTSI, *A survey on datasets for fairness-*
- 455 *aware machine learning*, WIREs Data Mining and Knowledge Discovery, 12 (2022), [https://doi.org/](https://doi.org/10.1002/widm.1452)
- 456 10.1002/widm.1452, <https://doi.org/10.1002/widm.1452>.
- 457 [12] L. F. WIGHTMAN, *Lsac national longitudinal bar passage study. lsac research report series.*, 1998, [https:](https://api.semanticscholar.org/CorpusID:151073942)
- 458 [/api.semanticscholar.org/CorpusID:151073942](https://api.semanticscholar.org/CorpusID:151073942).
- 459 [13] H. XU, X. LIU, Y. LI, A. K. JAIN, AND J. TANG, *To be robust or to be fair: Towards fairness in*
- 460 *adversarial training*, 2021, <https://arxiv.org/abs/2010.06121>.