
Understanding Why Generalized Reweighting Does Not Improve Over ERM

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Empirical risk minimization (ERM) is known to be non-robust in practice to
2 distributional shift where the training and the test distributions are different. A suite
3 of approaches, such as importance weighting, and variants of distributionally robust
4 optimization (DRO), have been proposed to solve this problem. But a line of recent
5 work has empirically shown that these approaches do not significantly improve
6 over ERM in real applications with distribution shift. The goal of this work is to
7 obtain a comprehensive theoretical understanding of this intriguing phenomenon.
8 We first posit the class of Generalized Reweighting (GRW) algorithms, as a broad
9 category of approaches that iteratively update model parameters based on iterative
10 reweighting of the training samples. We show that when overparameterized models
11 are trained under GRW, the resulting models are close to that obtained by ERM.
12 We also show that adding small regularization which does not greatly affect the
13 empirical training accuracy does not help. Together, our results show that a broad
14 category of what we term GRW approaches are not able to achieve distributionally
15 robust generalization. Our work thus has the following sobering takeaway: to
16 make progress towards distributionally robust generalization, we either have to
17 develop non-GRW approaches, or perhaps devise novel classification/regression
18 loss functions that are adapted to the class of GRW approaches.

19 1 Introduction

20 It has now been well established that empirical risk minimization (ERM) can empirically achieve high
21 test performance on a variety of tasks, particularly with modern overparameterized models where the
22 number of parameters is much larger than the number of training samples. This strong performance
23 of ERM however has been shown to degrade under *distributional shift*, where the training and test
24 distributions are different [HS15, BGO16, Tat17]. There are two broad categories of distribution
25 shift: *domain generalization* where the test distribution contains new *environments* not in the training
26 distribution like in domain adaptation, and *subpopulation shift* where the two distributions have the
27 same set of subpopulations but their mixture weights differ like in algorithmic fairness applications.

28 People have proposed various approaches to learn models that are robust to distributional shift. The
29 most classical approach is importance weighting (IW) [Shi00], which reweights training samples; in
30 the context of subpopulation shift these weights are typically set so that each subpopulation/group
31 has the same overall weight in the training objective. The approach most widely used today is
32 Distributional Robust Optimization (DRO) [DN18, HSNL18], in which we assume that the test
33 distribution belongs to a certain set of distributions that are close to the training distribution (called
34 the *uncertainty set*), and train the model on the worst distribution in that set. Many variants of DRO
35 have been proposed and are used in practice [HNSS18, SKHL20, XDKR20, ZDKR21, ZDS⁺21].

36 While these approaches have been developed for the express purpose of improving ERM for distri-
 37 bution shift, a line of recent work has empirically shown the negative result that when used to train
 38 overparameterized models, these methods do not improve over ERM. For IW, [BL19] observed that
 39 its effect under stochastic gradient descent (SGD) diminishes over training epochs, and finally does
 40 not improve over ERM. For variants of DRO, [SKHL20] found that these methods overfit very easily,
 41 i.e. their test performances will drop to the same low level as ERM after sufficiently many epochs if
 42 no regularization is applied. [GLP21, KSM⁺21] compared these methods with ERM on a number of
 43 real-world applications, and found that in most cases none of these methods improves over ERM.

44 This line of empirical results has also been bolstered by some recent theoretical results. [SRKL20]
 45 constructed a synthetic dataset where a linear model trained with IW is provably not robust to
 46 subpopulation shift. [XYR21] further proved that under gradient descent (GD) with a sufficiently
 47 small learning rate, a linear classifier trained with either IW or ERM converges to the same max-
 48 margin classifier, and thus upon convergence, are no different. These previous theoretical results are
 49 limited to linear models and specific approaches such as IW where sample weights are fixed during
 50 training. They are not applicable to more complex models, and more general approaches where the
 51 sample weights could iteratively change, including most DRO variants.

52 Towards placing the empirical results on a stronger theoretical footing, we define the class of
 53 *generalized reweighting* (GRW), which dynamically assigns weights to the training samples, and
 54 iteratively minimizes the weighted average of the sample losses. By allowing the weights to vary
 55 with iterations, we cover not just static importance weighting, but also DRO approaches outlined
 56 earlier; though of course, the GRW class is much broader than just these instances.

57 In this work, we prove the comprehensive result that in both regression and classification, and for
 58 both overparameterized linear models and wide neural networks, the models learnt via any GRW
 59 approach and ERM are similar, in the sense that their implicit biases are (almost) equivalent. We note
 60 that extending the analysis from linear models to wide neural networks is non-trivial since it requires
 61 the result that wide neural networks can be approximated by their linearized counterparts to hold
 62 *uniformly throughout the iterative process* of GRW algorithms. Our results extend the analysis in
 63 [LXS⁺19], but as we show, the proof in the original paper had some flaws, and due to which we have
 64 to fix the proof by changing the network initialization (Eqn. (9), see Appendix E).

65 Overall, the important takeaway is that *distributionally robust generalization* cannot be directly
 66 achieved by the broad class of GRW algorithms (which includes popular approaches such as impor-
 67 tance weighting and most DRO variants). Progress towards this important goal thus requires either
 68 going beyond GRW algorithms, or devising novel loss functions that are adapted to GRW approaches.
 69 In Section 6 we will discuss some promising future directions as well as the limitations of this work.

70 2 Preliminaries

71 Let the input space be $\mathcal{X} \subseteq \mathbb{R}^d$ and the output space be $\mathcal{Y} \subseteq \mathbb{R}$.¹ We assume that \mathcal{X} is a subset of the
 72 unit L_2 ball of \mathbb{R}^d , so that any $\mathbf{x} \in \mathcal{X}$ satisfies $\|\mathbf{x}\|_2 \leq 1$. We have a training set $\{\mathbf{z}_i = (\mathbf{x}_i, y_i)\}_{i=1}^n$
 73 *i.i.d.* sampled from an underlying distribution P over $\mathcal{X} \times \mathcal{Y}$. Denote $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$,
 74 and $\mathbf{Y} = (y_1, \dots, y_n) \in \mathbb{R}^n$. For any function $g : \mathcal{X} \mapsto \mathbb{R}^m$, we overload notation and use
 75 $g(\mathbf{X}) = (g(\mathbf{x}_1), \dots, g(\mathbf{x}_n)) \in \mathbb{R}^{m \times n}$ (except when $m = 1$, $g(\mathbf{X})$ is defined as a column vector).
 76 Let the loss function be $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, 1]$. ERM trains a model by minimizing its *expected risk*
 77 $\mathcal{R}(f; P) = \mathbb{E}_{\mathbf{z} \sim P}[\ell(f(\mathbf{x}), y)]$ via minimizing the *empirical risk* $\hat{\mathcal{R}}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), y_i)$.

78 In distributional shift, the model is evaluated not on the training distribution P , but a different test
 79 distribution P_{test} , so that we care about the expected risk $\mathcal{R}(f; P_{\text{test}})$. A large family of methods
 80 designed for such distributional shift is *distributionally robust optimization* (DRO), which minimizes
 81 the expected risk over the worst-case distribution $Q \ll P^2$ in a ball w.r.t. divergence D around the
 82 training distribution P . Specifically, DRO minimizes the *expected DRO risk* defined as:

$$\mathcal{R}_{D,\rho}(f; P) = \sup_{Q \ll P} \{\mathbb{E}_Q[\ell(f(\mathbf{x}), y)] : D(Q \| P) \leq \rho\} \quad (1)$$

83 for $\rho > 0$. Examples include CVaR, χ^2 -DRO [HSNL18], and DORO [ZDKR21], among others.

¹Our results can be easily extended to the multi-class scenario (see Appendix B).

²For distributions P and Q , Q is *absolute continuous* to P , or $Q \ll P$, means that for any event A , $P(A) = 0$ implies $Q(A) = 0$.

84 A common category of distribution shift is known as subpopulation shift. Let the data domain contain
 85 K groups $\mathcal{D}_1, \dots, \mathcal{D}_K$. The training distribution P is the distribution over all groups, and the test
 86 distribution P_{test} is the distribution over one of the groups. Let $P_k(\mathbf{z}) = P(\mathbf{z} \mid \mathbf{z} \in \mathcal{D}_k)$ be the
 87 conditional distribution over group k , then P_{test} can be any one of P_1, \dots, P_K . The goal is to train a
 88 model f that performs well over every group. There are two common ways to achieve this goal: one
 89 is minimizing the *balanced empirical risk* which is an unweighted average of the empirical risk over
 90 each group, and the other is minimizing the *worst-group risk* defined as

$$\mathcal{R}_{\max}(f; P) = \max_{k=1, \dots, K} \mathcal{R}(f; P_k) = \max_{k=1, \dots, K} \mathbb{E}_{\mathbf{z} \sim P} [\ell(f(\mathbf{x}), y) \mid \mathbf{z} \in \mathcal{D}_k] \quad (2)$$

91 3 Generalized Reweighting (GRW)

92 Various methods have been proposed towards learning models that are robust to distributional shift.
 93 In contrast to analyzing each of these individually, we instead consider a large class of what we call
 94 Generalized Reweighting (GRW) algorithms that includes the ones mentioned earlier, but potentially
 95 many others more. Loosely, GRW algorithms iteratively assign each sample a weight during training
 96 (that could vary with the iteration) and iteratively minimize the weighted average risk. Specifically, at
 97 iteration t , GRW assigns a weight $q_i^{(t)}$ to sample \mathbf{z}_i , and minimizes the weighted empirical risk:

$$\hat{\mathcal{R}}_{\mathbf{q}^{(t)}}(f) = \sum_{i=1}^n q_i^{(t)} \ell(f(\mathbf{x}_i), y_i) \quad (3)$$

98 where $\mathbf{q}^{(t)} = (q_1^{(t)}, \dots, q_n^{(t)})$ and $q_1^{(t)} + \dots + q_n^{(t)} = 1$.

99 *Static GRW* assigns to each $\mathbf{z}_i = (\mathbf{x}_i, y_i)$ a fixed weight q_i that does not change during training, i.e.
 100 $q_i^{(t)} \equiv q_i$. A classical method is *importance weighting* [Shi00], where if $\mathbf{z}_i \in \mathcal{D}_k$ and the size of \mathcal{D}_k
 101 is n_k , then $q_i = (Kn_k)^{-1}$. Under importance weighting, (3) becomes the balanced empirical risk in
 102 which each group has the same weight. Note that ERM is also a special case of static GRW.

103 On the other hand, in *dynamic GRW*, $\mathbf{q}^{(t)}$ changes with t . For instance, any approach that iteratively
 104 upweights samples with high losses in order to help the model learn “hard” samples, such as DRO,
 105 is an instance of GRW. When estimating the population DRO risk $\mathcal{R}_{D, \rho}(f; P)$ in Eqn. (1), if P
 106 is set to the empirical distribution over the training samples, then $Q \ll P$ implies that Q is also
 107 a distribution over the training samples. Thus, DRO methods belong to the broad class of GRW
 108 algorithms. There are two common ways to implement DRO. One uses Danskin’s theorem and
 109 chooses Q as the maximizer of $\mathbb{E}_Q[\ell(f(\mathbf{x}), y)]$ in each epoch. The other one formulates DRO as a
 110 bi-level optimization problem, where the lower level updates the model to minimize the expected risk
 111 over Q , and the upper level updates Q to maximize it. Both can be seen as instances of GRW. As one
 112 popular instance of the latter, *Group DRO* was proposed by [SKHL20] to minimize (2). Denote the
 113 empirical risk over group k by $\hat{\mathcal{R}}_k(f)$, and the model at time t by $f^{(t)}$. Group DRO iteratively sets
 114 $q_i^{(t)} = g_k^{(t)} / n_k$ for all $\mathbf{z}_i \in \mathcal{D}_k$ where $g_k^{(t)}$ is the group weight that is updated as

$$g_k^{(t)} \propto g_k^{(t-1)} \exp\left(\nu \hat{\mathcal{R}}_k(f^{(t-1)})\right) \quad (\forall k = 1, \dots, K) \quad (4)$$

115 for some $\nu > 0$, and then normalized so that $q_1^{(t)} + \dots + q_n^{(t)} = 1$. [SKHL20] then showed (in
 116 their Proposition 2) that for convex settings, the Group DRO risk of iterates converges to the global
 117 minimum with the rate $O(t^{-1/2})$ if ν is sufficiently small.

118 4 Theoretical Results for Regression

119 In this section, we will study GRW for regression tasks that use the squared loss

$$\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2. \quad (5)$$

120 We will prove that for both linear models and sufficiently wide fully-connected neural networks, the
 121 implicit bias of GRW is equivalent to ERM, so that starting from the same initial point, GRW and
 122 ERM will converge to the same point when trained for an infinitely long time, which explains why
 123 GRW does not improve over ERM without regularization and early stopping. We will further show
 124 that while regularization can affect this implicit bias, it must be large enough to *significantly lower*
 125 *the training performance*, or the final model will still be similar to the unregularized ERM model.

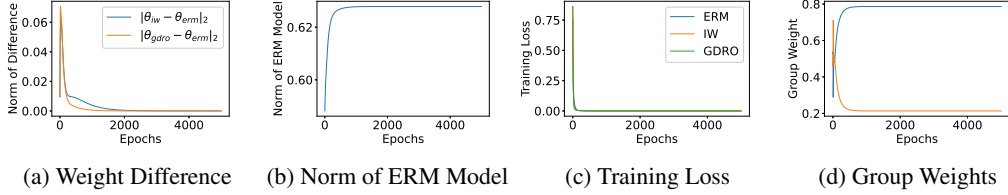


Figure 1: Experimental results of ERM, importance weighting (IW) and Group DRO (GDRO) with the squared loss on six MNIST images with a linear model. All norms are L_2 norms.

126 4.1 Linear Models

127 We first demonstrate our result on simple linear models to provide our readers with a key intuition;
 128 later, we will apply this same intuition to neural networks. This key intuition draws from results
 129 of [GLSS18]. Let the linear model be denoted by $f(\mathbf{x}) = \langle \theta, \mathbf{x} \rangle$, where $\theta \in \mathbb{R}^d$. We consider the
 130 overparameterized setting where $d > n$. The weight update rule of GRW under GD is the following:

$$\theta^{(t+1)} = \theta^{(t)} - \eta \sum_{i=1}^n q_i^{(t)} \nabla_{\theta} \ell(f^{(t)}(\mathbf{x}_i), y_i) \quad (6)$$

131 where $\eta > 0$ is the learning rate. For a linear model with the squared loss, the update rule is

$$\theta^{(t+1)} = \theta^{(t)} - \eta \sum_{i=1}^n q_i^{(t)} \mathbf{x}_i (f^{(t)}(\mathbf{x}_i) - y_i) \quad (7)$$

132 For this training scheme, we can prove that if the training error converges to zero, then the model
 133 converges to an interpolator θ^* (s.t. $\forall i, \langle \theta^*, \mathbf{x}_i \rangle = y_i$) independent of $q_i^{(t)}$ (proofs in Appendix D):

134 **Theorem 1.** *If $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly independent, then under the squared loss, for any GRW such*
 135 *that the empirical training risk $\hat{\mathcal{R}}(f^{(t)}) \rightarrow 0$ as $t \rightarrow \infty$, it holds that $\theta^{(t)}$ converges to an interpolator*
 136 *θ^* that only depends on $\theta^{(0)}$ and $\mathbf{x}_1, \dots, \mathbf{x}_n$, but does not depend on $q_i^{(t)}$.*

137 The proof is based on the following key intuition regarding the update rule (7): $\theta^{(t+1)} - \theta^{(t)}$ is
 138 a linear combination of $\mathbf{x}_1, \dots, \mathbf{x}_n$ for all t , so $\theta^{(t)} - \theta^{(0)}$ always lies in the linear subspace
 139 $\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, which is an n -dimensional linear subspace if $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly independent.
 140 By Cramer’s rule, there is exactly one $\tilde{\theta}$ in this subspace such that we get interpolation of all the
 141 data $\langle \tilde{\theta} + \theta^{(0)}, \mathbf{x}_i \rangle = y_i$ for all $i \in \{1, \dots, n\}$. In other words, the parameter $\theta^* = \tilde{\theta} + \theta^{(0)}$ in this
 142 subspace that interpolates all the data is unique. Thus the proof would follow if we were to show that
 143 $\theta^{(t)} - \theta^{(0)}$, which lies in the subspace, also converges to interpolating the data.

144 We have essentially proved the following sobering result: *the implicit bias of any GRW that achieves*
 145 *zero training error is equivalent to ERM, so GRW does not improve over ERM.* While the various
 146 distributional shift methods discussed in the introduction have been shown to satisfy the precondition
 147 of convergence to zero training error with overparameterized models and linearly independent
 148 inputs [SKHL20], we provide the following theorem that shows this for the broad class of GRW
 149 methods. Specifically, we show this result for any GRW satisfying the following assumption with a
 150 sufficiently small learning rate:

151 **Assumption 1.** There are constants q_1, \dots, q_n s.t. $\forall i, q_i^{(t)} \rightarrow q_i$ as $t \rightarrow \infty$. And $\min_i q_i = q^* > 0$.

152 **Theorem 2.** *If $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly independent, then there exists $\eta_0 > 0$ such that for any*
 153 *GRW satisfying Assumption 1 with the squared loss, and any $\eta \leq \eta_0$, the empirical training risk*
 154 *$\hat{\mathcal{R}}(f^{(t)}) \rightarrow 0$ as $t \rightarrow \infty$.*

155 Finally, we use a simple experiment to demonstrate the correctness of this result. The experiment is
 156 conducted on a training set of six MNIST images, five of which are digit 0 and one is digit 1. We use
 157 a 784-dimensional linear model and run ERM, importance weighting and group DRO. The results are
 158 presented in Figure 1, and they show that the training loss of each method converges to 0, and the gap
 159 between the model weights of importance weighting, Group DRO and ERM converges to 0, meaning
 160 that all three model weights converge to the same point, whose L_2 norm is about 0.63. Figure 1d also
 161 shows that the group weights in Group DRO empirically satisfy Assumption 1.

162 **4.2 Wide Neural Networks (Wide NNs)**

163 Now we study *sufficiently wide fully-connected neural networks*. We extend the analysis in [LXS⁺19]
 164 in the neural tangent kernel (NTK) regime [JGH18]. In particular we study the following network:

$$\mathbf{h}^{l+1} = \frac{W^l}{\sqrt{d_l}} \mathbf{x}^l + \beta \mathbf{b}^l \quad \text{and} \quad \mathbf{x}^{l+1} = \sigma(\mathbf{h}^{l+1}) \quad (l = 0, \dots, L) \quad (8)$$

165 where σ is a non-linear activation function, $W^l \in \mathbb{R}^{d_{l+1} \times d_l}$ and $W^L \in \mathbb{R}^{1 \times d_L}$. Here $d_0 = d$. The
 166 parameter vector θ consists of W^0, \dots, W^L and b^0, \dots, b^L (θ is the concatenation of all flattened
 167 weights and biases). The final output is $f(\mathbf{x}) = \mathbf{h}^{L+1}$. And let the neural network be initialized as

$$\begin{cases} W_{i,j}^{l(0)} \sim \mathcal{N}(0, 1) \\ \mathbf{b}_j^{l(0)} \sim \mathcal{N}(0, 1) \end{cases} \quad (l = 0, \dots, L-1) \quad \text{and} \quad \begin{cases} W_{i,j}^{L(0)} = 0 \\ \mathbf{b}_j^{L(0)} \sim \mathcal{N}(0, 1) \end{cases} \quad (9)$$

168 We also need the following assumption on the wide neural network:

169 **Assumption 2.** σ is differentiable everywhere. Both σ and its first-order derivative $\dot{\sigma}$ are Lipschitz.³

170 **Difference from [JGH18].** Our initialization (9) differs from the original one in [JGH18] in the last
 171 (output) layer, where we use the zero initialization $W_{i,j}^{L(0)} = 0$ instead of the Gaussian initialization
 172 $W_{i,j}^{L(0)} \sim \mathcal{N}(0, 1)$. This modification permits us to accurately approximate the neural network with
 173 its linearized counterpart (11), as we notice that the proofs in [LXS⁺19] (particularly the proofs of
 174 their Theorem 2.1 and their Lemma 1 in Appendix G) are flawed. In Appendix E we will explain
 175 what goes wrong in their proofs and how we manage to fix the proofs with our modification.

176 Denote the neural network at time t by $f^{(t)}(\mathbf{x}) = f(\mathbf{x}; \theta^{(t)})$ which is parameterized by $\theta^{(t)} \in \mathbb{R}^p$
 177 where p is the number of parameters. We use the shorthand $\nabla_{\theta} f^{(0)}(\mathbf{x}) := \nabla_{\theta} f(\mathbf{x}; \theta)|_{\theta=\theta_0}$. The
 178 *neural tangent kernel* (NTK) of this model is $\Theta^{(0)}(\mathbf{x}, \mathbf{x}') = \nabla_{\theta} f^{(0)}(\mathbf{x})^{\top} \nabla_{\theta} f^{(0)}(\mathbf{x}')$, and the *Gram*
 179 *matrix* is $\Theta^{(0)} = \Theta^{(0)}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{n \times n}$. For this wide NN, we still have the following NTK theorem:

180 **Lemma 3.** *If σ is Lipschitz and $d_l \rightarrow \infty$ for $l = 1, \dots, L$ sequentially, then $\Theta^{(0)}(\mathbf{x}, \mathbf{x}')$ converges*
 181 *in probability to a non-degenerate⁴ deterministic limiting kernel $\Theta(\mathbf{x}, \mathbf{x}')$.*

182 The *kernel Gram matrix* $\Theta = \Theta(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{n \times n}$ is a positive semi-definite symmetric matrix.
 183 Denote its largest and smallest eigenvalues by λ^{\max} and λ^{\min} . Note that Θ is non-degenerate, so we
 184 can assume that $\lambda^{\min} > 0$ (which is almost surely true when $d_L \gg n$). Then we have:

185 **Theorem 4.** *Let $f^{(t)}$ be a wide fully-connected neural network that satisfies Assumption 2 and is*
 186 *trained by any GRW satisfying Assumption 1 with the squared loss. Let $f_{\text{ERM}}^{(t)}$ be the same model*
 187 *trained by ERM from the same initial point. If $d_1 = \dots = d_L = \tilde{d}$, $\nabla_{\theta} f^{(0)}(\mathbf{x}_1), \dots, \nabla_{\theta} f^{(0)}(\mathbf{x}_n)$*
 188 *are linearly independent, and $\lambda^{\min} > 0$, then there exists a constant $\eta_1 > 0$ such that: if $\eta \leq \eta_1$ ⁵,*
 189 *then for any $\delta > 0$, there exists $\tilde{D} > 0$ such that as long as $\tilde{d} \geq \tilde{D}$, with probability at least $(1 - \delta)$*
 190 *over random initialization we have: for any test point $\mathbf{x} \in \mathbb{R}^d$ such that $\|\mathbf{x}\|_2 \leq 1$, as $\tilde{d} \rightarrow \infty$,*

$$\limsup_{t \rightarrow \infty} \left| f^{(t)}(\mathbf{x}) - f_{\text{ERM}}^{(t)}(\mathbf{x}) \right| = O(\tilde{d}^{-1/4}) \rightarrow 0 \quad (10)$$

191 Note that for simplicity, in the theorem we only consider the case where $d_1 = \dots = d_L = \tilde{d} \rightarrow \infty$,
 192 but in fact the result can be very easily extended to the case where $d_l/d_1 \rightarrow \alpha_l$ for $l = 2, \dots, L$ for
 193 some constants $\alpha_2, \dots, \alpha_L$, and $d_1 \rightarrow \infty$. Here we provide a proof sketch for this theorem. The key
 194 is to consider the *linearized neural network* of $f^{(t)}(\mathbf{x})$:

$$f_{\text{lin}}^{(t)}(\mathbf{x}) = f^{(0)}(\mathbf{x}) + \langle \theta^{(t)} - \theta^{(0)}, \nabla_{\theta} f^{(0)}(\mathbf{x}) \rangle \quad (11)$$

195 which is a linear model with features $\nabla_{\theta} f^{(0)}(\mathbf{x})$. Thus if $\nabla_{\theta} f^{(0)}(\mathbf{x}_1), \dots, \nabla_{\theta} f^{(0)}(\mathbf{x}_n)$ are linearly
 196 independent, then the linearized NN converges to the unique interpolator. Then we show that the

³ f is Lipschitz if there exists a constant $L > 0$ such that for any $\mathbf{x}_1, \mathbf{x}_2$, $|f(\mathbf{x}_1) - f(\mathbf{x}_2)| \leq L \|\mathbf{x}_1 - \mathbf{x}_2\|_2$.

⁴ *Non-degenerate* means that $\Theta(\mathbf{x}, \mathbf{x}')$ depends on \mathbf{x} and \mathbf{x}' and is not a constant.

⁵ For ease of understanding, later we will write this condition as “with a sufficiently small learning rate”.

197 wide neural network can be approximated by its linearized counterpart *uniformly throughout training*,
 198 which is considerably more subtle in our case due to the GRW dynamics. Here we prove that the gap
 199 is bounded by $O(\tilde{d}^{-1/4})$, but in fact we can prove that it is bounded by $O(\tilde{d}^{-1/2+\epsilon})$ for any $\epsilon > 0$:

200 **Lemma 5** (Approximation Theorem). *For a wide fully-connected neural network $f^{(t)}$ satisfying*
 201 *Assumption 2 and is trained by any GRW satisfying Assumption 1 with the squared loss, let $f_{\text{lin}}^{(t)}$ be its*
 202 *linearized neural network trained by the same GRW (i.e. $q_i^{(t)}$ are the same for both networks for any*
 203 *i and t). Under the conditions of Theorem 4, with a sufficiently small learning rate, for any $\delta > 0$,*
 204 *there exist constants $\tilde{D} > 0$ and $C > 0$ such that as long as $\tilde{d} \geq \tilde{D}$, with probability at least $(1 - \delta)$*
 205 *over random initialization we have: for any test point $\mathbf{x} \in \mathbb{R}^d$ such that $\|\mathbf{x}\|_2 \leq 1$,*

$$\sup_{t \geq 0} \left| f_{\text{lin}}^{(t)}(\mathbf{x}) - f^{(t)}(\mathbf{x}) \right| \leq C \tilde{d}^{-1/4} \quad (12)$$

206 Theorem 4 shows that at *any test point \mathbf{x}* within the unit ball, the gap between the outputs of wide
 207 NNs trained by GRW and ERM from the same initial point is arbitrarily close to 0. So we have shown
 208 that for regression, with both linear and wide NNs, GRW does not improve over ERM.

209 4.3 Wide Neural Networks, with L_2 Regularization

210 Previous work such as [SKHL20] proposed to improve DRO algorithms by adding L_2 penalty to the
 211 objective function. In this section, we thus study adding L_2 regularization to GRW algorithms:

$$\hat{\mathcal{R}}_{q^{(t)}}^\mu(f) = \sum_{i=1}^n q_i^{(t)} \ell(f(\mathbf{x}_i), y_i) + \frac{\mu}{2} \|\theta - \theta^{(0)}\|_2^2 \quad (13)$$

212 From the outset, it is easy to see that under L_2 regularization, GRW methods have different implicit
 213 biases than ERM. For example, when f is a linear model, ℓ is convex and smooth, then $\hat{\mathcal{R}}_{q^{(t)}}^\mu(f)$ with
 214 static GRW is a convex smooth objective function, so under GD with a sufficiently small learning rate,
 215 the model will converge to the global minimizer (see Appendix D.1). Moreover, the global optimum
 216 θ^* satisfies $\nabla_{\theta} \hat{\mathcal{R}}_{q^{(t)}}^\mu(f(\mathbf{x}; \theta^*)) = 0$, solving which yields $\theta^* = \theta^{(0)} + (\mathbf{X}\mathbf{Q}\mathbf{X}^\top + \mu\mathbf{I})^{-1} \mathbf{X}\mathbf{Q}(\mathbf{Y} -$
 217 $f^{(0)}(\mathbf{X}))$, which depends on $\mathbf{Q} = \text{diag}(q_1, \dots, q_n)$, so adding L_2 regularization at least seems to
 218 yield different results from ERM (whether it improves over ERM might depend on q_1, \dots, q_n).

219 However, the following result shows that this regularization must be large enough to *significantly*
 220 *lower the training performance*, or the resulting model would still be close to the unregularized ERM
 221 model. We still denote the largest and smallest eigenvalues of the kernel Gram matrix Θ by λ^{\max} and
 222 λ^{\min} . We use the subscript “reg” to refer to a regularized model (trained by minimizing (13)).

223 **Theorem 6.** *Suppose there exists $M_0 > 0$ s.t. $\|\nabla_{\theta} f^{(0)}(\mathbf{x})\|_2 \leq M_0$ for all $\|\mathbf{x}\|_2 \leq 1$. If $\lambda^{\min} > 0$*
 224 *and $\mu > 0$, then for a wide NN satisfying Assumption 2, and any GRW minimizing the squared loss*
 225 *with a sufficiently small learning rate η , if $d_1 = d_2 = \dots = d_L = \tilde{d}$, $\nabla_{\theta} f^{(0)}(\mathbf{x}_1), \dots, \nabla_{\theta} f^{(0)}(\mathbf{x}_n)$*
 226 *are linearly independent, and the empirical training risk of $f_{\text{reg}}^{(t)}$ satisfies*

$$\limsup_{t \rightarrow \infty} \hat{\mathcal{R}}(f_{\text{reg}}^{(t)}) < \epsilon \quad (14)$$

227 *for some $\epsilon > 0$, then with a sufficiently small learning rate, as $\tilde{d} \rightarrow \infty$, with probability close to 1*
 228 *over random initialization, for any \mathbf{x} such that $\|\mathbf{x}\|_2 \leq 1$ we have*

$$\limsup_{t \rightarrow \infty} \left| f_{\text{reg}}^{(t)}(\mathbf{x}) - f_{\text{ERM}}^{(t)}(\mathbf{x}) \right| = O(\tilde{d}^{-1/4} + \sqrt{\epsilon}) \rightarrow O(\sqrt{\epsilon}) \quad (15)$$

229 *where $f_{\text{reg}}^{(t)}$ is trained by regularized GRW and $f_{\text{ERM}}^{(t)}$ by unregularized ERM from same initial points.*

230 The proof again starts from analyzing linearized neural networks, and showing that regularization
 231 does not help there (Appendix D.4.2). Then, we need to prove a new approximation theorem for L_2
 232 regularized GRW connecting wide NNs to their linearized counterparts uniformly through the GRW
 233 training process (Appendix D.4.1). Note that with regularization, we no longer need Assumption
 234 1 to prove the new approximation theorem, because previously Assumption 1 is used to prove the
 235 convergence of GRW, but with regularization GRW naturally converges.

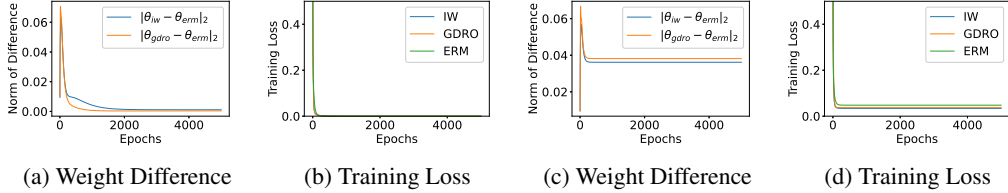


Figure 2: Experimental results of ERM, importance weighting (IW) and Group DRO (GDRO) with L_2 regularization with the squared loss. Left two: $\mu = 0.1$; Right two: $\mu = 10$.

Theorem 6 shows that if the training error can go below ϵ , then the gap between the outputs of the two models *on any test point* \mathbf{x} within the unit ball will be at most $O(\sqrt{\epsilon})$. Thus, if ϵ is very small, regularized GRW yields a very similar model to unregularized ERM, and thus makes improvement.

To empirically demonstrate this result, we run the same experiment as in Section 4.1 but with L_2 regularization. The results are presented in Figure 2. We can see that when the regularization is small, the training losses still converge to 0, and the three model weights still converge to the same point. On the contrary, with a large regularization, the training loss does not converge to 0, and the three model weights no longer converge to the same point. This shows that the regularization must be large enough to lower the training performance in order to make a significant difference to the implicit bias.

5 Theoretical Results for Classification

Now we consider classification where $\mathcal{Y} = \{+1, -1\}$. The big difference is that *classification losses don't have finite minimizers*. A classification loss converging to zero means that the model weight “explodes” to infinity instead of converging to a finite point. We focus on the canonical logistic loss:

$$\ell(\hat{y}, y) = \log(1 + \exp(-\hat{y}y)) \quad (16)$$

5.1 Linear Models

We first consider training the linear model $f(\mathbf{x}) = \langle \theta, \mathbf{x} \rangle$ with GRW under gradient descent with the logistic loss. As noted earlier, in this setting, [BL19] made the empirical observation that importance weighting does not improve over ERM. Then, [XYR21] proved that for importance weighting algorithms, as $t \rightarrow \infty$, $\|\theta^{(t)}\|_2 \rightarrow \infty$ and $\theta^{(t)} / \|\theta^{(t)}\|_2$ converges to a unit vector that does not depend on the sample weights, so it does not improve over ERM. To extend this theoretical result to the broad class of GRW algorithms, we will prove two results. First, in Theorem 7 we will show that under the logistic loss, any GRW algorithm satisfying the following weaker assumption:

Assumption 3. For all i , $\liminf_{t \rightarrow \infty} q_i^{(t)} > 0$,

if the training error converges to 0, and the direction of the model weight converges to a fixed unit vector, then this unit vector must be the *max-margin classifier* defined as

$$\hat{\theta}_{\text{MM}} = \arg \max_{\theta: \|\theta\|_2=1} \left\{ \min_{i=1, \dots, n} y_i \cdot \langle \theta, \mathbf{x}_i \rangle \right\} \quad (17)$$

Second, Theorem 8 shows that for any GRW satisfying Assumption 1, the training error converges to 0 and the direction of the model weight converges, so it does not improve over ERM.

Theorem 7. *If $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly independent, then for the logistic loss, we have: for any GRW satisfying Assumption 3, if as $t \rightarrow \infty$ the empirical training risk $\hat{\mathcal{R}}(f^{(t)})$ converges to 0 and $\theta^{(t)} / \|\theta^{(t)}\|_2 \rightarrow \mathbf{u}$ for some unit vector \mathbf{u} , then $\mathbf{u} = \hat{\theta}_{\text{MM}}$.*

This result is an extension of [SHN⁺18]. Note that $\hat{\theta}_{\text{MM}}$ does not depend on $q_i^{(t)}$, so this result shows that the sample weights have no effect on the implicit bias. Thus, for any GRW method that only satisfies the weak Assumption 3, as long as the training error converges to 0 and the model weight direction converges, GRW does not improve over ERM. We next show that any GRW satisfying Assumption 1 does have its model weight direction converge, and its training error converge to 0.

Theorem 8. *For any loss ℓ that is convex, L -smooth in \hat{y} and strictly monotonically decreasing to zero as $y\hat{y} \rightarrow +\infty$, and GRW satisfying Assumption 1, denote $F(\theta) = \sum_{i=1}^n q_i \ell(\langle \theta, \mathbf{x}_i \rangle, y_i)$. If $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly independent, then with a sufficiently small learning rate η , we have:*

- 273 (i) $F(\theta^{(t)}) \rightarrow 0$ as $t \rightarrow \infty$. (ii) $\|\theta^{(t)}\|_2 \rightarrow \infty$ as $t \rightarrow \infty$.
- 274 (iii) Let $\theta_R = \arg \min_{\theta} \{F(\theta) : \|\theta\|_2 \leq R\}$. θ_R is unique for any R such that $\min_{\|\theta\|_2 \leq R} F(\theta) < \min_i q_i \ell(0, y_i)$. And if $\lim_{R \rightarrow \infty} \frac{\theta_R}{R}$ exists, then $\lim_{t \rightarrow \infty} \frac{\theta^{(t)}}{\|\theta^{(t)}\|_2}$ also exists and they are equal.

275 This result is an extension of Theorem 1 of [JDST20]. For the logistic loss, it is easy to show that
 276 it satisfies the conditions of the above theorem and $\lim_{R \rightarrow \infty} \frac{\theta_R}{R} = \hat{\theta}_{\text{MM}}$. Thus, Theorems 8 and 7
 277 together imply that all GRW satisfying Assumption 1 (including ERM) have the same implicit bias
 278 (see Appendix D.5.3). We also have empirical verification for these results (see Appendix C).

279 **Remark.** It is impossible to extend these results to wide NNs like Theorem 4 because for a neural
 280 network, if $\|\theta^{(t)}\|_2$ goes to infinity, then $\|\nabla_{\theta} f\|_2$ will also go to infinity. However, for a linear model,
 281 the gradient is a constant. Consequently, the gap between the neural networks and its linearized
 282 counterpart will “explode” under gradient descent, so there can be no approximation theorem like
 283 Lemma 5 that can connect wide NNs to their linearized counterparts. Thus, we consider regularized
 284 GRW, for which $\theta^{(t)}$ converges to a finite point and there is an approximation theorem.

285 5.2 Wide Neural Networks, with L_2 Regularization

286 Consider minimizing the regularized weighted empirical risk (13) with ℓ being the logistic loss. As in
 287 the regression case, with L_2 regularization, GRW methods have different implicit biases than ERM
 288 for the same reasons as in Section 4.3. And similarly, we can show that in order for GRW methods to
 289 be sufficiently different from ERM, the regularization needs to be large enough to significantly lower
 290 the training performance. Specifically, in the following theorem we show that if the regularization
 291 is too small to lower the training performance, then a wide neural network trained with regularized
 292 GRW and the logistic loss will still be very close to the *max-margin linearized neural network*:

$$f_{\text{MM}}(\mathbf{x}) = \langle \hat{\theta}_{\text{MM}}, \nabla_{\theta} f^{(0)}(\mathbf{x}) \rangle \quad \text{where} \quad \hat{\theta}_{\text{MM}} = \arg \max_{\|\theta\|_2=1} \left\{ \min_{i=1, \dots, n} y_i \cdot \langle \theta, \nabla_{\theta} f^{(0)}(\mathbf{x}_i) \rangle \right\} \quad (18)$$

293 Note that f_{MM} does not depend on $q_i^{(t)}$. Moreover, using the result in the previous section we can
 294 show that a linearized neural network trained with unregularized ERM will converge to f_{MM} :

295 **Theorem 9.** Suppose there exists $M_0 > 0$ such that $\|\nabla_{\theta} f^{(0)}(\mathbf{x})\|_2 \leq M_0$ for all test point \mathbf{x} . For a
 296 wide NN satisfying Assumption 2, and for any GRW satisfying Assumption 1 with the logistic loss,
 297 if $d_1 = d_2 = \dots = d_L = \tilde{d}$ and $\nabla_{\theta} f^{(0)}(\mathbf{x}_1), \dots, \nabla_{\theta} f^{(0)}(\mathbf{x}_n)$ are linearly independent and the
 298 learning rate is sufficiently small, then for any $\delta > 0$ there exists a constant $C > 0$ such that: with
 299 probability at least $(1 - \delta)$ over random initialization, as $\tilde{d} \rightarrow \infty$ we have: for any $\epsilon \in (0, \frac{1}{4})$, if
 300 the empirical training error satisfies $\limsup_{t \rightarrow \infty} \hat{\mathcal{R}}(f_{\text{reg}}^{(t)}) < \epsilon$, then for any test point \mathbf{x} such that
 301 $|f_{\text{MM}}(\mathbf{x})| > C \cdot (-\log 2\epsilon)^{-1/2}$, $f_{\text{reg}}^{(t)}(\mathbf{x})$ has the same sign as $f_{\text{MM}}(\mathbf{x})$ when t is sufficiently large.

302 This result says that at any test point \mathbf{x} on which the max-margin linear classifier classifies with a
 303 margin of $\Omega((-\log 2\epsilon)^{-1/2})$, the neural network has the same prediction. And as ϵ decreases, the
 304 confidence threshold also becomes lower. Similar to Theorem 6, this theorem provides the scaling of
 305 the gap between the regularized GRW model and the unregularized ERM model *w.r.t.* ϵ .

306 This result justifies the empirical observation in [SKHL20] that with large regularization, some GRW
 307 algorithms can maintain a high worst-group test performance, with the cost of suffering a significant
 308 drop in training accuracy. On the other hand, if the regularization is small and the model can achieve
 309 nearly perfect training accuracy, then its worst-group test performance will still significantly drop.

310 6 Discussion

311 6.1 Distributionally Robust Generalization and Future Directions

312 A large body of prior work focused on distributionally robust optimization, but we show that these
 313 methods have (almost) equivalent implicit biases as ERM. In other words, *distributionally robust*
 314 *optimization* (DRO) does not necessarily have better *distributionally robust generalization* (DRG).

315 Therefore, we argue that it is necessary to design principled ways to improve DRG, which is what
 316 people really want in the first place. Here we discuss three promising approaches to improving DRG.

317 The first approach is data augmentation and pretraining on large datasets. Our theoretical findings
 318 suggest that the implicit bias of GRW is determined by the training samples and the initial point, but
 319 not the sample weights. Thus, to improve DRG, we can either obtain more training samples, or start
 320 from a better initial point, as demonstrated in two recent papers [WGS⁺22, SKL⁺22].

321 The second approach (for classification) is to go beyond the class of (iterative) sample reweighting
 322 based GRW algorithms, for instance via *logit adjustment* [MJR⁺21], which makes a classifier *have*
 323 *larger margins on smaller groups* to improve its generalization on smaller groups. An early approach
 324 by [CWG⁺19] proposed to add an $O(n_k^{-1/4})$ additive adjustment term to the logits output by the
 325 classifier. Following this spirit, [MJR⁺21] proposed the LA-loss which also adds an additive adjust-
 326 ment term to the logits. [YCZC20] proposed the CDT-loss which adds a multiplicative adjustment
 327 term to the logits by dividing the logits of different classes with different temperatures. [KPOT21]
 328 proposed the VS-loss which includes both additive and multiplicative adjustment terms, and they
 329 showed that only the multiplicative adjustment term affects the implicit bias, while the additive term
 330 only affects optimization, a fact that can be easily derived from our Theorem 8. Finally, [LZT⁺21]
 331 proposed AutoBalance which optimizes the adjustment terms with a bi-level optimization framework.

332 The third approach is to stay within the class of GRW algorithms, but to change the classifica-
 333 tion/regression loss function to be suited to GRW. A recent paper [WCHH22] showed that for linear
 334 classifiers, one can make the implicit bias of GRW dependent on the sample weights by replacing the
 335 exponentially-tailed logistic loss with the following *polynomially-tailed loss*:

$$\ell_{\alpha,\beta}(\hat{y}, y) = \begin{cases} \ell_{\text{left}}(\hat{y}y) & , \text{ if } \hat{y}y < \beta \\ \frac{1}{[\hat{y}y - (\beta - 1)]^\alpha} & , \text{ if } \hat{y}y \geq \beta \end{cases} \quad (19)$$

336 And this result can be extended to GRW satisfying Assumption 1 using our Theorem 8. The reason
 337 why loss (19) works is that it changes $\lim_{R \rightarrow \infty} \frac{\theta_R}{R}$, and the new limit depends on the sample weights.

338 6.2 Limitations

339 Like most theory papers, our work makes some strong assumptions. The two main assumptions are:

- 340 (i) The model is a linear model or a sufficiently wide fully-connected neural network.
- 341 (ii) The model is trained for sufficiently long time, *i.e.* without early stopping.

342 Regarding (i), [COB19] argued that NTK neural networks fall in the “lazy training” regime and
 343 results might not be transferable to general neural networks. However, this class of neural networks
 344 has been widely studied in recent years and has provided considerable insights into the behavior
 345 of general neural networks, which is hard to analyze otherwise. Regarding (ii), in some easy tasks,
 346 when early stopping is applied, existing algorithms for distributional shift can do better than ERM
 347 [SKHL20]. However, as demonstrated in [GLP21, KSM⁺21], in real applications these methods still
 348 cannot significantly improve over ERM even with early stopping, so early stopping is not the ultimate
 349 universal solution. Thus, though inevitably our results rely on some strong assumptions, we believe
 350 that they provide important insights into the problems of existing methods and directions for future
 351 work, which are significant contributions to the study of distributional shift problems.

352 7 Conclusion

353 In this work, we posit a broad class of what we call Generalized Reweighting (GRW) algorithms that
 354 include popular approaches such as importance weighting, and Distributionally Robust Optimization
 355 (DRO) variants, that were designed towards the task of learning models that are robust to distributional
 356 shift. We show that when used to train overparameterized linear models or wide NN models, even this
 357 very broad class of GRW algorithms does not improve over ERM, because they have the same implicit
 358 biases. We also showed that regularization does not help if it is not large enough to significantly
 359 lower the average training performance. Our results thus suggest to make progress towards learning
 360 models that are robust to distributional shift, we have to either go beyond this broad class of GRW
 361 algorithms, or design new losses specifically targeted to this class.

362 **References**

- 363 [BGO16] Su Lin Blodgett, Lisa Green, and Brendan O’Connor. Demographic dialectal variation
364 in social media: A case study of African-American English. In *Proceedings of the 2016*
365 *Conference on Empirical Methods in Natural Language Processing*, pages 1119–1130,
366 Austin, Texas, November 2016. Association for Computational Linguistics.
- 367 [BL19] Jonathon Byrd and Zachary Lipton. What is the effect of importance weighting in deep
368 learning? In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of*
369 *the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of*
370 *Machine Learning Research*, pages 872–881. PMLR, 09–15 Jun 2019.
- 371 [COB19] Lénaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable
372 programming. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox,
373 and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32.
374 Curran Associates, Inc., 2019.
- 375 [CWG⁺19] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning
376 imbalanced datasets with label-distribution-aware margin loss. *Advances in Neural*
377 *Information Processing Systems*, 32:1567–1578, 2019.
- 378 [DN18] John Duchi and Hongseok Namkoong. Learning models with uniform performance via
379 distributionally robust optimization. *arXiv preprint arXiv:1810.08750*, 2018.
- 380 [GLP21] Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In
381 *International Conference on Learning Representations*, 2021.
- 382 [GLSS18] Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit
383 bias in terms of optimization geometry. In Jennifer Dy and Andreas Krause, editors,
384 *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of
385 *Proceedings of Machine Learning Research*, pages 1832–1841. PMLR, 10–15 Jul 2018.
- 386 [HNSS18] Weihua Hu, Gang Niu, Issei Sato, and Masashi Sugiyama. Does distributionally robust
387 supervised learning give robust classifiers? In *International Conference on Machine*
388 *Learning*, pages 2029–2037. PMLR, 2018.
- 389 [HS15] Dirk Hovy and Anders Søgaard. Tagging performance correlates with author age. In
390 *Proceedings of the 53rd annual meeting of the Association for Computational Linguistics*
391 *and the 7th international joint conference on natural language processing (volume 2:*
392 *Short papers)*, pages 483–488, 2015.
- 393 [HSNL18] Tatsunori Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. Fair-
394 ness without demographics in repeated loss minimization. In Jennifer Dy and Andreas
395 Krause, editors, *International Conference on Machine Learning*, volume 80 of *Proceed-*
396 *ings of Machine Learning Research*, pages 1929–1938, Stockholm, Stockholm
397 Sweden, 10–15 Jul 2018. PMLR.
- 398 [JDST20] Ziwei Ji, Miroslav Dudík, Robert E. Schapire, and Matus Telgarsky. Gradient descent
399 follows the regularization path for general losses. In Jacob Abernethy and Shivani
400 Agarwal, editors, *Proceedings of Thirty Third Conference on Learning Theory*, volume
401 125 of *Proceedings of Machine Learning Research*, pages 2109–2136. PMLR, 09–12
402 Jul 2020.
- 403 [JGH18] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Conver-
404 gence and generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle,
405 K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information*
406 *Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- 407 [KPOT21] Ganesh Ramachandra Kini, Orestis Paraskevas, Samet Oymak, and Christos Thramp-
408 poulidis. Label-imbalanced and group-sensitive classification under overparameteriza-
409 tion. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.

- 410 [KSM⁺21] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang,
411 Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena
412 Gao, Tony Lee, Etienne David, Ian Stavness, Wei Guo, Berton Earnshaw, Imran Haque,
413 Sara M Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea
414 Finn, and Percy Liang. Wilds: A benchmark of in-the-wild distribution shifts. In Marina
415 Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on*
416 *Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages
417 5637–5664. PMLR, 18–24 Jul 2021.
- 418 [LXS⁺19] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha
419 Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve
420 as linear models under gradient descent. *Advances in neural information processing*
421 *systems*, 32:8572–8583, 2019.
- 422 [LZT⁺21] Mingchen Li, Xuechen Zhang, Christos Thrampoulidis, Jiasi Chen, and Samet Oymak.
423 Autobalance: Optimized loss functions for imbalanced data. In *Thirty-Fifth Conference*
424 *on Neural Information Processing Systems*, 2021.
- 425 [MJR⁺21] Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, An-
426 dreas Veit, and Sanjiv Kumar. Long-tail learning via logit adjustment. In *International*
427 *Conference on Learning Representations*, 2021.
- 428 [Shi00] Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting
429 the log-likelihood function. *Journal of statistical planning and inference*, 90(2):227–244,
430 2000.
- 431 [SHN⁺18] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro.
432 The implicit bias of gradient descent on separable data. *The Journal of Machine Learning*
433 *Research*, 19(1):2822–2878, 2018.
- 434 [SKHL20] Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distribu-
435 tionally robust neural networks for group shifts: On the importance of regularization for
436 worst-case generalization. In *International Conference on Learning Representations*,
437 2020.
- 438 [SKL⁺22] Shiori Sagawa, Pang Wei Koh, Tony Lee, Irena Gao, Sang Michael Xie, Kendrick Shen,
439 Ananya Kumar, Weihua Hu, Michihiro Yasunaga, Henrik Marklund, Sara Beery, Etienne
440 David, Ian Stavness, Wei Guo, Jure Leskovec, Kate Saenko, Tatsunori Hashimoto,
441 Sergey Levine, Chelsea Finn, and Percy Liang. Extending the WILDS benchmark for
442 unsupervised adaptation. In *International Conference on Learning Representations*,
443 2022.
- 444 [SRKL20] Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. An investigation of
445 why overparameterization exacerbates spurious correlations. In Hal Daumé III and Aarti
446 Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*,
447 volume 119 of *Proceedings of Machine Learning Research*, pages 8346–8356. PMLR,
448 13–18 Jul 2020.
- 449 [Tat17] Rachael Tatman. Gender and dialect bias in youtube’s automatic captions. In *Proceed-*
450 *ings of the First ACL Workshop on Ethics in Natural Language Processing*, pages 53–59,
451 2017.
- 452 [WCHH22] Ke Alexander Wang, Niladri Shekhar Chatterji, Saminul Haque, and Tatsunori
453 Hashimoto. Is importance weighting incompatible with interpolating classifiers? In
454 *International Conference on Learning Representations*, 2022.
- 455 [WGS⁺22] Olivia Wiles, Sven Gowal, Florian Stimberg, Sylvestre-Alvise Rebuffi, Ira Ktena, Krish-
456 namurthy Dj Dvijotham, and Ali Taylan Cemgil. A fine-grained analysis on distribution
457 shift. In *International Conference on Learning Representations*, 2022.
- 458 [XDKR20] Ziyu Xu, Chen Dan, Justin Khim, and Pradeep Ravikumar. Class-weighted classifi-
459 cation: Trade-offs and robust approaches. In Hal Daumé III and Aarti Singh, editors,
460 *Proceedings of the 37th International Conference on Machine Learning*, volume 119

- 461 of *Proceedings of Machine Learning Research*, pages 10544–10554. PMLR, 13–18 Jul
462 2020.
- 463 [XYR21] Da Xu, Yuting Ye, and Chuanwei Ruan. Understanding the role of importance weighting
464 for deep learning. In *International Conference on Learning Representations*, 2021.
- 465 [YCZC20] Han-Jia Ye, Hong-You Chen, De-Chuan Zhan, and Wei-Lun Chao. Identifying and
466 compensating for feature deviation in imbalanced deep learning. *arXiv preprint*
467 *arXiv:2001.01385*, 2020.
- 468 [ZDKR21] Runtian Zhai, Chen Dan, Zico Kolter, and Pradeep Ravikumar. Doro: Distributional
469 and outlier robust optimization. In Marina Meila and Tong Zhang, editors, *Proceedings*
470 *of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings*
471 *of Machine Learning Research*, pages 12345–12355. PMLR, 18–24 Jul 2021.
- 472 [ZDS⁺21] Runtian Zhai, Chen Dan, Arun Suggala, J Zico Kolter, and Pradeep Kumar Raviku-
473 mar. Boosted CVar classification. In *Thirty-Fifth Conference on Neural Information*
474 *Processing Systems*, 2021.

475 Checklist

476 The checklist follows the references. Please read the checklist guidelines carefully for information on
477 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
478 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
479 the appropriate section of your paper or providing a brief inline description. For example:

- 480 • Did you include the license to the code and datasets? **[Yes]** See Section.
- 481 • Did you include the license to the code and datasets? **[No]** The code and the data are
482 proprietary.
- 483 • Did you include the license to the code and datasets? **[N/A]**

484 Please do not modify the questions and only use the provided macros for your answers. Note that the
485 Checklist section does not count towards the page limit. In your paper, please delete this instructions
486 block and only keep the Checklist section heading above along with the questions/answers below.

- 487 1. For all authors...
 - 488 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
489 contributions and scope? **[Yes]**
 - 490 (b) Did you describe the limitations of your work? **[Yes]** See Section 6.2.
 - 491 (c) Did you discuss any potential negative societal impacts of your work? **[No]** Not
492 relevant.
 - 493 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
494 them? **[Yes]**
- 495 2. If you are including theoretical results...
 - 496 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - 497 (b) Did you include complete proofs of all theoretical results? **[Yes]** See Appendix D.
- 498 3. If you ran experiments...
 - 499 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
500 mental results (either in the supplemental material or as a URL)? **[Yes]**
 - 501 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
502 were chosen)? **[Yes]**
 - 503 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
504 ments multiple times)? **[No]** The experiments are only for demonstration.
 - 505 (d) Did you include the total amount of compute and the type of resources used (e.g., type
506 of GPUs, internal cluster, or cloud provider)? **[Yes]** See the code.
- 507 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 508 (a) If your work uses existing assets, did you cite the creators? [N/A]
509 (b) Did you mention the license of the assets? [N/A]
510 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
511
512 (d) Did you discuss whether and how consent was obtained from people whose data you're
513 using/curating? [N/A]
514 (e) Did you discuss whether the data you are using/curating contains personally identifiable
515 information or offensive content? [N/A]
516 5. If you used crowdsourcing or conducted research with human subjects...
517 (a) Did you include the full text of instructions given to participants and screenshots, if
518 applicable? [N/A]
519 (b) Did you describe any potential participant risks, with links to Institutional Review
520 Board (IRB) approvals, if applicable? [N/A]
521 (c) Did you include the estimated hourly wage paid to participants and the total amount
522 spent on participant compensation? [N/A]