# Just Train Twice: Improving Group Robustness without Training Group Information

Evan Zheran Liu [* 1]   Behzad Haghgoo [* 1]   Annie S. Chen [* 1]   Aditi Raghunathan [1]   Pang Wei Koh [1]
Shiori Sagawa [1]   Percy Liang [1]   Chelsea Finn [1]

## Abstract

Standard training via empirical risk minimization (ERM) can produce models that achieve high accuracy on average but low accuracy on certain groups, especially in the presence of spurious correlations between the input and label. Prior approaches that achieve high worst-group accuracy, like group distributionally robust optimization (group DRO) require expensive group annotations for each training point, whereas approaches that do not use such group annotations typically achieve unsatisfactory worst-group accuracy. In this paper, we propose a simple two-stage approach, JTT, that first trains a standard ERM model for several epochs, and then trains a second model that upweights the training examples that the first model misclassified. Intuitively, this upweights examples from groups on which standard ERM models perform poorly, leading to improved worst-group performance. Averaged over four image classification and natural language processing tasks with spurious correlations, JTT closes 75% of the gap in worst-group accuracy between standard ERM and group DRO, while only requiring group annotations on a small validation set in order to tune hyperparameters.

## 1. Introduction

The standard approach of empirical risk minimization (ERM)—training machine learning models to minimize average training loss—can produce models that achieve low test error on average but still incur high error on certain groups of examples (Hovy & Søgaard, 2015; Blodgett et al., 2016; Tatman, 2017; Hashimoto et al., 2018; Duchi et al., 2019). These performance disparities across groups can be especially pronounced in the presence of *spurious correla-*

*tions*. For example, in the task of classifying whether an online comment is toxic, the training data is often biased so that mentions of particular demographics (e.g., certain races or religions) are correlated with toxicity. Models trained via ERM then associate these demographics with toxicity and thus perform poorly on groups of examples in which the correlation does not hold, such as non-toxic comments mentioning a particular demographic (Borkan et al., 2019). Similar performance disparities due to spurious correlations occur in many other applications, including other language tasks, facial recognition, and medical imaging (Gururangan et al., 2018; McCoy et al., 2019; Badgeley et al., 2019; Sagawa et al., 2020a; Oakden-Rayner et al., 2020).

Following prior work, we formalize this setting by considering a set of pre-defined groups (e.g., corresponding to different demographics) and seeking models that have low worst-group error (Sagawa et al., 2020a). Previous approaches typically require annotations of the group membership of each training example (Sagawa et al., 2020a; Goel et al., 2020; Zhang et al., 2020). While these approaches have been successful at improving worst-group performance, the required training group annotations are often expensive to obtain; for example, in the toxicity classification task mentioned above, each comment has to be annotated with all the demographic identities that are mentioned.

In this paper, we propose a simple algorithm, JTT (Just Train Twice), for improving the worst-group error *without training group annotations*, instead only requiring group annotations on a much smaller validation set to tune hyperparameters. JTT is composed of two stages: we first identify training examples that are misclassified by a standard ERM model, and then we train the final model by upweighting the examples identified in the first stage. Intuitively, this procedure exploits the observation that sufficiently-regularized ERM models tend to incur high worst-group training error (and subsequently high worst-group test error). This makes selecting misclassified examples an effective heuristic for identifying examples from groups that ERM models fail on, such as minority groups. Since the final classifier upweights such examples, it performs better on such groups and achieves better minority group performance.

---

[*]Equal contribution   [1]Stanford University. Correspondence to: Evan Zheran Liu <evanliu@cs.stanford.edu>.

We evaluate JTT on two image classification datasets with spurious correlations, Waterbirds (Wah et al., 2011; Sagawa et al., 2020a) and CelebA (Liu et al., 2015) and two natural language processing datasets, MultiNLI (Williams et al., 2018) and CivilComments-WILDS (Borkan et al., 2019; Koh et al., 2021). We use the versions of Waterbirds, CelebA, and MultiNLI from Sagawa et al. (2020a), where in Waterbirds, the label *waterbird* or *landbird* spuriously correlates with water in the background; in CelebA, the label *blond* or *non-blond* spuriously correlates with binary gender; and in MultiNLI, the label spuriously correlates with the presence of negation words. In CivilComments-WILDS, where the input is online comments, the label *toxic, non-toxic* spuriously correlates with the mention of particular demographics, as discussed above. Our method outperforms ERM on all four datasets, with an average worst-group accuracy improvement of 16.2%, while maintaining competitive average accuracy (only 4.2% worse on average). Furthermore, despite having no group annotations during training, JTT closes 75% of the gap between ERM and group DRO, which uses complete group information on the training data.

We then empirically analyze JTT. First, we analyze the examples identified by JTT and show that JTT upweights groups on which standard ERM models perform poorly, e.g., minority groups that do not have the spurious correlation (such as waterbirds on land in the Waterbirds dataset). Second, we show that having validation group annotations is essential for hyperparameter tuning for JTT and other related algorithms.

Finally, we compare JTT with the distributionally robust optimization (DRO) algorithm that minimizes the conditional value at risk (CVaR). CVaR DRO aims to train models that are robust to a wide range of potential distribution shifts by minimizing the worst-case loss over all subsets of the training set of a certain size (Duchi et al., 2019). This objective does not require training group annotations, and it can be optimized by dynamically upweighting training examples with the highest losses in each minibatch (Levy et al., 2020). Though CVaR DRO and JTT share conceptual similarities—they both upweight high loss training points and do not require training group information—the difference is that JTT upweights a static set of examples, while CVaR DRO dynamically re-computes which examples to update. Empirically, we find that JTT empirically substantially outperforms CVaR DRO on worst-group accuracy in the datasets we tested.

## 2. Related Work

In this paper, we focus on group robustness (i.e., training models that obtain good performance on each of a set of predefined groups in the dataset), though other notions of robustness are also studied, such as adversarial examples (Big-

gio et al., 2013; Szegedy et al., 2014) or domain generalization (Blanchard et al., 2011; Muandet et al., 2013). Approaches for group robustness fall into the two main categories we discuss below.

**Robustness using group information.** Several approaches leverage group information during training, either to combat spurious correlations or handle shifts in group proportions between train and test distributions. For example, Mohri et al. (2019); Sagawa et al. (2020a); Zhang et al. (2020) minimize the worst-group loss during training; Goel et al. (2020) synthetically expand the minority groups via generative modeling; Shimodaira (2000); Byrd & Lipton (2019); Sagawa et al. (2020b) reweight or subsample to artificially balance the majority and minority groups; Cao et al. (2019; 2020) impose heavy Lipschitz regularization around minority points. These approaches substantially reduce worst-group error, but obtaining group annotations for the entire training set can be extremely expensive.

Another line of work studies worst-group performance in the context of fairness (Hardt et al., 2016; Woodworth et al., 2017; Pleiss et al., 2017; Agarwal et al., 2018; Khani et al., 2019). While these works also aim to improve the worst-group loss, they explicitly focus on equalizing the loss across all groups.

**Robustness without group information.** We focus on the setting where group annotations are expensive and unavailable on the training data, and potentially only available on a much smaller validation set. Many approaches for this setting fall under the general DRO framework, where models are trained to minimize the worst-case loss across all distributions in a ball around the empirical distribution (Ben-Tal et al., 2013; Lam & Zhou, 2015; Duchi et al., 2016; Namkoong & Duchi, 2017; Oren et al., 2019). Pezeshki et al. (2020) modify the dynamics of stochastic gradient descent to avoid learning spurious correlations. Sohoni et al. (2020) automatically identify groups based by clustering the data points. Kim et al. (2019) propose an auditing scheme that searches for high-loss groups defined by a function within a pre-specified complexity class and postprocess the model to minimize discrepancies identified by the auditor. Khani et al. (2019) minimize the variance in the loss across all data points to encourage lower discrepancy in the losses across all possible groups. Another approach is to directly learn how to reweight the training examples either using small amounts of metadata (Shu et al., 2019) or automatically via meta-learning (Ren et al., 2018).

Most closely related to JTT are several approaches that also train a pair of models, where the performance of the first model is used to help train the second model (Yaghoobzadeh et al., 2019; Utama et al., 2020; Nam et al., 2020). We compare JTT with one such approach, called Learning from

Failure (LfF) (Nam et al., 2020). In LfF, the first model is intentionally biased and tries to identify minority examples where the spurious correlation does not hold. The identified examples are then upweighted while training the second model. This approach interleaves the updates of both models and requires an intentional biasing with the the first model. In contrast, our approach of JTT is simpler, though conceptually similar: we only identify points to upweight once (i.e. no interleaved updates which generally destabilize training), and we just perfom standard ERM with regularization to identify points without any artificial biasing. Empirically, despite its simplicity, JTT performs better than LfF.

Concurrently, Creager et al. (2021) also proposed a method that leverages similar intuition to JTT. This work first uses the errors of a standard ERM model to infer group labels, similar to JTT. Then, they learn a model that is invariant to the predicted labels.

# 3. Preliminaries

## 3.1. Problem Setup

We consider the setting of classifying an input $x \in \mathcal{X}$ as a label $y \in \mathcal{Y}$. We are given $n$ training points $\{(x_1, y_1), \ldots, (x_n, y_n)\}$. Our goal is to learn a model $f_\theta : \mathcal{X} \to \mathcal{Y}$, parameterized by $\theta \in \Theta$. We measure performance across a set of pre-defined groups $\mathcal{G}$. Each point $(x, y)$ belongs to some group $g \in \mathcal{G}$ and we evaluate classifiers on their worst-group error defined as follows:

$$\max_{g \in \mathcal{G}} \mathbb{E}\left[\ell_{0-1}(x, y; \theta) \mid g\right], \quad (1)$$

where $l_{0-1}(x, y; \theta) = \mathbf{1}[f_\theta(x) \neq y]$ is the 0-1 loss.

We are interested in the setting where we do not have group annotations on training points because they are expensive to obtain. Our goal is to achieve good worst-group error at test time without training group annotations. However, we are given a small validation set of $m$ points with group annotations $\{(x_1, y_1, g_1), \ldots, (x_m, y_m, g_m)\}$. These group annotations allow us to compute the worst-group validation error, which we use to tune hyperparameters.

**Groups based on spurious correlations.** In our experiments, we primarily consider the setting where each group $g = (a, y) \in \mathcal{G}$ is defined by the label $y$ and a spurious attribute $a \in \mathcal{A}$ that spuriously correlates with the label (i.e., $\mathcal{G} = \mathcal{A} \times \mathcal{Y}$). Figure 1 illustrates the four groups on the Waterbirds dataset, where the background spuriously correlates with the label.

## 3.2. Comparisons

Here, we describe four other algorithms that we use as comparisons in this paper: (i) Empirical risk minimization
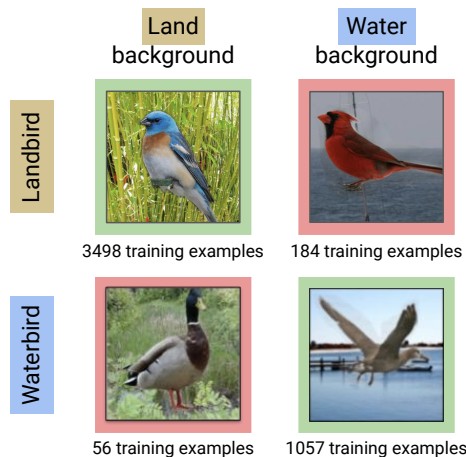


*Figure 1.* The four groups on the Waterbirds dataset, formed by the background spurious attribute and bird type label. Most training examples belong to the groups where the background matches the bird type (highlighted in green), while only a small fraction belong to the groups where the background does not match the bird type (highlighted in red).

(ERM), which is the standard approach for training machine learning models by minimizing the average training loss (Section 3.2.1). (ii) A distributionally robust optimization (DRO) method for minimizing the conditional value at risk (CVaR), which seeks to minimize error over all groups above a certain size (Duchi et al., 2019), and is a natural approach to training models with low worst-group error without group annotations (Section 3.2.2). (iii) Learning from Failure (LfF) (Nam et al., 2020), a recent approach that is conceptually similar to JTT (Section 3.2.3). (iv) Group DRO (Sagawa et al., 2020a), which—unlike all of the preceding methods—uses training group annotations, and can therefore be considered as an oracle method that upper bounds the performance we might expect from methods that do not use training group annotations (Section 3.2.4).

### 3.2.1. EMPIRICAL RISK MINIMIZATION (ERM)

Empirical risk minimization minimizes the average training loss across training points. Given a loss function $\ell(x, y; \theta) : \mathcal{X} \times \mathcal{Y} \times \Theta \to \mathbb{R}_+$ (e.g. cross-entropy loss), ERM minimizes the following objective:

$$J_{\text{ERM}}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(x_i, y_i; \theta). \quad (2)$$

### 3.2.2. DISTRIBUTIONALLY ROBUST OPTIMIZATION OF THE CONDITIONAL VALUE AT RISK (CVAR DRO)

Instead of minimizing the expected loss over the empirical training distribution, distributionally robust learning algorithms define an uncertainty set over distributions that are within some distance of the empirical training distribution,

and then minimize the expected loss over the worst-case distribution in this uncertainty set (Duchi et al., 2019).

In this paper, we study a classic instance of this type of worst-case loss known as the conditional value at risk (CVaR) at level $\alpha \in (0, 1]$, which corresponds to an uncertainty set that contains all $\alpha$-sized subpopulations of the training distribution (Rockafellar & Uryasev, 2000). The idea is that the worst loss over $\alpha$-sized subpopulations upper bounds the worst-group loss over the (unknown) groups in $\mathcal{G}$ when $\alpha$ is close to the size of the smallest group in $\mathcal{G}$.

In practice, we treat $\alpha$ as a hyperparameter. Concretely, for some loss function $\ell(x, y; \theta)$, the CVaR objective can be written as

$$J_{\text{CVaR}}(\theta, \alpha) = \sup_{q \in \Delta^n} \left\{ \sum_{i=1}^{n} q_i \ell(x_i, y_i; \theta) \text{ s.t. } \|q\|_\infty \leq \frac{1}{\alpha n} \right\},$$
(3)

where $\Delta^n$ is the probability simplex in $\mathbb{R}^n$.

Note that the CVaR objective is equivalent to the average loss incurred by the $\alpha$-fraction of training points that have the highest loss.

### 3.2.3. LEARNING FROM FAILURE (LFF)

LfF attempts to automatically upweight examples from challenging groups, such as those where the spurious correlation does not hold. It does this by learning two models $f_B(y \mid x; \theta_B)$ and $f_D(y \mid x; \theta_D)$, parameterized by $\theta_B$ and $\theta_D$.

The first model $f_B$ is trained with ERM using generalized cross-entropy (GCE) loss (Zhang & Sabuncu, 2018):

$$\ell_{\text{GCE}}(x_i, y_i; \theta_B, q) = \frac{1 - f_B(y_i \mid x_i; \theta_B)^q}{q},$$
(4)

where $q \in [0, 1)$ is a hyperparameter. Compared to standard cross-entropy loss, the gradient of GCE loss upweights examples where $f_B(y_i \mid x_i; \theta_B)$ is large, which intentionally biases $f_B$ to perform better on easier examples and poorly on examples group challenging groups.

The second model is also trained with ERM, using cross-entropy loss, where each example $(x_i, y_i)$ is reweighted by a factor of:

$$\mathcal{W}(x_i, y_i) = \frac{\log f_B(y_i \mid x_i)}{\log f_B(y_i \mid x_i) + \log f_D(y_i \mid x_i)}.$$
(5)

The hope is that early in training, $\log f_B(y_i \mid x_i)$ will be smaller than $\log f_D(y_i \mid x_i)$ on the easier examples, which leads to smaller weights on the easier examples and larger weights on the challenging examples.

### 3.2.4. GROUP DISTRIBUTIONALLY ROBUST OPTIMIZATION (GROUP DRO)

Group DRO uses training group annotations to directly minimize the worst-group error on the training set. Our primary focus in this paper is the setting where we do not have access to training group annotations. However, we use group DRO as an oracle method that upper bounds the performance we can expect without any training group annotations. Assume we have access to group annotations on the training data such that the $n$ training points are $\{(x_1, y_1, g_1), \ldots (x_n, y_n, g_n)\}$. For some loss function $\ell(x, y; \theta)$, the group DRO objective can then be written as:

$$J_{\text{group-DRO}}(\theta) = \max_{g \in \mathcal{G}} \frac{1}{n_g} \sum_{i \mid g_i = g} \ell(x_i, y_i; \theta)$$
(6)

where $n_g$ is the number of training points with group $g_i = g$.

## 4. JTT: Just Train Twice

We now present JTT, a simple two-stage approach that does not require group annotations at training time. In the first stage, we train an identification model and select examples with high training loss. Then, in the second stage, we train a final model while upweighting the selected examples.

**Stage 1 (identification).** The key empirical observation that JTT builds on is that sufficiently low complexity ERM models tend to fit groups with easy-to-learn spurious correlations (e.g., landbirds on land and waterbirds on water in the Waterbirds dataset), but not groups that do not exhibit the same correlation (e.g., waterbirds on land) (Sagawa et al., 2020a). We therefore use the simple heuristic of first training an *identification model* $\hat{f}_{\text{id}}$ via ERM and then identifying an *error set* $E$ of training examples that $\hat{f}_{\text{id}}$ misclassifies:

$$E = \{(x_i, y_i) \text{ s.t. } \hat{f}_{\text{id}}(x_i) \neq y_i\}.$$
(7)

**Stage 2 (upweighting).** Next, we train a final model $\hat{f}_{\text{final}}$ by upweighting the points in the error set $E$ identified in step one:

$$J_{\text{up-ERM}}(\theta, E) = \left( \lambda_{\text{up}} \sum_{(x,y) \in E} \ell(x, y; \theta) + \sum_{(x,y) \notin E} \ell(x, y; \theta) \right),$$
(8)

where $\lambda_{\text{up}} \in \mathbb{R}_+$ is a hyperparameter. The hope is that if the examples in the error sets come from challenging groups, such as those where the spurious correlation does not hold, then upweighting them will lead to better worst-group performance.

**Practical implementation.** Overall, training JTT is summarized in Algorithm 1. In practice, to restrict the capacity

---

**Algorithm 1** JTT training

**Input:** Training set $\mathcal{D}$ and hyperparameters $T$ and $\lambda_{\mathrm{up}}$.

**Stage one: identification**

1. Train $\hat{f}_{\mathrm{id}}$ on $\mathcal{D}$ via ERM for $T$ steps (Equation 2).

2. Construct the error set $E$ of training examples misclassified by $\hat{f}_{\mathrm{id}}$ (Equation 7).

**Stage two: upweighting identified points**

3. Construct upsampled dataset $\mathcal{D}_{\mathrm{up}}$ containing examples in the error set $\lambda_{\mathrm{up}}$ times and all other examples once.

4. Train final model $\hat{f}_{\mathrm{final}}$ on $\mathcal{D}_{\mathrm{up}}$ via ERM (Equation 2).

---

of the identification model, we only train it for $T$ steps, where $T$ is a hyperparameter (line 1). This prevents it from potentially overfitting the training data and yielding an empty error set. To implement the upweighted objective (8), we simply upsample the examples from the error set by $\lambda_{\mathrm{up}}$ (line 3) and train the final model on the upsampled data (line 4). Specifically, in each epoch of training, we sample each example from the error set $\lambda_{\mathrm{up}}$ times and all other examples only once.

We tune the algorithm hyperparameters (the number of training epochs $T$ for the identification model $\hat{f}_{\mathrm{id}}$, and the upweight factor $\lambda_{\mathrm{up}}$) and both identification and final model hyperparameters (e.g., the learning rate and $\ell_2$ regularization strength) based on the worst-group error of the final model $\hat{f}_{\mathrm{final}}$ on the validation set. In our experiments, we share the same hyperparameters and architecture between the identification and final models, outside of the early stopping $T$ of the identification model, and we sometimes find it helpful to learn them with different optimizers. Note that setting the upweight factor $\lambda_{\mathrm{up}}$ to 1 recovers ERM, so JTT should perform at least as well as ERM, given a sufficiently large validation set. We describe full training details in Appendix A.

## 5. Experiments

In our experiments, we first demonstrate that JTT substantially improves worst-group performance compared to standard ERM models (Section 5.2). We also show that it recovers a significant fraction of the performance gains yielded by group DRO, which, as discussed in Section 3, is an oracle that relies on group annotations on training examples. We then present empirical analysis of JTT, including the analysis of the error set (Section 5.3), exploration on the role of the validation set (Section 5.4), and comparison with CVaR DRO (Section 5.5).

### 5.1. Setup

We study four datasets in which prior work has observed poor worst-group performance due to spurious correlations

(Figure 2). Full details about these datasets are in Appendix B.

- **Waterbirds** (Wah et al., 2011; Sagawa et al., 2020a): The task is to classify images of birds as "waterbird" or "landbird", and the label is spuriously correlated with the image background, which is either "land" or "water."

- **CelebA** (Liu et al., 2015): We consider the task from Sagawa et al. (2020a) of classifying the hair color of celebrities as "blond" or "not blond." The label is spuriously correlated with gender, which is either "male" or "female."

- **MultiNLI** (Williams et al., 2018): Given a pair of sentences, the task is to classify whether the second sentence is entailed by, neutral with, or contradicts the first sentence. We use the spurious attribute from Sagawa et al. (2020a), which is the presence of negation words in the second sentence; due to the artifacts from the data collection process, contradiction examples often include negation words.

- **CivilComments-WILDS** (Borkan et al., 2019; Koh et al., 2021): The task is to classify whether an online comment is toxic or non-toxic, and the label is spuriously correlated with mentions of certain demographic identities (male, female, White, Black, LGBTQ, Muslim, Christian, and other religion). We use the evaluation metric from Koh et al. (2021), which defines 16 overlapping groups $(a, \textit{toxic})$ and $(a, \textit{non-toxic})$ for each of the above 8 demographic identities $a$, and report the worst-group performance over these groups.

**Points of comparison.** We aim to answer two main questions: (1) How does JTT compare with other approaches that also do not use training group information? (2) How does JTT compare with approaches that *do* use training group information?

To answer the first question, we compare JTT with ERM, CVaR DRO, and a recently proposed approach called Learning from Failure (LfF) (Nam et al., 2020). To answer the second question, we compare JTT with group DRO (Sagawa et al., 2020a), an oracle that uses training group annotations. For details about these approaches, see Section 3. Note that on CivilComments, group DRO cannot be directly applied on the 16 defined groups, since it is not designed for overlapping groups. Instead, our group DRO minimizes worst-group loss over 4 groups $(y, a)$, where the spurious attribute $a$ is a binary indicator of whether any demographic identity is mentioned and the label $y$ is *toxic* or *non-toxic*. We tune the hyperparameters of all approaches based on worst-group performance on a small validation set with group annotations.

*Figure 2.* Examples from the tasks we evaluate on. The spurious attribute $a$ is correlated with the label $y$ on the training data.

| Method | Group labels in train set? | Waterbirds | | CelebA | | MultiNLI | | CivilComments-WILDS | |
|---|---|---|---|---|---|---|---|---|---|
| | | Avg Acc. | Worst-group Acc. | Avg Acc. | Worst-group Acc. | Avg Acc. | Worst-group Acc. | Avg Acc. | Worst-group Acc. |
| ERM | No | 97.3% | 72.6% | 95.6% | 47.2% | 82.4% | 67.9% | 92.6% | 57.4% |
| CVaR DRO (Levy et al., 2020) | No | 96.0% | 75.9% | 82.5% | 64.4% | 82.0% | 68.0% | 92.5% | 60.5% |
| LfF (Nam et al., 2020) | No | 91.2% | 78.0% | 85.1% | 77.2% | 80.8% | 70.2% | 92.5% | 58.8% |
| JTT (Ours) | No | 93.3% | **86.7%** | 88.0% | **81.1%** | 78.6% | **72.6%** | 91.1% | **69.3%** |
| Group DRO (Sagawa et al., 2020a) | Yes | 93.5% | 91.4% | 92.9% | 88.9% | 81.4% | 77.7% | 88.9% | 69.9% |

*Table 1.* Average and worst-group test accuracies of models trained via JTT and baselines. JTT substantially improves worst-group accuracy relative to ERM and CVaR DRO and outperforms LfF (Nam et al., 2020), a recently proposed algorithm for improving worst-group accuracy without group annotations. We also compare with group DRO, an oracle that assumes group annotations. JTT recovers a significant fraction of the gap in worst-group accuracy between ERM and group DRO.

## 5.2. Main Results

Table 4 reports the average and worst-group accuracies of all approaches. Compared to other approaches that do not use training group information, JTT consistently achieves higher worst-group accuracy on all 4 datasets. Additionally, JTT performs well even relative to approaches that use training group information. In particular, JTT recovers a significant portion of the gap in worst-group accuracy between ERM and group DRO, closing 75% of the gap on average. As a caveat, we note that simple label balancing also achieves comparable worst-group accuracy to group DRO on CivilComments.

JTT's worst-group accuracy improvements come at only a modest drop in average accuracy, averaging only 4.2% worse than the highest average accuracy on each dataset. This drop is consistent with Sagawa et al. (2020a), which observes a tradeoff between average and worst-group accuracies.

## 5.3. Error set analysis

We find it surprising that just a small amount of group information on the validation set can allow JTT to achieve high worst-group accuracy with no knowledge of the groups on the training set. We now probe into how JTT achieves such high worst-group accuracy. In order to perform this analysis, we use the group annotations on the training data to closely examine what examples are upweighted in the

| Dataset | Worst-group Recall | Worst-group Precision | Worst-group Empirical Rate |
|---|---|---|---|
| Waterbirds | 87.5% | 19.1% | 1.2% |
| CelebA | 94.7% | 9.4% | 0.9% |
| MultiNLI | 67.1% | 2.2% | 1.0% |
| CivilComments | 96.9% | 7.8% | 0.9% |

*Table 2.* The precision and recall of the worst-group examples (i.e., the group with lowest validation accuracy) belonging to JTT's error set. The error set includes a high fraction of the worst-group examples and includes them at a much higher rate than they occur in the training data.

| Group | Enrichment | ERM test acc. |
|---|---|---|
| (land background, waterbird) | 15.92x | 72.6% |
| (water background, landbird) | 6.97x | 73.3% |
| (water background, waterbird) | 2.40x | 96.3% |
| (land background, landbird) | 0.02x | 99.3% |

*Table 3.* Waterbirds error set breakdowns.

error set identified in the first step of JTT, though we don't use such training group annotations for training JTT.

To start, we define the worst group as the group on which the standard ERM model achieves the lowest test accuracy, when tuned for worst-group validation accuracy. We analyze how well the error set captures this worst group. To do this, we measure *precision*, the fraction of examples in the error set that belong to the worst group, *recall*, the fraction of the worst group examples that are included in the error set,

| Group | Enrichment | ERM test acc. |
|---|---|---|
| (male, blond) | 10.44x | 47.2% |
| (female, blond) | 5.42x | 89.1% |
| (male, non-blond) | 0.32x | 99.3% |
| (female, non-blond) | 0.01x | 95.1% |

*Table 4.* CelebA error set breakdowns.

| Group | Enrichment | ERM test acc. |
|---|---|---|
| (negation, neutral) | 2.2x | 67.9% |
| (no negation, contradiction) | 1.35x | 77.0% |
| (negation, entailment) | 1.14x | 80.4% |
| (no negation, neutral) | 1.07x | 81.8% |
| (no negation, entailment) | 0.73x | 86.1% |
| (negation, contradiction) | 0.19x | 94.5% |

*Table 5.* MultiNLI error set breakdowns.

and the *empirical rate*, the rate at which the worst group examples appear in the training data.

As reported in Table 2, we observe that the error set contains worst-group examples at a much higher rate (precision) than they appear in the training dataset (empirical rate). Worst-group examples appear in the error set 2.2x to 15.9x more frequently in the error set than in the training data, across the 4 datasets. In other words, the worst group is significantly *enriched* in the error set compared to the training dataset, which may explain why JTT has much better worst-group performance over ERM. Additionally, the error set has high worst-group recall, ranging from 67.1% to 96.9% and averaging to 86.4% across the 4 datasets. Together, these results indicate that the worst group is included in the error set at relatively high both precision and recall.

Empirically, ERM performs poorly on several groups, not just on a single worst group. We therefore next examine what other groups the examples in the error set belong to, beyond the worst group. For each group, we compute two metrics: (i) *enrichment*, defined as how much more frequently examples from a group appear in the error set than in the training data (i.e., the precision of the group divided by the empirical rate of the group); (ii) the test accuracy that ERM achieves on this group, when tuned for worst-group validation accuracy.

Tables 3 to 5 and Table 12 in Appendix B.4 report these results for Waterbirds, CelebA, MultiNLI, and CivilComments respectively. We observe that the enrichment roughly inversely correlates with ERM's test accuracy on that group: examples from low performance groups are included at high rates in the error set relative to the empirical rate. This may help JTT perform better across all groups that ERM performs poorly on, which in turn improves worst-group accuracy.

| | Worst-group test acc. |
|---|---|
| Standard error set | 86.7% |
| No waterbirds on water backgrounds | 80.7% |
| Swap error set examples | 86% |

*Table 6.* Worst-group test accuracies with 3 variants of the error set on Waterbirds: (i) standard unchanged error set; (ii) removing all waterbird on water background examples; (iii) swapping each error set example with a random example from the same group.

Finally, we note that while the groups with high enrichment often correspond to groups where the spurious correlation does not hold, this is not always the case. In particular, the waterbird on water background group in Waterbirds and the blonde female group in CelebA have high enrichments, even though the spurious correlation holds in these groups. We hypothesize that this occurs due to label imbalance, since the waterbird label and blonde label are relatively rare and appear only in 23% and 15% of the training examples, respectively. Empirically, upweighting examples from these groups is indeed important for JTT's worst-group test accuracy. When we remove all waterbird on water background examples from the error set, JTT's worst-group test accuracy drops by 6%. However, while the group composition of the error set (i.e., the fraction of the error set in each group) is important, the exact examples inside the error set do not seem to matter. When we replace each error set example with another example from the same group on Waterbirds, worst-group test accuracy drops by only 0.7%. These two results are shown in Table 6.

### 5.4. Hyperparameter Tuning and the Role of the Validation Set

In all of our experiments, we tune the algorithm and model hyperparameters based on the worst-group accuracy on the validation set. In general, across all methods, we found hyperparameter tuning in this fashion to be critical. Table 7 shows that even for CVaR DRO, LfF, and JTT, which all try to improve worst-group accuracy without relying on training group annotations, the worst-group test accuracies on Waterbirds and CelebA plummet when the hyperparameters are tuned for average accuracy on the validation set, instead of worst-group accuracy on the validation set. Importantly, this means that even though these methods do not require training group annotations, they still require *validation group annotations* in order to have high worst-group test accuracy. Existing methods for improving worst-group accuracy generally rely on some form of this assumption (either by explicitly tuning for validation group accuracy, or by assuming access to a validation set that is balanced by groups). Removing this reliance is an important direction for future work; we discuss this further in Section 6.

| | Waterbirds worst-group test acc. | | CelebA worst-group test acc. | |
|---|---|---|---|---|
| | Tuned for average | Tuned for worst-group | Tuned for average | Tuned for worst-group |
| CVaR DRO (Levy et al., 2020) | 62.0% | 75.9% | 36.1% | 64.4% |
| LfF (Nam et al., 2020) | 44.1% | 78.0% | 24.4% | 77.2% |
| JTT (Ours) | 62.5% | 86.7% | 40.6% | 81.1% |

*Table 7.* Across the methods that do not use training group annotations, worst-group test performance is significantly higher when hyperparameters are tuned for worst-group validation accuracy instead of average validation accuracy. This shows that for these methods, it is still critical to have validation group annotations.

| | Waterbirds | | | | CelebA | | |
|---|---|---|---|---|---|---|---|
| 1x | $\frac{1}{5}$x | $\frac{1}{10}$x | $\frac{1}{20}$x | 1x | $\frac{1}{5}$x | $\frac{1}{10}$x | $\frac{1}{20}$x |
| **86.7%** | **84.0%** | **86.9%** | 76.0% | **81.1%** | **81.1%** | **81.1%** | **82.2%** |

*Table 8.* JTT retains high test worst-group accuracy on Waterbirds and CelebA when the validation set size is reduced to $\frac{1}{10}$x, though performance drops at $\frac{1}{20}$x on CelebA.

The sensitivity of CVaR DRO, LfF, and JTT to hyperparameters is consistent with prior work showing that reweighting methods like these require appropriate capacity control via techniques like early stopping or strong $\ell_2$ regularization (Byrd & Lipton, 2019; Sagawa et al., 2020a). These parameters—e.g., how many epochs to train for, or what $\ell_2$ regularization strength to use—generally need to be set by tuning directly for worst-group accuracy; we found that tuning for average accuracy typically results in models with lower $\ell_2$ regularization strength and that were trained for more epochs.

Compared to ERM, JTT has two additional hyperparameters: the number of epochs to train the identification model $T$ and the upweight factor $\lambda_{up}$. As an illustration of the sensitivity to hyperparameters, Figure 3 shows how the worst-group accuracy of JTT's final model changes as we vary $T$ between 20 and 100 epochs on Waterbirds. Worst-group accuracy is high when $T$ is between 40 and 60, but drops when $T$ is too small or too large.

**Reducing the size of the validation set.** In the main experiments in Section 5.2, we use the default validation sets provided with each of the datasets. Using group annotations on these validation sets is already cheaper than training group annotations, as these default validation sets are 2–10x smaller than their corresponding training sets. However, we additionally test if JTT can continue to achieve high worst-group performance using even smaller validation sets to further reduce the cost of obtaining group annotations on those sets. On Waterbirds and CelebA, we reduce the validation set size by a factor of 1x (no reduction), $\frac{1}{5}$x, $\frac{1}{10}$x, and $\frac{1}{20}$x and tune JTT's hyperparameters based on worst-group accuracy on the reduced validation set. We find that JTT continues to achieve high worst-group accuracy, even when reducing the validation set size by $\frac{1}{10}$x and $\frac{1}{20}$x, amount-
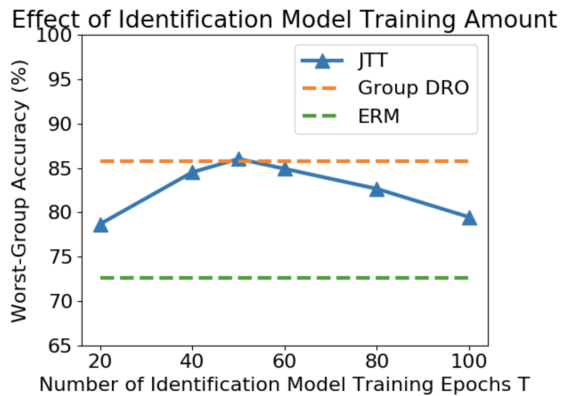


*Figure 3.* Effect of number of epochs of identification model training in JTT on Waterbirds. Worst-group test accuracy is high for $T$ between 40 and 60 epochs, but degrades when $T$ becomes too small or large, which results in less informative error sets.

ing to only 119 and 993 total examples, on Waterbirds and CelebA respectively.

### 5.5. Comparison with CVaR DRO

In this section, we explore the relation between JTT and CVaR DRO. Recall that the CVaR objective in Equation 3 is the average loss incurred by the $\alpha$-fraction of training examples with the highest loss. We can view minimizing this objective as upweighting this $\alpha$-fraction of examples while ignoring the remaining examples. In this way, JTT is conceptually similar to CVaR DRO: both upweight training points with high loss, without requiring group annotations of training points. However, their empirical performance is widely different: CVaR DRO offers only small gains in worst-group accuracy over ERM, while JTT offers substantial gains. One key difference is that in JTT, the set of points that get upweighted $E$ is computed once during stage 1, and then held fixed. In contrast, minimizing the CVaR objective involves dynamically computing the $\alpha$-subset of points with the highest loss at each step, upweighting them and updating the model, and then repeating to update the $\alpha$-subset. As we show next, ablating this key difference from JTT substantially degrades worst-group performance.

| | Waterbirds | |
|---|---|---|
| Epochs per update ($K$) | Average Acc. | Worst-group Acc. |
| 10 | 89.2% | 77.1% |
| 20 | 86.3% | 72.1% |
| 30 | 91.4% | 86.8% |
| 50 | 93.1% | 88.6% |

*Table 9.* Effect of dynamically computing JTT's error set on Waterbirds. We first train the identification model for $T = 50$ epochs, as usual. Then, we dynamically update the error set using the final model after every $K$ epochs of training the final model. Lower values of $K$ have significantly lower accuracies.

**Dynamically computing the error set in JTT lowers accuracy.** We start by observing that the performance of JTT drops when we dynamically recompute the error set $E$, instead of only computing $E$ once using the identification model. Concretely, we study a variant of JTT on the Waterbirds dataset: as usual, we first train an identification model for $T = 50$ epochs, but then while training the final model, every $K$ epochs, we dynamically update the error set $E$ as the errors of the final model over the training set. Setting $K$ to be $\infty$—which means that we only compute the error set $E$ once after training the identification model for $T$ epochs—recovers standard JTT. On the other hand, lowering $K$ makes the algorithm more similar to minimizing CVaR, since this more frequently updates the upweighted set to be the examples with higher loss under the current model, instead of the examples with higher loss under the static identification model.

Table 9 shows the results as we vary $K$ between 10, 20, 30, and 50 epochs on Waterbirds, re-tuning all hyperparameters for each value of $K$. At high values of $K$, where the error set remains fixed for many epochs, both average and worst-group accuracies are high. However, as $K$ decreases, the average and worst-group accuracies drop. These results show that at least on Waterbirds, holding the error set fixed appears to be critical for JTT.

**Which examples does CVaR upweight?** The analysis above suggests that the relatively poor worst-group performance of CVaR DRO might stem from how it dynamically computes which examples to upweight. We further study the behavior of CVaR DRO by analyzing the examples that it upweights. Concretely, throughout CVaR DRO training, we periodically identify the $\alpha$-fraction of training examples with the highest loss and measure the worst-group precision and recall, where the worst group is defined as the group on which ERM achieves the lowest test accuracy.

Figure 4 shows the results using the value of $\alpha$ achieving the highest worst-group validation accuracy: $\alpha = 0.2$ on Waterbirds, $\alpha = 0.00852$ on CelebA, and $\alpha = 0.5$ on MultiNLI. On Waterbirds, the worst-group examples (which comprise

approximately 1% of the training set) make up 19% of the error set for JTT, whereas they oscillate between 1% and 10% of the worst-$\alpha$ fraction for CVaR DRO. As a result, JTT consistently upweights nearly 90% of the worst-group examples, whereas CVaR DRO oscillates between upweighting the worst group and the other groups, upweighting as little as 20% of the examples at some points during training. On CelebA, CVaR DRO upweights the worst-group examples with slightly higher precision than JTT, but $\alpha$ is much smaller than the size of the error set; as a result, JTT upweights nearly 95% of the worst-group examples, whereas CVaR DRO only upweights 13% of them. On MultiNLI, the worst group steadily gets less and less upweighted for CVaR DRO, whereas JTT upweights it at a higher rate, though it still only comprises 2% of the error set for JTT.

These results suggest that the CVaR objective might be overly conservative where the $\alpha$-fraction of examples with highest loss often include many examples from other groups. Furthermore, the set of examples varies widely across different iterations of training. In contrast, JTT upweights a fixed set of points. Empirically, we find that this allows JTT to successfully use the worst-group accuracy on a small validation set to identify error sets that improve accuracy on groups we care about.

## 6. Discussion

In this work, we presented Just Train Twice (JTT), a simple algorithm that substantially improves worst-group performance without requiring expensive group annotations during training. We conclude by discussing several directions for future work.

First, a better theoretical understanding of when and why JTT works would help us to refine and further develop methods for training models that are less susceptible to spurious correlations. For example, it would be useful to understand why early-stopped ERM models (as in the identification models used by JTT) seem to consistently latch onto the spurious correlations in our datasets, and why it seems to be important to fix the upweighted set instead of dynamically recomputing it, as in CVaR DRO.

Second, JTT and many prior methods on robustness without group information all rely on a validation set that is representative of the distribution shift or annotated with group information. While these annotations are significantly cheaper that labeling the entire training set, it still requires the practitioner to be aware of any spurious correlations and define groups accordingly. Doing so may be notably difficult in real-world applications. Therefore, this leaves open the question of whether methods can perform well with mis-specified groups or no group annotations whatsoever.

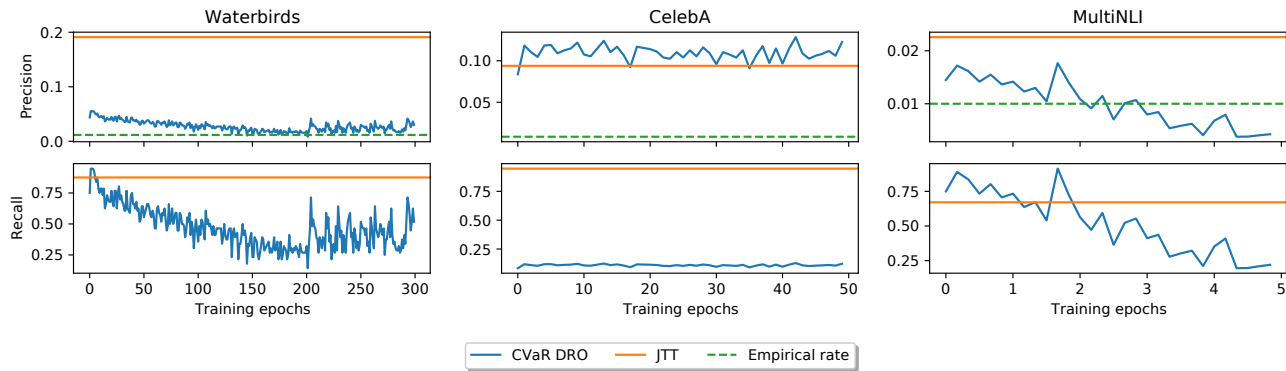Finally, while our experiments focus on group robustness in

*Figure 4.* The composition of the CVaR set (the $\alpha$-fraction of training examples with the highest loss) as training progresses for CVaR DRO models. In these plots, the worst group is defined as the group with the lowest test accuracy under the ERM model. For each dataset, the top plot shows the worst-group precision: the fraction of the CVaR set that belongs to the worst group (blue), with the analogous fraction of the JTT error set (orange) and the overall training data (green) provided for comparison. The bottom plot shows the worst-group recall: the fraction of total worst-group examples that are in the respective sets. For Waterbirds and MultiNLI, the CVaR set is less enriched for the worst group compared to JTT. For CelebA, it is slightly more enriched, but $\alpha$ is much smaller than the size of the JTT error set, so it only upweights a small fraction of the worst group.

the presence of spurious correlations, JTT is not specifically tailored to spurious correlations. Given JTT's simplicity, it would be straightforward to experiment with JTT to see if it might improve performance under different types of distribution shifts, such as in domain generalization settings (Blanchard et al., 2011; Muandet et al., 2013).

**Reproducibility.** Our code is publicly available at https://github.com/anniesch/jtt.

## References

Agarwal, A., Beygelzimer, A., Dudik, M., Langford, J., and Wallach, H. A reductions approach to fair classification. In *International Conference on Machine Learning (ICML)*, pp. 60–69, 2018.

Badgeley, M. A., Zech, J. R., Oakden-Rayner, L., Glicksberg, B. S., Liu, M., Gale, W., McConnell, M. V., Percha, B., Snyder, T. M., and Dudley, J. T. Deep learning predicts hip fracture using confounding patient and healthcare variables. *npj Digital Medicine*, 2, 2019.

Ben-Tal, A., den Hertog, D., Waegenaere, A. D., Melenberg,

B., and Rennen, G. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59:341–357, 2013.

Biggio, B., Corona, I., Maiorca, D., Nelson, B., Šrndić, N., Laskov, P., Giacinto, G., and Roli, F. Evasion attacks against machine learning at test time. In *Joint European conference on machine learning and knowledge discovery in databases*, pp. 387–402, 2013.

Blanchard, G., Lee, G., and Scott, C. Generalizing from several related classification tasks to a new unlabeled sample. In *Advances in neural information processing systems*, pp. 2178–2186, 2011.

Blodgett, S. L., Green, L., and O'Connor, B. Demographic dialectal variation in social media: A case study of African-American English. In *Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1119–1130, 2016.

Borkan, D., Dixon, L., Sorensen, J., Thain, N., and Vasserman, L. Nuanced metrics for measuring unintended bias with real data for text classification. In *World Wide Web (WWW)*, pp. 491–500, 2019.

Byrd, J. and Lipton, Z. What is the effect of importance weighting in deep learning? In *International Conference on Machine Learning (ICML)*, pp. 872–881, 2019.

Cao, K., Wei, C., Gaidon, A., Arechiga, N., and Ma, T. Learning imbalanced datasets with label-distribution-aware margin loss. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

Cao, K., Chen, Y., Lu, J., Arechiga, N., Gaidon, A., and Ma, T. Heteroskedastic and imbalanced deep learning with adaptive regularization. *arXiv preprint arXiv:2006.15766*, 2020.

Creager, E., Jacobsen, J.-H., and Zemel, R. Environment inference for invariant learning. In *International Conference on Machine Learning (ICML)*, pp. 2189–2200, 2021.

Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Association for Computational Linguistics (ACL)*, pp. 4171–4186, 2019.

Duchi, J., Glynn, P., and Namkoong, H. Statistics of robust optimization: A generalized empirical likelihood approach. *arXiv*, 2016.

Duchi, J., Hashimoto, T., and Namkoong, H. Distributionally robust losses against mixture covariate shifts. https://cs.stanford.edu/~thashim/assets/publications/condrisk.pdf, 2019.

Goel, K., Gu, A., Li, Y., and Ré, C. Model patching: Closing the subgroup performance gap with data augmentation. *arXiv preprint arXiv:2008.06775*, 2020.

Gururangan, S., Swayamdipta, S., Levy, O., Schwartz, R., Bowman, S., and Smith, N. A. Annotation artifacts in natural language inference data. In *Association for Computational Linguistics (ACL)*, pp. 107–112, 2018.

Hardt, M., Price, E., and Srebo, N. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 3315–3323, 2016.

Hashimoto, T. B., Srivastava, M., Namkoong, H., and Liang, P. Fairness without demographics in repeated loss minimization. In *International Conference on Machine Learning (ICML)*, 2018.

He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Computer Vision and Pattern Recognition (CVPR)*, 2016.

Hovy, D. and Søgaard, A. Tagging performance correlates with age. In *Association for Computational Linguistics (ACL)*, pp. 483–488, 2015.

Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML)*, pp. 448–456, 2015.

Khani, F., Raghunathan, A., and Liang, P. Maximum weighted loss discrepancy. *arXiv preprint arXiv:1906.03518*, 2019.

Kim, M. P., Ghorbani, A., and Zou, J. Multiaccuracy: Black-box post-processing for fairness in classification. In *Association for the Advancement of Artificial Intelligence (AAAI)*, pp. 247–254, 2019.

Koh, P. W., Sagawa, S., Marklund, H., Xie, S. M., Zhang, M., Balsubramani, A., Hu, W., Yasunaga, M., Phillips, R. L., Gao, I., Lee, T., David, E., Stavness, I., Guo, W., Earnshaw, B. A., Haque, I. S., Beery, S., Leskovec, J., Kundaje, A., Pierson, E., Levine, S., Finn, C., and Liang, P. WILDS: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning (ICML)*, 2021.

Lam, H. and Zhou, E. Quantifying input uncertainty in stochastic optimization. In *2015 Winter Simulation Conference*, 2015.

Levy, D., Carmon, Y., Duchi, J. C., and Sidford, A. Large-scale methods for distributionally robust optimization. *arXiv preprint arXiv:2010.05893*, 2020.

Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning face attributes in the wild. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3730–3738, 2015.

Loshchilov, I. and Hutter, F. Sgdr: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations (ICLR)*, 2017.

McCoy, R. T., Pavlick, E., and Linzen, T. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. In *Association for Computational Linguistics (ACL)*, 2019.

Mohri, M., Sivek, G., and Suresh, A. T. Agnostic federated learning. In *International Conference on Machine Learning (ICML)*, pp. 4615–4625, 2019.

Muandet, K., Balduzzi, D., and Schölkopf, B. Domain generalization via invariant feature representation. In *International Conference on Machine Learning (ICML)*, pp. 10–18, 2013.

Nam, J., Cha, H., Ahn, S., Lee, J., and Shin, J. Learning from failure: Training debiased classifier from biased classifier. *arXiv preprint arXiv:2007.02561*, 2020.

Namkoong, H. and Duchi, J. Variance regularization with convex objectives. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

Oakden-Rayner, L., Dunnmon, J., Carneiro, G., and Ré, C. Hidden stratification causes clinically meaningful failures in machine learning for medical imaging. In *Proceedings of the ACM Conference on Health, Inference, and Learning*, pp. 151–159, 2020.

Oren, Y., Sagawa, S., Hashimoto, T., and Liang, P. Distributionally robust language modeling. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2019.

Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. Automatic differentiation in pytorch, 2017.

Pezeshki, M., Kaba, S.-O., Bengio, Y., Courville, A., Precup, D., and Lajoie, G. Gradient starvation: A learning proclivity in neural networks. *arXiv preprint arXiv:2011.09468*, 2020.

Pleiss, G., Raghavan, M., Wu, F., Kleinberg, J., and Weinberger, K. Q. On fairness and calibration. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 5684–5693, 2017.

Ren, M., Zeng, W., Yang, B., and Urtasun, R. Learning to reweight examples for robust deep learning. In *International Conference on Machine Learning (ICML)*, 2018.

Rockafellar, R. T. and Uryasev, S. Optimization of conditional value-at-risk. *Journal of Risk*, 2:21–41, 2000.

Sagawa, S., Koh, P. W., Hashimoto, T. B., and Liang, P. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In *International Conference on Learning Representations (ICLR)*, 2020a.

Sagawa, S., Raghunathan, A., Koh, P. W., and Liang, P. An investigation of why overparameterization exacerbates spurious correlations. In *International Conference on Machine Learning (ICML)*, 2020b.

Shimodaira, H. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*, 90:227–244, 2000.

Shu, J., Xie, Q., Yi, L., Zhao, Q., Zhou, S., Xu, Z., and Meng, D. Meta-Weight-Net: Learning an explicit mapping for sample weighting. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 1919–1930, 2019.

Sohoni, N. S., Dunnmon, J. A., Angus, G., Gu, A., and Ré, C. No subclass left behind: Fine-grained robustness in coarse-grained classification problems. *arXiv preprint arXiv:2011.12945*, 2020.

Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. In *International Conference on Learning Representations (ICLR)*, 2014.

Tatman, R. Gender and dialect bias in YouTube's automatic captions. In *Workshop on Ethics in Natural Langauge Processing*, volume 1, pp. 53–59, 2017.

Utama, P. A., Moosavi, N. S., and Gurevych, I. Towards debiasing NLU models from unknown biases. *arXiv preprint arXiv:2009.12303*, 2020.

Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. The Caltech-UCSD Birds-200-2011 dataset. Technical report, California Institute of Technology, 2011.

Williams, A., Nangia, N., and Bowman, S. A broad-coverage challenge corpus for sentence understanding through inference. In *Association for Computational Linguistics (ACL)*, pp. 1112–1122, 2018.

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., and Brew, J. HuggingFace's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.

Woodworth, B., Gunasekar, S., Ohannessian, M. I., and Srebro, N. Learning non-discriminatory predictors. In *Conference on Learning Theory (COLT)*, pp. 1920–1953, 2017.

Yaghoobzadeh, Y., Mehri, S., Tachet, R., Hazen, T. J., and Sordoni, A. Increasing robustness to spurious correlations using forgettable examples. *arXiv preprint arXiv:1911.03861*, 2019.

Zhang, J., Menon, A., Veit, A., Bhojanapalli, S., Kumar, S., and Sra, S. Coping with label shift via distributionally robust optimisation. *arXiv preprint arXiv:2010.12230*, 2020.

Zhang, Z. and Sabuncu, M. R. Generalized cross entropy loss for training deep neural networks with noisy labels. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.

Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., and Torralba, A. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(6):1452–1464, 2017.

# A. Training Details

In this section, we detail the model architectures and hyperparameters used by each approach. Within each dataset, we used the same model architecture across all approaches: ResNet-50 (He et al., 2016) for Waterbirds and CelebA, and BERT for MultiNLI and CivilComments (Devlin et al., 2019). For ResNet-50, we used the PyTorch (Paszke et al., 2017) implementation of ResNet-50, starting from ImageNet-pretrained weights. For BERT, we used the the HuggingFace implementation (Wolf et al., 2019) of BERT, also starting from pretrained weights.

We use the LfF implementation released by Nam et al. (2020). We use the group DRO and ERM implementations released by Sagawa et al. (2020a) and also implement CVaR DRO and JTT on top of this code base, with the CVaR DRO implementation adapted from Levy et al. (2020). For the group DRO experiments on Waterbirds, CelebA, and MultiNLI, we directly use the reported performance numbers from Sagawa et al. (2020a). We note that these numbers utilize group-specific loss adjustments that encourage the model to attain lower training losses on smaller groups, which was shown to improve worst-group generalization. We train our own group DRO model on CivilComments-WILDS as it was not included in Sagawa et al. (2020a); for this, we did not implement these group adjustments. We train our own models for all other algorithms.

For all approaches, we tune all hyperparameters as well as early stop based on highest worst-group accuracy on the validation set. On top of the hyperparameters shared between all algorithms (e.g., learning rate, $\ell_2$ regularization), which we detail for each dataset below, CVaR DRO, JTT, and LfF have additional hyperparameters that we tuned separately for each dataset:

- For **CVaR DRO**, we tune the size of the worst-case subpopulation $\alpha \in \{0.1, 0.2, 0.5\}$. For CelebA, we additionally tried $\alpha = \frac{\text{\# smallest group examples}}{\text{\# training examples}} = 0.00852$.

- For **LfF**, we tune the hyperparameter $q$ by grid searching over $q \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$. For Civil-Comments, we additionally sample two values log-uniformly from $(0, 0.1]$. This hyperparameter was not tuned in the experiments in Nam et al. (2020).

- For **JTT**, we additionally tune the number of epochs of training the identification model $T$ and the upsampling factor $\lambda_{\text{up}}$. While developing JTT, we tried the following values of $T$ and $\lambda_{\text{up}}$ without looking at test results, though not necessarily all combinations: $T \in \{1, 2, 40, 50, 60\}$ and $\lambda_{\text{up}} \in \{5, 10, 20, 30, 40, 50, 100, \frac{|\text{training set}|}{|\text{error set}|}\}$. For the final experiments we tune over $\lambda_{\text{up}} \in \{20, 50, 100\}$ for the vision datasets (Waterbirds and CelebA) and $\lambda_{\text{up}} \in$

| Learning rate | $\ell_2$ regularization strength |
|---|---|
| 1e-3 | 1e-4 |
| 1e-4 | 1e-1 |
| 1e-5 | 1 |

*Table 10.* Learning rates and $\ell_2$ regularization strengths for Waterbirds.

| Learning rate | $\ell_2$ regularization strength |
|---|---|
| 1e-4 | 1e-4 |
| 1e-4 | 1e-2 |
| 1e-5 | 1e-1 |

*Table 11.* Learning rates and $\ell_2$ regularization strengths for CelebA.

$\{4, 5, 6\}$ for the NLP datasets (MultiNLI and Civil-Comments). Additionally, Waterbirds requires more training epochs than the others, due to its much smaller training set size, so we tune over $T \in \{40, 50, 60\}$ for Waterbirds and $T \in \{1, 2\}$ for all other datasets.

ERM has no additional algorithm-specific hyperparameters. For group DRO, we fixed the step size $\eta_q$ for updating group weights to its default value of 0.01 from Sagawa et al. (2020a), without tuning.

In general, for JTT, we fixed the initialization model and the final model to share the same hyperparameters, with two exceptions. First, the initialization model is trained only for $T$ epochs, whereas the final model is trained for longer; exact values vary by dataset. Second, for BERT, we found that it was helpful for JTT to be able to choose different optimizers for the initialization model and final model. Specifically, for the ResNet-50 models, we used SGD with momentum 0.9 and no learning rate scheduler or gradient clipping. For BERT, we additionally considered the standard AdamW optimizer Loshchilov & Hutter (2017) with a linearly-decaying learning rate and gradient clipping (setting the max $\ell_2$-norm of the gradients to 1), but we set this as a hyperparameter and allowed the initialization and final models to independently be optimized by SGD or AdamW.

**Waterbirds.** All approaches are optimized for up to 300 epochs with batch size 64, using batch normalization (Ioffe & Szegedy, 2015), and no data augmentation. We chose this smaller batch size (compared to the batch size of 128 used in Sagawa et al. (2020a)) for computational convenience. All approaches are optimized with stochastic gradient descent (SGD) with momentum 0.9.

For all approaches, we tune over the 3 pairs of learning rate and $\ell_2$ regularization strength used by Sagawa et al. (2020a) detailed in Table 10. For ERM and LfF, this yields learning rate 1e-3 and $\ell_2$ regularization 1e-4. For CVaR DRO, this

yields learning rate 1e-4 and $\ell_2$ regularization 1e-1. For JTT and group DRO, this yields learning rate 1e-5 and $\ell_2$ regularization 1.

Our grid search over $\alpha$ for CVaR DRO yields $\alpha = 0.2$. Our grid search over $q$ for LfF yields $q = 0.5$. Our grid search over $T$ and $\lambda_{\text{up}}$ for JTT yields $T = 60$ epochs and $\lambda_{\text{up}} = 100$.

**CelebA.** We train all approaches for up to 50 epochs with batch size 128, using batch normalization and no data augmentation. Like in Waterbirds, we optimize all approaches with SGD with momentum 0.9.

Also like in Waterbirds, we tune over the 3 pairs of learning rate and $\ell_2$ regularization strength used by Sagawa et al. (2020a), detailed in Table 11. For ERM, this yields learning rate 1e-4 and $\ell_2$ regularization 1e-4. For LfF, this yields learning rate 1e-4 and $\ell_2$ regularization 1e-2. For all other approaches, this yields learning rate 1e-5 and $\ell_2$ regularization 1e-1.

Our grid search over $\alpha$ for CVaR DRO yields $\alpha = 0.00852$. Our grid search over $q$ for LfF yields $q = 0.5$. Our grid search over $T$ and $\lambda_{\text{up}}$ for JTT yields $T = 1$ epoch and $\lambda_{\text{up}} = 50$.

**MultiNLI.** We train each approach for up to 5 epochs with default tokenization, dropout, batch size 32, no $\ell_2$-regularization, and an initial learning rate of 0.00002.

JTT achieves the highest validation worst-group accuracy using SGD optimization without clipping for the initial model, and using the AdamW optimizer with clipping for the final model. All other approaches achieve highest validation worst-group accuracy using AdamW with clipping. Our grid search over $\alpha$ for CVaR DRO yields $\alpha = 0.5$. Our grid search over $q$ for LfF yields $q = 0.1$. Our grid search over $T$ and $\lambda_{\text{up}}$ for JTT yields $T = 2$ epochs and $\lambda_{\text{up}} = 6$.

**CivilComments-WILDS.** All approaches use the details from Koh et al. (2021). We capped the number of tokens per example at 300 and used an initial learning rate of 0.00001. We train all approaches for up to 5 epochs with batch size 16 and $\ell_2$ regularization strength of 0.01.

JTT achieves the highest validation worst-group accuracy using SGD optimization without clipping for the initial model, and using the AdamW optimizer with clipping for the final model. All other approaches achieve highest validation worst-group accuracy using AdamW with clipping. Our grid search over $\alpha$ for CVaR DRO yielded $\alpha = 0.5$. Our grid search over $q$ for LfF yielded $q = 0.00001$. Our grid search over $T$ and $\lambda_{\text{up}}$ for JTT yields $T = 2$ epochs and $\lambda_{\text{up}} = 6$.

We also note that our group DRO approach uses a differ-ent spurious attribute compared to the group DRO results reported in Koh et al. (2021). Our group DRO uses the spurious attribute of any demographic identity being mentioned, while the one in Koh et al. (2021) uses only mentions of the Black demographic. Both perform similarly: ours achieves 0.3% lower worst-group accuracy, but 0.5% higher average accuracy.

## B. Dataset Details

### B.1. Waterbirds

We use the Waterbirds dataset introduced by Sagawa et al. (2020a), which is constructed by cropping out images of birds from the CUB dataset (Wah et al., 2011) and pasting them on backgrounds from the Places dataset (Zhou et al., 2017). In this dataset, images of seabirds (albatross, auklet, cormorant, frigatebird, fulmar, gull, jaeger, kittiwake, pelican, puffin, or tern) and waterfowl (gadwall, grebe, mallard, merganser, guillemot, or Pacific loon) are labeled as *waterbirds*, and all other birds are labeled as *landbirds*.

Backgrounds from the *ocean* and *natural lake* categories in the Places dataset are considered to have spurious attribute $a = water\ background$, while backgrounds from the *bamboo forest* or *broadleaf forest* categories are considered to have spurious attribute $a = land\ background$.

There are two minority groups: (land background, waterbird) and (water background, landbird); and two majority groups: (land background, landbird) and (water background, waterbird). We use the same training / valid / test splits from Sagawa et al. (2020a). In the training data, 95% of the waterbirds appear on water backgrounds, and 95% of the landbirds appear on land backgrounds, so the minority groups contain far fewer examples than the majority groups. In the validation and test sets, both the landbirds and waterbirds are evenly split between the water and land backgrounds.

### B.2. CelebA

We use the task setup from Sagawa et al. (2020a) on the CelebA celebrity face dataset (Liu et al., 2015). The label $y$ is set to be the *Blond_Hair* attribute, and the spurious attribute $a$ is set to be the *Male* attribute: being female spurious correlates with having blond hair. The minority groups are (blond, male) and (not blond, female), although the (blond, male) group is significantly smaller than the (not blond, female) group. The majority groups are (blond, female) and (not blond, male). We use the standard train / valid / test splits from Sagawa et al. (2020a).

| Group | Enrichment | ERM test acc. |
|---|---|---|
| (muslim, toxic) | 8.58x | 62.1% |
| (christian, toxic) | 8.58x | 57.4% |
| (LGBTQ, toxic) | 8.58x | 72.1% |
| (other religion, toxic) | 8.56x | 62.9% |
| (black, toxic) | 8.53x | 74.6% |
| (white, toxic) | 8.49x | 68.6% |
| (female, toxic) | 8.49x | 64.7% |
| (male, toxic) | 8.48x | 66.6% |
| (white, non-toxic) | 0.09x | 84.7% |
| (LGBTQ, non-toxic) | 0.07x | 82.2% |
| (black, non-toxic) | 0.07x | 74.6% |
| (male, non-toxic) | 0.06x | 93.0% |
| (muslim, non-toxic) | 0.05x | 89.4% |
| (female, non-toxic) | 0.05x | 94.7% |
| (other religion, non-toxic) | 0.04x | 93.4% |
| (christian, non-toxic) | 0.03x | 96.3% |

*Table 12.* CivilComments error set breakdowns.



*Figure 5.* Effect on worst-group accuracy of removing the $y = a$ and $y \neq a$ examples from JTT's error set. Both upsampling $y = a$ examples and upsampling $y \neq a$ examples substantially contribute to improving worst-group accuracy.

### B.3. MultiNLI

We use the task setup from Sagawa et al. (2020a) on the MultiNLI natural language inference dataset (Williams et al., 2018). Given two sentences, a premise and a hypothesis, the task is to predict whether the hypothesis is *entailed by*, *neutral with*, or *contradicted by* the premise. The spurious attribute $a$ is a binary indicator for when any of the negation words *nobody*, *no*, *never*, or *nothing* appear in the second sentence (the hypothesis), which spuriously correlates with the *contradiction* label. We use the standard train / valid / test splits from Sagawa et al. (2020a).

### B.4. CivilComments-WILDS

We use the CivilComments-WILDS dataset from Koh et al. (2021), which is derived from the Jigsaw dataset (Borkan et al., 2019). Given a real online comment, the task is to predict whether the comment is *toxic* or *not toxic*. The spurious attribute $a$ is an 8-dimensional binary vector, where each entry is a binary indicator of whether the following 8 demographic identities are mentioned in the online comment: *male*, *female*, *LGBTQ*, *Christian*, *Muslim*, *other religion*, *Black*, and *White*.

Following Koh et al. (2021), we consider the 16 *potentially overlapping* groups equal to (*identity*, *toxic*) and (*identity*, *not toxic*) for all 8 identities. We use the standard train / valid / test splits from Koh et al. (2021).

## C. Additional Experimental Results

### C.1. Error set analysis for CivilComments

We include the CivilComments error set analysis in Table 12 for space constraints.
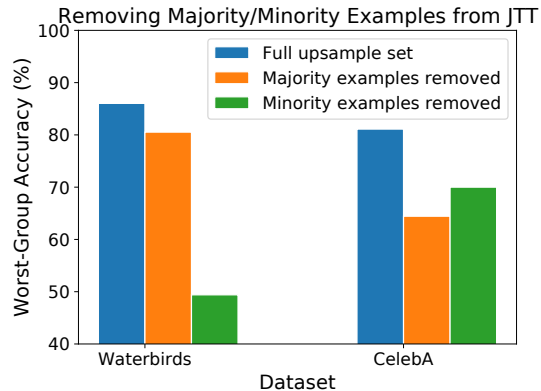
### C.2. Additional analysis

Below, we present a series of analyses that involve partitioning the dataset into two groups: groups in which spurious correlation holds with $y = a$, and groups in which spurious correlation does not hold with $9 \neq a$. For this investigation, we focus on Waterbirds and CelebA, where all groups can be clearly partitioned as above because we consider binary classification tasks with binary spurious attributes. In Waterbirds, the $y = a$ groups are waterbirds on water backgrounds and landbirds on land backgrounds; the $y \neq a$ groups are waterbirds on land backgrounds and landbirds on water backgrounds. In CelebA, the $y = a$ groups are blond females and non-blond males; the $y \neq a$ groups are non-blond females and blond males. In contrast, it is unclear how to partition the groups as above in MultiNLI, in which we consider a multi-class classification problem, and in CivilComments-WILDS, in which we have multiple spurious attributes corresponding to different demographic identities.

**Impact of $y = a$ and $y \neq$ examples in the error set.** We first study how worst-group accuracy changes when we remove $y = a$ examples or $y \neq a$ examples from the error set, as summarized in Figure 5. In both datasets, removing either the $y = a$ or $y \neq a$ examples from the error set significantly decreases worst-group accuracy. While this reduction in worst-group accuracy could stem from the fact that we consider a fixed set of hyperparameters including the upweight factor (which was tuned for JTT with the full error set), it is possible that both $y = a$ and $y \neq a$ examples contribute to the improvement in worst-group accuracy. In particular, because both datasets have substantial label imbalance, it is expected that upweighting groups from rare labels is important to perform well on all groups, and in
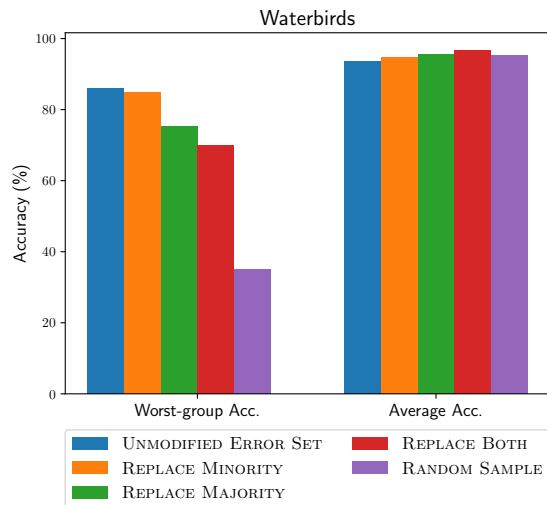
*Figure 6.* Effect of replacing the $y \neq a$ and $y = a$ examples in JTT's error set with randomly-selected $y \neq a$ and $y = a$ examples on Waterbirds. Worst-group accuracy decreases when replacing either the $y = a$ or $y \neq a$ examples, suggesting that JTT successfully automatically identifies informative examples that improve worst-group accuracy when upsampled.

fact, groups with $y \neq a$ and with rare $y$ are upweighted as discussed in Section 5.3.

Next, we explore if the particular $y = a$ or $y \neq a$ examples that JTT upsamples is important, or if upsampling *any* collection of examples in these groups yields high worst-group accuracy. To do this, we study how average and worst-group accuracies change when we replace examples in JTT's error set with randomly selected examples from specific groups. Specifically, we study what happens when we upsample the following four variants of JTT's error set:

- REPLACE $y \neq a$: We replace the $y \neq a$ examples in the error set with an equal number of randomly selected $y \neq a$ examples, leaving the $y = a$ examples in the error set uncha nged.

- REPLACE $y = a$: We replace the $y = a$ examples in the error set with an equal number of randomly selected $y = a$ examples, leaving the $y \neq a$ examples in the error set unchanged.

- REPLACE both: We replace both the $y \neq a$ examples in the error set with an equal number of randomly selected $y \neq a$ examples, and the $y = a$ examples in the error set with an e qual number of randomly selected $y = a$ examples.

- RANDOM sample: We replace all examples in the error set with an equal number of randomly selected examples. This yields a different fraction of $y \neq a$ examples in the error set compared to REPLACE both.

| | Waterbirds | | CelebA | |
|---|---|---|---|---|
| | Avg. | Worst-group | Avg. | Worst-group |
| UPSAMPLE MINORITY | 96.7% | 75.9% | 93.4% | 57.2% |
| JTT | 93.3% | **86.7%** | 88.0% | **81.1%** |

*Table 13.* Average and worst-group test accuracies. UPSAMPLE MINORITY, which upsamples $y \neq a$ examples, achieves higher worst-group accuracy than ERM, but lower than JTT.

Figure 6 compares upsampling these variants of the error set with the original unmodified error set (UNMODIFIED ERROR SET). Compared to upsampling the original error set, upsampling REPLACE $y \neq a$ slightly decreases worst-group accuracy and leaves average accuracy unchanged. This suggests that upsampling most $y \neq a$ examples helps improve worst-group accuracy, though the particular $y \neq a$ examples JTT identifies in the error set are still slightly better than rand om. On the other hand, upsampling REPLACE $y = a$ significantly decreases worst-group accuracy, although it slightly improves average accuracy compared to the original error set. This suggests that the particular $y = a$ examples JTT identifies in the error set are important for improving worst-group accuracy. This could be because the label balance within upsampled $y = a$ changes, or for other reasons. Finally, the low worst-group and average accuracies of both REPLACE both and RANDOM sample show that merely upsampling random $y = a$ and $y \neq a$ examples is insufficient to achieve high worst-group accuracy.

**Upsampling $y \neq a$ groups.** We present the performance of a simple baseline, in which we upweight $y \neq a$ examples using ground-truth group annotations, in Table 13. While UPSAMPLE MINORITY improves worst-group error over ERM, this baseline is limited in a few ways. First, $y \neq a$ groups are not necessarily groups with the worst accuracies or smallest number of examples, for example due to label imbalance. So while we $y \neq a$ examples are counter-examples to the spurious correlations, it's not necessarily expected that they improve the worst-group performance well. Secondly, in the presence of ground-truth examples, it is possible to reweight each of the groups independently, rather than reweighting $y = a$ and $y \neq a$ groups. Prior work has observed much higher worst-group performance by reweighting the groups than UPSAMPLE MINORITY (Sagawa et al., 2020a).