

An Investigation of Critical Issues in Bias Mitigation Techniques

Robik Shrestha¹ Kushal Kafle² Christopher Kanan^{1,3,4}
 Rochester Institute of Technology¹ Adobe Research² Paige³ Cornell Tech⁴
 {rss9369, kk6055, kanan}@rit.edu

Abstract

A critical problem in deep learning is that systems learn inappropriate biases, resulting in their inability to perform well on minority groups. This has led to the creation of multiple algorithms that endeavor to mitigate bias. However, it is not clear how effective these methods are. This is because study protocols differ among papers, systems are tested on datasets that fail to test many forms of bias, and systems have access to hidden knowledge or are tuned specifically to the test set. To address this, we introduce an improved evaluation protocol, sensible metrics, and a new dataset, which enables us to ask and answer critical questions about bias mitigation algorithms. We evaluate seven state-of-the-art algorithms using the same network architecture and hyper-parameter selection policy across three benchmark datasets. We introduce a new dataset called *Biased MNIST* that enables assessment of robustness to multiple bias sources. We use *Biased MNIST* and a visual question answering (VQA) benchmark to assess robustness to hidden biases. Rather than only tuning to the test set distribution, we study robustness across different tuning distributions, which is critical because for many applications the test distribution may not be known during development. We find that algorithms exploit hidden biases, are unable to scale to multiple forms of bias, and are highly sensitive to the choice of tuning set. Based on our findings, we implore the community to adopt more rigorous assessment of future bias mitigation methods. All data, code, and results are publicly available¹.

1. Introduction

Deep learning systems are trained to minimize their loss on a training dataset. However, datasets often contain spurious correlations and hidden biases which result in systems that have low loss on the training data distribution, but then fail to work appropriately on minority groups because they exploit and even amplify these spurious correlations [71, 35]. For example, in systems trained to infer hair color on the

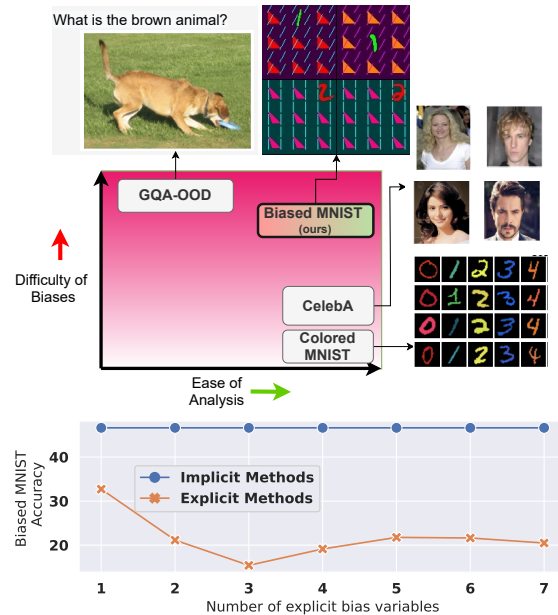


Figure 1: Current bias mitigation systems are tested on simple datasets that are easy to analyze, but do not offer challenges present in realistic cases. Addressing this, we propose the *Biased MNIST* dataset which is easy to analyze, yet is reflective of real world challenges since it contains multiple sources of biases. We find that methods fail on *Biased MNIST* even when all the biases are explicitly labeled. We also test on *GQA-OOD*, where the sources of biases are not very obvious and, thus the methods fail to generalize.

CelebA dataset [43], the majority group of non-blond males occurs 50 times more than the minority group of blond males, resulting in systems incorrectly predicting non-blond as hair color for the minority group.

While this is a toy problem, in the real world, hidden minority patterns are common and failing on them can have dire consequences. Systems designed to aid human resources, help with medical diagnosis, determine probation, or loan qualification could be biased against minority groups based on age, gender, religion, sexual orientation, ethnicity, or race [54, 8, 16, 13, 48]. Systems can exploit correlated variables even if they are not directly a part of the input e.g.,

¹<https://github.com/erobic/bias-mitigators>

through inferred zip codes [21], failing to work effectively on minority groups.

Recently, many methods have been proposed to make neural networks bias resistant. These methods can be grouped into two types: 1) those that assume the bias variables e.g., the gender label in CelebA, are explicitly annotated and can be accessed during training [55, 55, 69, 37] and, 2) those that do not require explicit access [46, 50]. Assuming explicit access requires extra annotations in addition to the actual target, and for many tasks it may not be immediately clear what the bias variables are e.g., biases may only be discovered years later [51, 50]. Methods that do not assume access to these bias variables have only recently been proposed [46, 65, 50].

So far, there is no study comparing methods from either group comprehensively. Often papers fail to compare against recent methods and vary widely in the protocols, datasets, architectures, and optimizers used. For instance, the widely used Colored MNIST dataset, where colors and digits are spuriously correlated with each other, is setup differently across papers. Some use it as a binary classification task (class 0: digits 0-4, class 1: digits: 5-9) [5, 50], whereas others use a multi-class setting (10 classes) [37, 40]. For CelebA, [46] uses ResNet-18 whereas [50] uses ResNet-50, but the comparison was done without taking this architectural change into account. These discrepancies make it difficult to judge the methods on an even ground.

Methods are typically highly sensitive to hyperparameter choices, and papers report numbers on systems in which the hyperparameters were tuned using the test set distribution [18, 50, 64]. In the real world, biases may stem from multiple factors and may change in different environments, making this setup unrealistic. Furthermore, tuning on the test distribution can lead to methods that are right for the wrong reasons. When this is done, systems can perform well just by exploiting the biases they are supposed to overcome [62, 64], and they will then fail once deployed because they have not really learned to solve the task.

In addition, we posit that the commonly used benchmarks are not challenging enough to test generalization to realistic scenarios. For example CelebA and Colored MNIST, two of the most widely used benchmarks, contain a single bias variable to mitigate: gender and color respectively. It is unclear how well methods would fare in presence of multiple types of bias, e.g., position or co-occurring objects/patterns, which are commonly present in real-world datasets. For some tasks it can be impossible to exhaustively enumerate all bias variables. For example, in visual question answering (VQA), where a system answers questions about images, biases can stem from: object-context co-occurrences, visual concept/language correlations, question type/answer distributions, and more. Annotating all such sources of bias is unrealistic. Even when the bias variables are explicitly la-

beled, it is still unclear if the methods can remain robust to all of the bias sources, since this entails generalization to a large number of dataset groups e.g., hundred thousand groups for GQA-OOD [36].

We address the above issues via these contributions:

1. We describe our new Biased MNIST dataset and corresponding evaluation protocol for measuring resistance to multiple forms of bias. It measures resistance to spuriously correlated background/foreground color, texture, co-occurring distractors, position, and more.
2. We compare seven state-of-the-art bias mitigation methods on classification tasks using Biased MNIST and CelebA, measuring generalization to minority patterns, scalability to multiple sources of biases, sensitivity to hyperparameters, etc. We ensure fair comparisons by using the same architecture, optimizer, and performing grid searches for hyperparameters.
3. To go beyond image classification, we measure the performance of these methods on the biased GQA benchmark for VQA.
4. We provide concrete recommendations for future studies, so that comparisons among algorithms are meaningful and reflective of real-world challenges.

2. Problem Statement

To properly study bias mitigation, it is necessary to provide a definition of biased data and biased behavior in a model. We study bias in supervised classification i.e., the goal is to learn a function $f : X \rightarrow Y$ which outputs a categorical target $y \in Y$ given $x \in X$. Each x is itself a mixture of a signal s that we wish the system to use for inference and bias b that is spuriously correlated with y . Since the spurious correlations between y and b do not always hold, systems exploiting b to infer y fail to generalize.

We can measure the robustness to such tendencies by intentionally introducing covariate shift e.g., with a test dataset distribution that differs from training or a metric that balances performance across groups. For our study, we use the mean per group accuracy/unbiased accuracy, which weighs all the groups equally. Furthermore, we focus on the cross-bias setting defined by [6] where the same set of bias variables are present in both train and test sets.

3. Bias Mitigation Strategies

Without bias mitigation mechanisms, **standard models (StdM)** often use spurious bias variables for inference, rather than developing invariance to them, which often results in their inability to perform well on minority patterns [27, 11, 3, 61]. To address this, several bias mitigation mechanisms have been proposed, and they can be categorized into two groups: 1) methods that access explicit bias labels during

training, and 2) methods that do not assume such access. We briefly review methods from these categories, with an emphasis on the methods assessed in our studies.

3.1. Explicit Bias Mitigation Methods

Explicit bias mitigation techniques directly access the bias variables: b_{expl} . during training to develop invariance to them. Based on the way these variables are utilized during training, we choose five different explicit methods for our study. We refer to them as **explicit methods** for conciseness.

Re-sampling/Re-weighting: These approaches balance out the spurious correlations. The classical approach is to re-balance the class distribution by adjusting the sampling probability/ loss weight for majority/minority samples [14, 26, 41, 72, 20]. This includes synthesizing minority instances too [14, 26]. Moving beyond class imbalances, REPAIR [40] proposed learning dynamic weights to mitigate representation bias [39]. However, [55] have shown promising results by using static weights to upweight minority patterns. We choose this method due to its simplicity.

Group Upweighting (Up Wt) [55] attempts to mitigate the correlations between y and b_{expl} . by upweighting the minority patterns. Specifically, each sample (x, y) is assigned to a group: $g = (y, b_1, b_2, \dots, b_E)$, where E is the total number of variables contained in b_{expl} . and the loss is scaled by $\frac{1}{N_g}$, where N_g is the number of instances in group g . Up Wt requires the models to be sufficiently regularized, i.e., be trained with low learning rates and/or high weight decays to be robust to the minority groups.

Distributionally Robust Optimization (DRO): DRO [22] minimizes the worst-case expected loss over potential test distributions. Often, such distributions are approximated by sampling from a uniform divergence ball around the train distribution [10, 23, 47]. However, this lacks structured priors about the potential shifts, and instead hurts generalization [32].

Group DRO (GDRO) [55] provides DRO with the necessary prior that it must generalize to all groups. Similar to Up Wt, GDRO also uses y and b_{expl} . to create groups and has been shown to work well with sufficiently regularized models. However, unlike Up Wt, it performs weighted sampling from each group and has an optimization procedure to minimize the loss over the worst-case group.

Ensembling Approaches: Ensembling approaches [28, 17, 12] have a two-branch setup: a) a bias-only branch f_b that predicts y from b alone to identify the bias-prone samples, and b) a de-biased branch $f_d(\cdot)$ that is trained to focus on samples that f_b finds difficult so that it learns richer features that work on difficult samples too. The two branches can be ensembled in different ways. DRiFt [28] uses product-of-experts [31] and LearnedMixIn [17] extends this through learned weights and entropy constraints that control f_b .

Reduction of Unimodal Biases (RUBi) [12] multiplies

the outputs from $f_d(\cdot)$ with sigmoided outputs from $f_b(\cdot)$, thereby assigning higher loss weights to samples that cannot be predicted through biases alone. RUBi was the previous state-of-the-art on VQA-CP [3], a testbed for measuring robustness to biases in VQA. For the bimodal problem of VQA, the original implementation focused on linguistic biases, training f_b on question features only. For our studies, we instead train f_b on b_{expl} . directly, to control the type of biases captured by f_b . We assess RUBi [12] over others since it performed better in the preliminary studies.

Adversarial Debiasing: These techniques impair the ability of the representation learner to encode biases [69, 1, 52, 25]. Like ensembling methods, they also employ a two-branch setup, with the representation encoder in the main branch being penalized if the bias-only branch: $f_b(\cdot)$ is successful at predicting biases from them [69]. Alternately, $f_b(\cdot)$ may be trained to predict the class label from the biased features [52, 25], but in either case, the gradient from $f_b(\cdot)$ is reversed during backpropagation for debiasing.

Learning Not to Learn (LNL) [37] uses an adversarial setup derived from minimization of mutual information between representation and bias. In addition to the gradient reversal, the mutual information formulation introduces an entropy regularization on the bias predictions.

Invariant Risk Minimization (IRM): The goal of IRM is to extract representations that are invariant across environments: $\mathcal{E} = \{e_1, e_2, \dots, e_E\}$, each encoding different spurious correlations [5, 63, 15]. Such representations enable the same classifier to be simultaneously optimal over all \mathcal{E} . For this, [5] propose to regularize the gradient norm of a fixed linear classifier. More recent variants include regularization of variance of risks [38, 67]. However, [53] have shown that such objectives can fail to recover the invariant features in practice. Despite this negative result, we still compare against IRM since it is a promising research direction.

IRMv1 [5] is an efficient approximation of an otherwise computationally expensive bi-level IRM objective. It consists of a regularization constraint on the gradient norm with respect to a fixed scalar $\theta_c = 1.0$:

$$\min_{\theta} \sum_{e \in \mathcal{E}_{tr}} l^e(\hat{y}) + \lambda \|\nabla_{\theta_c} l^e(\theta_c \cdot \hat{y})\|^2,$$

where, l^e is loss on environment e , \hat{y} is the logit vector yielded by the model parameterized by θ and λ balances between the empirical risk and invariance. In our experiments, we use the previously defined explicit data groups as the training environments for IRM.

3.2. Implicit Bias Mitigation Methods

Since explicit access to bias variables is an undesirable requirement in practice, some recent methods have proposed to mitigate biases without such assumptions. We call them **implicit methods** for conciseness and describe them below.

Limited Capacity Models: Most implicit methods assume that easy-to-learn biases can be captured by limiting the capacity of the models [65, 57, 46]. The capacity of such bias-prone models: $f_b()$ can be limited by using a small subset of train instances for a few epochs [65], using fewer model parameters [57], attaching a classifier to intermediate layers, instead of using the final representation layers [18], using bias-prone architectures [7] or amplifying biases [46]. Main network: f_d is then debiased by assigning higher weights to the harder samples, so that it generalizes to samples that cannot be predicted through biases alone.

Learning From Failure (LFF) [46] amplifies the bias in f_b using the generalized cross entropy loss [70]:

$$GCE(s(x; \theta), y) = \frac{1 - s_y(x; \theta)^\gamma}{\gamma},$$

where s_y is the softmax score for the ground truth class and γ determines the degree of bias amplification. Samples with high f_b loss are then assigned higher weights while training f_d . While γ seems critical, the original paper does not discuss a way to tune it and instead fixes it to a default value of $\gamma = 0.7$. In fact, the paper does not provide details on model selection at all; however, this is an important question, which we discuss in Sec. 6.4.

Gradient Starvation Mitigation: Different from the limited capacity methods, spectral decoupling [50] aims to overcome the issue of gradient starvation [19], which is the tendency to only rely on statistically dominant features. This is related to the simplicity bias exploitation, where models exploit the simplest features despite having access to more predictive features [59, 30], which are more complex.

Spectral Decoupling (SD) [50] aims to decouple the learning dynamics between features. The authors show that regularizing the network outputs (\hat{y}) as:

$$\frac{\lambda}{2} \|\hat{y} - \gamma\|_2^2,$$

where, λ and γ are hyperparameters, provably decouples the learning dynamics, enabling learning of better features.

4. Datasets

We use datasets that enable probing existing methods with critical questions regarding their robustness. We test on datasets with varying scales and types of biases, allowing us to perform highly controlled studies that analyze scalability to a large number of hidden groups.

4.1. Biased MNIST

Existing datasets for assessing bias mitigation methods do not enable analysis of multiple bias sources, e.g., Colored MNIST only tests for color versus class bias. To address this, we created the Biased MNIST dataset, which requires

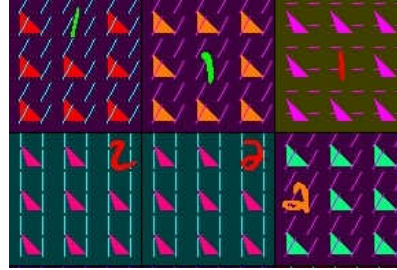


Figure 2: Biased MNIST requires the methods to classify the target digit while remaining invariant to multiple sources of biases.

recognizing digits while remaining robust to multiple sources of biases i.e., other factors which are also correlated with the target variable. Specifically, the dataset consists of 3×3 grids of cells, where the target digit is placed at one of the grid cells and is correlated with multiple bias variables $\{b_j\}$ including: a) background color, b) digit color, c) digit position in the grid, d) distracting shapes present in other cells, e) color of the distracting shapes, f) texture type and g) texture color (see Fig. 2). Each variable can take one of ten discrete values: each variable takes the majority/biased value corresponding to the target digit with a probability of p_{bias} , otherwise takes one of the remaining nine values with uniform probability. For instance, if $p_{bias} = 0.7$, then 70% of 1's are colored as green, 70% of 2's are colored red and so on. This applies to all the variables e.g., 70% of 1's placed on a purple background, while 70% of 2's are placed on a green background. The bias level: p_{bias} can be specified for each variable to control the types and degrees of biases. For convenience, p_{bias} is set to 0.7 in the train set, unless otherwise specified. The test set is unbiased i.e., has $p_{bias} = 0.1$.

4.2. CelebA

The CelebA dataset [43] of celebrity faces is widely used to assess bias mitigation techniques [55, 56, 46, 50]. Following earlier work, it is used for binary hair color classification (blond or non-blond), which is correlated with gender. There are two major bias sources: a) class imbalance, with non-blond occurring 5.7 times more than blond hair color, and b) presence of a rare group, i.e., blond male celebrities only account for 0.86% of the training instances.

4.3. GQA-OOD

We use the GQA visual question answering dataset [33] to highlight the challenges of using bias mitigation methods on real-world tasks. It has multiple sources of biases including imbalances in answer distribution, visual concept co-occurrences, question word correlations, and question type/answer distribution. It is unclear how the explicit bias variables should be defined so that the methods can gener-

alize to all minority groups. GQA-OOD [36] divides the evaluation and test sets into majority (head) and minority (tail) groups based on the answer frequency within each ‘local group’ (e.g., colors of bags), which is a unique combination of ‘global group’ or answer type (e.g., objects or colors) and the main concept asked in the question (e.g., ‘bag’, ‘chair’, etc.). The head/tail categorization makes analysis easier; however, it is unclear how one should specify the explicit biases so that the models generalize even to the rarest of local groups. Therefore, we explore multiple ways of defining the explicit bias variable in separate experiments: a) majority/minority group label (2 groups), b) answer class (1833 groups), c) global group (115 groups) and d) local group (133328 groups). It is unknown if bias mitigation methods can scale to hundreds and thousands of groups in GQA, yet natural tasks require such an ability.

5. Network Architecture & Tuning Procedure

For each dataset, we assess all bias mitigation methods with the same neural network architecture. For CelebA, we use ResNet-18 [29]. For Biased MNIST, we use a convolutional neural network with four ReLU layers consisting of a max pooling layer attached after the first convolutional layer. For GQA-OOD, we employ the UpDn architecture [4], which is widely used for VQA [58, 36, 66].

For each dataset, we use the class label y and the explicit bias variables $b_{expl.}$ to define explicit groups for Up Wt, GDRO and IRMv1. For instance, for CelebA, hair color and gender result in four explicit groups while for Biased MNIST, the number of groups: $|G|$ depends on the number of explicit bias variables: $|b_{expl.}|$, with $|G| = 10^{|b_{expl.}|}$. We will specify the exact $b_{expl.}$ for each experiment in Sec. 6. For GQA, we use head/tail, answer class, global and local groups as explicit variables. For all datasets, RUBi uses $b_{expl.}$ to predict y , whereas LNL trains the adversarial branch to predict $b_{expl.}$ from representations. Of course the implicit methods: StdM, LFF and SD are invariant to the choice of the explicit biases during training. Unless otherwise specified, results from Biased MNIST are averaged across 3 random seeds, but due to computational constraints, we ran models on CelebA and GQA-OOD only once.

Hyperparameters for each method were chosen using a grid search with unbiased accuracy on each dataset’s validation set. To make this tractable, we first ran a grid search for the learning rate over $\{10^{-3}, 10^{-4}, 10^{-5}\}$ and weight decay over $\{0.1, 10^{-3}, 10^{-5}, 0\}$. After the best values were chosen, we searched for method-specific hyperparameters. Due to the size of GQA-OOD, hyperparameter search was performed by training on only 10% of instances, and then the best selected hyperparameters were used with the full training dataset. The exact values for the hyperparameters are specified in the Appendix.

Table 1: Unbiased accuracies $Acc(\alpha = 0)$ on all datasets for all methods. We format the *first*, *second* and *third* best results. Methods that do not access explicit biases have gray background.

Methods/ Datasets	CelebA	Biased MNIST	GQA
StdM	80.3	42.0	<u>44.8</u>
Up Wt [56]	<u>87.4</u>	30.1	30.0
GDRO [55]	88.5	27.2	26.4
RUBi [12]	87.2	38.9	24.1
LNL [37]	79.2	40.6	28.6
IRMv1 [5]	79.8	38.7	39.3
LFF [46]	77.8	<u>56.6</u>	45.1
SD [50]	<u>88.6</u>	<u>41.3</u>	<u>46.9</u>

6. Questions Posed and Answered

In this section, we probe the existing methods with critical questions regarding their robustness. For each question, we first describe the empirical setup to explore the question, and then present the results.

6.1. Head-to-Head Comparisons

Question 1: Are there clear winners in a head-to-head comparisons across datasets?

We first present the mean per group accuracy for all eight methods on all three datasets in Table. 1 to see if any method does consistently well across benchmarks. For this, we used class and gender labels as explicit biases for CelebA. For Biased MNIST, there are multiple ways to define explicit biases, but for this section, we simply use each of the seven variables as explicit biases in different runs and average across the runs. We study combinations of multiple explicit variables in Sec. 6.3. We set $p_{bias} = 0.7$ for this section, and present results across different p_{bias} in the Appendix. Similarly for GQA, we consider each of the four variables described in Sec. 4.3 as explicit bias in separate runs and present the average.

Results. As shown in Table. 1, no method performs universally well across datasets; however, the implicit methods LFF and SD obtain high unbiased accuracies on most datasets. This shows that implicit methods can deal with multiple bias sources without explicit access. Explicit methods work well on CelebA but fail on Biased MNIST and GQA. Specifically, Up Wt, GDRO and RUBi obtain 7-8% improvements over StdM on CelebA, which requires generalization to only 4 groups. However, all explicit methods perform worse than StdM on Biased MNIST and GQA, signifying their inability to deal with multiple bias sources. LNL and IRMv1 were comparable to StdM even on CelebA, demonstrating lack of generalization even on simple settings. Despite being a simpler method, Up Wt outperformed GDRO on both Biased MNIST and GQA, but both were

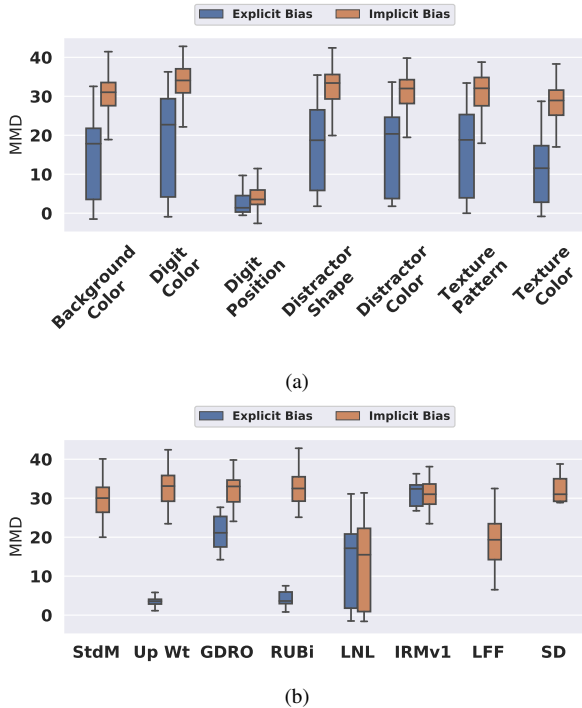


Figure 3: Boxplots of majority/minority difference (MMD) on Biased MNIST over: a) bias variables, and b) different methods.

worse than StdM. These results show that implicit methods can outperform explicit methods.

6.2. Bias Exploitation

Question 2: Do methods show robustness to both explicit and implicit biases?

In this set of experiments, we compare the resistance to explicit and implicit biases. We primarily focus on the Biased MNIST dataset, reserving each individual variable as the explicit bias in separate runs of the explicit methods, while treating the remaining variables as implicit biases. To ease analysis, we compute the accuracy gap between the majority and minority groups i.e., majority/minority difference (MMD). Majority/minority groups are defined per variable e.g., for foreground color, green 1’s, red 2’s etc are placed in the majority group and the rest in the minority group and MMD simply computes the accuracy difference between the two groups for each variable. High MMDs indicate that the methods rely heavily on spurious patterns favoring the majority groups and thus fail on the minority groups.

Results. In Fig. 3a, we present the MMD boxplots for all bias variables, comparing cases when the label of the variable is either explicitly specified (explicit bias), or kept hidden (implicit bias) from the methods. Barring digit position, we observe that the MMD values are higher when the variables are not explicitly labeled for the methods, indicating that the

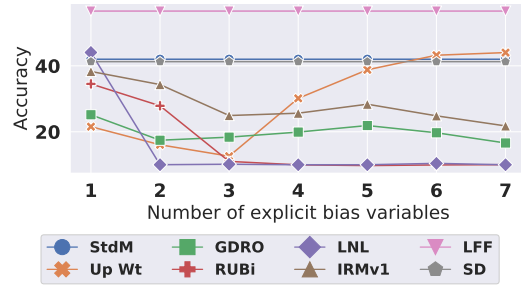


Figure 4: Unbiased accuracy as a function of number of explicit bias variables. StdM, LFF and SD are invariant to the choice of explicit bias variables.

Table 2: Mean of head and tail accuracies on GQA, considering different variables as explicit biases.

Methods	Head/Tail (2 groups)	Answer Class (1833 groups)	Global Group (115 groups)	Local Group (133328 groups)
StdM		44.8		
Up Wt [56]	43.3	26.0	26.4	24.2
GDRO [55]	46.9	28.6	10.8	19.4
RUBi [12]	44.1	N/A	5.6	22.6
LNL [37]	42.9	N/A	32.4	10.7
IRMv1 [5]	47.2	35.8	40.4	33.8
LFF [46]		45.1		
SD [50]		46.9		

explicit methods in general fail to mitigate implicit biases. Fig. 3b breaks down exploitation of explicit and implicit biases for each method. UpWt, GDRO and RUBi have low MMD values for explicit biases, but high MMD values for implicit biases, showing that they mitigate the explicit biases to some extent, but are not robust to the implicit biases. LNL and IRMv1 seem to be equanimously affected by both explicit and implicit biases, and thus fail to improve upon the baseline as previously shown in Table 1. LFF has a relatively low range of MMDs and as shown by the improvements in Table 1, the method outperforms others on Biased MNIST.

Interestingly, MMD was low for digit position. We hypothesize this is because CNNs are unable to use position information for inference [42]. To confirm this, we add CoordConv layers [42] before and after the maxpooling layer in CNN to enable usage of position information. This resulted in methods exploiting digit position too, showing larger MMD values of 11.1%-25.6% as compared to the 2.2%-8.7% without the CoordConv layers. Such inductive biases affect whether or not methods exploit certain dataset biases, and we discuss this in Sec. 7.

6.3. Scalability of Methods

Question 3: Do methods scale up to multiple types of biases and a large number of dataset groups?

It is unknown how well the methods scale up to multiple sources of biases and large number of groups, even when they are explicitly annotated. To study this, we train the explicit methods with multiple explicit variables for Biased MNIST and individual variables that lead to hundreds and thousands of groups for GQA and compare them with the implicit methods. For Biased MNIST, we first sort the seven total variables in the descending order of MMD (obtained by StdM) and then conduct a series of experiments. In the first experiment, the most exploited variable, distractor shape, is used as the explicit bias. In the second experiment, the two most exploited variables, distractor shape and texture, are used as explicit biases. This is repeated until all seven variables are used². Note that conducting the seventh experiment entails annotating each instance with every possible source of bias. While this may not be realistic in practice, such a controlled setup will reveal if the explicit methods can generalize when they have complete information about every bias source.

To test scalability on a natural dataset, we conduct four experiments per explicit method on GQA-OOD with the explicit bias variables: a) head/tail (2 groups), b) answer class (1833 groups), c) global group (115 groups), and d) local group (133328 groups). Unlike Biased MNIST, we do not test with combinations of these variables since the last three variables already entail generalization to many groups.

Results. We find that implicit methods either improve or are comparable with StdM, but most explicit methods fail when asked to generalize to multiple bias variables and a large number of groups, even when the bias variables are explicitly provided. As shown in Fig. 4, all explicit methods are below StdM on Biased MNIST. Barring LNL and Up Wt, other explicit methods exhibit degraded accuracy as the number of explicit bias variables increases. Because the implicit methods do not rely on the choice of explicit biases, we simply repeat the same accuracy across x-axis. Among the implicit methods, LFF obtains the highest improvement, whereas SD is close to StdM.

Results for GQA-OOD are similar, with explicit methods failing to scale up to a large number of groups, while implicit methods showing some improvements over StdM. As shown in Table 2, when the number of groups is small, i.e., when using a head/tail binary indicator as the explicit bias, explicit methods remain comparable or even outperform StdM, but when the number of groups grow to hundreds and thousands, they fail. IRMv1 and GDRO obtain the highest improvements of 2.4% and 2.1% over StdM, respectively, with the binary head/tail bias, but they show large drops when using

²The exact order is given in the Appendix.

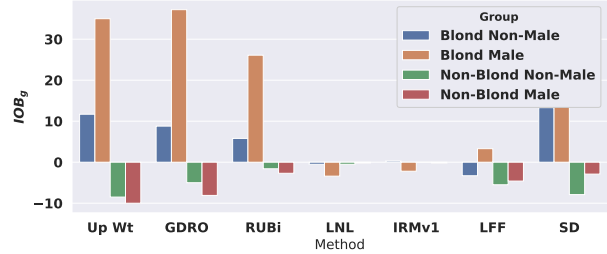


Figure 5: Improvement Over the Standard Model (IOSM) for each group of CelebA.

answer class, global group or local group as explicit bias variables. Some drops are extreme, e.g., RUBi drops 39% when using global group as the explicit bias variable.

Results on a simpler setting. We further study bias exploitation on CelebA. For this, we plot improvement over the standard model (*IOSM*) in Fig. 5, which is the accuracy gain over the standard model on each dataset group. The improvements in blond (minority group) incur degradation in non-blond (majority group). The methods tilt predictions either in the favor of minority or majority groups, which shows the inability to learn the signal even on simple settings.

6.4. Robustness to Model Selection Criteria

Question 4: Can the methods generalize to the test set without being tuned on the test distribution? Do they exhibit robustness across a wide range of hyperparameters?

Assuming access to the test distribution for model selection is unrealistic and can result in models being right for the wrong reasons [64]. Rather, it is ideal if the methods can generalize without being tuned on the test distribution and we study this ability by comparing models selected through varying tuning distributions. To control the tuning distribution, we define a generalization of the mean per group accuracy (MPG) metric, that can interpolate within as well as extrapolate beyond the train and test distributions:

$$Acc(\alpha) = \frac{\sum_{g=1}^{|G|} p_g^\alpha Acc_g}{\sum_{g=1}^{|G|} p_g^\alpha}.$$

Here, p_g denotes the ratio of samples present in group g , $|G|$ is the total number of groups and α is used to control the group prior. When $\alpha = 0$, $Acc(\alpha = 0)$ yields the MPG i.e., a balanced distribution where all groups are weighed equally. When $\alpha = 1$, then the weights reflect the train priors/biases. When $0 < \alpha < 1$, it interpolates between biased (train) priors and unbiased group weights. When $\alpha < 0$, minority groups are weighed more and when $\alpha > 1$, majority groups are weighed more i.e., it amplifies the train bias.

Ideally, methods should yield robust models regardless of the tuning distribution i.e., the value of α . To test this ability, we train a set of candidate models for model selection, with

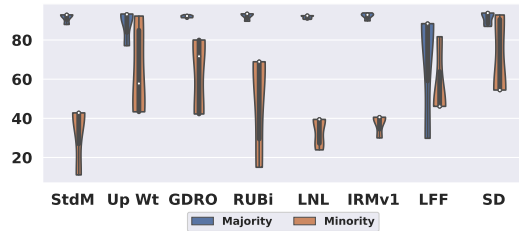


Figure 6: Range of minority (blond haired males) and majority (mean over rest of the groups) test accuracies on CelebA when varying the validation distribution from $\alpha = -1$ (inverted train bias) to $\alpha = 2$ (increased train bias).

different sets of hyperparameters. Specifically, we train nine different models with $learning\ rate \in \{10^{-3}, 10^{-4}, 10^{-5}\}$ and $weight\ decay \in \{10^{-1}, 10^{-3}, 10^{-5}\}$ for CelebA and Biased MNIST and, then perform model selection by computing $Acc(\alpha)$ at $\alpha \in \{-1.0, -0.5, 0.0, 0.5, 1.0, 1.5, 2.0\}$ on the validation sets.

Results. For CelebA, methods generally show large variance on the minority patterns (blond haired male celebrities), and lower variance on the majority patterns (mean over rest of the groups), whereas for Biased MNIST, we find that methods only work for certain set of hyperparameters and show degraded results on both majority/minority patterns if the hyperparameters change. As illustrated by the violin plots of CelebA’s unbiased test accuracies in Fig. 6, LNL and IRMv1 have the lowest variance, but neither improves over StdM. Up Wt, GDRO and RUBi show the largest variances for the minority group, indicating they are highly sensitive to the choice of tuning distribution. For LFF, we found high variance for both majority and minority patterns. In fact, we were unable to replicate the published LFF results, with $\gamma = 0.7$ yielding a high accuracy (86%) on the rarest group, but low accuracies on the rest. After tuning it, we found that $\gamma = 0.1$ gave the best unbiased test accuracy.

Interestingly, for Biased MNIST we found that $learning\ rate = 10^{-3}$ and $weight\ decay = 10^{-5}$ worked best for all methods. Even though Up Wt and GDRO are known to generalize to minority groups when using a low learning rate and high weight decay [56], we did not observe this for Biased MNIST. We hypothesize that when multiple sources of bias are present, as in Biased MNIST, methods have multiple ways of predicting the classes, some of which maybe easier to learn than the others. When the hyperparameters are suitable to exploit these biases, methods obtain their best accuracies, which are still lower than StdM.

7. Discussion

Our study demonstrates that systems are highly sensitive to the tuning distribution, that explicit methods cannot handle multiple bias sources, and that more rigorous analysis is

critical for bias mitigation algorithms for future progress. Based on our results, we argue that the community should focus on implicit methods, rather than explicit, not only because explicit methods require additional annotations, but also because they perform worse on realistic settings.

We make the following recommendations:

1. Compare against multiple state-of-the-art methods under fair settings.
2. Test on datasets that enable control over the number and degrees of biases, including realistic datasets.
3. Analyze generalization to both explicit and implicit sources of bias.
4. Be forthcoming about whether test distribution was used for model selection and compare robustness to tuning distributions that differ from the test.

If these guidelines are adopted, we believe significant progress can be made so that bias mitigation algorithms can have real-world benefit for deployed systems.

An interesting observation was that a weaker architecture, CNNs, were able to ignore position bias, whereas a more powerful architecture, CoordConv, resorted to exploiting this bias resulting in worse performance. While the community has largely focused on training procedures for bias mitigation, an exciting avenue for future work is to incorporate appropriate inductive biases into the architectures, perhaps endowing them with the ability to choose the the minimal computational power to do a task so that they are less sensitive to unwanted biases. This will essentially enable the algorithms to use Occam’s razor to determine the minimal capabilities required to do a task to reduce their ability to utilize biases.

We have pointed to issues with the existing bias mitigation approaches, which alter the loss or use resampling. An orthogonal avenue for attacking bias mitigation is to use alternative architectures. Neuro-symbolic and graph-based systems could be created that focus on learning and grounding predictions on structured concepts, which have shown promising generalization capabilities [68, 44, 34, 24, 60]. Causality is another relevant line of research, where the goal is to uncover the underlying causal mechanisms [49, 45, 9, 2]. Discovery and usage of causal concepts is a promising direction for building robust systems. These areas have not been explicitly studied for their ability to overcome dataset bias.

Acknowledgements. This work was supported in part by the DARPA/SRI Lifelong Learning Machines program [HR0011-18-C-0051], AFOSR grant [FA9550-18-1-0121], and NSF award #1909696. The views and conclusions contained herein are those of the authors and should not be interpreted as representing the official policies or endorsements of any sponsor.

References

- [1] E. Adeli, Qingyu Zhao, A. Pfefferbaum, E. Sullivan, Li Fei-Fei, Juan Carlos Niebles, and K. Pohl. Bias-resilient neural network. *ArXiv*, abs/1910.03676, 2019. [3](#)
- [2] Vedika Agarwal, Rakshith Shetty, and Mario Fritz. Towards causal VQA: revealing and reducing spurious correlations by invariant and covariant semantic editing. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 9687–9695. IEEE, 2020. [8](#)
- [3] Aishwarya Agrawal, Dhruv Batra, Devi Parikh, and Anirudha Kembhavi. Don’t just assume; look and answer: Overcoming priors for visual question answering. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 4971–4980. IEEE Computer Society, 2018. [2](#), [3](#)
- [4] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 6077–6086. IEEE Computer Society, 2018. [5](#)
- [5] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019. [2](#), [3](#), [5](#), [6](#)
- [6] Hyojin Bahng, Sanghyuk Chun, Sangdoon Yun, Jaegul Choo, and Seong Joon Oh. Learning de-biased representations with biased representations. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 528–539. PMLR, 2020. [2](#)
- [7] Hyojin Bahng, Sanghyuk Chun, Sangdoon Yun, Jaegul Choo, and Seong Joon Oh. Learning de-biased representations with biased representations. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 528–539. PMLR, 2020. [4](#)
- [8] Solon Barocas and Andrew D Selbst. Big data’s disparate impact. *Calif. L. Rev.*, 104:671, 2016. [1](#)
- [9] Alexis Bellot and Mihaela van der Schaar. Accounting for unobserved confounding in domain generalization. *arXiv preprint arXiv:2007.10653*, 2020. [8](#)
- [10] Aharon Ben-Tal, Dick Den Hertog, Anja De Waegenare, Bertrand Melenberg, and Gijs Rennen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59(2):341–357, 2013. [3](#)
- [11] Tolga Bolukbasi, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam Tauman Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 4349–4357, 2016. [2](#)
- [12] Rémi Cadène, Corentin Dancette, Hedi Ben-younes, Matthieu Cord, and Devi Parikh. RUBi: Reducing Unimodal Biases for Visual Question Answering. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 839–850, 2019. [3](#), [5](#), [6](#)
- [13] Robert Challen, Joshua Denny, Martin Pitt, Luke Gompels, Tom Edwards, and Krasimira Tsaneva-Atanasova. Artificial intelligence, bias and clinical safety. *BMJ Quality & Safety*, 28(3):231–237, 2019. [1](#)
- [14] Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002. [3](#)
- [15] Yo Joong Choe, Jiyeon Ham, and Kyubyong Park. An empirical study of invariant risk minimization. *ICML 2020 Workshop on Uncertainty and Robustness in Deep Learning*, 2020. [3](#)
- [16] Alexandra Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2):153–163, 2017. [1](#)
- [17] Christopher Clark, Mark Yatskar, and Luke Zettlemoyer. Don’t take the easy way out: Ensemble based methods for avoiding known dataset biases. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4069–4082, Hong Kong, China, 2019. Association for Computational Linguistics. [3](#)
- [18] Christopher Clark, Mark Yatskar, and Luke Zettlemoyer. Learning to model and ignore dataset bias with mixed capacity ensembles. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3031–3045, Online, 2020. Association for Computational Linguistics. [2](#), [4](#)
- [19] Remi Tachet des Combes, Mohammad Pezeshki, Samira Shabanian, Aaron Courville, and Yoshua Bengio. On the learning dynamics of deep neural networks. *arXiv preprint arXiv:1809.06848*, 2018. [4](#)
- [20] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge J. Belongie. Class-balanced loss based on effective number of samples. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 9268–9277. Computer Vision Foundation / IEEE, 2019. [3](#)
- [21] Anupam Datta, Matthew Fredrikson, Gihyuk Ko, Piotr Mardziel, and Shayak Sen. Use privacy in data-driven systems: Theory and experiments with machine learnt programs. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pages 1193–1210, 2017. [2](#)
- [22] Erick Delage and Yinyu Ye. Distributionally robust optimization under moment uncertainty with application to data-driven problems. *Operations research*, 58(3):595–612, 2010. [3](#)
- [23] John C Duchi, Peter W Glynn, and Hongseok Namkoong. Statistics of robust optimization: A generalized empirical

- likelihood approach. *Mathematics of Operations Research*, 2021. 3
- [24] Difei Gao, Ke Li, Ruiping Wang, Shiguang Shan, and Xilin Chen. Multi-modal graph neural network for joint reasoning on vision and scene text. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 12743–12753. IEEE, 2020. 8
- [25] Gabriel Grand and Yonatan Belinkov. Adversarial regularization for visual question answering: Strengths, shortcomings, and side effects. In *Proceedings of the Second Workshop on Shortcomings in Vision and Language*, pages 1–13, Minneapolis, Minnesota, 2019. Association for Computational Linguistics. 3
- [26] Haibo He, Yang Bai, Eduardo A Garcia, and Shutao Li. Adasyn: Adaptive synthetic sampling approach for imbalanced learning. In *2008 IEEE international joint conference on neural networks (IEEE world congress on computational intelligence)*, pages 1322–1328. IEEE, 2008. 3
- [27] Haibo He and Eduardo A Garcia. Learning from imbalanced data. *IEEE Transactions on knowledge and data engineering*, 21(9):1263–1284, 2009. 2
- [28] He He, Sheng Zha, and Haohan Wang. Unlearn dataset bias in natural language inference by fitting the residual. In *Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP (DeepLo 2019)*, pages 132–142, Hong Kong, China, 2019. Association for Computational Linguistics. 3
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778. IEEE Computer Society, 2016. 5
- [30] Katherine L Hermann and Andrew K Lampinen. What shapes feature representations? exploring datasets, architectures, and training. *NeurIPS*, 2020. 4
- [31] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002. 3
- [32] Weihua Hu, Gang Niu, Issei Sato, and Masashi Sugiyama. Does distributionally robust supervised learning give robust classifiers? In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 2034–2042. PMLR, 2018. 3
- [33] Drew A. Hudson and Christopher D. Manning. GQA: A new dataset for real-world visual reasoning and compositional question answering. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 6700–6709. Computer Vision Foundation / IEEE, 2019. 4
- [34] Drew A. Hudson and Christopher D. Manning. Learning by abstraction: The neural state machine. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 5901–5914, 2019. 8
- [35] Kushal Kafle, Robik Shrestha, and Christopher Kanan. Challenges and prospects in vision and language research. *Frontiers in Artificial Intelligence*, 2:28, 2019. 1
- [36] Corentin Kervadec, Grigory Antipov, Moez Baccouche, and Christian Wolf. Roses are red, violets are blue... but should vqa expect them to? *arXiv preprint arXiv:2006.05121*, 2020. 2, 5
- [37] Byungju Kim, Hyunwoo Kim, Kyungsu Kim, Sungjin Kim, and Junmo Kim. Learning not to learn: Training deep neural networks with biased data. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 9012–9020. Computer Vision Foundation / IEEE, 2019. 2, 3, 5, 6
- [38] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghui Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). *arXiv preprint arXiv:2003.00688*, 2020. 3
- [39] Yingwei Li, Yi Li, and Nuno Vasconcelos. Resound: Towards action recognition without representation bias. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 513–528, 2018. 3
- [40] Yi Li and Nuno Vasconcelos. REPAIR: removing representation bias by dataset resampling. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 9572–9581. Computer Vision Foundation / IEEE, 2019. 2, 3
- [41] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2999–3007. IEEE Computer Society, 2017. 3
- [42] Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 9628–9639, 2018. 6
- [43] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Large-scale celebfaces attributes (celeba) dataset. *Retrieved August, 15:2018*, 2018. 1, 4
- [44] Jiayuan Mao, Chuang Gan, Pushmeet Kohli, Joshua B. Tenenbaum, and Jiajun Wu. The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. 8
- [45] Nicolai Meinshausen. Causality from a distributional robustness point of view. In *2018 IEEE Data Science Workshop (DSW)*, pages 6–10. IEEE, 2018. 8
- [46] Junhyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. Learning from failure: Training debiased classi-

- fier from biased classifier. In *Advances in Neural Information Processing Systems*, 2020. [2](#), [4](#), [5](#), [6](#)
- [47] Hongseok Namkoong and John C. Duchi. Stochastic gradient methods for distributionally robust optimization with f-divergences. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 2208–2216, 2016. [3](#)
- [48] Peter A Noseworthy, Zachi I Attia, LaPrincess C Brewer, Sharonne N Hayes, Xiaoxi Yao, Suraj Kapa, Paul A Friedman, and Francisco Lopez-Jimenez. Assessing and mitigating bias in medical artificial intelligence: the effects of race and ethnicity on a deep learning model for ecg analysis. *Circulation: Arrhythmia and Electrophysiology*, 13(3):e007988, 2020. [1](#)
- [49] Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, pages 947–1012, 2016. [8](#)
- [50] Mohammad Pezeshki, Sékou-Oumar Kaba, Yoshua Bengio, Aaron Courville, Doina Precup, and Guillaume Lajoie. Gradient starvation: A learning proclivity in neural networks. *arXiv preprint arXiv:2011.09468*, 2020. [2](#), [4](#), [5](#), [6](#)
- [51] Oskar Pfüngst. *Clever Hans:(the horse of Mr. Von Osten.) a contribution to experimental animal and human psychology*. Holt, Rinehart and Winston, 1911. [2](#)
- [52] Sainandan Ramakrishnan, Aishwarya Agrawal, and Stefan Lee. Overcoming language priors in visual question answering with adversarial regularization. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 1548–1558, 2018. [3](#)
- [53] Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. The risks of invariant risk minimization. *arXiv preprint arXiv:2010.05761*, 2020. [3](#)
- [54] Swati Sachan, Jian-Bo Yang, Dong-Ling Xu, David Eraso Benavides, and Yang Li. An explainable ai decision-support-system to automate loan underwriting. *Expert Systems with Applications*, 144:113100, 2020. [1](#)
- [55] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*, 2019. [2](#), [3](#), [4](#), [5](#), [6](#)
- [56] Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. An investigation of why overparameterization exacerbates spurious correlations. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 8346–8356. PMLR, 2020. [4](#), [5](#), [6](#), [8](#)
- [57] Victor Sanh, Thomas Wolf, Yonatan Belinkov, and Alexander M Rush. Learning from others’ mistakes: Avoiding dataset biases without modeling them. *arXiv preprint arXiv:2012.01300*, 2020. [4](#)
- [58] Ramprasaath Ramasamy Selvaraju, Stefan Lee, Yilin Shen, Hongxia Jin, Shalini Ghosh, Larry P. Heck, Dhruv Batra, and Devi Parikh. Taking a HINT: leveraging explanations to make vision and language models more grounded. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages 2591–2600. IEEE, 2019. [5](#)
- [59] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The pitfalls of simplicity bias in neural networks. *arXiv preprint arXiv:2006.07710*, 2020. [4](#)
- [60] Jiaxin Shi, Hanwang Zhang, and Juanzi Li. Explainable and explicit visual reasoning over scene graphs. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 8376–8384. Computer Vision Foundation / IEEE, 2019. [8](#)
- [61] Robik Shrestha, Kushal Kafle, and Christopher Kanan. Answer them all! toward universal visual question answering models. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 10472–10481. Computer Vision Foundation / IEEE, 2019. [2](#)
- [62] Robik Shrestha, Kushal Kafle, and Christopher Kanan. A negative case analysis of visual grounding methods for VQA. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8172–8181, Online, July 2020. Association for Computational Linguistics. [2](#)
- [63] Damien Teney, Ehsan Abbasnejad, and Anton van den Hengel. Unshuffling data for improved generalization. *arXiv preprint arXiv:2002.11894*, 2020. [3](#)
- [64] Damien Teney, Kushal Kafle, Robik Shrestha, Ehsan Abbasnejad, Christopher Kanan, and Anton van den Hengel. On the value of out-of-distribution testing: An example of goodhart’s law. 2020. [2](#), [7](#)
- [65] Prasetya Ajie Utama, Nafise Sadat Moosavi, and Iryna Gurevych. Towards debiasing NLU models from unknown biases. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7597–7610, Online, 2020. Association for Computational Linguistics. [2](#), [4](#)
- [66] Jialin Wu, Liyan Chen, and Raymond J Mooney. Improving vqa and its explanations by comparing competing explanations. *arXiv preprint arXiv:2006.15631*, 2020. [5](#)
- [67] Chuanlong Xie, Fei Chen, Yue Liu, and Zhenguo Li. Risk variance penalization: From distributional robustness to causality. *arXiv preprint arXiv:2006.07544*, 2020. [3](#)
- [68] Kexin Yi, Jiajun Wu, Chuang Gan, Antonio Torralba, Pushmeet Kohli, and Josh Tenenbaum. Neural-symbolic VQA: disentangling reasoning from vision and language understanding. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing*

Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pages 1039–1050, 2018. [8](#)

- [69] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 335–340, 2018. [2](#), [3](#)
- [70] Zhilu Zhang and Mert R. Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 8792–8802, 2018. [4](#)
- [71] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2979–2989, Copenhagen, Denmark, 2017. Association for Computational Linguistics. [1](#)
- [72] Yang Zou, Zhiding Yu, BVK Vijaya Kumar, and Jinsong Wang. Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In *Proceedings of the European conference on computer vision (ECCV)*, pages 289–305, 2018. [3](#)