

On The Impact of Machine Learning Randomness on Group Fairness

Prakhar Ganesh

National University of Singapore
Singapore
pganesh@comp.nus.edu.sg

Martin Strobel

National University of Singapore
Singapore
martin.r.strobel@gmail.com

Hongyan Chang

National University of Singapore
Singapore
hongyan@comp.nus.edu.sg

Reza Shokri

National University of Singapore
Singapore
reza@comp.nus.edu.sg

ABSTRACT

Statistical measures for group fairness in machine learning reflect the gap in performance of algorithms across different groups. These measures, however, exhibit a high variance between different training instances, which makes them unreliable for empirical evaluation of fairness. What causes this high variance? We investigate the impact on group fairness of different sources of randomness in training neural networks. We show that the variance in group fairness measures is rooted in the high volatility of the learning process on *under-represented groups*. Further, we recognize the dominant source of randomness as the stochasticity of *data order* during training. Based on these findings, we show how one can control group-level accuracy (i.e., model fairness), with high efficiency and negligible impact on the model's overall performance, by simply changing the data order for a single epoch.

CCS CONCEPTS

• **Computing methodologies** → **Machine learning**; • **General and reference** → *Evaluation*.

KEYWORDS

neural networks, fairness, randomness in training, evaluation

ACM Reference Format:

Prakhar Ganesh, Hongyan Chang, Martin Strobel, and Reza Shokri. 2023. On The Impact of Machine Learning Randomness on Group Fairness. In *2023 ACM Conference on Fairness, Accountability, and Transparency (FAccT '23)*, June 12–15, 2023, Chicago, IL, USA. ACM, New York, NY, USA, 31 pages. <https://doi.org/10.1145/3593013.3594116>

1 INTRODUCTION

Machine learning models are shown to manifest and escalate historical biases present in their training data [1, 4, 16, 59]. Understanding these biases and the resulting ethical obligations have led to the rise of fair machine learning research [13, 15, 37]. However, recent work

has observed high variance in fairness measures across multiple training runs, usually attributed to non-determinism in training (e.g., weight initialization, data reshuffling, etc.). These findings challenge the effectiveness of many bias mitigation algorithms [3, 48], and even the legitimacy of several fairness trends present in literature [51]. Thus, a reliable extraction of fairness measures requires accounting for the high variance due to randomness in the learning process to avoid lottery winners (see Fig. 1).

The standard solution to this concern is executing a large number of training runs with different randomness. However, such a solution creates huge computational demands when examining biases in neural networks. For instance, it costs about \$450K to train a model of similar quality as GPT-3 [57], and thus executing multiple training runs of such a model is not practical. But, are multiple identical runs essential? Can we instead find an efficient alternative to measure this variance? Our paper answers this critical yet unsolved question.

In this work, we perform an empirical investigation into the high fairness variance due to randomness in neural network training, with a diverse set of experiments on a multitude of settings, including different datasets across modalities, various fairness metrics, and changing hyperparameters and model architecture. More specifically, our empirical analysis answers the following questions

- **Is there a dominant source of randomness?** We show that the fairness variance observed in the literature is dominated by randomness due to data reshuffling during training. Reshuffling causes large changes in fairness even between consecutive epochs within a single run, while other forms of randomness have minimal influence.
- **Why are fairness measures highly sensitive to data reshuffling?** We show a higher vulnerability of minorities to changing model behavior, i.e., a higher prediction uncertainty for under-represented groups. This disparate prediction uncertainty between groups is reflected in any statistical fairness measure defined on model predictions.
- **How does data order impact fairness?** We demonstrate an immediate impact of the data order on fairness. That is, we show that a model's fairness score is heavily influenced by the most recent gradient updates, irrespective of the preceding training. We also demonstrate how to create custom data orders that can efficiently control group-level performances (and thus in turn, model fairness), with a minor impact on the overall accuracy.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

FAccT '23, June 12–15, 2023, Chicago, IL, USA

© 2023 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-0192-4/23/06.

<https://doi.org/10.1145/3593013.3594116>

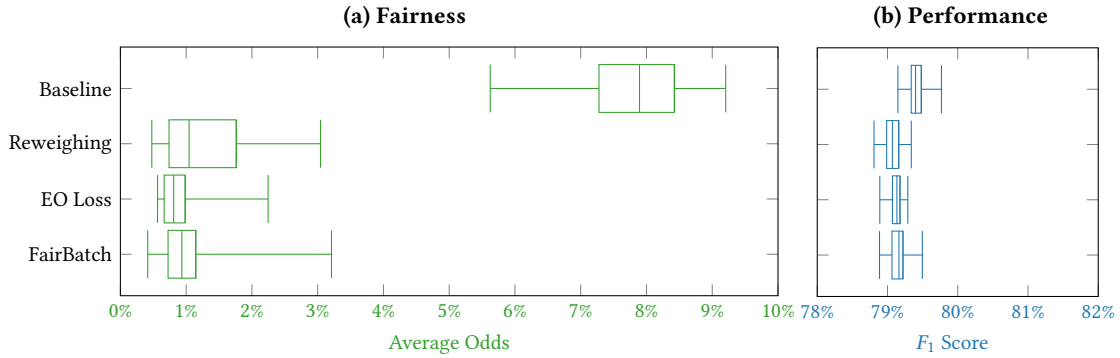


Figure 1: Variance of fairness and accuracy: (a) Fairness (average odds) has a high variance across multiple runs due to non-determinism in training. This variance persists even with state-of-the-art bias mitigation algorithms (Reweighing [30]; Equalized Odds Loss [22]; FairBatch [46]) (b) The overall performance (F_1 score), however, has a significantly smaller range of variance.

• What are the practical implications?

- Given the immediate impact of data order on model fairness and the nature of data order reshuffling in neural network training, we propose that using fairness variance across epochs in a *single* training run is a good proxy to study fairness variance across multiple runs, thus reducing the computational requirements by a significant factor.
- We also propose a custom data order that can improve model fairness within a single training epoch, and compete with existing bias mitigation algorithms. Interestingly, we show that similar custom data orders can also be created by adversaries to *freely control* fairness gaps in only a single epoch of training, even under explicit bias mitigation.

2 BACKGROUND AND RELATED WORK

In this section, we first introduce the relevant background in fair machine learning and randomness in neural network training. We then discuss the related work on the impact of randomness on fairness evaluation in deep learning.

2.1 Fairness in Machine Learning

Fair machine learning can be broadly divided into two categories, (i) group fairness [15], and (ii) individual fairness [20]. Group fairness relies on measuring the disparity between the average performance of protected groups against other privileged groups, and thus focuses on highlighting systematic bias against certain groups. Individual fairness instead relies on some form of similarity between individuals and requires consistency in the decision-making, i.e., similar individuals should be treated similarly.

In this work, we focus specifically on group fairness. Group fairness has a diverse set of definitions in the literature, usually chosen based on the stakeholders involved, known even to have opposing behavior in specific settings [40, 47]. We will rely on three commonly used group fairness metrics, i.e., demographic parity, average odds, and equal opportunity [26]. Demographic parity is the measure of disparity between the percentage of positive outcomes for each group, i.e., it does not allow model predictions to depend on sensitive attributes. Average odds (and its relaxed version, equal

opportunity) is instead a measure of disparity between predictions for each group conditioned on the true labels, i.e., it does allow overall predictions to depend on sensitive attributes, but does not allow predictions for certain ground-truth labels to depend on sensitive attributes. Bias calculation and mitigation for group fairness have accumulated extensive literature in recent years, along with many open-source benchmarks [5, 8, 22, 45].

2.2 Randomness in Neural Network Training

Deep learning involves various forms of randomness that impact a neural network’s path to convergence. This randomness during training can introduce noise into the optimization objective and works as a regularizer for the learning algorithm [42]. It makes the model prioritize generalization, avoid overfitting, escape local minima, and even speed up convergence [12]. Thus, randomness during training is integral to the success of neural networks, but its impact on model behavior needs to be carefully examined [7, 10].

Broadly, randomness in neural networks can be studied in the context of the following categories (see Fig. 2),

- **Data Splitting:** For any experimental setup in machine learning, the dataset under consideration is randomly divided into train-val-test (or just train-test) splits, to avoid information leakage and perform a fair evaluation.
- **Weight Initialization:** Weight initialization refers to the initial parameter vector that is the starting point for the gradient descent. Randomness in weight initialization is crucial for breaking the symmetry between model parameters and allows the neural network to learn complex functions [25].
- **Random Reshuffling:** Neural network training relies on gradient descent to optimize a chosen objective iteratively. Calculating the gradient over the entire dataset for every optimization step is expensive. A commonly adopted alternative is uniformly sampling a subset of the dataset to approximate the gradient, known as stochastic gradient descent (SGD). In practice, it has been shown that instead of uniform sampling, SGD can also be implemented by simply traversing a random order, i.e., random reshuffling, of the training data [12, 38, 41].

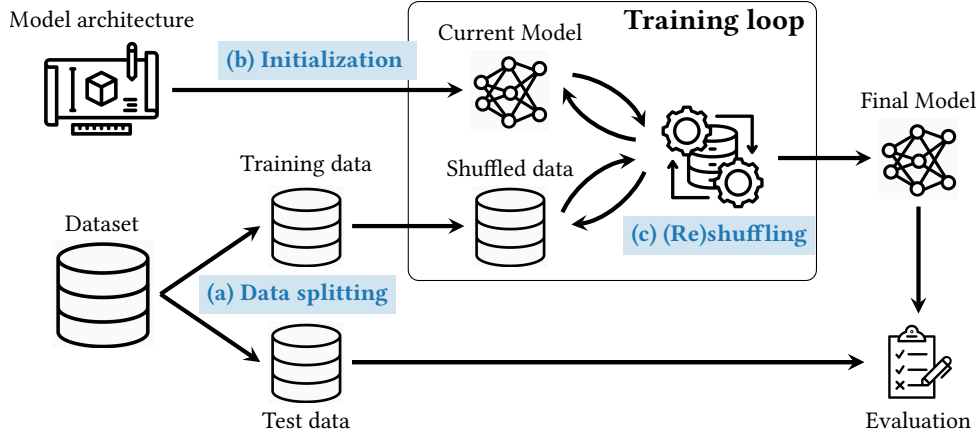


Figure 2: Sources of randomness during training: Randomness is introduced at several points during the training of a neural network. (a) Data splitting before training is important to avoid information leakage during evaluation, and involves randomness for the same. (b) Randomness in weight initialization at the start of the training is necessary to break symmetry and learn complex representations. (c) Random reshuffling at every epoch implements stochastic gradient descent (SGD) in practice. There are many other forms of randomness in neural network training. Yet, most are only used for certain specific settings.

- **Other Sources of Randomness:** Several additional components in the learning algorithm can introduce further randomness. These components are not standard, but are introduced in special cases to achieve specific objectives. For example, data augmentation to increase dataset size [31], dropout to perform regularization [58], gradient noise for private training [14], etc.

2.3 High Variance in Fair Deep Learning

There has been a growing awareness of high variance in fair deep learning, associated with non-determinism in model training or the underlying implementation [3, 21, 43, 48, 51], and the uncertainty of existing results in the literature.

Soares et al. [51] investigate the relationship between various algorithmic choices and the corresponding fairness variance, in large language models. They study the correlation of fairness with model size and found no obvious trends, as opposed to various claims previously made in literature [6, 28]. They also found that fairness is heavily affected by the random seed, i.e. simply changing the randomness can cause a huge variance in fairness.

Sellam et al. [48] have trained and released 25 pre-trained BERT checkpoints, each trained from scratch under identical settings but with a different random seed. They also analyze the variance of model fairness and the impact of commonly used bias mitigation algorithms on downstream tasks when starting with different pre-trained models. They show significant variance across changing random seeds and question the value of such mitigation techniques.

Amir et al. [3] revisit bias mitigation techniques in clinical texts and show a lack of statistically significant improvement after accounting for non-determinism in training. Friedler et al. [21] explore the stability of fairness under a rarely studied source of randomness, i.e. data splitting, and show notable impact on fairness evaluation.

While existing literature focuses on exploring the impact of high fairness variance in bias evaluation, we instead focus on investigating its source. Furthermore, we propose to move away from

the practice of simply executing multiple runs to capture fairness variance and instead provide a computationally efficient proxy.

3 PROBLEM STATEMENT

We start by formally defining the problem statement and detailing our experiment setting for the rest of the paper.

3.1 Neural Network Training

Most machine learning algorithms can be abstracted down to an optimization problem for a given objective, usually a loss function. More specifically, for a training dataset $(x, y) \in \mathcal{D}$, a family of hypothesis functions \mathcal{F} , and a loss function \mathcal{L} , the optimization goal for the learning algorithm can be defined as,

$$f^* \leftarrow \operatorname{argmin}_{f \in \mathcal{F}} \sum_{(x, y) \in \mathcal{D}} \mathcal{L}(f(x), y) \quad (1)$$

The above formulation of the learning objective is known as empirical risk minimization (ERM) [56]. However, finding a global optimum for ERM in deep learning is typically intractable, due to the high dimensional, non-convex formulation of neural networks. Neural networks are instead trained iteratively, starting with a randomly sampled function f_0 , refining the model with a learning algorithm \mathcal{A} for T epochs, to finally output the trained model f_T . The learning algorithm at every epoch t takes in the current model, complete training data \mathcal{D} , and a number of hyperparameters ξ (e.g., batch size, learning rate, etc.), to progressively improve the model by a single epoch of training. The learning algorithm can contain various sources of randomness, as discussed above. In our work, we will focus on two standard forms of randomness found in every neural network training, i.e., weight initialization and random reshuffling of data order at every epoch. More specifically, neural network training can be defined as,

$$f_t := \mathcal{A}(f_{t-1}, \mathcal{D}, \xi, r_s, t) \quad f_0 \sim \mathcal{F}; r_s \sim R \quad (2)$$

The function f_0 , i.e., the weight vector initialization in a parameterized neural network, is randomly sampled from pre-defined distribution \mathcal{F} , and the random seed for reshuffling r_s is sampled from a uniform distribution R . Note that both random seed r_s and epoch number t are together responsible for the data shuffling of epoch t . Thus, for fixed reshuffling (i.e., fixed r_s), the data order is still shuffled at every epoch during a single training run but is the same at any epoch t across two different training runs. More details on the random seed setup can be found in Appendix A.

3.2 Metrics and Variance

A model f 's performance can be evaluated using its outputs on the test dataset. We will stick to the commonly used binary classification and binary sensitive attribute $a \in \{0, 1\}$ setting in fairness literature for the rest of the discussion. In the main paper, we will rely on F_1 Score and average odds (AO) to measure the model's average performance and group fairness respectively. AO can be empirically interpreted as the average disparity between separately calculated *true positive rates* (TPR) and *false positive rates* (FPR) of various groups [26]. The metrics are defined as

$$F_1(f, \mathcal{D}) := \frac{2TP}{P + PP} = \frac{2 \sum_{\mathcal{D}^e} \mathbb{1}[f(x)=y \wedge y=1]}{\sum_{\mathcal{D}^e} \mathbb{1}[y=1] + \sum_{\mathcal{D}^e} \mathbb{1}[f(x)=1]} \quad (3)$$

$$\begin{aligned} AO(f, \mathcal{D}) &:= \frac{(\Delta TPR + \Delta FPR)}{2} \\ &= \frac{1}{2} \sum_{\substack{h=\{0,1\} \\ y=h}} \left| \frac{\sum_{\mathcal{D}^e} \mathbb{1}[f(x)=1 \wedge a=0]}{\sum_{\mathcal{D}^e} \mathbb{1}[a=0]} - \frac{\sum_{\mathcal{D}^e} \mathbb{1}[f(x)=1 \wedge a=1]}{\sum_{\mathcal{D}^e} \mathbb{1}[a=1]} \right| \end{aligned} \quad (4)$$

where $\sum_{\mathcal{D}^e}$ is the sum over all data points in the test set \mathcal{D}^e , i.e., $\sum_{(x,y,a) \in \mathcal{D}^e}$, and $\mathbb{1}[z]$ is an indicator function which is 1 when the boolean expression z is true, and 0 otherwise. We also include additional experiments for two more fairness measures, equal opportunity (EOpp) and demographic parity (DP) in Appendix I. Moreover, we will show that the non-determinism in fairness originates from high prediction uncertainty for minority (Section 5), and thus will be reflected in any fairness metric defined on these predictions. We report all metrics in percentage.

At the heart of our work is the study of fairness variance across model checkpoints. We define variance across multiple runs and variance across epochs in a single run as,

$$Var_{F_1}^{runs}(\mathcal{A}, T) := \underset{f_0 \sim \mathcal{F}; r_s \sim R}{Var} (F_1(f_T)), \quad (5)$$

$$Var_{F_1}^{epochs}(\mathcal{A}, f_0, r_s, T_1, T_2) := \underset{t \in [T_1, T_2]}{Var} (F_1(f_t)), \quad (6)$$

$$Var_{AO}^{runs}(\mathcal{A}, T) := \underset{f_0 \sim \mathcal{F}; r_s \sim R}{Var} (AO(f_T)), \quad (7)$$

$$Var_{AO}^{epochs}(\mathcal{A}, f_0, r_s, T_1, T_2) := \underset{t \in [T_1, T_2]}{Var} (AO(f_t)), \quad (8)$$

Existing work in the literature has shown high variance in fairness scores across multiple runs $Var_{AO}^{runs}(\mathcal{A}, T)$. In our work, we first decouple the impact of two standard sources of randomness, i.e., study $Var_{f_0 \sim \mathcal{F}}(F_1(f_T))$ and $Var_{r_s \sim R}(F_1(f_T))$ separately. In doing so, we find high variance in fairness scores even across epochs in a single training run (Section 4), and thus further study variance across epochs in fairness scores $Var_{AO}^{epochs}(\mathcal{A}, f_0, r_s, T_1, T_2)$. Note

that for F_1 score, variance across multiple runs $Var_{F_1}^{runs}(\mathcal{A}, T)$ and across epochs in a single run $Var_{F_1}^{epochs}(\mathcal{A}, f_0, r_s, T_1, T_2)$ are both relatively stable. Unless otherwise specified, we train our models for a total of $T = 300$ epochs, and we measure variance across epochs from $T_1 = 100$ to $T_2 = 300$. We make this choice because the models have converged to stable accuracy before epoch 100 (refer to the training curve in Appendix C for more details).

3.3 Datasets and Models

We will conduct our investigation on ACSIncome and ACSEmployment tasks of the Folktables dataset [17], and binary classification of the 'smiling' label in CelebA dataset [33], with perceived gender (Female vs. Male) as the sensitive attribute for all datasets. For CelebA, input features are obtained by passing the image through a pre-trained ResNet-50 backbone and extracting the output feature vector. More details on the datasets are provided in Appendix B.

We train a feed-forward network with a single hidden layer of 64 neurons and *ReLU* activation, and train the model with *cross-entropy* (CE) loss for $T = 300$ epochs at batch size 128 and learning rate $1e-3$, in all our experiments unless specified otherwise. Note that while we measure fairness scores in our experiments, we do not explicitly train the models with any fairness constraints (except in Section 7.2 when training baseline bias mitigation algorithms). We also include additional experiments by changing training hyperparameters, i.e., batch size, learning rate, and model architecture, in Section 4.2 (and Appendix G). We use a train-val-test split of 0.7 : 0.1 : 0.2, and maintain the same split throughout all our experiments, i.e. we do not consider potential randomness due to data splitting. All our experiments and evaluations are performed only on the test split. We will focus primarily on the ACSIncome task in the main text, while additional experiments on CelebA and the ACSEmployment task are included in Appendix F.

4 THE DOMINANT SOURCE OF RANDOMNESS

In this section, we move past the observation that different training runs lead to different outcomes, and investigate the high fairness variance by studying and contrasting the two canonical sources of randomness.

4.1 Impact of Weight Initialization and Random Reshuffling on Fairness Variance

We start by decoupling the two sources of randomness inherent to the widely adapted SGD, i.e., weight initialization and random reshuffling, and study their impact on fairness variance separately in Fig. 3. We collect average odds (AO) and F_1 score at epoch 300 for 50 unique training runs each while, (i) allowing for both sources of randomness, (ii) changing only the weight initialization while keeping the random reshuffling fixed, and (iii) changing only the random reshuffling while keeping the weight initialization fixed, respectively. The large range of fairness scores reported by allowing for both sources of randomness in Fig. 3(a) represents the variance observed in the existing literature. Interestingly, when these sources are examined separately, the variance under fixed initialization in Fig. 3(c) is equivalently large but the variance under fixed reshuffling in Fig. 3(b) drops significantly. It is clear that fairness variance

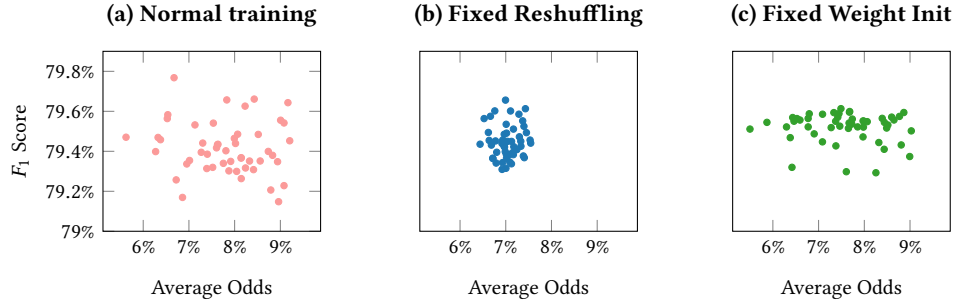


Figure 3: Decoupling the effect of randomness in weight initialization and reshuffling: (a) Variance in average odds (AO) by allowing both sources of randomness simultaneously represents the fairness variance in existing literature. (b) We see a significant drop in variance if we change only the weight initialization while keeping the reshuffling fixed. (c) However, we observe high range of variance by changing only the reshuffling, even for a fixed weight initialization. These results suggest reshuffling of the data order as the dominant source of fairness variance, with little influence from weight initialization.

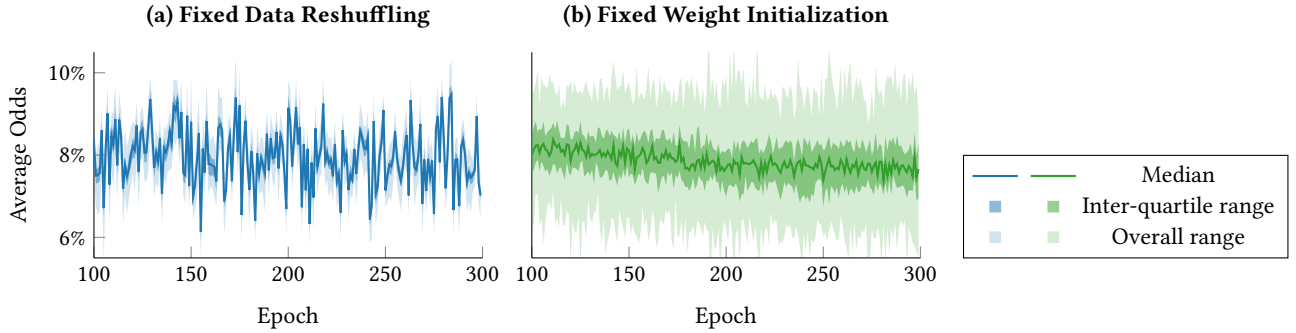


Figure 4: Training dynamics under fixed weight initialization and reshuffling: Median, inter-quartile range, and overall range of average odds across 50 training runs while keeping the data reshuffling or the weight initialization fixed respectively. (a) Despite different initializations, models with fixed data reshuffling have very little variance across training runs, but high variance across epochs. This highlights the dominant impact of random reshuffling on model fairness. (b) High variance even across training runs at the same epoch under fixed weight initialization further supports our claim.

originates from data order shuffling, while randomness in weight initialization has minimal impact.

To further probe the difference between these two sources of randomness, we study the training dynamics of the previous set of models across epochs, instead of just the final model checkpoint, in Fig. 4. We plot the median, inter-quartile range, and overall range of average odds (AO) across the complete set of 50 training runs from epoch 100 to 300 in the two isolated settings from the previous experiment. We find a high correlation in fairness scores across training runs with fixed data reshuffling (average pairwise pearson coefficient ≈ 0.94), which supports our observations of low variance in fairness scores at the final epoch in Fig. 3(b). Furthermore, there is a lack of any reasonable correlation between fairness scores of training runs with fixed weight initialization (average pairwise pearson coefficient ≈ 0.04). Interestingly, the high fairness variance across epochs inside a single training run in Fig. 4(a) closely matches the variance that we observe across multiple training runs in Fig. 4(b). In other words, for fixed data reshuffling the average odds value at any epoch is almost the same between different training runs, but the average change even between consecutive epochs

is large, while for fixed weight initialization, even the variance between runs is quite high.

4.2 Dominance of Random Reshuffling across Datasets, Metrics and Hyperparameters

We extend our previous experiment and calculate correlation across multiple runs for additional datasets, fairness measures, as well as hyperparameter choices of batch size, learning rate, model architecture, and dropout regularization with different dropout rates in Table 1. Here we measure the correlation (i.e., average pairwise pearson coefficient) across 50 training runs in each setting for fixed data shuffling and fixed weight initialization. It is clear from the results that even under diverse settings, the correlation between multiple runs with fixed data reshuffling is significantly high, while the correlation with fixed weight initialization is close to zero.

In addition to the overall trends supporting our initial claim, individual trends in Table 1 under various settings are also quite interesting. The correlation score for fixed weight initialization under hyperparameters that induce noisier training (i.e., smaller

Table 1: Average pairwise pearson coefficient for correlation across multiple runs: Fixed random reshuffling (RR) (i.e., changing only the weight initialization) has a high correlation score across multiple runs, while fixed weight initialization (WI) (i.e., changing only the random reshuffling) has a correlation score close to zero, which establishes the dominance of data reshuffling on fairness. These trends exist across different datasets, fairness metrics, and hyperparameter choices.

(a) Different Datasets		
	Fixed RR	Fixed WI
ACSIIncome	.94	.04
ACSEmployment	.89	.01
CelebA	.92	.01
(b) Different Fairness Metrics		
	Fixed RR	Fixed WI
Average Odds	.94	.04
Equal Opportunity	.94	.05
Demographic Parity	.96	.03
(c) Changing Hyperparameters		
	Fixed RR	Fixed WI
Default Hyperparameters	.94	.04
Batch Size = 16	.95	.00
Learning Rate = 0.01	.93	.00
Arch = {2048, 64}	.92	.00
No Dropout	.94	.04
Dropout Rate = 10%	.88	.03
Dropout Rate = 20%	.85	.04
Dropout Rate = 30%	.80	.13

batch size, larger learning rate, etc.) drops even lower than the default setup. Note that we do not report results for a bigger batch size or a smaller learning rate, as these models face issues with convergence even after 300 epoch of training (see Appendix G for more details), which further indicates the need of randomness and noise in the learning algorithm to facilitate better and faster convergence. However, this randomness will also create a dominant dependence of model fairness on data reshuffling, as we study in our work. We also see the trends diminishing (although still clear) with higher dropout rates. This suggests that appropriate regularization during training can indeed, to some extent, reduce the impact of randomness in training (or more specifically, data reshuffling) on fairness scores. We also provide detailed results for each hyperparameter setting above in Appendix G, H.

Takeaway 1: *Random reshuffling of data order during training is the dominant cause of high fairness variance as seen in the literature, while randomness in weight initialization has minimal influence.*

5 WHY IS FAIRNESS HIGHLY SENSITIVE TO RANDOMNESS?

In the previous section, we showed the dominant impact of data reshuffling on model fairness. In this section, we show that its in fact the imbalance in the underlying data distribution for training which creates high volatility in predictions for the minority. Thus, groups with smaller representations are more significantly influenced by the randomness in reshuffling, resulting in high fairness variance.

5.1 Changing Predictions Across Epochs

We observed high variance in model fairness even between consecutive epochs during training (see Fig. 4). The changing predictive behavior of neural networks beyond training loss convergence is not surprising, and has been studied extensively in literature [29, 32, 54, 55]. As we are concerned with the fairness of the final decisions made by the model, we will focus on a change in the model’s discrete output class when discussing changing predictions. More specifically, a model is said to have undergone a change

in prediction for some input x during epoch t , if $f_t(x) \neq f_{t-1}(x)$, where $f_t(x)$ is the output class when passing the input x through the model checkpoint at the end of epoch t . While these changing predictions maintain an overall stable average performance, they can still have a disparate impact on individual groups, the exact characteristics of which are less known.

We study this instability by investigating individual data points which change their predictions. We plot the dataset distribution across groups in Fig. 5(a) and the percentage of data points from each group that changed their prediction at least once between epochs 100 and k , where we gradually increase the value of k , in Fig. 5(b). Clearly, the trends in the percentage of unique data points with changing predictions mirror the representation of each group in the original dataset, i.e. the groups which are represented the least are the most vulnerable to changing model behavior. For example, positive labels from the group Female are severely under-represented and consequently have almost twice the percentage of unique examples with changing predictions than any other groups.

5.2 Disparate Prediction Uncertainty

Higher vulnerability to changing discrete predictions for minorities can be interpreted as an indication of higher uncertainty in the underlying model predictions. To further probe the disparate model behavior across different groups, we record the cumulative distribution of prediction uncertainty for each group separately in Fig. 6. We rely on two commonly used methods to measure prediction uncertainty, i.e., (i) monte-carlo dropout [23], and (ii) training a bayesian neural network [11]. We execute 1000 forward passes for each method and record the standard deviation in outputs as the prediction uncertainty. Despite different distribution trends, both methods highlight the higher prediction uncertainty of minorities. As expected, the order of prediction uncertainty across groups follows the training data distribution (Fig. 5(a)), i.e. groups with a larger representation in the training data have smaller number of examples with large prediction uncertainty, which is quite intuitive.

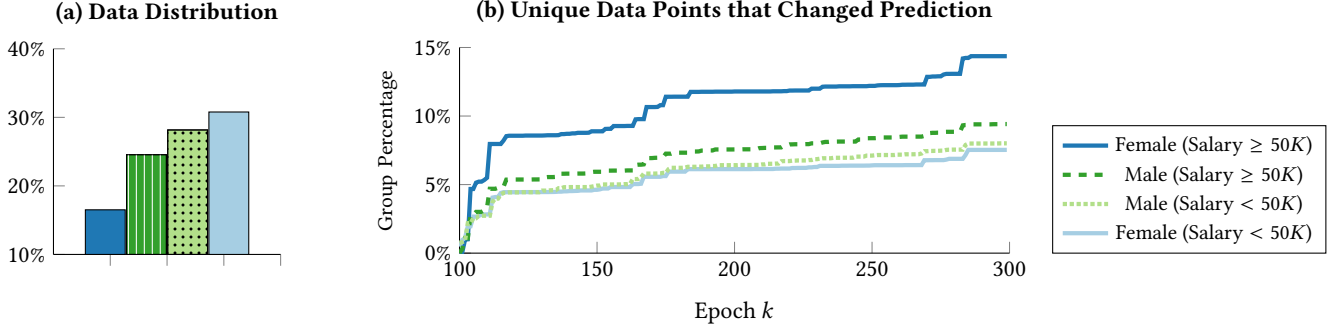


Figure 5: Disparate percentage of changing predictions across epochs: (a) The underlying data distribution of ACSIncome shows positive labels from group Female as an under-represented minority. (b) Total percentage of unique data points from each subgroup that change prediction across epochs follow the opposite order of their representation in the training data. These results highlight subgroups with least representation being the most vulnerable to changing predictions.

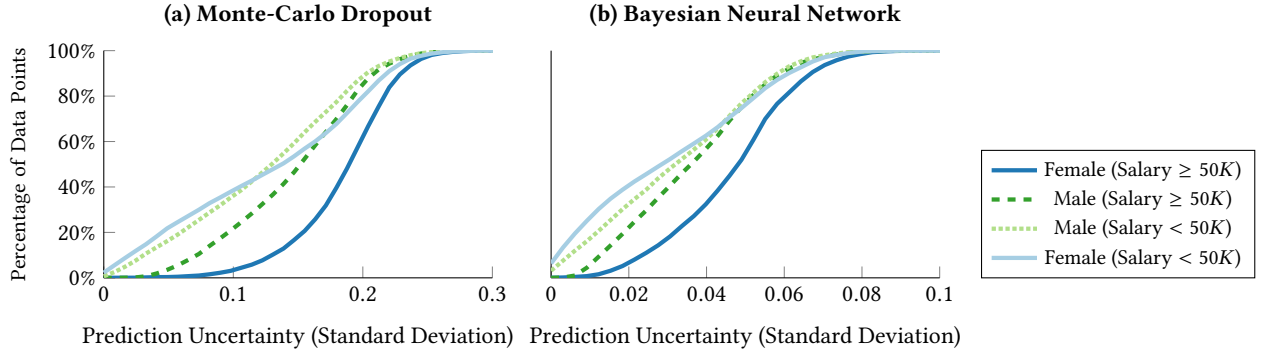


Figure 6: Normalized cumulative distribution of prediction uncertainty for various groups: The distribution of the minority group (Female with Salary $\geq 50K$) is significantly more skewed towards higher uncertainty than any other group, i.e. the minority contains far more percentage of data points with high prediction uncertainty than the majority.

As the model has skewed cumulative distribution towards uncertainty for the minority, any fairness metric defined on the output of such a model will also reflect this instability, and thus manifests as fairness variance in existing literature [3, 48, 51].

Takeaway 2: Under-represented groups have higher prediction uncertainty in the final trained model, and thus predictions for data points from such minorities are more sensitive to the randomness.

6 IMPACT OF DATA ORDER ON FAIRNESS

In Section 4, we observed the dominance of reshuffling on fairness variance. Data order during training governs gradient updates and thus its impact on fairness is unsurprising [34, 39, 44, 49, 50, 52]. Even under randomly shuffled data order, neural networks are known to undergo changes in predictive behavior during training [29, 32, 54, 55]. However, it is the immediacy of the impact of data order that is surprising and a novel observation of our work. We now study the impact of data order in a single epoch on fairness.

6.1 Data Order’s Immediate Impact

To study the immediacy and characteristics of the impact of data order on model fairness, we fine-tune a set of already converged

model checkpoints for a common sequence of b batches and record the fairness variance across checkpoints for different values of b in Fig. 7. This allows us to measure fairness variance across models which have experienced the same most recent b gradient updates. The choice of the common sequence of b batches fed to the model was done by separately training a model and choosing the suffix of data order corresponding to epochs with the best and worst fairness scores on the validation set. As the number of fixed batches b increases, the fairness variance decreases, and it is clear from the results that the impact of data order on fairness is quite immediate (an epoch of ACSIncome dataset is $b = 1070$ batches). Moreover, the resulting fairness is also characteristically stable for a specific data order, i.e. batches taken from the suffix of the data order corresponding to the best fairness epoch of an individual training run also helps fine-tune all other checkpoints towards the same best fairness, and vice-versa for the worst fairness epoch.

The set of checkpoints were chosen by sampling 1000 different checkpoints from epochs 100 to 300 of 50 different training runs while allowing for both forms of randomness simultaneously. The checkpoints were chosen in this manner to create a diverse training history, and show that these models achieve the same fairness

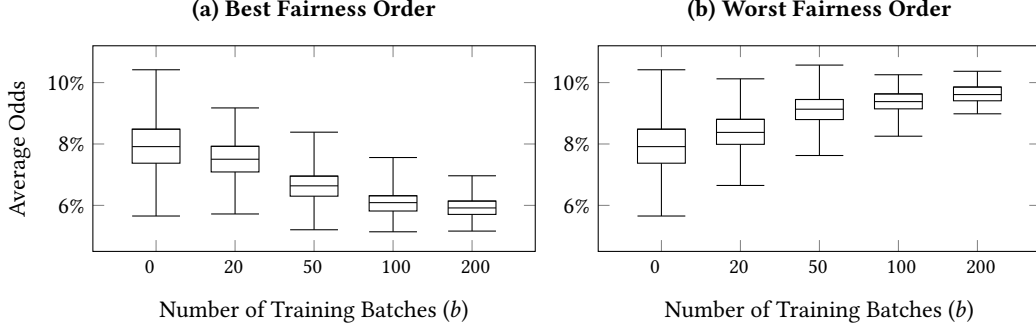


Figure 7: Fairness variance under common b most recent gradients: Average odds stabilize as the number of most recent common training batches b increases, highlighting the immediate impact of these gradient updates on model fairness. Moreover, it even predictably stabilizes to low or high fairness based on the corresponding data order from the epoch with best or worst fairness.

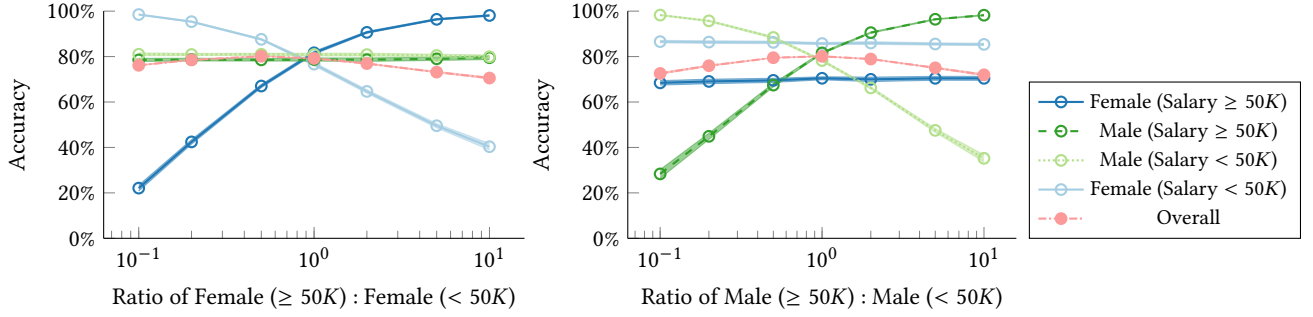


Figure 8: Manipulating group level accuracy with data order: We show how to control group level accuracy by changing data distribution of the most recent gradient updates, tested separately for ratio between positive and negative labels for group Female, and group Male, respectively, while keeping other ratios fixed. In only a single epoch of training, we are able to manipulate the group level accuracy trade-off, with relatively small impact on overall accuracy.

scores when fine-tuned on a common set of batches, irrespective of training prior to those updates. We also extend our experiment to show the same behavior even for batches from any random data order, instead of deliberately chosen data orders as above, in Appendix D. One possible explanation of this immediate impact is the presence of no energy valleys in deep learning loss landscape between minima of separately trained models [18, 24]. Another possible explanation can be built on a recent line of work that shows there is only one functionally unique minima in loss landscape of neural networks, while all other minima simply contain permutation or scale symmetries of the same set of models [2].

6.2 Manipulating Group Accuracy Distribution with Data Order

In the previous section, we saw a stable relationship between data order and the resulting fairness score. However, the data orders in the experiment above were sampled from a set of random data orders. We now show that it is possible to create our own custom data order to achieve any target fairness score. We hypothesize that since the data order in the most recent batches has an immediate impact on model fairness, the distribution in these batches must be temporarily changing the loss landscape and nudging the

group-level accuracy. This could allow us to manipulate group-level accuracy in only a single epoch of fine-tuning.

To test this hypothesis, we fine-tune a set of 50 already converged models for exactly a single epoch on custom data orders with chosen distributions and record group-level accuracy in each setting in Fig. 8. To create this custom data order, we start by fixing the ratio between different groups and then form batches for the data order suffix in this exact ratio until we run out of data points (which will happen for any ratio that is not the dataset distribution). The excess data points are then shuffled randomly and placed at the prefix of data order. We do two sets of isolated experiments, by changing the ratio between positive and negative labels for group Female (Fig. 8(a)) and group Male (Fig. 8(b)) respectively, while keeping other ratios fixed. It is clear that by manipulating the data distribution in the most recent gradient updates, we can control the group-level accuracy of the model. While the overall accuracy drops noticeably for extreme ratios, it does not change much in the middle despite significant variance in group-level performances.

Note that our custom data order still has the same distribution as the original dataset, i.e., we have not changed the distribution of the complete data order, but only moved around the data points to change the distribution of the data order suffix. These results further

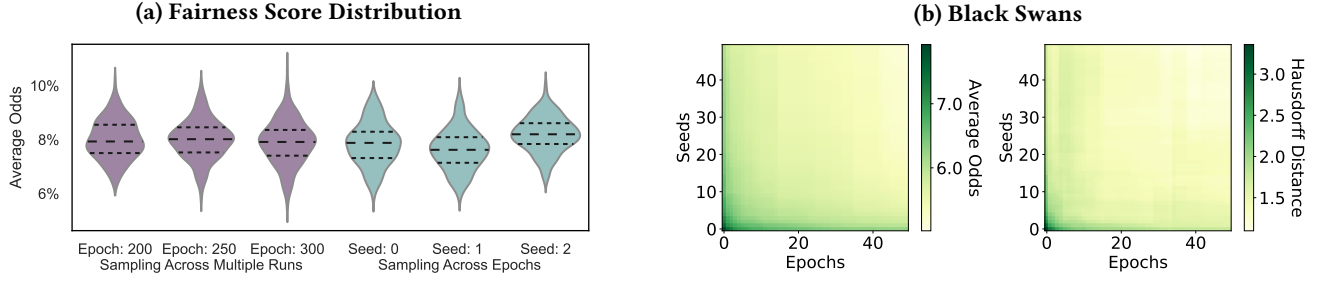


Figure 9: Fairness variance across multiple runs vs across epochs in a single run: (a) Fairness scores (average odds) across multiple training runs and across epochs in a single training run have similar empirical distributions. (b) Quality of black swans (i.e. extremely rare checkpoints) improves with more checkpoints collected, either in terms of fairness (the lowest achievable average odds score) or the trade-off between overall performance and fairness (the Hausdorff distance to the best pareto front). This improvement occurs at the same rate, irrespective of sampling checkpoints across multiple random seeds (x-axis) or multiple epochs in a single seed (y-axis). But, sampling across epochs is significantly cheaper, providing a highly efficient alternative to executing multiple runs.

strengthen our claim on the immediate impact of the most recent gradient updates on model fairness and group-level accuracy. *Fairbatch* [46], a recently proposed bias mitigation algorithm, follows a similar formulation (although it additionally changes the overall distribution). Fairbatch creates batches with a fixed ratio between groups, and this ratio is continuously optimized to counter the existing bias in the model. Our results not only explain their success, but also state that instead of regularly adapting the distribution to compensate for the model bias (and oversampling/undersampling certain groups), one can directly use the desired distribution to create a custom data order and the model will adapt to it immediately.

Takeaway 3: *The training data order has an immediate impact on the model’s fairness scores. That is, the data distribution in the most recent gradient updates can control the model’s group level accuracy in only a single epoch of fine-tuning.*

7 APPLICATIONS OF THE IMPACT OF DATA ORDER ON FAIRNESS

With a better understanding of the impact of data order on model fairness and how to control it, we now explore some practical applications of our observations.

7.1 Capturing Fairness Variance in a Single Run

We now return to our original problem of capturing fairness variance without wasting computing resources on a large number of training runs. We saw an immediate impact of data order on model fairness, which shows that fairness variance across multiple training runs can instead be studied as simply the randomness in data order at their last epochs. As these orders randomly reshuffled, their distribution across multiple trainings should be the same as their distribution across epochs in a single training run. Thus, we propose evaluating fairness of intermediate checkpoints in a single training run as a proxy for multiple runs.

To test the similarity in both distributions, we simply plot the distribution of fairness scores for checkpoints across 200 training runs (for three different stopping epochs) and across epochs 100 to 300 in a single training run (for three different training runs)

in Fig. 9(a). Clearly, the empirical distribution of fairness across multiple training runs closely matches the distribution across a single training run. We also perform the Kolmogorov–Smirnov (KS) test [35] to match sampling across multiple runs (at epoch 300) and sampling across epochs (for a single training run with seed 0). The maximum difference in empirical CDF of the two distributions was only 0.07, and the KS test gave the p-value of 0.712, i.e. the probability of the hypothesis that both set of fairness scores were sampled from the same underlying distribution.

Furthermore, we also test the quality of black swans, i.e. the best models under certain quality measure, as a function of number of unique training runs and number of epochs evaluated per training run, in Fig. 9(b). For all $(t, s) \in [1, 50]$, we perform s unique training runs (while changing both forms of randomness, i.e., weight initialization and random reshuffling), and evaluate the model for last t epochs in each training run, thus accumulating a total of $t * s$ checkpoints. We then calculate two different quality measures for these set of checkpoints, i.e., (i) the best fairness achieved across all checkpoints, and (ii) the Hausdorff distance [9] of the Pareto-front (including both fairness and F_1 scores) from the best achievable Pareto-front, i.e., for $t = 50; s = 50$. Finally, the experiment is repeated and averaged 50 times to compensate for randomness in the s training runs. Interestingly, the black swans for both quality measures show no significant distinction between increasing the number of training runs or evaluating multiple epochs per training run, i.e., sampling more checkpoints in either direction gives us similar improvements.

It is clear from our results that the commonly used method to capture fairness variance in literature ($t = 1; s = 50$) is highly inefficient use of computing resources, and one can extract the same quality of black swans (and overall fairness variance) by simply observing fairness across multiple checkpoints in a single training run ($t = 50; s = 1$), which would require 50 times less computation. With these experiments, we showed direct benefits of evaluating multiple epochs in a single training run, saving huge amounts of resources and time in capturing the overall fairness variance.

Takeaway 4: *Fairness distribution across multiple runs is empirically the same as that across epochs within a single run.*

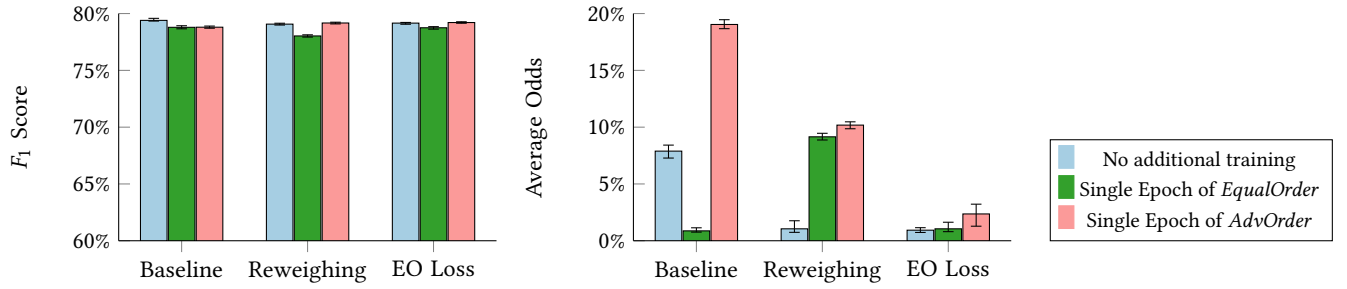


Figure 10: Comparing data order manipulation with other bias mitigation algorithms: Under baseline training setup, by changing the data order for just a single epoch of fine-tuning, *EqualOrder* gets competitive performance to commonly used bias mitigation methods. Similarly, *AdvOrder* gets significantly worse fairness than even the baseline. For Reweighing, we see an increase in bias even with *EqualOrder*, as the beneficial weighing of minority creates bias towards the majority. On the other hand, even reweighing is not enough to counter the effects of *AdvOrder*. Finally, equalized odds loss is capable of dynamically adapting to the model’s changing predictive behavior, yet we still observe increase in bias under *AdvOrder*.

7.2 Bias Mitigation via Data Order Manipulation

To measure the effectiveness of our group accuracy manipulation, we extend the discussion from Section 6.2 for two special ratios, 1 : 1 and 1 : 3 between positive and negative labels of group Female (for ACSIncome dataset), and call them *EqualOrder* and *AdvOrder* respectively. More specifically, we fine-tune converged models with a single epoch of *EqualOrder* (and *AdvOrder*), and record F_1 score and average odds in Fig. 10. We perform experiments with three unique setups, (i) Baseline training, (ii) Reweighing [30], a data pre-processing which weighs every label-group pair based on its representation in the overall dataset, and (iii) Equalized Odds Loss [22], an in-processing loss function to nudge the model towards fair predictions. By training with *EqualOrder* for a single epoch, the baseline model achieves competitive fairness scores to other bias mitigation algorithms. On the other hand, using *AdvOrder* can further push the model bias, emphasizing the adversarial power of data ordering, even in presence of explicit mitigation algorithms.

Notably, reweighing suffers from an unexpected high bias even under *EqualOrder*, as the combination of ideally distributed data order suffix along with increase in the minority data weights pushes the model towards significant unfair behavior against the majority (as opposed to against the minority in all other results). Moreover, equalized odds loss shows controlled damage under *AdvOrder* due to the loss function regularly adapting to the degrading behavior, but the unfairness still increases. *AdvOrder* is dangerous as it still maintains the overall accuracy, but favors the majority. We can force even worse fairness gaps by pushing the ratio to its extreme, however that will impact the model’s overall accuracy. These results cement the effectiveness of manipulating group level accuracy by controlling the data order for just a single epoch of fine-tuning.

Takeaway 5: A data order with a balanced suffix can significantly improve in fairness scores. Similarly, even bias mitigation algorithms can fail when trained with an adversarial data order.

8 CONCLUSION

Fairness variance due to changing randomness in deep learning has raised concerns regarding the reliability of existing results in

literature [3, 48, 51]. In our work, we took a close look at various sources of randomness, and found a dominant impact of data order on model fairness, which we showed was in turn due to a higher prediction uncertainty of the trained model on under-represented groups in the dataset. We further demonstrated that the distribution seen by the model in the most recent gradient updates can be easily exploited to achieve desirable group-level accuracy behavior, and proposed several practical applications of this immediate impact of data order on model fairness, including a highly efficient alternative to executing multiple training runs when studying fairness variance due to randomness in training.

In our work, we focused only on the discrete decisions made by the model, as we were investigating the impact of non-determinism in model training on its fairness. However, further extensions of this discussion to trends in the internal state of the learned model can reveal even granular characteristics, and has potential application in similar fields of research, for example, understanding high variance in out-of-distribution generalization [36], exploiting model multiplicity under various settings [10], and many more.

ACKNOWLEDGMENTS

This research is supported by Google PDPO faculty research award, Intel within the www.private-ai.org center, Meta faculty research award, the NUS Early Career Research Award (NUS ECRA award number NUS ECRA FY19 P16), and the National Research Foundation, Singapore under its Strategic Capability Research Centres Funding Initiative. Any opinions, findings, conclusions, or recommendations expressed in this material are those of the author(s) and do not reflect the views of the National Research Foundation, Singapore.

REFERENCES

- [1] Mohsen Abbasi, Sorelle A Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. 2019. Fairness in representation: quantifying stereotyping as a representational harm. In *Proceedings of the 2019 SIAM International Conference on Data Mining*. SIAM, 801–809.
- [2] Samuel K Ainsworth, Jonathan Hayase, and Siddhartha Srinivasa. 2022. Git re-basin: Merging models modulo permutation symmetries. *arXiv preprint arXiv:2209.04836* (2022).

- [3] Silvio Amir, Jan-Willem van de Meent, and Byron C Wallace. 2021. On the Impact of Random Seeds on the Fairness of Clinical Classifiers. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 3808–3823.
- [4] Solon Barocas and Andrew D Selbst. 2016. Big data's disparate impact. *Calif. L. Rev.* 104 (2016), 671.
- [5] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilovic, et al. 2018. AI Fairness 360: An Extensible Toolkit for Detecting, Understanding, and Mitigating Unwanted Algorithmic Bias. *arXiv preprint arXiv:1810.01943* (2018).
- [6] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. 610–623.
- [7] Steven Bethard. 2022. We need to talk about random seeds. *arXiv preprint arXiv:2210.13393* (2022).
- [8] Sarah Bird, Miro Dudík, Richard Edgar, Brandon Horn, Roman Lutz, Vanessa Milan, Mehrnoosh Sameki, Hanna Wallach, and Kathleen Walker. 2020. Fairlearn: A toolkit for assessing and improving fairness in AI. *Microsoft, Tech. Rep. MSR-TR-2020-32* (2020).
- [9] T Birsan and Dan Tiba. 2005. One hundred years since the introduction of the set distance by Dimitrie Pompeiu. In *IFIP Conference on System Modeling and Optimization*. Springer, 35–39.
- [10] Emily Black, Manish Raghavan, and Solon Barocas. 2022. Model multiplicity: Opportunities, concerns, and solutions. In *2022 ACM Conference on Fairness, Accountability, and Transparency*. 850–863.
- [11] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. 2015. Weight uncertainty in neural network. In *International conference on machine learning*. PMLR, 1613–1622.
- [12] Léon Bottou. 2012. Stochastic gradient descent tricks. In *Neural networks: Tricks of the trade*. Springer, 421–436.
- [13] Simon Caton and Christian Haas. 2020. Fairness in machine learning: A survey. *arXiv preprint arXiv:2010.04053* (2020).
- [14] Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. 2011. Differentially private empirical risk minimization. *Journal of Machine Learning Research* 12, 3 (2011).
- [15] Alexandra Chouldechova and Aaron Roth. 2018. The frontiers of fairness in machine learning. *arXiv preprint arXiv:1810.08810* (2018).
- [16] Kate Crawford. 2013. The hidden biases in big data. *Harvard business review* 1, 4 (2013).
- [17] Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. 2021. Retiring adult: New datasets for fair machine learning. *Advances in Neural Information Processing Systems* 34 (2021).
- [18] Felix Draxler, Kambis Veschgini, Manfred Salmhofer, and Fred Hamprecht. 2018. Essentially no barriers in neural network energy landscape. In *International conference on machine learning*. PMLR, 1309–1318.
- [19] Dheeru Dua and Casey Graff. 2017. UCI Machine Learning Repository. <http://archive.ics.uci.edu/ml>
- [20] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*. 214–226.
- [21] Sorelle A Friedler, Carlos Scheidegger, Suresh Venkatasubramanian, Sonam Choudhary, Evan P Hamilton, and Derek Roth. 2019. A comparative study of fairness-enhancing interventions in machine learning. In *Proceedings of the conference on fairness, accountability, and transparency*. 329–338.
- [22] Akihiko Fukuchi, Yoko Yabe, and Masashi Sode. 2020. FairTorch. <https://github.com/wbawakate/fairtorch>.
- [23] Yarín Gal and Zoubin Ghahramani. 2016. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*. PMLR, 1050–1059.
- [24] Timur Gariipov, Pavel Izmailov, Dmitrii Podoprikin, Dmitry P Vetrov, and Andrew G Wilson. 2018. Loss surfaces, model connectivity, and fast ensembling of dnns. *Advances in neural information processing systems* 31 (2018).
- [25] Abhijit Ghatak and Abhijit Ghatak. 2019. Initialization of Network Parameters. *Deep Learning with R* (2019), 87–102.
- [26] Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of opportunity in supervised learning. *Advances in neural information processing systems* 29 (2016).
- [27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [28] Sara Hooker, Nyalleng Moorosi, Gregory Clark, Samy Bengio, and Emily Denton. 2020. Characterising bias in compressed models. *arXiv preprint arXiv:2010.03058* (2020).
- [29] Matthew Jagielski, Om Thakkar, Florian Tramèr, Daphne Ippolito, Katherine Lee, Nicholas Carlini, Eric Wallace, Shuang Song, Abhradeep Thakurta, Nicolas Papernot, et al. 2022. Measuring Forgetting of Memorized Training Examples. *arXiv preprint arXiv:2207.00099* (2022).
- [30] Faisal Kamiran and Toon Calders. 2012. Data preprocessing techniques for classification without discrimination. *Knowledge and information systems* 33, 1 (2012), 1–33.
- [31] Cherry Khosla and Baljit Singh Saini. 2020. Enhancing performance of deep learning models with different data augmentation techniques: A survey. In *2020 International Conference on Intelligent Engineering and Management (ICIEM)*. IEEE, 79–85.
- [32] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences* 114, 13 (2017), 3521–3526.
- [33] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep Learning Face Attributes in the Wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.
- [34] Yucheng Lu, Wentao Guo, and Christopher M De Sa. 2022. GraB: Finding Provably Better Data Permutations than Random Reshuffling. *Advances in Neural Information Processing Systems* 35 (2022), 8969–8981.
- [35] Frank J Massey Jr. 1951. The Kolmogorov-Smirnov test for goodness of fit. *Journal of the American statistical Association* 46, 253 (1951), 68–78.
- [36] R Thomas McCoy, Junghyun Min, and Tal Linzen. 2020. BERTs of a feather do not generalize together: Large variability in generalization across models with similar test set performance. In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*. 217–227.
- [37] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)* 54, 6 (2021), 1–35.
- [38] Konstantin Mishchenko, Ahmed Khaled, and Peter Richtárik. 2020. Random reshuffling: Simple analysis with vast improvements. *Advances in Neural Information Processing Systems* 33 (2020), 17309–17320.
- [39] Amirkeivan Mohtashami, Sebastian Stich, and Martin Jaggi. 2022. Characterizing & Finding Good Data Orderings for Fast Convergence of Sequential Gradient Methods. *arXiv preprint arXiv:2202.01838* (2022).
- [40] Arvind Narayanan. 2018. Translation tutorial: 21 fairness definitions and their politics. In *Proc. Conf. Fairness Accountability Transp., New York, USA*, Vol. 1170. 3.
- [41] Lam M Nguyen, Quoc Tran-Dinh, Dzong T Phan, Phuong Ha Nguyen, and Marten Van Dijk. 2021. A unified convergence analysis for shuffling-type gradient methods. *The Journal of Machine Learning Research* 22, 1 (2021), 9397–9440.
- [42] Hyeonwoo Noh, Tackgeun You, Jonghwan Mun, and Bohyung Han. 2017. Regularizing deep neural networks by noise: Its interpretation and optimization. *Advances in Neural Information Processing Systems* 30 (2017).
- [43] Shangshu Qian, Viet Hung Pham, Thibaud Lutellier, Zeou Hu, Jungwon Kim, Lin Tan, Yaoliang Yu, Jiahao Chen, and Sameena Shah. 2021. Are my deep learning systems fair? An empirical study of fixed-seed training. *Advances in Neural Information Processing Systems* 34 (2021).
- [44] Shashank Rajput, Kangwook Lee, and Dimitris Papailiopoulos. 2022. Permutation-Based SGD: Is Random Optimal?. In *International Conference on Learning Representations*.
- [45] Charan Reddy, Deepak Sharma, Soroush Mehri, Adriana Romero-Soriano, Samira Shabanian, and Sina Honari. 2021. Benchmarking Bias Mitigation Algorithms in Representation Learning through Fairness Metrics. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*.
- [46] Yuji Roh, Kangwook Lee, Steven Euijong Whang, and Changho Suh. 2020. FairBatch: Batch Selection for Model Fairness. In *International Conference on Learning Representations*.
- [47] Nripsuta Ani Saxena. 2019. Perceptions of fairness. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. 537–538.
- [48] Thibault Sellam, Steve Yadlowsky, Ian Tenney, Jason Wei, Naomi Saphra, Alexander D'Amour, Tal Linzen, Jasmijn Bastings, Iulia Raluca Turc, Jacob Eisenstein, et al. 2021. The MultiBERTs: BERT Reproductions for Robustness Analysis. In *International Conference on Learning Representations*.
- [49] Vatsal Shah, Xiaoxia Wu, and Sujay Sanghavi. 2020. Choosing the sample with lowest loss makes sgd robust. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 2120–2130.
- [50] Ilya Shumailov, Zakhar Shumaylov, Dmitry Kazhdan, Yiren Zhao, Nicolas Papernot, Murat A Erdogdu, and Ross J Anderson. 2021. Manipulating sgd with data ordering attacks. *Advances in Neural Information Processing Systems* 34 (2021), 18021–18032.
- [51] Ioana Baldini Soares, Dennis Wei, Karthikeyan Natesan Ramamurthy, Moninder Singh, and Mikhail Yurochkin. 2022. Your Fairness May Vary: Pretrained Language Model Fairness in Toxic Text Classification. In *Annual Meeting of the Association for Computational Linguistics*.
- [52] Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe. 2022. Curriculum learning: A survey. *International Journal of Computer Vision* (2022), 1–40.
- [53] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research* 15, 1 (2014), 1929–1958.

- [54] Kushal Tirumala, Aram H Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. 2022. Memorization Without Overfitting: Analyzing the Training Dynamics of Large Language Models. *arXiv preprint arXiv:2205.10770* (2022).
- [55] Mariya Toneva, Alessandro Sordani, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. 2018. An Empirical Study of Example Forgetting during Deep Neural Network Learning. In *International Conference on Learning Representations*.
- [56] Vladimir Vapnik. 1991. Principles of risk minimization for learning theory. *Advances in neural information processing systems* 4 (1991).
- [57] Abhinav Venigalla and Linden Li. 2022. Mosaic LLMs (Part 2): GPT-3 quality for < \$500k. <https://www.mosaicml.com/blog/gpt-3-quality-for-500k>
- [58] Stefan Wager, Sida Wang, and Percy S Liang. 2013. Dropout training as adaptive regularization. *Advances in neural information processing systems* 26 (2013).
- [59] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. 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*. 2979–2989.

IMAGE CREDITS

- Blueprint by Berkah Icon from Noun Project
- Neural Network by Ian Rahmadi Kurniawan from Noun Project
- Data by shashank singh from Noun Project
- Data processing by Eko Purnomo from Noun Project
- Evaluate by Justin Blake from Noun Project

A NON-DETERMINISM IN MODEL TRAINING

In our paper, we focus on fairness variance due to randomness in the training algorithm (i.e., weight initialization and random reshuffling). To control the randomness, we set manual seeds at various intermediate locations in our code. We refer to the seed set right before building the neural network as the weight initialization seed, which influences the randomness in sampling the weight values. Similarly, we refer to the seed set right before the first training data shuffling as random reshuffling seed, which influences the data order that will be used as reference for the rest of the training. During training, we simply set the epoch number as seed right before reshuffling the reference data order at every epoch. As we can change the reference data order by changing the random reshuffling seed, this setup allows us to control non-determinism in data order throughout training with a single random seed.

We also perform experiments with dropout regularization during training. The epoch number seed set at the start of every epoch serves dual functionality, and also controls the randomness from dropout regularization. That is, fixing both weight initialization and random reshuffling seeds allows us to deterministically replicate model training, while changing both seeds simultaneously is similar to the discussion of non-determinism currently present in fair deep learning literature. We provide a minimal pseudo-code to explain our control over training non-determinism and define weight initialization and random reshuffling seeds.

```
....
torch.manual_seed(weight_initialization_seed)
model = MLPModel() # Initialize Neural Model
...
train_data, .. = load_dataset() # load data with
↳ deterministic order
train_data = shuffle(train_data,
↳ seed=random_resuffling_seed)
...
trainloader = torch.utils.data.DataLoader(train_data,
↳ ..., shuffle=True)
...
for epoch in range(total_epochs):
    ...
    torch.manual_seed(epoch)
    for batch in trainloader:
        ....
```

B DATASETS

ACSIIncome. ACSIIncome is one of the five pre-defined tasks in the Folktables dataset [17], which was recently collected to improve the older and commonly used UCI Adult Income dataset [19]. More specifically, we use the subset of Folktables dataset from the state of California, USA for 2018. The dataset contains a total of 195,665 data points with 10 features each, where each data points represents an individual. The task is a binary classification to predict whether the individual's income is above \$50,000. For fairness measures, we use perceived gender (*Sex*) as the sensitive attribute, which is also one of the 10 input features.

List of features : *Age, Class of worker, Educational attainment, Marital status, Occupation, Place of birth, Relationship, Usual hours worked per week in past 12 months, Sex, Recoded detailed race code*

ACSEmployment. ACSEmployment is another one of the five pre-defined tasks in the Folktables dataset [17]. We use the same subset from the state of California, USA for 2018 as above. The dataset contains a total of 378,817 data points with 16 features each, and the task is a binary classification to predict whether the individual is employed or not. For fairness measures, same as above, we use perceived gender (*Sex*) as the sensitive attribute.

List of features : *Age, Educational attainment, Marital status, Sex, Disability recode, Employment status of parents, Mobility status, Citizenship status, Military service, Ancestry recode, Nativity, Relationship, Hearing difficulty, Vision difficulty, Cognitive difficulty, Recoded detailed race code, Grandparents living with grandchildren*

CelebA. CelebA dataset is a large scale celebrity face attributes dataset [33], which contains a total of 202,599 images of celebrities with 40 different binary labels each. We focus on the binary classification task of the 'smiling' label in the dataset, while we use the 'gender' label as the sensitive attribute for fairness evaluation. Moreover, we do not directly use the images of CelebA dataset, but instead pass them through a pre-trained and frozen ResNet-50 backbone [27] to extract image representations, which are treated as inputs to our model.

C TRAINING CURVE AND CONVERGENCE

To understand why we choose to study the model behavior between epochs 100 and 300, we point the reader towards the overall training curve of the model plotted in Fig. 11. It is clear that the model has converged by epoch 100, and maintains stable accuracy scores for the last 200 epochs.

Even though the accuracy scores have converged, we still find high variance in fairness scores as discussed in the main text of the paper. One might suspect this implies that fairness scores could take longer to converge. To check this, we allow a single training to run for a total of 3000 epochs (as opposed to the standard 300 epochs used in all other experiments in our paper) and collect the fairness scores in Fig. 12. It is clear that the model fairness does not stabilize by simply increasing the number of epochs, which further supports our hypothesis of never-ending local oscillations due to SGD noise that cause high fairness variance.

D RANDOM ORDER FOR IMMEDIATE IMPACT

We used carefully chosen data order to show the immediate impact of data order in Section 6. Here, we provide additional results on randomly chosen data order to show that the property does hold even for a random data order. Sae=me as in the original experiment, we sample 1000 unique checkpoints randomly from last 200 epochs of 50 different training runs while allowing both forms of training non-determinism simultaneously. We then fine-tune each of these checkpoints for exactly one epoch on a common, ad this time randomly chosen, data order. We collect fairness variance across checkpoints before and after this single epoch of training in Fig. 13.

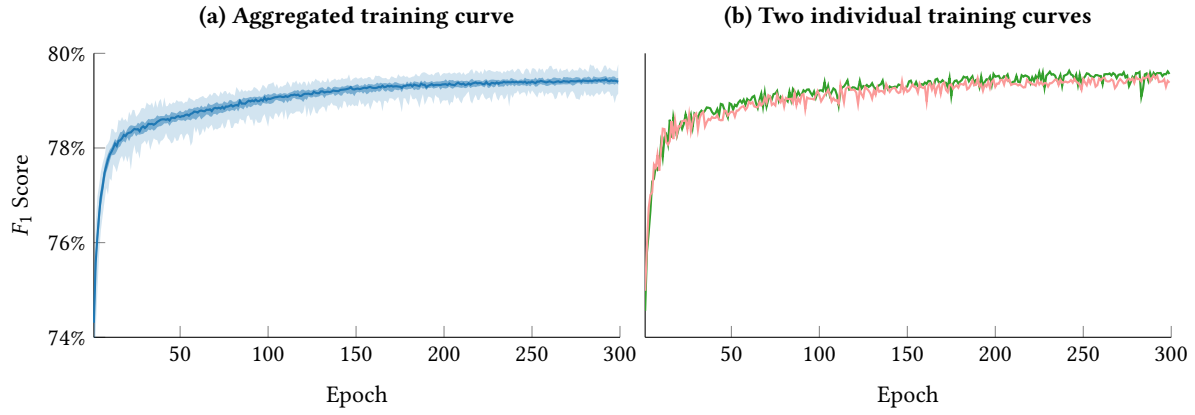


Figure 11: (Left) Aggregated training curve for all 50 training runs with both changing weight initialization and random reshuffling. (Right) Example of two individual training runs.

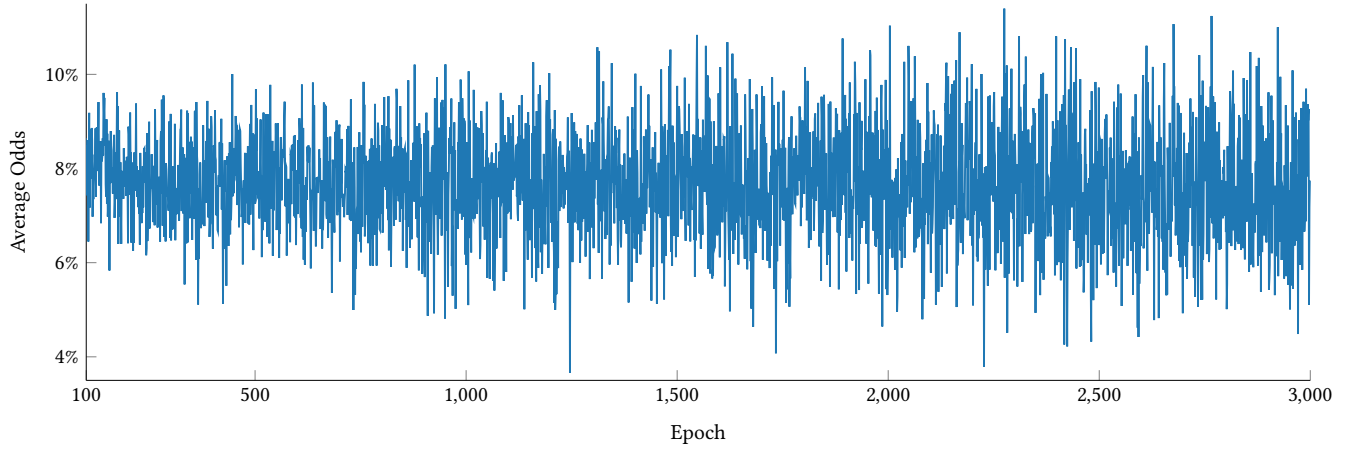


Figure 12: A single training run extended to 3000 epochs.

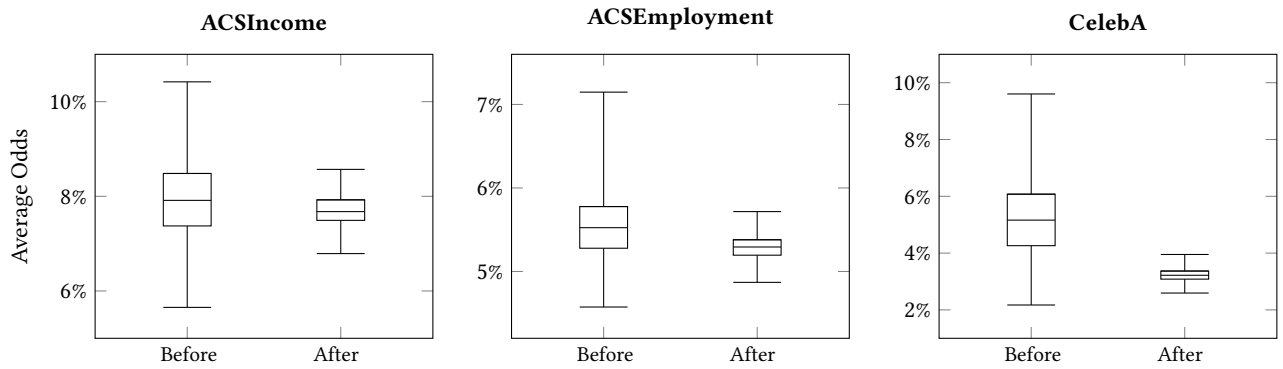


Figure 13: (Before) Multiple checkpoints taken from training runs with changing weight initialization, training data order, and even number of training epochs, show a high range of fairness variance (i.e., average odds variance), as expected. (After) However, these fairness scores are stabilized significantly by training these checkpoints for only a single epoch on a common randomly chosen data order. This shows the immediate impact of data order on model fairness.

Before training on the same data order for a single epoch, these checkpoints represent the complete range of fairness variance previously noted. However, after only a single epoch of training on a common data order, these models have all moved towards the same fairness score with significantly lower variance. This highlights the immediate impact of data order on model fairness, which is stable based on only the data order of the most recent epoch.

E SANITY CHECK FOR DATA ORDER SUFFIX

We showed that the fairness scores are dominated by the most recent gradient updates as seen by the model. As a sanity check, we also provide ablation study for experiments in Fig. 7, but choose the b batches randomly instead of from the suffix. The results are collected in Fig. 14. The results show while the fairness variance does get smaller with more common batches (as expected, see also Fig. 13), its predictability is lost when choosing the b batches randomly. Thus, it is indeed the batches from data order suffix that govern the predictability of model fairness.

F ADDITIONAL EXPERIMENTS FOR ALL DATASETS

F.1 Weight Initialization and Random Reshuffling

We provide additional results in Fig. 15, 16 on CelebA and ACSEmployment to show the dominance of reshuffling.

F.2 Changing Predictions and Data Distribution

We show the relationship between data distribution and changing predictions for CelebA and ACSEmployment in Fig. 17. The group with least representation maintains being the most vulnerable to changing predictions.

F.3 Capturing Variance in a Single Training Run

We show the empirical similarity of distribution across multiple training runs and multiple epochs in a single training run for additional datasets CelebA and ACSEmployment in Fig. 18.

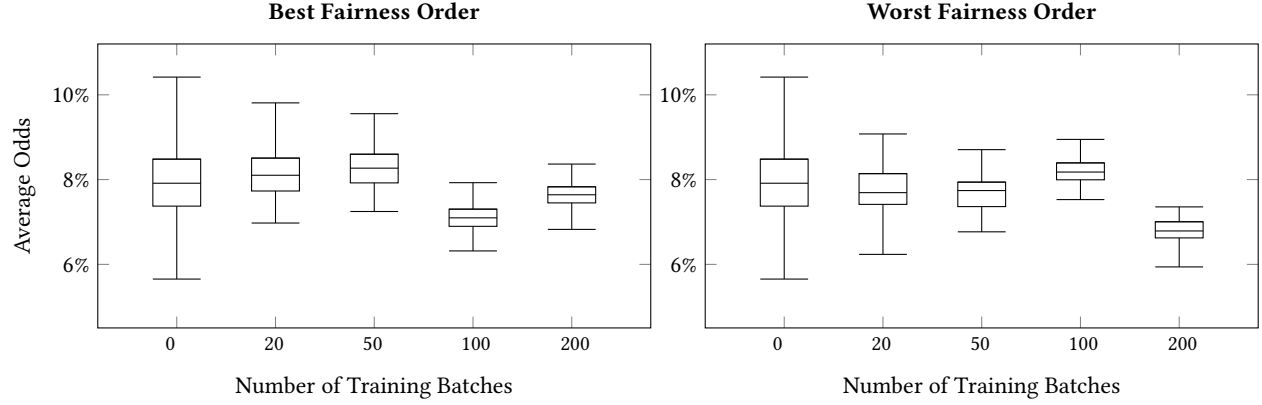


Figure 14: Additional experiments by choosing b batches, similar to the experiments in Figure 7, but choosing batches randomly instead of the suffix. The results still maintain stability as b increases but now they stabilize to random fairness scores instead of best and worst fairness scores as seen in Figure 7. Clearly, it's the most recent batches in that order which truly governs the model fairness.

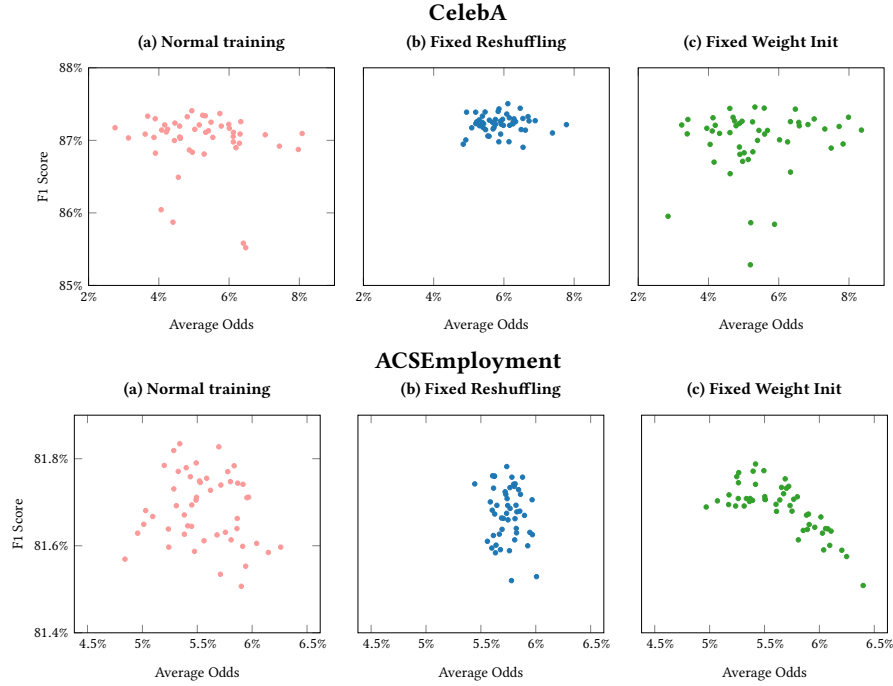


Figure 15: Additional experiments on CelebA and ACSEmployment datasets reveal similar trends as seen in Figure 3. Random reshuffling of data order is the dominant source of variance in fairness scores, with very little influence from the weight initialization.

F.4 Manipulating Group Level Accuracy with Data Order

We provide additional results in Figure 19, 20, highlighting the predictability of model fairness based on the data order.

G ADDITIONAL EXPERIMENTS FOR CHANGING HYPERPARAMETERS

We provide additional experiments for (i) batch size 16 and 1024, deviating from the default batch size of 128 in the main text, (ii) learning rate 0.01 and 0.0001, deviating from the default learning rate of 0.001 in the main text, and (iii) model architecture with two hidden layers containing 2048 and 64 neurons respectively,

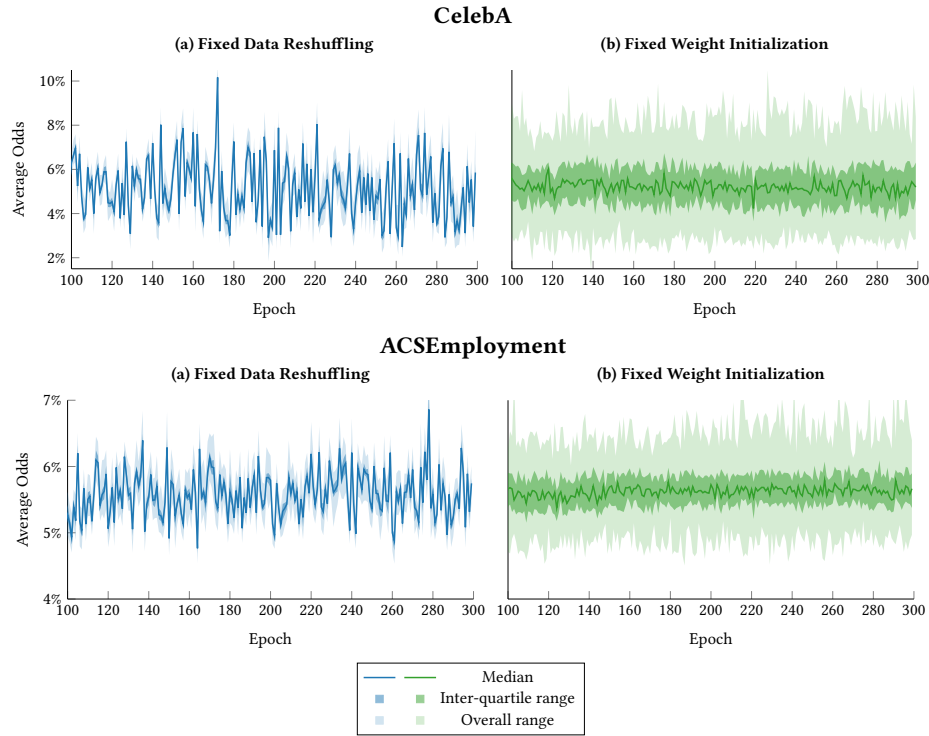


Figure 16: Additional experiments on CelebA and ACSEmployment datasets reveal similar trends as in Figure 4. These results further highlight the dominant impact of random reshuffling on fairness.

deviating from the single hidden layer architecture used in the main text.

stability (i.e. large batch size or lower learning rate) is an inefficient solution to solving the problem of fairness variance.

G.1 Training Curves and Convergence

We first observe the training curve in all the settings described above, along with the original setting, in Fig. 21. It is clear that models with a bigger batch size (or a smaller learning rate) do not converge to the same F1 scores as other models. This further indicates the need of randomness and noise in the learning algorithm to give neural networks more 'exploration energy' and facilitate better and faster convergence. On the other hand, a smaller batch size (or a larger learning rate) comes with excessive variance in the model's F1 score, which highlights the importance of achieving an appropriate balance between several hyperparameter choices.

To better understand the difficulty of reaching convergence, we continue training a single instance with batch size 1024 (and separately with learning rate 0.0001) for a total of 1000 epochs and compare it against the standard training setup (i.e., with batch size 128 and learning rate 0.001). The results are collected in Fig. 22, 23. It is clear that a larger batch size (or lower learning rate) slows down the convergence speed. Moreover, it also does not achieve the same F1 scores previously seen, possibly due to not being able to take complete advantage of the noisy nature of mini-batch gradient descent. Thus, training a model with hyperparameters that enforces

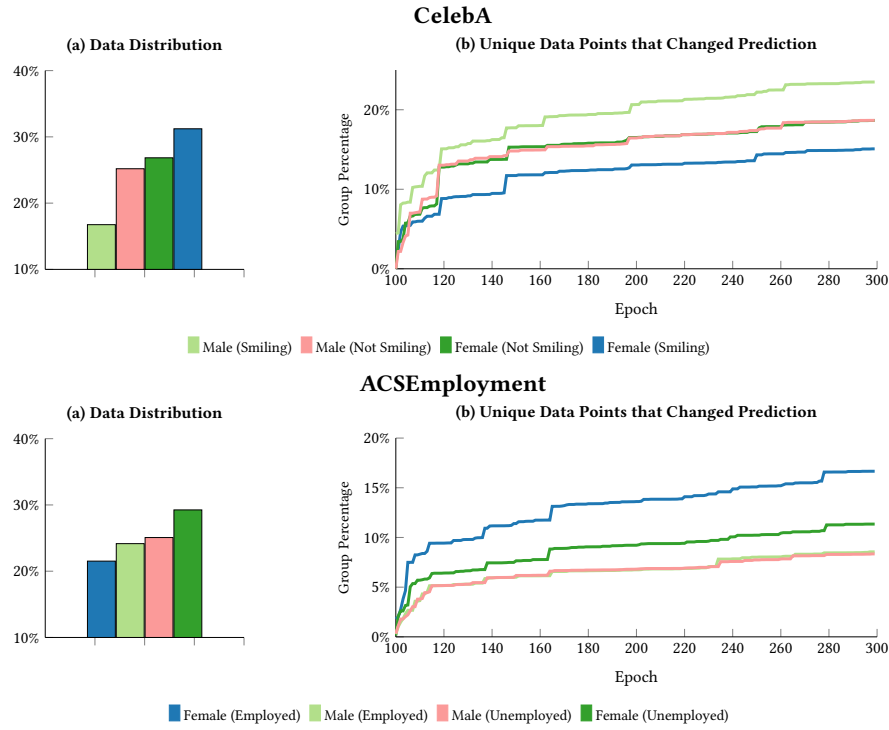


Figure 17: Additional experiments on CelebA and ACSEmployment datasets reveal similar trends as seen in Figure 5. These results highlight subgroups with least representation being the most vulnerable to changing predictions.

G.2 Weight Initialization and Random Reshuffling

We provide results in Fig. 24, 25 for changing batch size, learning rate, and architecture on ACSIncome dataset to show the dominance of random reshuffling on model fairness. With a decrease in batch size, the range of fairness variance increases significantly, but the overall expected trends follow the same behavior as noted in the main text, i.e. a sharp change in fairness scores across epochs under fixed data reshuffling, and high variance even for a single epoch across multiple runs under changing data reshuffling (fixed weight initialization). Similar trends can be observed when we increase the learning rate, or provide the training algorithm with a bigger neural model. The increase in variance suggests high instability with smaller batch size, higher learning rate, and bigger neural models, all of which is expected.

On the other hand, one would expect more stable model behavior with a bigger batch size or a smaller learning rate. While this is indeed the case, the comparison here is not fair because as noted earlier (Fig. 21, these models have not yet converged, also evident by the clear downward trend of fairness scores. It is clear that certain hyperparameter settings are not conducive to efficient convergence, even though they might eventually provide more stable fairness scores.

G.3 Changing Predictions and Data Distribution

We note changing predictions for all 5 settings described above in Fig. 26. Note that since we are focusing on changing training hyperparameters for ACSIncome, the data distribution remains the same as in Fig. 5(a) (also copied here in Fig. 26 for reference).

Same as before, model instability has increased with smaller batch size, higher learning rate, or bigger model architecture, but the trends of vulnerability for various subgroups remains the same. Even though we see same trends of higher vulnerability for a higher batch size or smaller learning rate, we know that these models have not converged and thus would recommend not making any inference from these two set of results.

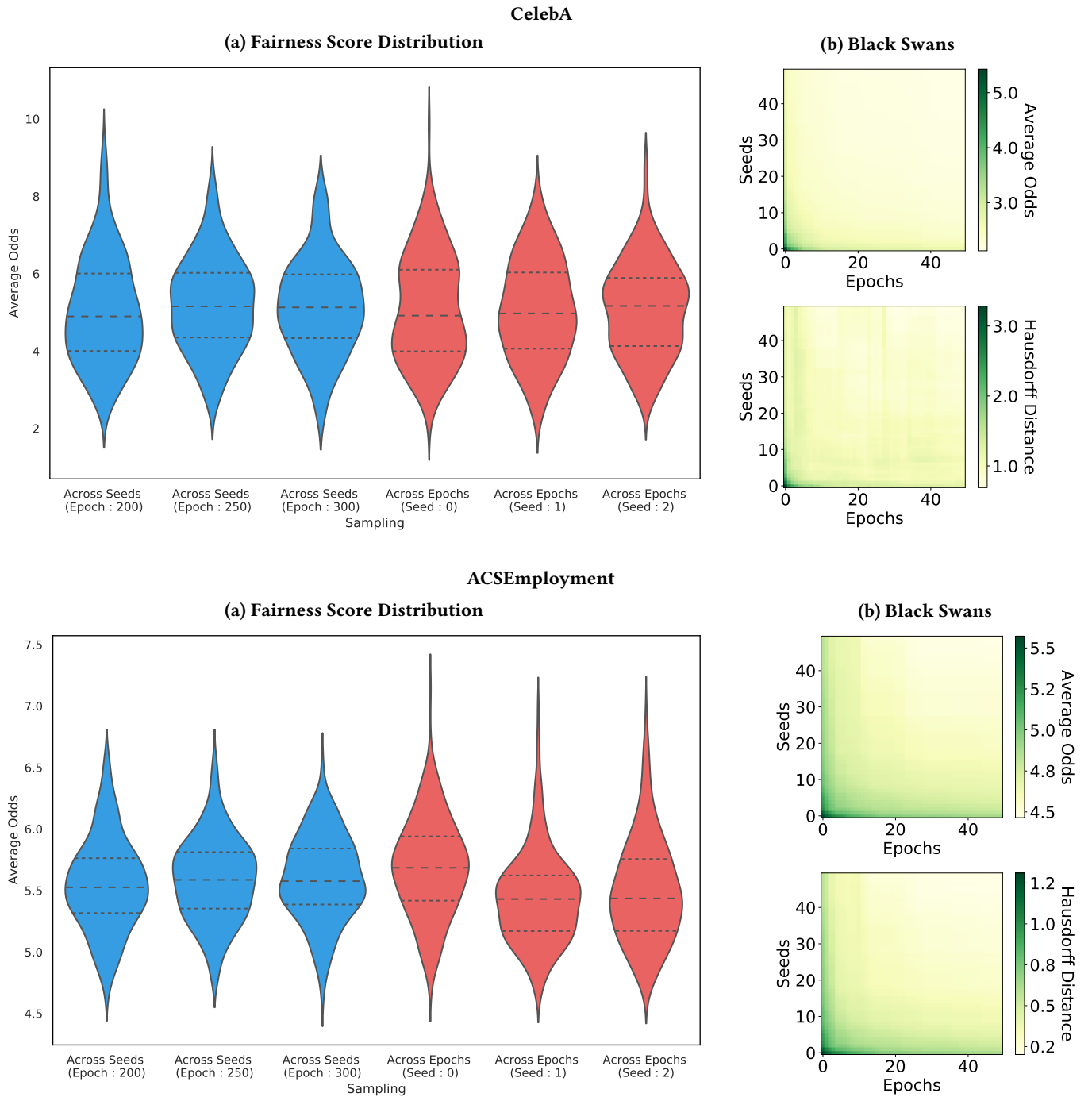


Figure 18: Additional experiments on CelebA and ACSEmployment datasets reveal similar trends as seen in Figure 9. Fairness scores (average odds) across multiple training runs and across epochs in a single training run have similar empirical distributions. Thus, studying this distribution across epochs provides a highly efficient alternative to executing multiple training runs.

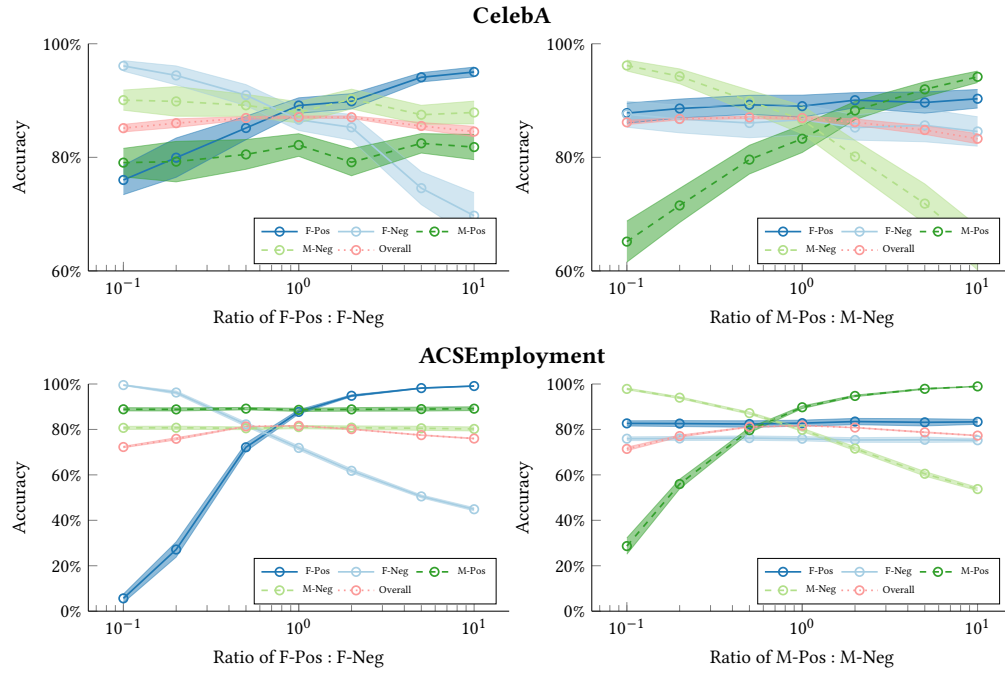


Figure 19: Additional experiments on CelebA and ACSEmployment datasets reveal similar trends as seen in Figure 8. In only a single epoch of training, we are able to manipulate the group level accuracy trade-off, with relatively small impact on overall accuracy.

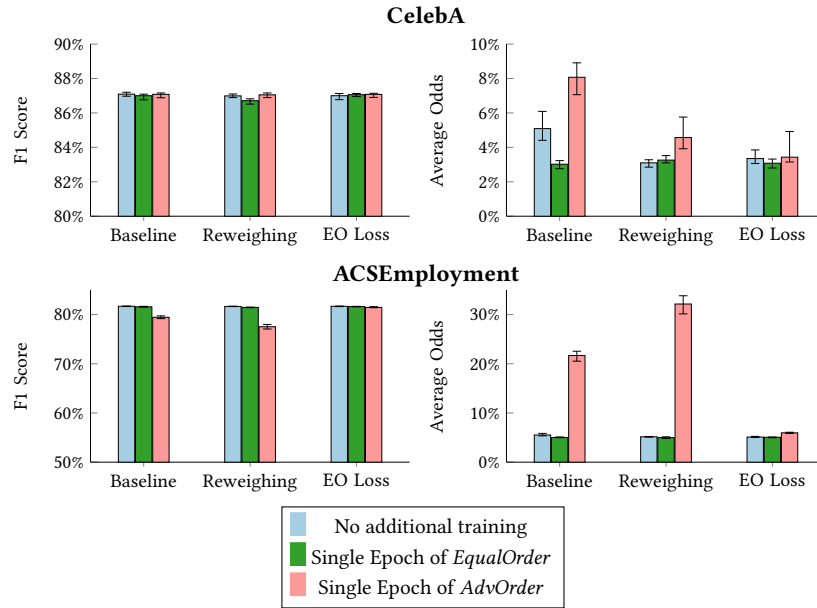


Figure 20: Additional experiments on CelebA and ACSEmployment reveal similar trends as seen in Figure 10. *EqualOrder* gets competitive performance to commonly used bias mitigation methods. Similarly, *AdvOrder* gets significantly worse fairness than even the baseline.

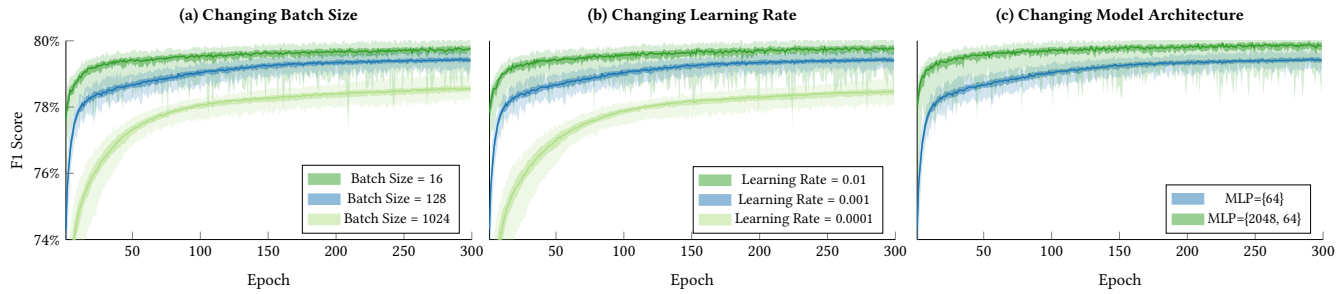


Figure 21: Training curve under changing batch size, learning rate, and model architecture. It is clear that decreasing the batch size, increasing the learning rate, or using a bigger model architecture results in a faster model convergence but with higher variance even in accuracy scores. On the other hand, models with high batch size or low learning rates tend to not achieve the same accuracy scores and have not reached convergence even at 300 epochs.

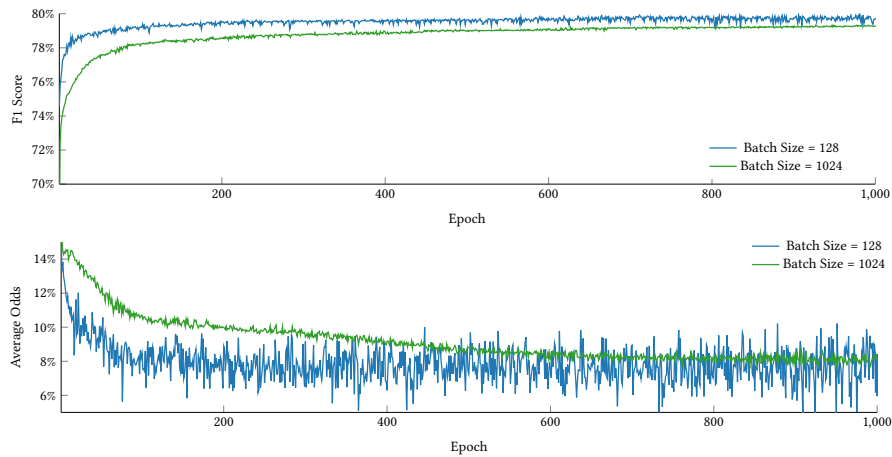


Figure 22: Using a higher batch size can over time achieve stable fairness scores, however the convergence speed is significantly slower. Moreover, it loses a noticeable margin of F1 score.

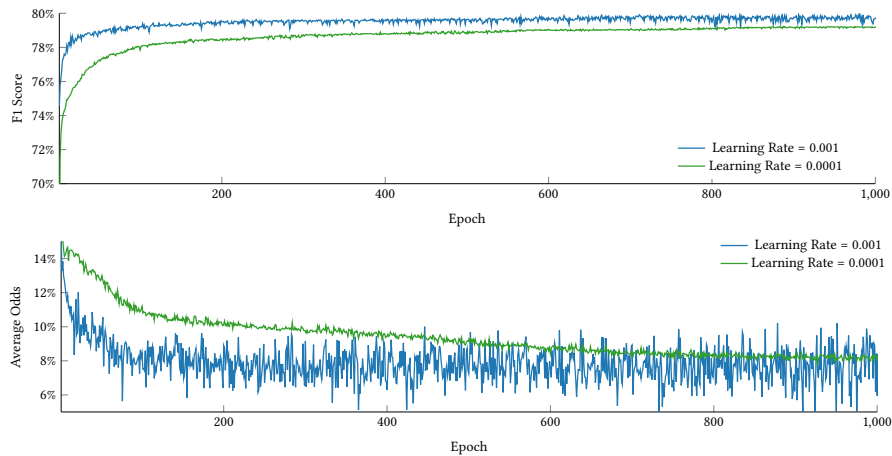


Figure 23: Using a lower learning rate can over time achieve stable fairness scores, however the convergence speed is significantly slower. Moreover, it loses a noticeable margin of F1 score.

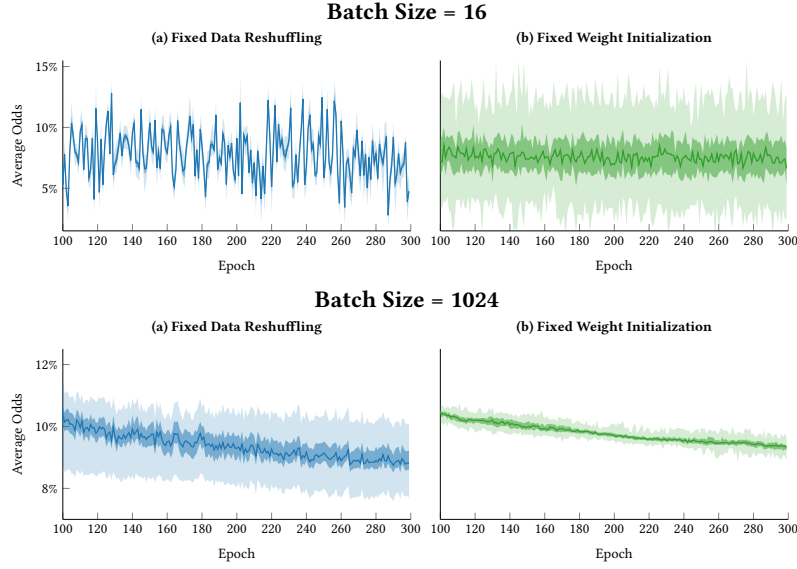


Figure 24: Additional experiments for changing batch size with experiment setting as in Figure 4. These results further highlight the dominant impact of random reshuffling on fairness.

G.4 Manipulating Group Level Accuracy with Data Order

We provide additional results in Fig. 27, highlighting the predictability of model fairness based on the data order. Note that the hyperparameter setting for an additional epoch of fine-tuning during group accuracy manipulation is the same as the setup used for training that particular model. For example, when manipulating models which were trained with batch size 16, the single epoch of fine-tuning is also done with batch size 16.

We also provide results for models with high batch size and low learning rate separately in Fig. 28. Again, note that not only are these models not converged, but they are also fine-tuned on the same inefficient hyperparameter settings. Thus, the trends here are not comparable, but added for completeness.

H ADDITIONAL EXPERIMENTS FOR DROPOUT REGULARIZATION

Another notable hyperparameter in neural model training is dropout regularization [53]. Dropout regularization randomly drops a certain percentage (known as dropout rate) of connections between consecutive layers in the model at every forward pass during training. We extend our discussion from Section 4 to study the impact of dropout on the trends of random reshuffling seen in Figure 4. More specifically, we repeat the experiments while introducing various rates of dropout in the training setup. The results are collected in Fig. 29.

Even with dropout regularization, the impact of data reshuffling on model fairness clearly dominates weight initialization. That is, the trends of high correlation between multiple runs with the same data order, despite starting from different initializations, and lack of any correlation between multiple runs with different data order, even after starting from the same initialization, are still evident even in presence of other sources of randomness.

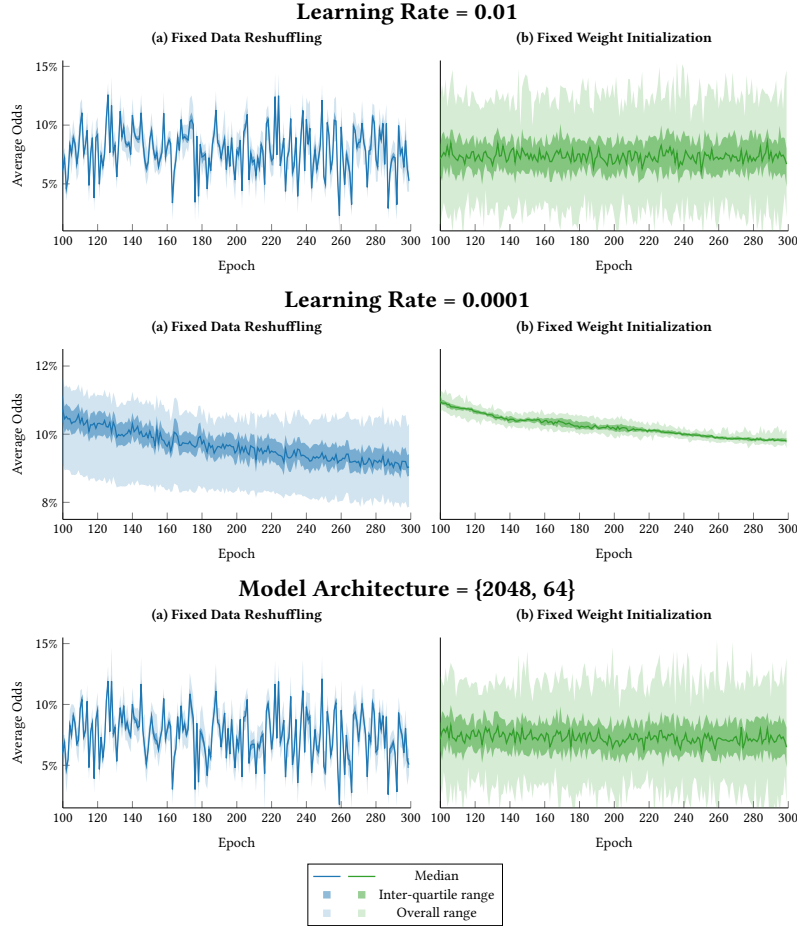


Figure 25: Additional experiments for changing learning rate and model architecture with experiment setting as in Figure 4. These results further highlight the dominant impact of random reshuffling on fairness.

I ADDITIONAL EXPERIMENTS FOR OTHER FAIRNESS METRICS

We repeat the experiment in the main text for two more fairness metrics, Equal Opportunity (EOpp) and Demographic Parity (DP) [26]

Using the same notations as in Section 3, EOpp can be defined as,

$$EOpp(f, \mathcal{D}) := \left| \frac{\sum_{\mathcal{D}^e} \mathbb{1}[f(x)=1 \wedge y=1 \wedge a=0]}{\sum_{\mathcal{D}^e} \mathbb{1}[y=1 \wedge a=0]} - \frac{\sum_{\mathcal{D}^e} \mathbb{1}[f(x)=1 \wedge y=1 \wedge a=1]}{\sum_{\mathcal{D}^e} \mathbb{1}[y=1 \wedge a=1]} \right|. \quad (9)$$

Similarly, DP can be defined as,

$$DP(f, \mathcal{D}) := 1 - \min \left(\frac{\sum_{\mathcal{D}^e} \mathbb{1}[f(x)=1 \wedge a=0]}{\sum_{\mathcal{D}^e} \mathbb{1}[f(x)=1 \wedge a=1]}, \frac{\sum_{\mathcal{D}^e} \mathbb{1}[f(x)=1 \wedge a=1]}{\sum_{\mathcal{D}^e} \mathbb{1}[f(x)=1 \wedge a=0]} \right). \quad (10)$$

I.1 High Variance in Fairness Scores

We start by repeating the experiment comparing various state-of-the-art bias mitigation techniques and intersecting range of fairness scores in Fig. 30.

I.2 Weight Initialization and Random Reshuffling

We provide additional results for EOpp and DP for the experiment conducted in Fig. 4 to show correlation between multiple runs on ACSIncome dataset. The results are collected in Fig. 31, and show similar trends as seen in the main text.

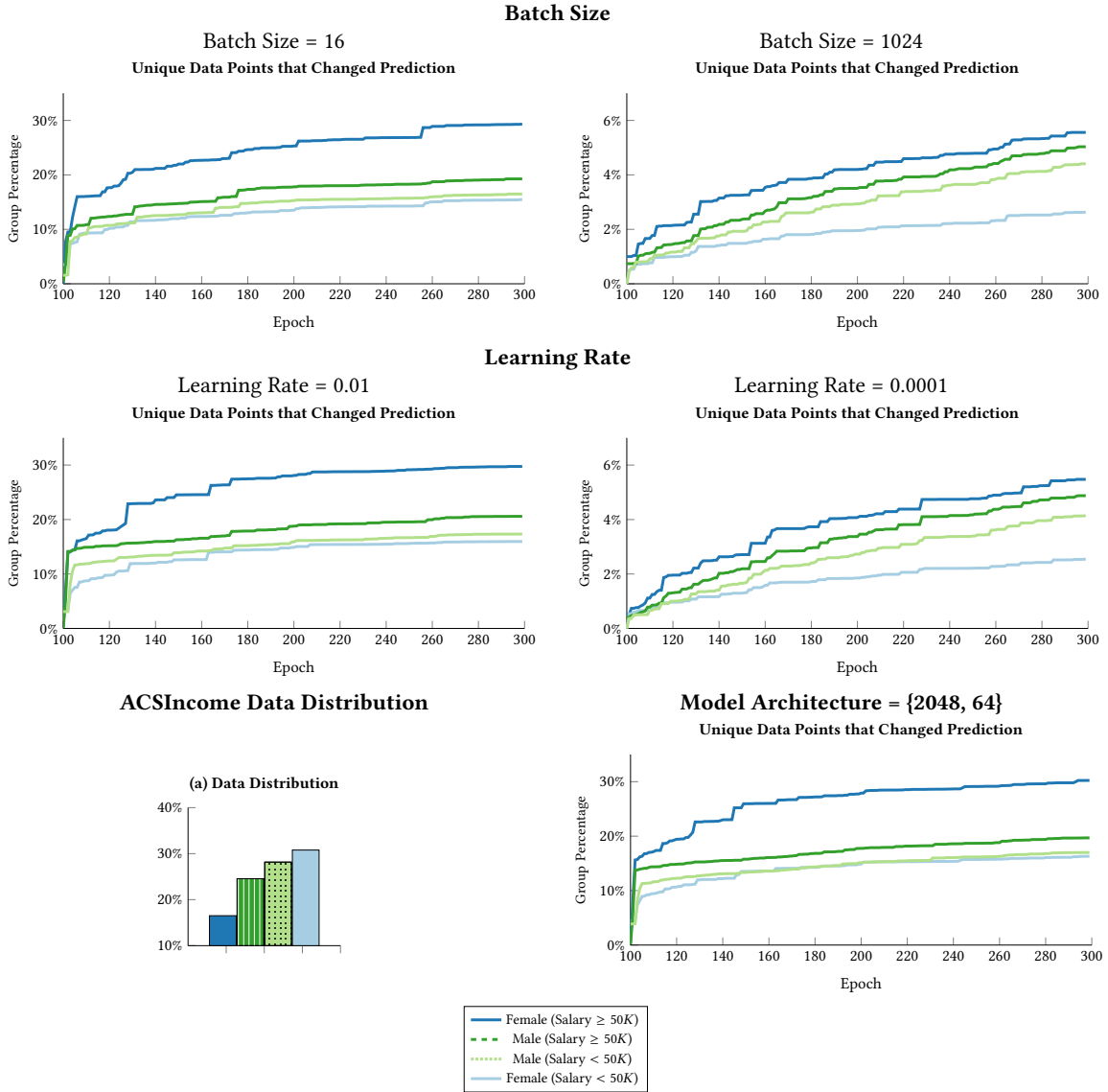


Figure 26: Additional experiments reveal similar trends as seen in Figure 5. These results highlight subgroups with least representation being the most vulnerable to changing predictions.

I.3 Capturing Variance in a Single Training Run

We show the empirical similarity of distribution across multiple training runs and multiple epochs in a single training run for fairness measures EO_{pp} and DP in Figure 32.

J 10 RAW TRAINING RUNS FROM FIG. 4

We plot 10 randomly chosen training runs each for fixed weight initialization and fixed random reshuffling (see Fig. 4) in Fig. 33 and Fig. 34, respectively. As expected, each individual training run in both settings has high variance across epochs even after convergence. More importantly, the trends of fairness matches closely across multiple runs for fixed random reshuffling, even though they started from different weight initialization.

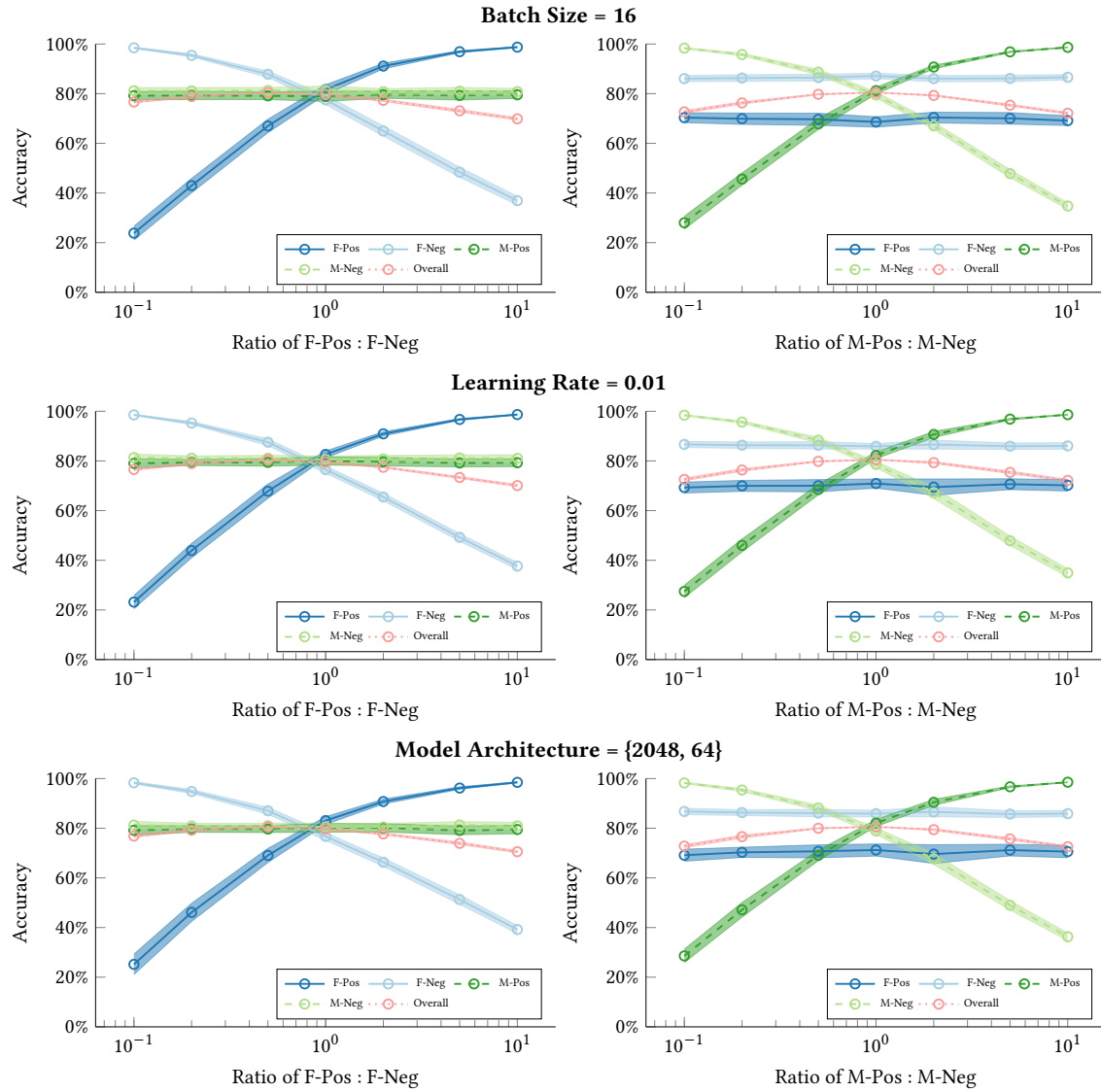


Figure 27: Additional experiments reveal similar trends as seen in Figure 8. In only a single epoch of training, we are able to manipulate the group level accuracy trade-off, with relatively small impact on overall accuracy.

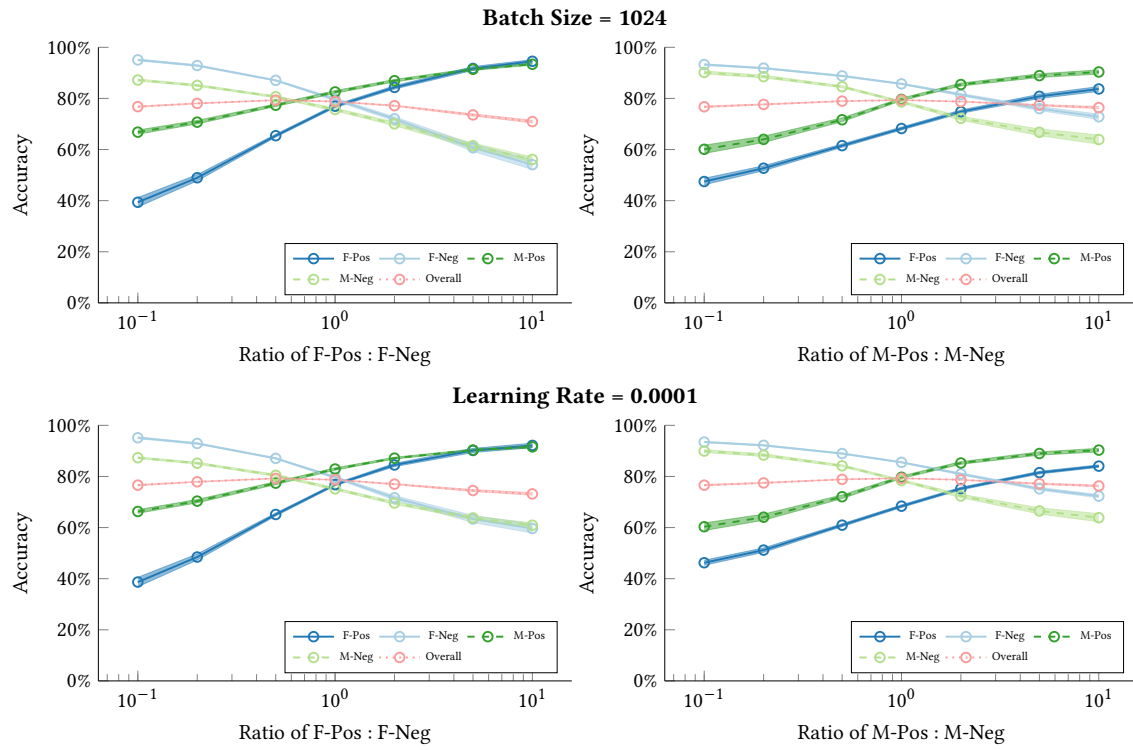


Figure 28: Additional experiments as seen in Figure 27, but for models that did not converge as noted earlier. The results are added for completeness.

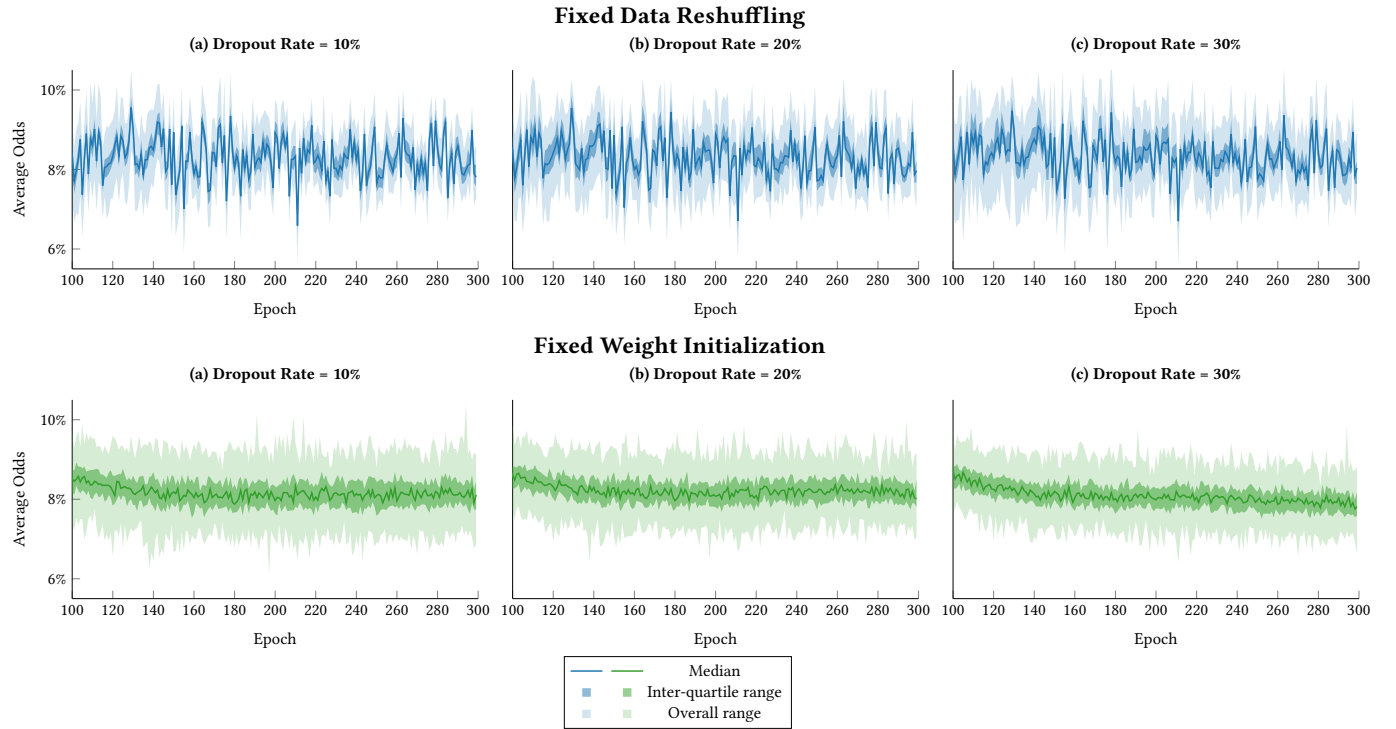


Figure 29: Median, inter-quartile range, and overall range of average odds across 50 different training runs with dropout regularization, while changing only the weight initialization or the random reshuffling, for various dropout rates. The overall range of fairness variance has decreased, and the range of variance across multiple training runs with fixed data reshuffling increases with higher dropout rate, when compared to training without dropout layers. Despite this, the trends of data order dominance over fairness variance are clearly visible even with dropout regularization.

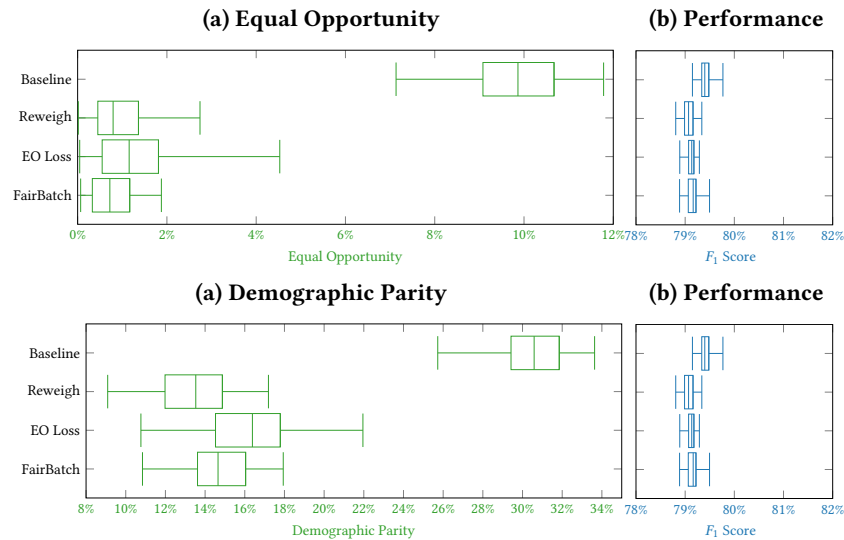


Figure 30: Additional experiments with different fairness metrics for setting in Figure 1. Fairness has a high variance across multiple runs. Note that the x-axis across fairness and F1 score is not similarly scaled for demographic parity to keep the results readable.

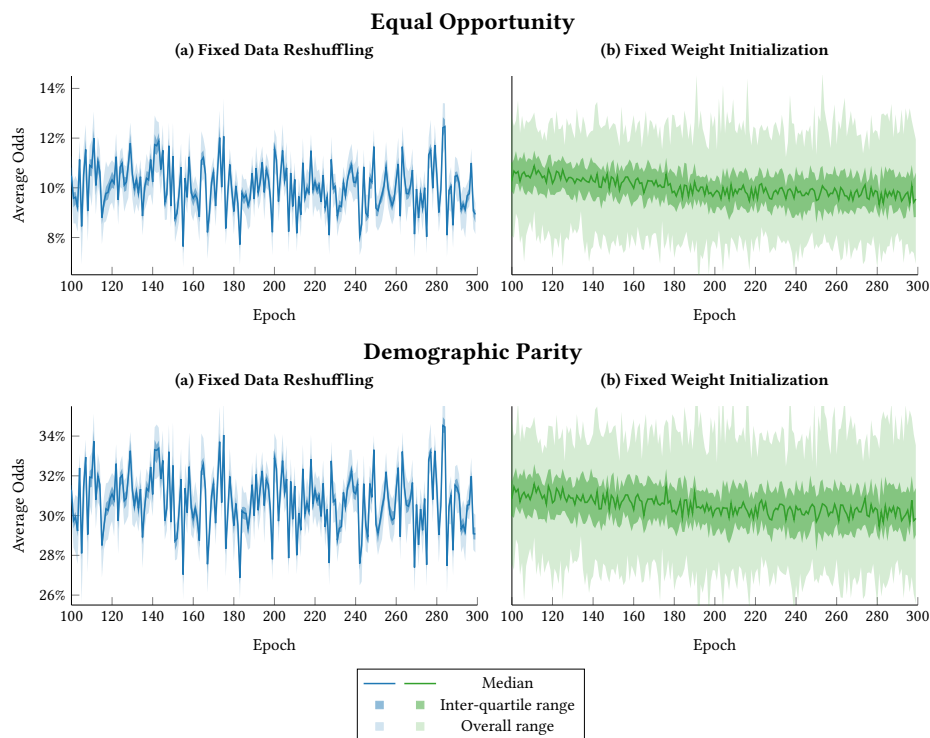


Figure 31: Additional experiments with fairness metric EOpp and DP reveal similar trends as in Figure 4. These results further highlight the dominant impact of random reshuffling on fairness.

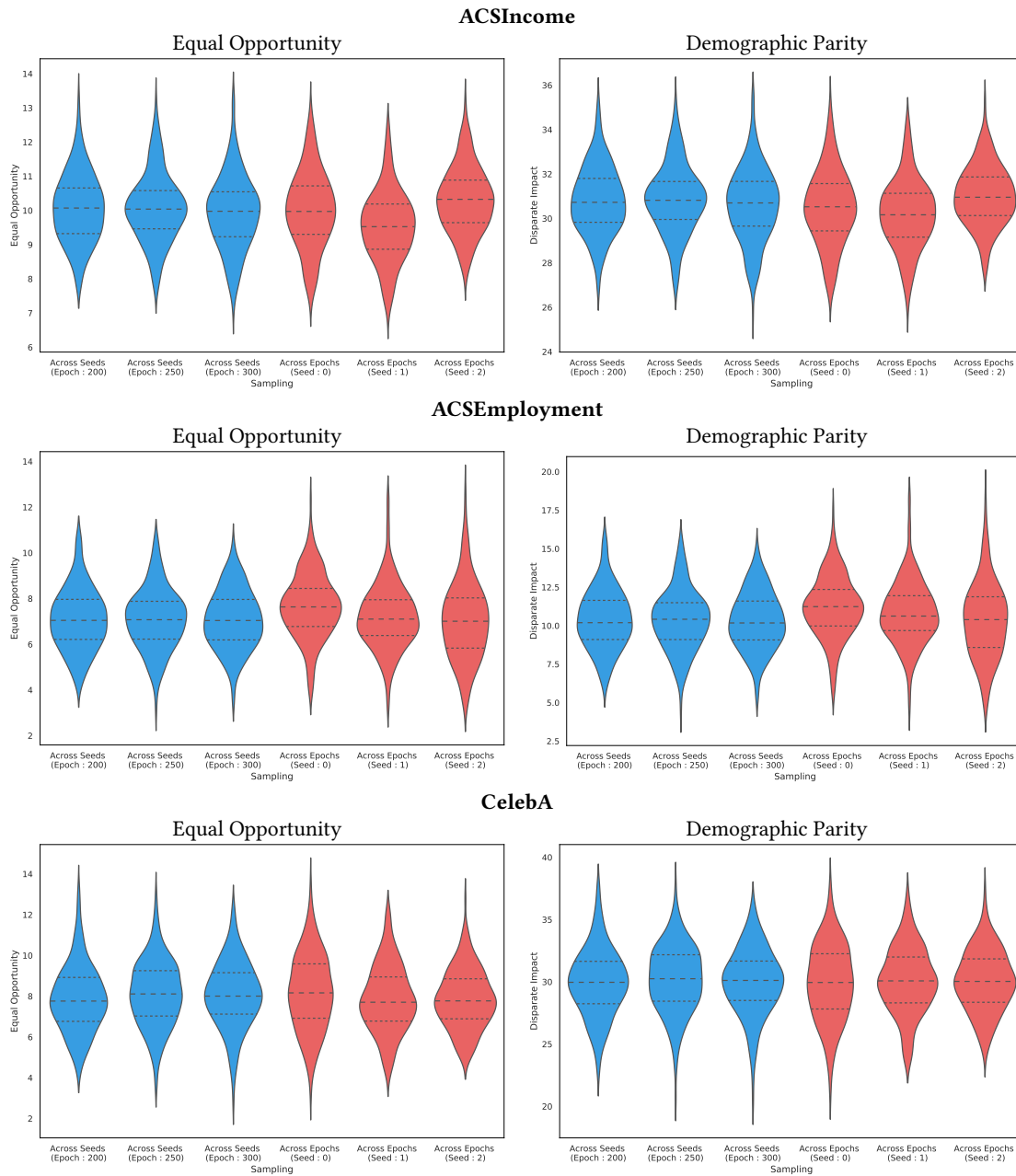


Figure 32: Additional experiments for fairness metrics EOpp and DP. Fairness scores across multiple training runs and across epochs in a single training run have similar empirical distributions. Thus, studying this distribution across epochs provides a highly efficient alternative.

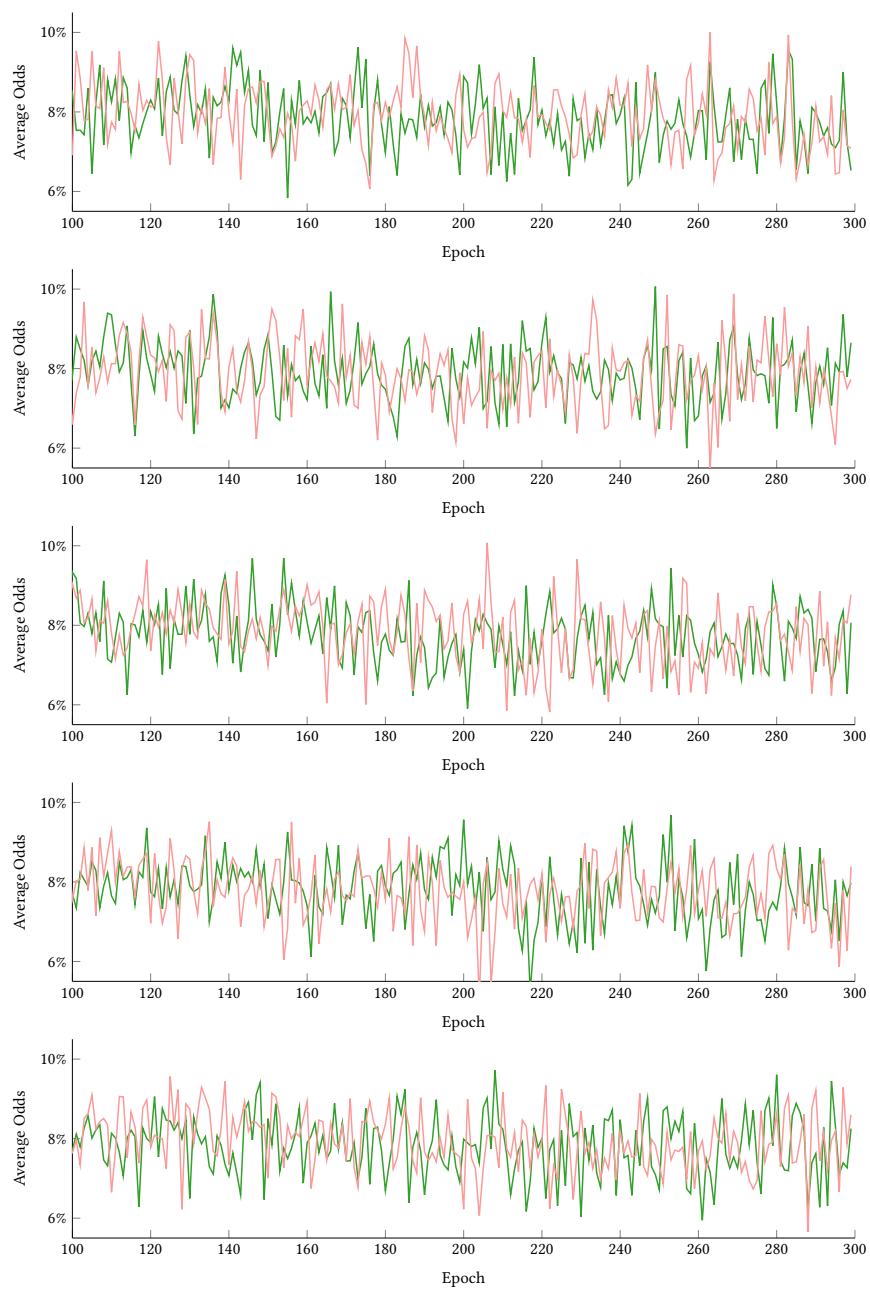


Figure 33: 10 randomly chosen raw training runs plotted in groups of 2 for fixed weight initialization, but changing random reshuffling.

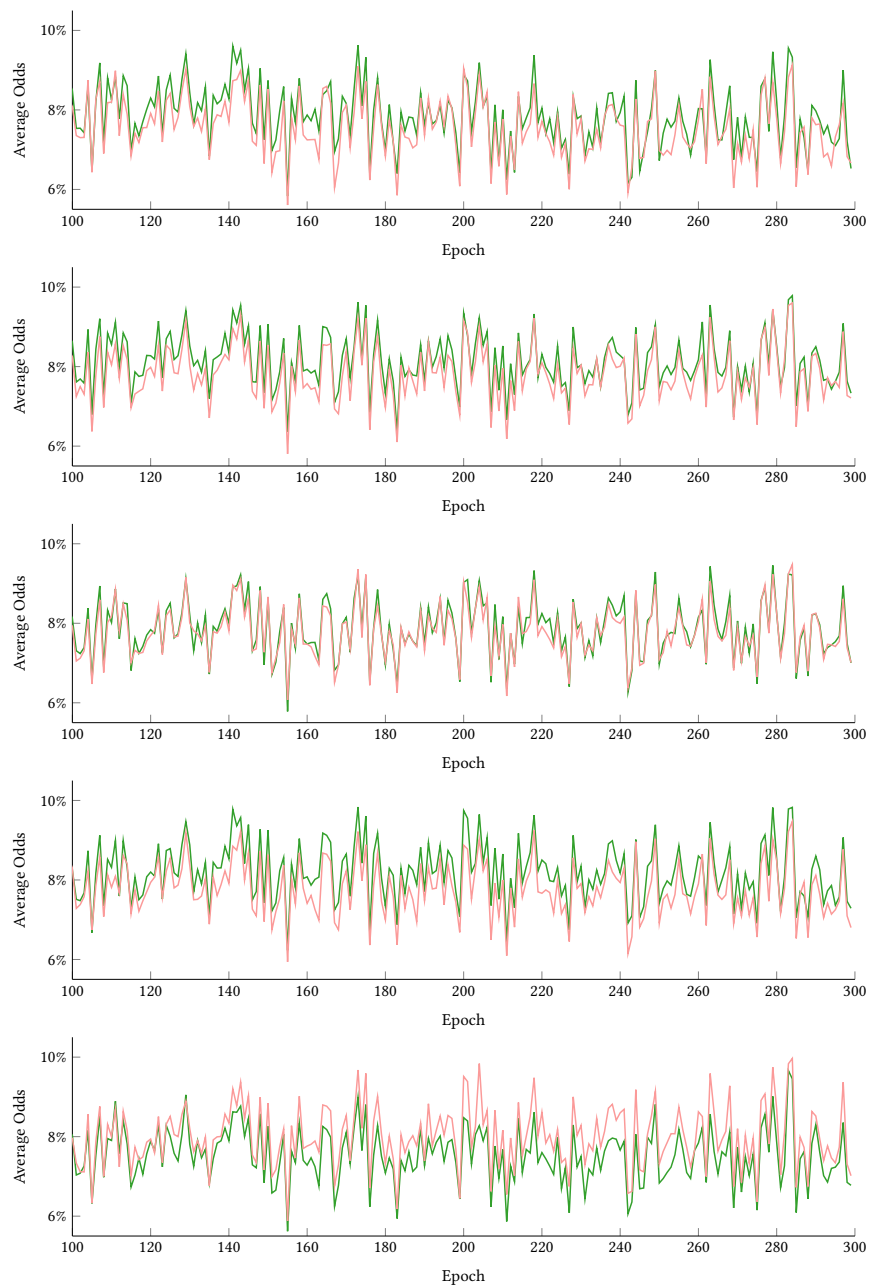


Figure 34: 10 randomly chosen raw training runs plotted in groups of 2 for fixed random reshuffling, but changing weight initialization.