# Pruning has a disparate impact on model accuracy

Anonymous Author(s) Affiliation Address email

## Abstract

1	Network pruning is a widely-used compression technique that is able to significantly
2	scale down overparameterized models with minimal loss of accuracy. This paper
3	shows that pruning may create or exacerbate disparate impacts. The paper sheds
4	light on the factors to cause such disparities, suggesting differences in gradient
5	norms and distance to decision boundary across groups to be responsible for this
6	critical issue. It analyzes these factors in detail, providing both theoretical and
7	empirical support, and proposes a simple, yet effective, solution that mitigates the
8	disparate impacts caused by pruning.

## 9 1 Introduction

As deep learning models evolve and become more powerful, they also become larger and more 10 costly to store and execute. The trend hinders their deployment in resource-constrained platforms, 11 such as embedded systems or edge devices, which require efficient models in time and space. 12 To address this challenge, studies have developed a variety of techniques to prune the relatively 13 insignificant or insensitive parameters from a neural network while ensuring competitive accuracy 14 11. 4. 5. 23. 24. 25. 30. When a model needs to be developed to fit given and certain requirements in 15 size and resource consumption, a pruned model which is derived from a large, rigorously-trained, 16 and (often) over-parameterized model, is regarded as a de-facto standard. That is because it performs 17 incomparably better than a same-size dense model which is trained from scratch, when the same 18 amount of effort and resources are invested. 19

In spite of the strengths of pruning, this paper shows that *pruning can induce or exacerbate disparate* 20 effects in the accuracy of the resulting reduced models. Intuitively, the removal of model weights 21 affects the process in which the network separates different classes, which can have contrasting 22 consequences for different groups of individuals. Specifically, this paper shows that the accuracy of 23 the pruned models tends to increase (decrease) more in classes that had already high (low) accuracy 24 in the original model, leading to a "the rich get richer" and "the poor get poorer" effect. This Matthew 25 effect is illustrated in Figure 1. The figure shows the accuracy of a facial recognition task on different 26 demographic groups for several pruning rates (indicating the percentage of parameters removed from 27 the original models). Notice how the accuracy of the majority group (White) tends to increase while 28 that of the minority groups tends to decrease as the pruning ratio increases. 29

Following these observations, the paper sheds light on the factors to cause such disparities. The theoretical findings suggest the presence of two key factors responsible for why accuracy disparities arise in pruned models: (1) disparity in *gradient norms* across groups, and (2) disparity in *Hessian matrices* associated with the loss function computed using a group's data. Informally, the former carries information about the groups' local optimality, while the latter relates to model separability. The paper analyzes these factors in detail, providing both theoretical and empirical support on a variety of settings, networks, and datasets.



Figure 1: Accuracy of each demographic group in the UTK-Face dataset using Resnet18 [14], at the increasing of the pruning rate.

<sup>37</sup> By recognizing these factors, the paper also develops a simple yet effective training technique that

<sup>38</sup> largely mitigates the disparate impacts caused by pruning. The method is based on an alteration of

<sup>39</sup> the loss function to include components that penalize disparity of the average gradient norms and

<sup>40</sup> distance to decision boundary across groups.

These findings are significant: Pruning is a key enabler for neural network models in embedded
systems with deployments in security cameras and sensors for autonomous devices for applications
where fairness is an essential need. (e.g., face recognition), Without careful consideration of the
fairness impact of these techniques, the resulting models can have profound effects on our society and
economy. To the best of the authors' knowledge, this work is the first to note, analyze, and mitigate

<sup>46</sup> the disparities arising due to network pruning, providing what the authors believe will be a useful

tool for researchers and practitioners in this field.

#### 48 Related work

<sup>49</sup> Fairness and network pruning have been long studied in isolation. The reader is referred to the related

<sup>50</sup> papers and surveys on fairness [3, 6, 8, 13, 17] and pruning [1, 4, 5, 23, 24, 25, 30] for a review on

<sup>51</sup> these areas. Their intersection, however, received little attention.

The recent interest in assessing societal values of machine learning models has seen an increase of studies at the intersection of different properties of a learning model and their effects on fairness. For example, Xu et al. [28] studies the setting of adversarial robustness and show that adversarial training introduces unfair outcomes in term of accuracy parity [31]. Zhu et al. [33] show that semisupervised settings can introduce unfair outcomes in the resulting accuracy of the learned models. Finally, several authors have also shown that private training can have unintended disparate impacts to the resulting models' outputs [2, 10, 26, 32] and downstream decisions [22, 27].

<sup>59</sup> Unfortunately, the literature on the fairness effects of pruning, or more generally, network compression,

has received very sparse attention. Hosseini et al. [15] observed empirically that knowledge distillation processes may produce unfair student models and Paganini [20] observed that a form of network

<sup>61</sup> processes may produce unfair student models and Paganini [20] observed th <sup>62</sup> compression can introduce accuracy disparity among different groups.

These observations are however poorly understood and have not received the attention they deserve given their broad impact on various population segments. It is the goal of this paper to address this critical knowledge gap and provide a step towards a deeper understanding of the fairness issues arising as a result of pruning.

#### 67 **2** Problem settings and goals

<sup>68</sup> The paper considers datasets *D* consisting of *n* datapoints  $(x_i, a_i, y_i)$ , with  $i \in [n]$ , drawn i.i.d. from <sup>69</sup> an unknown distribution  $\Pi$ . Therein,  $x_i \in X$  is a feature vector,  $a_i \in \mathcal{A}$  with  $\mathcal{A} = [m]$  (for some finite <sup>70</sup> *m*) is a demographic group attribute, and  $y_i \in \mathcal{Y}$  is a class label. For example, consider the case of a <sup>71</sup> face recognition task. The training example feature  $x_i$  may describe a headshot of an individual, the <sup>72</sup> protected attribute  $a_i$  may describe the individual's gender or ethnicity, and  $y_i$  represents the identity <sup>73</sup> of the individual. The goal is to learn a predictor  $f_{\theta} : X \to \mathcal{Y}$ , where  $\theta$  is a *k*-dimensional real-valued vector of parameters that minimizes the empirical risk function:

$$\overset{\star}{\boldsymbol{\theta}} = \operatorname*{argmin}_{\boldsymbol{\theta}} J(\boldsymbol{\theta}; D) = \frac{1}{n} \sum_{i=1}^{n} \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), y_i), \tag{1}$$

<sup>75</sup> where  $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_+$  is a non-negative *loss function* that measures the model quality.

The paper focuses on analyzing properties arising when extracting a small model  $f_{\bar{\theta}}$  with  $\bar{\theta} \subset \dot{\theta}$  of

<sup>77</sup> size  $|\theta| = \bar{k} \ll k$ . Model  $f_{\bar{\theta}}$  is constructed by pruning the least important values or filters from vector <sup>78</sup>  $\stackrel{*}{\theta}$  (i.e., those with smaller values in magnitude) according to a prescribed criterion, such as an  $\ell_p$ 

<sup>78</sup>  $\theta$  (i.e., those with smaller values in magnitude) according to a prescribed criterion, such as an  $\ell_p$ <sup>79</sup> norm [18] 24]. The paper focuses on understanding the fairness impacts (as defined next) arising

<sup>80</sup> when pruning general classifiers, such as neural networks.

Fairness The fairness analysis focuses on the notion of *excessive loss*, defined as the difference between the original and the pruned risk functions over some group  $a \in \mathcal{A}$ :

$$R(a) = J(\bar{\theta}; D_a) - J(\hat{\theta}; D_a), \tag{2}$$

where  $D_a$  denotes the subset of the dataset D containing samples  $(x_i, a_i, y_i)$  whose group membership

 $a_i = a$ . Intuitively, the excessive loss represents the change in loss (and thus, in accuracy) that a given

group experiences as a result of pruning. Fairness is measured in terms of the maximal excessive loss

86 *difference*, also referred to as *fairness violation*:

$$\xi(D) = \max_{a,a' \in \mathcal{A}} |R(a) - R(a')|,\tag{3}$$

<sup>87</sup> defining the largest excessive loss difference across all protected groups. (Pure) fairness is achieved <sup>88</sup> when  $\xi(D) = 0$ , and thus a fair pruning method aims at minimizing the excessive loss difference.

The goal of this paper is to shed light on why fairness issues arise (i.e., R(a) > 0) as a result of pruning, why some groups suffer more than others (i.e., R(a) > R(a')), and what mitigation measures could be taken to minimize unfairness due to pruning.

The paper uses the following notation: variables are denoted by calligraph symbols, vectors or matrices by bold symbols, and sets by uppercase symbols. Finally,  $\|\cdot\|$  denotes the Euclidean norm and the paper uses  $f_{\theta}(x)$  to refer to the model' *soft* outputs. All proofs are reported in Appendix A

#### 95 **3** Fairness analysis in pruning: Roadmap

To gain insights on how pruning may introduce unfairness, the paper starts with providing a useful upper bound for a group's excessive loss. Its goal is to isolate key aspects of model pruning that are responsible for the observed unfairness. The following discussion assumes the loss function  $\ell(\cdot)$  to be at least twice differentiable, which is the case for common ML loss functions, such as mean squared error or cross entropy loss.

101 **Theorem 1.** The excessive loss of a group  $a \in \mathcal{A}$  is upper bounded by

$$R(a) \le \left\| \boldsymbol{g}_{a}^{\ell} \right\| \times \left\| \bar{\boldsymbol{\theta}} - \overset{\star}{\boldsymbol{\theta}} \right\| + \frac{1}{2} \lambda \left( \boldsymbol{H}_{a}^{\ell} \right) \times \left\| \bar{\boldsymbol{\theta}} - \overset{\star}{\boldsymbol{\theta}} \right\|^{2} + O\left( \left\| \bar{\boldsymbol{\theta}} - \overset{\star}{\boldsymbol{\theta}} \right\|^{3} \right), \tag{4}$$

where  $g_a^{\ell} = \nabla_{\theta} J(\dot{\theta}; D_a)$  is the vector of gradients associated with the loss function  $\ell$  evaluated at  $\dot{\theta}$ and computed using group data  $D_a$ ,  $H_a^{\ell} = \nabla_{\theta}^2 J(\dot{\theta}; D_a)$  is the Hessian matrix of the loss function  $\ell$ , at the optimal parameters vector  $\dot{\theta}$ , computed using the group data  $D_a$  (henceforth simply referred to as group hessian), and  $\lambda(\Sigma)$  is the maximum eigenvalue of a matrix  $\Sigma$ .

The bound above follows from a second order Taylor expansion of the loss function, Cauchy-Schwarz
 inequality, and properties of the Rayleigh quotient.

Notice that, in addition to the difference in the original and pruned parameters vectors, two key terms appear in Equation (4): (1) The norms of the gradients  $g_a^{\ell}$  and (2) the maximum eigenvalue of the Hessian matrix  $H_a^{\ell}$  for a group *a*. Informally, the former is associated with the groups' local

<sup>&</sup>lt;sup>1</sup>With a slight abuse of notation, the results refer to  $\bar{\theta}$  as the homonymous vector which is extended with  $k - \bar{k}$  zeros.

optimality while the latter relates to the ability of the model to separate the groups data. As we will 111 show next these components represent the main sources of unfairness due to model pruning. 112

The following is an important corollary of Theorem **I**. It shows that the larger the pruning, the larger 113 will be the excessive loss for a given group. 114

**Corollary 1.** Let  $\bar{k}$  and  $\bar{k}'$  be the size of parameter vectors  $\bar{\theta}$  and  $\bar{\theta}'$ , respectively, resulting from 115 pruning model  $f_{\check{h}}$ , where  $\bar{k} < \bar{k}'$  (i.e., the former model prunes more weight than the latter one). Then, 116 for any group  $a \in \mathcal{A}$ , 117

$$\tilde{R}(a,\bar{\theta}) \ge \tilde{R}(a,\bar{\theta}'),\tag{5}$$

where  $\tilde{R}(a,\omega)$  is the excessive loss upper bound computed using pruned model parameters  $\omega$  (Eq. (4)). 118

A consequence of the corollary above is that as the pruning regime increases, the unfairness in 119 accuracy across groups may also become more significant, which the paper shows next. 120

The next sections analyze the effect of gradient norms and the Hessian to unfairness in the pruned 121 models. The theoretical claims are supported and complemented by analytical results. These results 122 use the UTKFace dataset [29] for a vision task whose goal is to classify ethnicity. The experiments 123 use a ResNet-18 architecture and the pruning counterparts remove the P% parameters with the 124 smallest absolute values for various P. All reported metrics are normalized and an average of 10 125 repetitions. While the theoretical analysis focuses on the notion of disparate impacts under the lens 126 of excessive loss, the empirical results report differences in accuracy of the resulting models. The 127 empirical results thus reflect the setting commonly adopted when measuring accuracy parity [31]. 128

The paper reports a glimpse of the empirical results, with the purpose of supporting the theoretical 129 claims, and extended experiments, as well as additional descriptions of the datasets and settings, are 130 reported in Appendix C 131

#### Why disparity in groups' gradients causes unfairness? 4 132

This section analyzes the effect of gradients norms on the unfairness observed in the pruned models. 133 In more detail, it shows that unbalanced datasets result in a model with large differences in gradient 134 norms between groups (Proposition I), it connects gradients norms for a group with the resulting 135 model errors in such a group (Proposition 2), and connects these concepts with the excessive loss 136 (Theorem  $\Pi$ ) to show that unfairness in model pruning is largely controlled by the difference in 137 gradient norms among groups. 138

Gradient norms and group sizes. The section first shows 139 that imbalanced datasets lead a model to have imbalanced 140 gradient norms across groups. The following result assumes 141 that the training converges to a local minima. 142





That is, groups with more data samples will result in smaller 145 gradients norms than groups with fewer data samples and 146 Figure 2: Group size vs. gradient norms. vice-versa. Figure 2 illustrates Proposition 1. The plot shows 147

the relation between groups sizes  $|D_a|$  and their associated gradient norms  $||g_a^{\ell}||$  on the UTK dataset 148 and settings described above. Notice the strong trend between decreasing group sizes and increasing 149 gradient norms for such groups. 150

Gradient norms and accuracy. Next, the section shows a strong connection between the gradient 151 norms of a group and its associated accuracy. The following assumes the models adopt a cross 152 entropy loss (or mean squared error for regression tasks, as shown Appendix A). 153

**Proposition 2.** For a given group  $a \in \mathcal{A}$ , gradient norms can be upper bounded as: 154

$$\|\boldsymbol{g}_{a}^{\ell}\| \in O\left(\sum_{(\boldsymbol{x}, \boldsymbol{y}) \in D_{a}} \underbrace{\|f_{\boldsymbol{\theta}}^{\star}(\boldsymbol{x}) - \boldsymbol{y}\|}_{Accuracy} \times \left\|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}}^{\star}(\boldsymbol{x})\right\|\right).$$



Figure 4: Accuracy, gradient norm, and group Hessian max eigenvalues of each ethnicity group, before and after increasing pruning ratios for UTK-Face dataset. The percentage of data samples across groups *White, Black, Asian, Indian, and Others* is  $\sim 0.42, 0.19, 0.15, 0.15, 0.07$ , respectively.

The above relates gradient norms with an error measure of the classifier to a target label multiplied by the gradient of the predictions. For example, in a classification task with cross entropy loss,  $\ell(f_{\theta}(x), y) = -\sum_{z \in \mathcal{Y}} f_{\theta}^{z}(x)y^{z}$ , where  $f_{\theta}^{z}(x)$  represents the *z*-th element of the output associated with the soft-max layer of model  $f_{\theta}$ , and y is a one-hot encoding of the true label *y*, with  $y^{z}$  representing its *z*-th element, then,

$$\begin{aligned} |\boldsymbol{g}_{a}|| &= \|\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}; \boldsymbol{D}_{a}, )\| = \left\| \frac{1}{|\boldsymbol{D}_{a}|} \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in \boldsymbol{D}_{a}} \nabla_{f} \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}), \boldsymbol{y}) \times \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}}(\boldsymbol{x}) \right\| \\ &= \left\| \frac{1}{|\boldsymbol{D}_{a}|} \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in \boldsymbol{D}_{a}} (f_{\boldsymbol{\theta}}(\boldsymbol{x}) - \boldsymbol{y}) \times \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}}(\boldsymbol{x}) \right\| \\ &\leq \frac{1}{|\boldsymbol{D}_{a}|} \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in \boldsymbol{D}_{a}} \|f_{\boldsymbol{\theta}}(\boldsymbol{x}) - \boldsymbol{y}\| \times \|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}}(\boldsymbol{x})\| . \end{aligned}$$



A similar observation holds for mean square error loss, as
 illustrated in Appendix A. The observation above sheds light
 on the correlation between the prediction error of a group and

163

164

its model gradients. This relation is emphasized in Figure 3: Figure 3: Accuracy vs. gradient norms. which illustrates that the gradient norm for a given group increases as its prediction accuracy decreases.

Proposition 2 allows us to link the gradient norms with the group accuracy of the resulting model, which, together with the result above will be useful to reason about the impact of gradient norms on the disparities in the group excessive losses.

**The role of gradient norms in pruning.** Having highlighted the connection between gradients norms of a group with the accuracy of the pruned model on such a group, this section provides theoretical intuitions on the role of gradient norms in the disparate group losses during pruning.

From Theorem 1, notice that the excessive loss is controlled by term  $\|g_a^{\ell}\| \times \|\bar{\theta} - \check{\theta}\|$ . As already 171 noted in Corollary 1, the term  $\|\bar{\theta} - \dot{\theta}\|$  regulates the impact of pruning on the excessive loss, as the 172 difference between the pruned and non-pruned parameters vectors directly depends on the pruning 173 rate. For a fixed pruning rate, however, notice that groups with different gradient norms will have a 174 disparate effect on the resulting term. In particular, groups with very small gradients norms (those 175 generally associated with highly accurate predictions) will be less sensitive to the effects of the 176 pruning rate. Conversely, groups with large gradient norms will be affected by the pruning rate to a 177 greater extent, with larger pruning rates, *typically* reflecting in larger excessive losses. 178

These observations of the factors of disparity, accuracy, and group size, can also be appreciated empirically in Figures 4a and 4b. The plots report accuracy (a) and gradient norms (b) on the UTKFace datasets for a variety of pruning rates. Consider group *White* (containing 42% of the total samples) and *Others* (containing 7% of the total samples). The unpruned model has high accuracy on the former group and small gradient norms. The accuracy of this group is insensitive to various pruning rates and even increases at large pruning regimes. In contrast, group *Others* has much lower accuracy and larger gradient norms in the unpruned model. As the pruning rate increase, their

accuracies drastically drop. As a result, in high pruning regimes, this minority group exhibits poor 186 accuracy and very high gradient norms. 187

Notice that the empirical results apply to much more complex settings than those which can be 188 analyzed formally, thus they complement the theoretical observations. 189

#### Why disparity in groups' Hessians causes unfairness? 5 190

Having examined the properties of the groups gradients and their relation to unfairness in pruning, 191 this section turns on analyzing how the Hessian associated with the loss function for a group is linked 192 to the unfairness observed during pruning. In more detail, it connects the groups' Hessian to the 193 distance to the decision boundary for the samples in that group and their resulting model errors 194 (Theorem 3), it illustrates a strong positive correlation between groups' Hessian and gradient norms, 195 196 and links these concepts with the excessive loss (Theorem 1) to show that unfairness in model pruning is controlled by the difference in maximum eigenvalues of the Hessians among groups. 197

**Group Hessians and accuracy.** The section first shows that groups presenting large Hessian values 198 may suffer larger disparate impacts due to pruning, when compared with groups that have smaller 199 Hessians. It does so by connecting the maximum eigenvalues of the groups Hessians with their 200 distance to decision boundary and the group accuracy. The following result sheds light on these 201 observations. It restricts its attention to models trained under binary cross entropy losses, for clarity 202 of explanation, although an extension to a multi-class case is directly attainable. 203

**Theorem 2.** Let  $f_{\theta}$  be a binary classifier trained using a binary cross entropy loss. For any group 204  $a \in \mathcal{A}$ , the maximum eigenvalue of the group Hessian  $\lambda(\mathbf{H}_a^{\ell})$  can be upper bounded by: 205

$$\lambda(\boldsymbol{H}_{a}^{\ell}) \leq \frac{1}{|\boldsymbol{D}_{a}|} \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in \boldsymbol{D}_{a}} \underbrace{\left(f_{\boldsymbol{\theta}}^{\star}(\boldsymbol{x})\right) \left(1 - f_{\boldsymbol{\theta}}^{\star}(\boldsymbol{x})\right)}_{Distance \ to \ decision \ boundary} \times \left\|\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}}^{\star}(\boldsymbol{x})\right\|^{2} + \underbrace{\left|f_{\boldsymbol{\theta}}^{\star}(\boldsymbol{x}) - \boldsymbol{y}\right|}_{Accuracy} \times \lambda\left(\nabla_{\boldsymbol{\theta}}^{2} f_{\boldsymbol{\theta}}^{\star}(\boldsymbol{x})\right). \tag{6}$$

The proof relies on derivations of the Hessian associated with model loss function and Weyl inequality. 206 In other words, Theorem 3 highlights a direct connection between the maximum eigenvalue of the 207 group Hessian and (1) the closeness to the decision boundary of the group samples, and (2) the 208 accuracy of the group. The distance to the decision boundary is derived from [7]. Intuitively this 209 term is maximized when the classifier is highly uncertain about the prediction:  $f^{*}_{\theta}(x) \to 0.5$ , and 210 minimized when it is highly certain  $f^*_{\theta}(x) \to 0$  or 1, as showed in the following proposition. 211

**Proposition 3.** Consider a binary classifier  $f_{\theta}(x)$ . For a given sample  $x \in D$ , the term  $f_{\theta}(x)(1 - f_{\theta}(x))$ 212  $f^{\star}_{\theta}(x)$ ) is maximized when  $f^{\star}_{\theta}(x) = 0.5$  and minimized when  $f^{\star}_{\theta}(x) \in \{0, 1\}$ . 213

Observe that a group consisting of samples that are far from 214 the decision boundary will have smaller Hessians and, thus, be 215 less subject to a drop in accuracy due to model pruning. These 216 results can be appreciated in Figure 5. Notice the inverse 217 relationship between maximum eigenvalues of the groups' 218 Hessians and their average distance to the decision boundary. 219 The same relation also holds for accuracy: the higher the 220 Hessians maximum eigenvalues, the smaller the accuracy. This 221 is intuitive as samples which are close to the decision boundary 222 will be more prone to errors due to small changes in the model 223



Figure 5: Group Hessians, distance to decision boundary, and accuracy.

due to pruning, when compared with samples lying far from the decision boundary. 224

Correlation between group Hessians and gradient norms. This section observes a positive 225 correlation between maximum eigenvalues of the Hessian of a group and their gradient norms. This 226 relation can be appreciated in Figure 6. While mainly empirical, this observation is important as it 227 illustrates that both the Hessian  $\lambda(H_a^l)$  and the gradient  $\|g_a^\ell\|$  terms appearing in the upper bound of 228 the excessive loss R(a) reported in Theorem T are in agreement. This relation was observed in all 229 our experiments and settings. Such observation allows us to infer that it is the combined effect of 230 gradient norms and group Hessians that is responsible for the excessive loss of a group and, in turn, 231 for the exacerbation of unfairness in the pruned models. 232

The role of the group Hessian in pruning. Having highlighted the connection between Hessian for a group with the resulting accuracy of the model on such a group, this section provides theoretical intuitions on the role of the Hessians in the disparate group losses during pruning.

In Theorem 1 notice that the excessive loss is controlled by term  $\|\boldsymbol{H}_{a}^{\ell}\| \times \|\boldsymbol{\bar{\theta}} - \boldsymbol{\dot{\theta}}\|^{2}$ . As also noted in the previous section, the term  $\|\boldsymbol{\bar{\theta}} - \boldsymbol{\dot{\theta}}\|$  regulates the impact of pruning on the excessive loss as the difference between the pruned and non-pruned parameters vectors directly depends on the pruning rate. Similar to the observation for gradient norms, with a fixed pruning rate, groups with



Figure 6: Group Hessians and gradient norms.

different Hessians will have a disparate effect on the resulting term. In particular, groups with small
Hessians eigenvalues (those generally distant from the decision boundary and highly accurate) will
be less sensitive to the effects of the pruning rate. Conversely, groups with large Hessians eigenvalues
will be affected by the pruning rate to a greater extent, *typically* resulting in larger excessive losses.
These observations can further be appreciated empirically in Figures 4a (for accuracy) and 4c (for
maximum group Hessian eigenvalues) on the UTKFace datasets for a variety of pruning rates.

# **6** Mitigation solution and evaluation

The previous sections highlighted the presence of two key factors playing a role in the observed model accuracy disparities due to pruning: the difference in gradient norms, and the difference in Hessians losses across groups. This section first shows how to leverage these findings to provide a simple, yet effective solution to reduce the disparate impacts of pruning. Then, the section illustrates the benefits of this mitigating solution on a variety of tasks, datasets, and network architectures.

#### 256 6.1 Mitigation solution

To achieve fairness, the aforementioned findings suggest to equalize the disparity associated with gradient norms  $\|g_a^{\ell}\|$  and Hessians  $\lambda(H_a^{\ell})$  across different groups  $a \in \mathcal{A}$ . For this goal, the paper adopts a constrained empirical risk minimization approach:

$$\underset{\boldsymbol{\theta}}{\text{minimize}} \quad J(\boldsymbol{\theta}; D) \quad \text{such that:} \quad \|\boldsymbol{g}_{a}^{\ell}\| = \|\boldsymbol{g}^{\ell}\|, \quad \lambda(\boldsymbol{H}_{a}^{\ell}) = \lambda(\boldsymbol{H}^{\ell}) \quad \forall a \in \mathcal{A}, \tag{7}$$

where  $g^{\ell} = \nabla_{\theta} J(\theta; D)$  and  $H^{\ell} = \nabla^2_{\theta} J(\theta; D)$  refer to the gradients and Hessian associated with loss function  $\ell$ , respectively, and are computed using the whole dataset D. The approach (7) is a common strategy adopted in fair learning tasks, and the paper uses the Lagrangian Dual method of Fioretto et al. (9) which exploits Lagrangian duality to extend the loss function with trainable and weighted regularization terms that encapsulate constraints violations (see Appendix C for additional details).

A shortcoming of this approach is, however, that requires computing the gradient norms and Hessian matrices of the group losses in each and every training iteration, rendering the process computationally unviable, especially for deep, overparametrized networks.

To overcome this computational burden, we will use two observations made earlier in the paper. 268 First, recall the strong relation between gradient norms for a group and their associated losses. This 269 aspect was noted in Proposition 2. That is, when the losses across the groups are similar, the gradient 270 norms across such groups will also tend to be similar. Next, Theorem 3 noted a positive correlation 271 272 between model errors (and thus loss values) for a group and its associated Hessian eigenvalues. Thus, when the losses across the groups are similar, the group Hessians will also tend to be similar. This 273 intuition is also complemented by the strong correlation between group Hessians and gradient norms 274 reported in Section 5. Based on the above observations, the paper proposes a simpler version of the 275 constrained minimizer (7) defined as 276

$$\underset{\boldsymbol{\theta}}{\text{minimize}} \quad J(\boldsymbol{\theta}; D) \quad \text{such that:} \quad J(\boldsymbol{\theta}; D_a) = J(\boldsymbol{\theta}; D) \quad \forall a \in \mathcal{A}, \tag{8}$$

that substitutes the gradient norms and max eigenvalues of group Hessians equality constraints with proxy terms capturing the group  $J(\theta; D_a)$  and population  $J(\theta; D)$  losses.



Figure 8: Accuracy and Fairness violations attained by all models on ResNet50, UTK-Face dataset with *ethnicity* (5 classes) as group attribute (and labels) [left] and *age* (9 classes) [right].



Figure 9: Accuracy and Fairness violations attained by all models on VGG-19, CIFAR-10 dataset (left) and SVHN (right) with 10 class labels also used as group attribute.

The impact of such proxy terms in the fairness 279 constrained program above can be appreciated, 280 empirically, in Figure 7. The plots, that use the 281 282 UTK-Face dataset, with Ethnicity as protected group, show an original unfair model (top) and 283 a fair counterpart obtained through Program (8) 284 (bottom). Notice how enforcing balance in the 285 group losses also helps reducing and balancing 286 the gradient norms and group's average distance 287 to the decision boundary. As a consequence, 288 the resulting model fairness is dramatically en-289 hanced (bottom-left subplot). 290

#### 291 6.2 Assessment of the mitigation solution

292

Datasets, models, and settings. This section



Decision boundar

Figure 7: Effects of fairness constraints in balancing not only group accuracy (left) but also gradient norms (middle) and group average distance to the decision boundary (right).

analyzes the results obtained using the proposed mitigation solution with ResNet50 and VGG19 on the UTKFace dataset [29], CIFAR-10 [16], and SVHN [19] for various protected attributes. The experiments compare the following four models:

• No Mitigation: it refers to the standard pruning approach which uses no fairness mitigation strategy.

• *Fair Bf Pruning*: it applies the fairness mitigation process (Problem (8)) exclusively to the original large network, thus *before* pruning.

• Fair Aft Pruning: it applies the mitigation exclusively to the pruned network, thus after pruning.

• *Fair Both*: it applies the mitigation both to the original large network and to the pruned network.

The experiments report the overall accuracy of resulting models as well as their fairness violations, defined here as the difference between the maximal and minimal group accuracy. The reported metrics are the average of 10 repetitions. Additional details on datasets, architectures, hyper-parameters adopted, as well as additional and extended results are reported in Appendix [C].

Effects on accuracy. The section first focuses on analyzing the effects of accuracy drop due to applying the proposed mitigation solution for fair pruning. Figure 8 compares the four models on the UTK-Face dataset using a ResNet50 architecture. The left subplots use *ethnicity* as protected group and class label, with  $|\mathcal{Y}| = 5$ , while the right subplots use *age* as protected group and class label, with  $|\mathcal{Y}| = 9$ . Notice that, as expected, all compared models present some drop in accuracy as the pruning rates increase. However, notably, the accuracy drops of the models that apply the fair mitigation steps are comparable to (or even improved) those of the "*No mitigation*" model, which applies standard pruning.

A similar trend can be seen in Figure that reports results on CIFAR (left) and SVHN (right). Both use the ten class labels as protected attributes. These results clearly illustrate the ability of the mitigating solution to preserve highly accurate models.

**Effects on fairness.** The section next illustrates the ability of the proposed solution to achieve fair 316 pruned models. The second and fourth subplots presented in Figures 8 and 9 illustrate the fairness 317 violations obtained by the four models analyzed on different datasets and settings. The paper makes 318 the following observations: First, all the plots exhibit a consistent trend in that the mitigation solution 319 produces models which improve the fairness of the baseline, "No mitigation" model. Observe that, as 320 already illustrated in Figure 7 the fair models tend to equalize the gradient norms and group Hessians 321 components (and thus the distance to the decision boundary across groups). Thus, the resulting 322 pruned models also attain better fairness, when compared to their standard counterparts. 323

Next, notice that "*Fair Aft Pruning*" often achieves better fairness violations than "*Fair Bf Pruning*", especially at high pruning regimes. This is because the former has the advantage to apply the mitigation solution directly to the pruned model to ensure that the resulting model has low differences in gradient norms and group Hessians. The presentation also illustrates the application of the mitigation strategies both before and after pruning (*Fair Both*) which shows once again the significance of applying the mitigation solution over the pruned network.

Finally, it is notable that "*Fair Aft Pruning*" achieves good reductions in fairness violation. Indeed, pre-trained large (non-pruned) fair models may not be available and the ability to retrain these large models prior to pruning may be hindered by their size and complexity.

# **7** Discussion and limitations

This section discusses three key messages found in this study. First, we notice that pruning affecting model separability and distance to the decision boundary is related to concepts also explored in robust machine learning [11, [21]]. Not surprisingly, some recent literature in network pruning has empirically observed that pruning may have a negative impact on adversarial robustness [12]. These observations raise questions about the connection between pruning, robustness, and fairness, which we believe is an important direction to further investigate.

Next, although the solution proposed in Problem (8) allows it to be adopted in large models, the size of modern ML models (together with the amount of hyperparameters searches) may hinder retraining such original massive models from incorporating fairness constraints. Notably, however, the proposed mitigation solution can be used as a post-processing step to be applied during the pruning operation directly. The previous section shows that the proposed method delivers desirable performance in terms of both accuracy and fairness.

Finally, we notice that the results analyzed in this paper pertain to losses that are twice differentiable. Lifting such an assumption will be an interesting and challenging future research avenue.

# 348 8 Conclusion

This work observed that pruning, while effective in compressing large models with minimal loss of accuracy, can result in substantial disparate accuracy impacts. The paper examined the factors causing such disparities both theoretically and empirically showing that: (1) disparity in gradient norms across groups and (2) disparity in Hessian matrices associated with the loss functions computed using a groups' data are two key factors responsible for such disparities to arise. By recognizing these factors, the paper also developed a simple yet effective retraining technique that largely mitigates the disparate impacts caused by pruning.

As reduced versions of large, overparametrized models become increasingly adopted in embedded systems to facilitate autonomous decisions, we believe that this work makes an important step toward *understanding* and *mitigating* the sources of disparate impacts observed in compressed learning models.

## 360 **References**

- [1] N. Aghli and E. Ribeiro. Combining weight pruning and knowledge distillation for cnn
   compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3191–3198, 2021.
- [2] E. Bagdasaryan, O. Poursaeed, and V. Shmatikov. Differential privacy has disparate impact on model accuracy. In *Advances in Neural Information Processing Systems*, pages 15479–15488, 2019.
- [3] S. Barocas, M. Hardt, and A. Narayanan. Fairness in machine learning. *Nips tutorial*, 1:2, 2017.
- [4] C. Baykal, L. Liebenwein, I. Gilitschenski, D. Feldman, and D. Rus. Sipping neural networks:
   Sensitivity-informed provable pruning of neural networks. *arXiv preprint arXiv:1910.05422*, 2019.
- [5] D. Blalock, J. J. G. Ortiz, J. Frankle, and J. Guttag. What is the state of neural network pruning? *arXiv preprint arXiv:2003.03033*, 2020.
- [6] S. Caton and C. Haas. Fairness in machine learning: A survey. *arXiv preprint arXiv:2010.04053*, 2020.
- [7] J. Cohen, E. Rosenfeld, and Z. Kolter. Certified adversarial robustness via randomized smoothing. In *International Conference on Machine Learning*, pages 1310–1320. PMLR, 2019.
- [8] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. Zemel. Fairness through awareness. In
   *Proceedings of the 3rd innovations in theoretical computer science conference*, pages 214–226,
   2012.
- [9] F. Fioretto, P. V. Hentenryck, T. W. Mak, C. Tran, F. Baldo, and M. Lombardi. Lagrangian
   duality for constrained deep learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 118–135. Springer, 2020.
- [10] F. Fioretto, C. Tran, P. V. Hentenryck, and K. Zhu. Differential privacy and fairness in decisions and learning tasks: A survey. *CoRR*, abs/2202.08187, 2022. URL https://arxiv.org/abs/ 2202.08187.
- [11] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples,
   2014. URL https://arxiv.org/abs/1412.6572.
- [12] Y. Guo, C. Zhang, C. Zhang, and Y. Chen. Sparse dnns with improved adversarial robustness.
   In *Proceedings of the International Conference on Neural Information Processing Systems* (*NeurIPS*), page 240–249, 2018.
- [13] M. Hardt, E. Price, E. Price, and N. Srebro. Equality of opportunity in supervised learning. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL https://proceedings.neurips.cc/paper/2016/file/
- 9d2682367c3935defcb1f9e247a97c0d-Paper.pdf
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In
   *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–
   778, 2016.
- [15] S. Hosseini, M. A. Shabani, M. M. Jahanara, and B. Salamatian. Learning fair from unfair teachers.
- [16] A. Krizhevsky, V. Nair, and G. Hinton. Cifar-10 (canadian institute for advanced research).
   URL http://www.cs.toronto.edu/~kriz/cifar.html.
- <sup>403</sup> [17] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan. A survey on bias and <sup>404</sup> fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6):1–35, 2021.
- [18] M. C. Mozer and P. Smolensky. Skeletonization: A technique for trimming the fat from a network via relevance assessment. In D. Touretzky, editor, *Advances in Neural Information*
- 407 Processing Systems, volume 1. Morgan-Kaufmann, 1988. URL https://proceedings.
   408 neurips.cc/paper/1988/file/07e1cd7dca89a1678042477183b7ac3f-Paper.pdf.
- [19] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Ng. Reading digits in natural images
   with unsupervised feature learning. 2011.
- [20] M. Paganini. Prune responsibly. *arXiv preprint arXiv:2009.09936*, 2020.

- [21] N. Papernot, P. McDaniel, and I. Goodfellow. Transferability in machine learning: from
   phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv: Arxiv- 1605.07277*, 2016.
- [22] D. Pujol, R. McKenna, S. Kuppam, M. Hay, A. Machanavajjhala, and G. Miklau. Fair decision
   making using privacy-protected data. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 189–199, 2020.
- [23] A. Renda, J. Frankle, and M. Carbin. Comparing rewinding and fine-tuning in neural network
   pruning. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=S1gSj0NKvB.
- [24] J. T. S. Han, J. Pool and W. J. Dally. Learning both weights and connections for efficient neural
   networks. In *NIPS*, 2015. URL https://arxiv.org/abs/1506.02626v3
- [25] V. Sehwag, S. Wang, P. Mittal, and S. Jana. Towards compact and robust deep neural networks.
   *preprint arXiv:1906.06110*, 2019.
- <sup>425</sup> [26] C. Tran, M. Dinh, and F. Fioretto. Differentially private empirical risk minimization under the fairness lens. In *Advances in Neural Information Processing Systems*, 2021.
- [27] C. Tran, F. Fioretto, P. V. Hentenryck, and Z. Yao. Decision making with differential privacy
   under a fairness lens. In Z. Zhou, editor, *International Joint Conference on Artificial Intelligence* (*IJCAI*), pages 560–566, 2021.
- [28] H. Xu, X. Liu, Y. Li, A. K. Jain, and J. Tang. To be robust or to be fair: Towards fairness in adversarial training, 2021.
- [29] S. Y. Zhang, Zhifei and H. Qi. Age progression/regression by conditional adversarial autoencoder. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017.
- [30] T. Zhang, S. Ye, K. Zhang, J. Tang, W. Wen, M. Fardad, and Y. Wang. A systematic dnn weight
   pruning framework using alternating direction method of multipliers. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 184–199, 2018.
- [31] H. Zhao and G. Gordon. Inherent tradeoffs in learning fair representations. *Advances in neural information processing systems*, 32:15675–15685, 2019.
- [32] K. Zhu, P. Van Hentenryck, and F. Fioretto. Bias and variance of post-processing in differential privacy. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 11177–11184, 2021.
- [33] Z. Zhu, T. Luo, and Y. Liu. The rich get richer: Disparate impact of semi-supervised learning.
   *arXiv preprint arXiv:2110.06282*, 2021.

# 445 NeurIPS 2022 Paper Checklist

446	1. For all authors
447	a) Do the main claims made in the abstract and introduction accurately reflect the paper's
448	contributions and scope?
449	Yes. The paper contributions are stated in the abstract and listed in the Introduction.
450	b) (b) Have you read the ethics review guidelines and ensured that your paper conforms
451	to them?
452	Yes.
453	c) Did you discuss any potential negative societal impacts of your work?
454	This work sheds light on the reasons behind the observed disparate impacts and fairness
455	violations through pruning. Pruning is a widely-used compression technique for large-
456	scale models which are deployed in settings with less resources. Thus, the insights
457	d) Did you describe the limitations of your work?
458	Vas Plaga, see section 7
459	2. If you are including theoretical results
460	2. If you are including theoretical results
461	(a) Did you state the full set of assumptions of all theoretical results?
462	Yes. The assumptions were stated in or before each Theorem and also reported in the
463	Appendix. (b) Did you include complete proofe of all theoretical results?
464	(b) Did you include complete proofs of all theoretical results?
465	proofs are reported in Appendix A
400	2 If you ran experiments
467	5. If you ran experiments
468	(a) Did you include the code, data, and instructions needed to reproduce the main experi-
469	Vas Code datasets and experiments were submitted in the supplemental material. We
470	also provide a link in Appendix $\square$
472	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
473	were chosen)?
474	Yes. See Appendix $\overline{C}$
475	(c) Did you report error bars (e.g., with respect to the random seed after running experi-
476	ments multiple times)?
477	The main evaluation metric adopted in this work is the excessive loss (see Equations
478	(2) and (3) which implicitly captures the randomness of the private mechanisms.
479	Providing error bars would be misleading.
480	(d) Did you include the amount of compute and the type of resources used (e.g., type of GPUs internal cluster or cloud provider)?
401	Yes See Appendix
402	A If you are using existing assets (e.g. code data models) or curating/releasing new assets
483	4. If you are using existing assets (e.g., code, data, models) of curating/releasing new assets
484	(a) If your work uses existing assets, did you cite the creators?
485	(b) Did you mention the license of the assets?
485	(b) Did you mention the needse of the assets: Yes, when available
407	(c) Did you include any new assets either in the supplemental material or as a URL?
489	No new asset was required to perform this research.
490	(d) Did you discuss whether and how consent was obtained from people whose data you're
491	using/curating?
492	Yes. The paper uses public datasets.
493	(e) Did you discuss whether the data you are using/curating contains personally identifiable
494	information or offensive content?
495	No. The adopted data is composed of standard benchmarks that have been used
496	extensively in the ML literature and we believe the above does not apply.
497	5. If you used crowdsourcing or conducted research with human subjects
498	No. This research did not use crowdsourcing.