

Motivation

1. Pruning is a widely used network compression technique in resource-constrained platforms.
2. It works by removing least important parameters/channels/filters in a network.
3. A majority of prior work only focus on accuracy, but..
4. We showed empirically and theoretically that model pruning can exacerbate the unfairness among demographic groups.

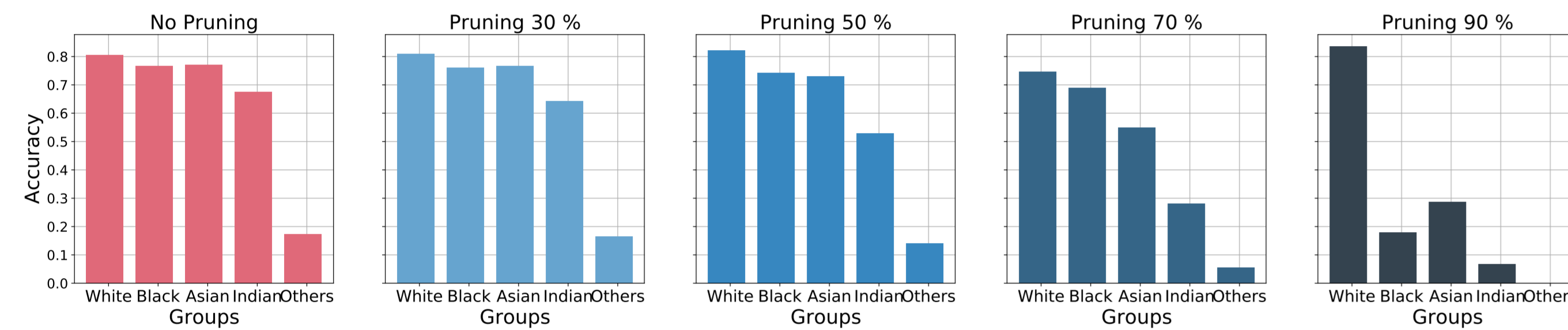


Figure 1. Accuracy of each demographic group in the UTK-Face dataset using Resnet18 [18], at the increasing of the pruning rate.

Settings

- Standard model training:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta; D) = \frac{1}{n} \sum_{i=1}^n \ell(f_{\theta}(x_i), y_i), \quad (1)$$

- Usually θ^* is over-parameterized which can costs memory and storage. Pruning is applied by removing $p\%$ least important parameters in θ^* to obtain $\hat{\theta}$.
- We define the excessive risk as the difference of loss before and after pruning:

$$R(a) = J(\hat{\theta}; D_a) - J(\hat{\theta}^*; D_a), \quad (2)$$

- Fairness is defined as the maximum difference of group excessive risks:

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

Fairness analysis in pruning: Roadmap

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

$$R(a) \leq \|g_a^{\ell}\| \times \|\hat{\theta} - \hat{\theta}^*\| + \frac{1}{2} \lambda(\mathbf{H}_a^{\ell}) \times \|\hat{\theta} - \hat{\theta}^*\|^2 + \mathcal{O}(\|\hat{\theta} - \hat{\theta}^*\|^3), \quad (4)$$

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

Why disparity in groups' gradients causes unfairness?

Proposition 1: Consider two groups a and b in \mathcal{A} with $|D_a| \geq |D_b|$. Then $\|g_a^{\ell}\| \leq \|g_b^{\ell}\|$.

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

$$\|g_a^{\ell}\| \in \mathcal{O} \left(\sum_{(x,y) \in D_a} \underbrace{\|f(x) - y\|}_{\text{Error}} \times \|\nabla_{\theta} f(x)\| \right).$$

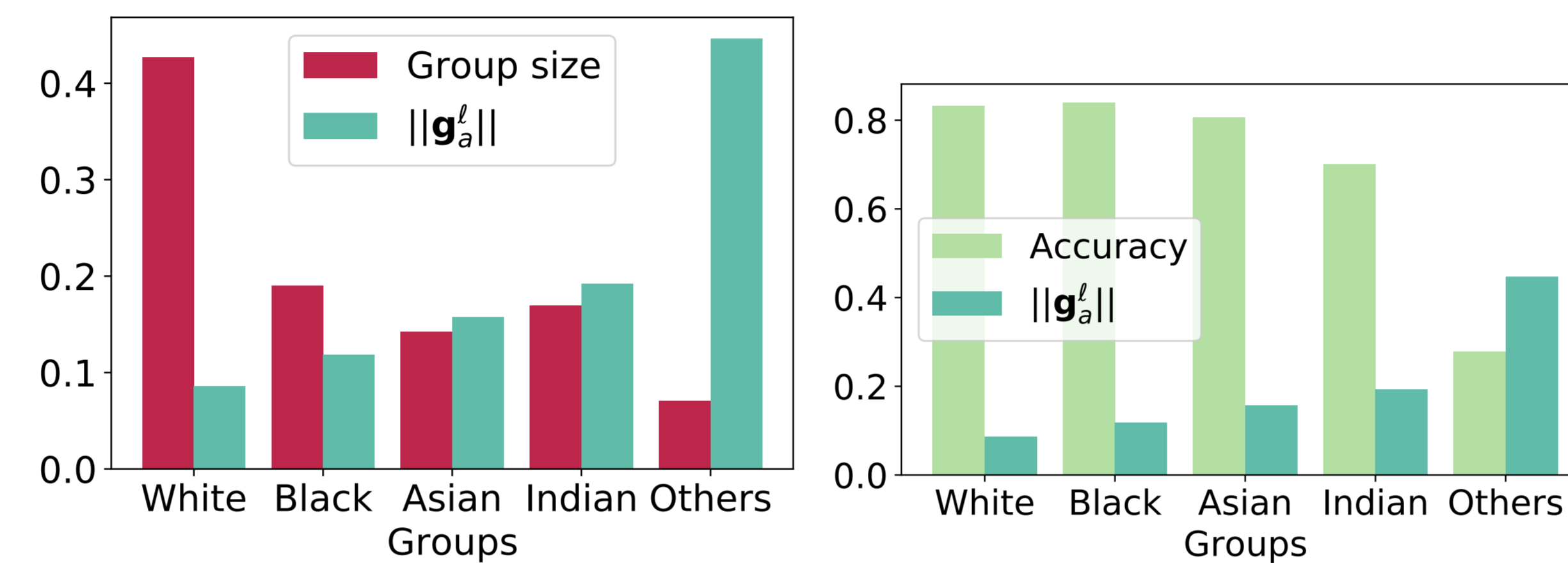


Figure 2. Impact of group sizes and accuracy towards group gradient norms

Why disparity in groups' Hessians causes unfairness?

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

$$\lambda(\mathbf{H}_a^{\ell}) \leq \frac{1}{|D_a|} \sum_{(x,y) \in D_a} \underbrace{(f(x))(1-f(x))}_{\text{Closeness to decision boundary}} \times \|\nabla_{\theta} f(x)\|^2 + \underbrace{\|f(x) - y\|}_{\text{Error}} \times \lambda(\nabla_{\theta}^2 f(x)). \quad (5)$$

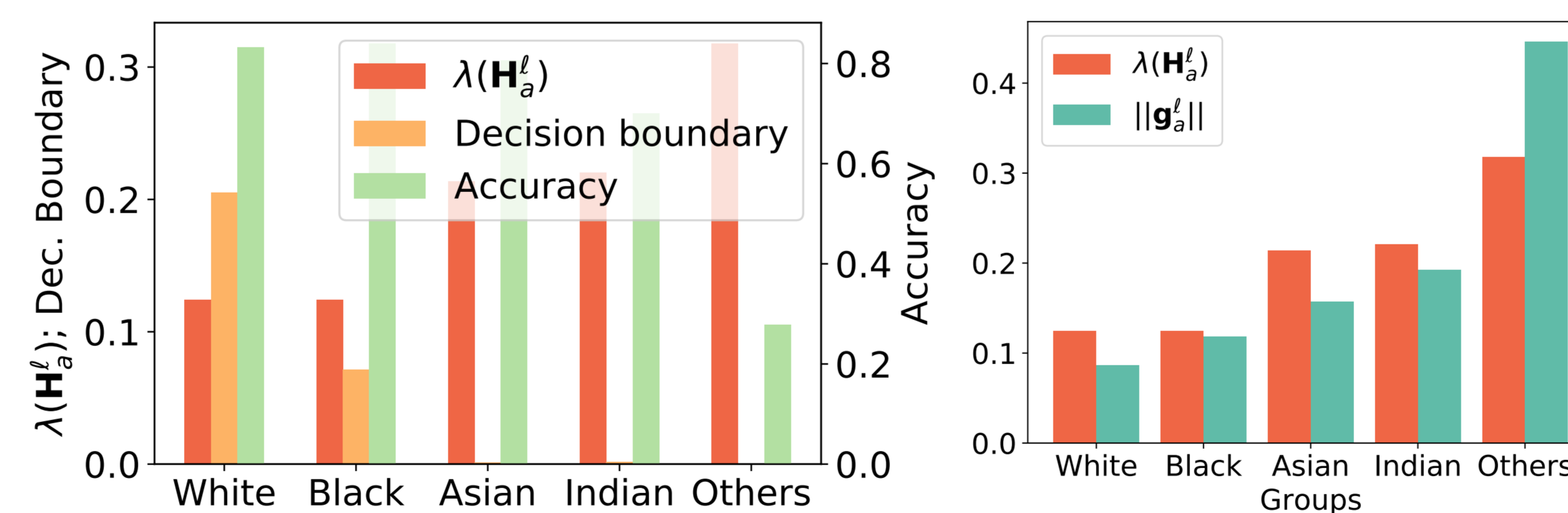


Figure 3. Left: Group Hessian vs decision boundary and accuracy. Right: Group Hessian vs gradient norm.

Mitigation solutions

- The disparity in gradient and Hessian norms are two direct factors for unfairness. The first straightforward mitigation solution is to equalize the gradient and Hessian norms:

$$\underset{\theta}{\operatorname{minimize}} J(\theta; D) \quad \text{such that: } \|g_a^{\ell}\| = \|g_b^{\ell}\|, \quad \lambda(\mathbf{H}_a^{\ell}) = \lambda(\mathbf{H}_b^{\ell}) \quad \forall a \in \mathcal{A}, \quad (6)$$

- However this solution can be computationally expensive.
- Hence, we propose to equalize the group losses due to their correlation with gradient and Hessian norms.

$$\underset{\theta}{\operatorname{minimize}} J(\theta; D) \quad \text{such that: } J(\theta; D_a) = J(\theta; D) \quad \forall a \in \mathcal{A}, \quad (7)$$

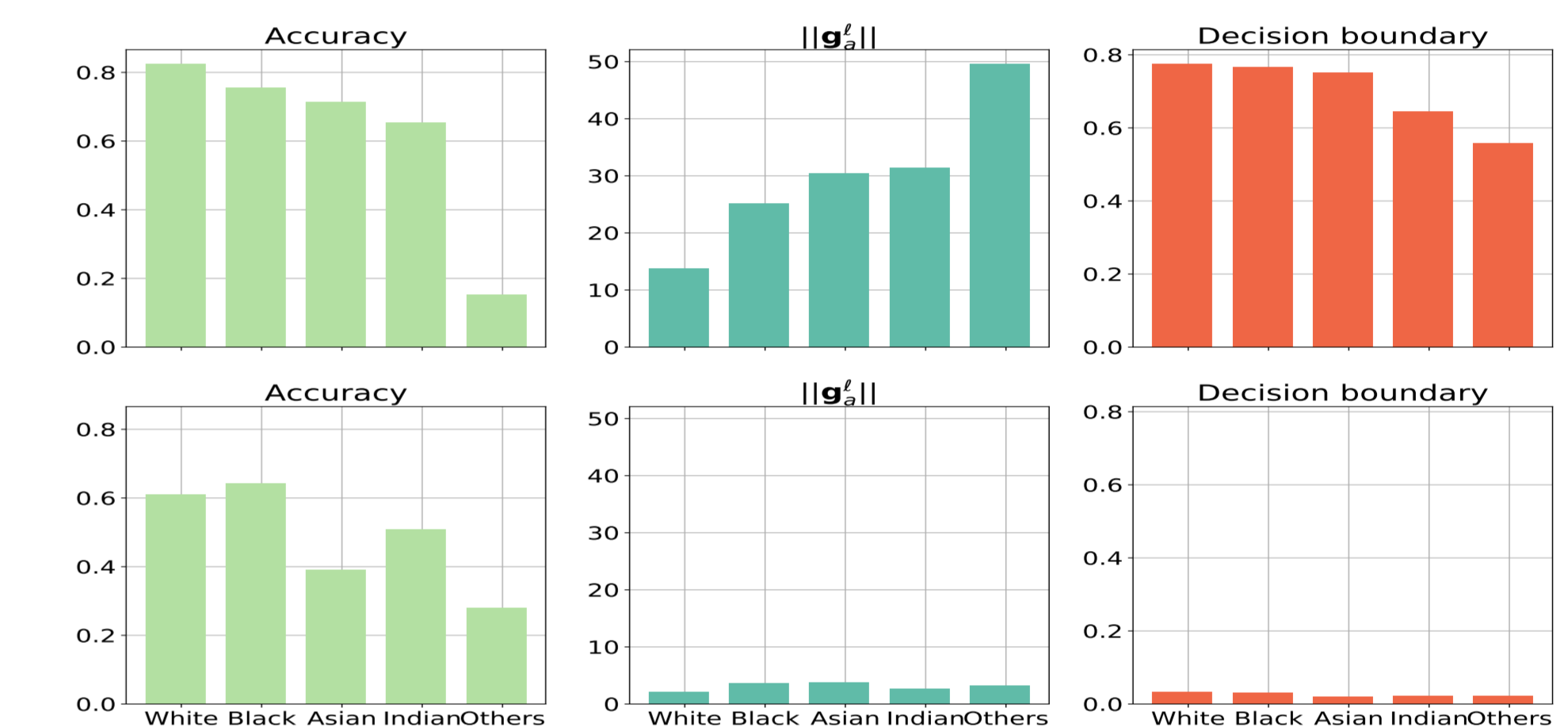


Figure 4. Effects of fairness constraints in balancing not only group accuracy (left) but also gradient norms (middle) and group average distance to the decision boundary (right)

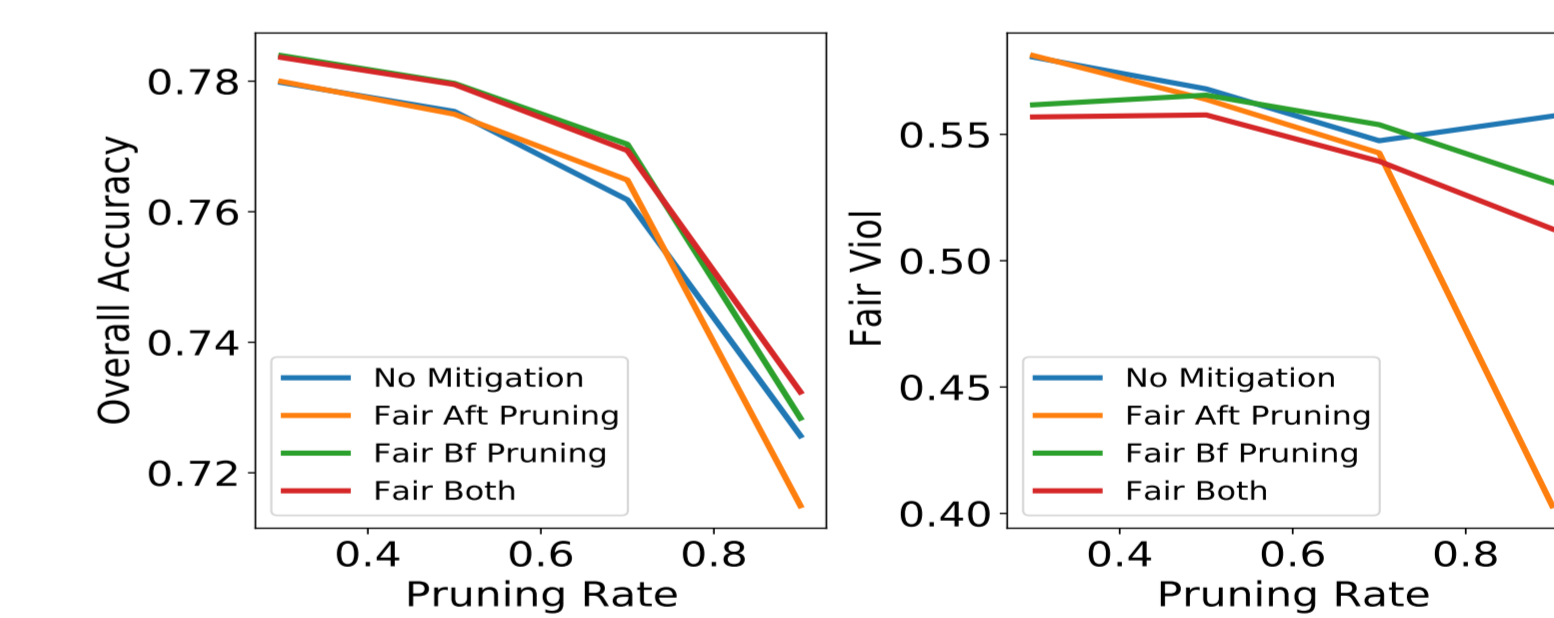


Figure 5. Accuracy and Fairness violations attained by all models on ResNet50, UTK-Face dataset with ethnicity (5 classes) as group attribute (and labels)

Discussion

- Pruning affects groups that are closer to the decision boundary. It also noted that these groups are more sensitive under adversarial attacks. Thus, understanding the interconnection among pruning, adversarial robustness and fairness is worth of investigation.
- Relax twice differentiable assumption over loss functions in fairness analysis can be a potential future work.

References

[1] Cuong Tran, Ferdinando Fioretto, Jung-Eun Kim, and Rakshit Naidu. Pruning has a disparate impact on model accuracy. *arXiv preprint arXiv:2205.13574*, 2022.