

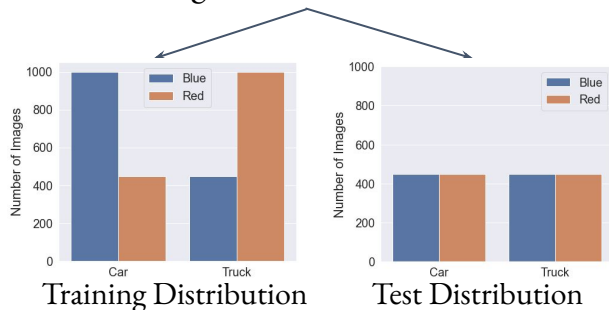
Bias Mimicking: A Simple Sampling Approach For Bias Mitigation

Maan Qraitem, Kate Saenko, Bryan A. Plummer

THU-AM-364

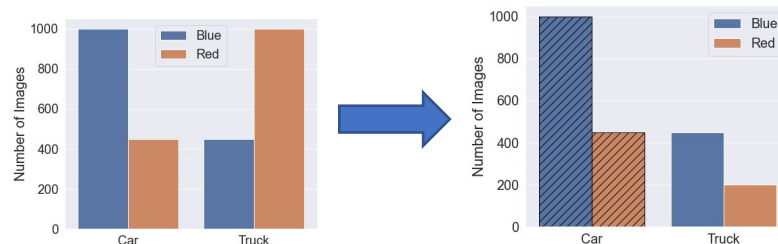
Task Definition

Y : Target Labels B : Bias Labels



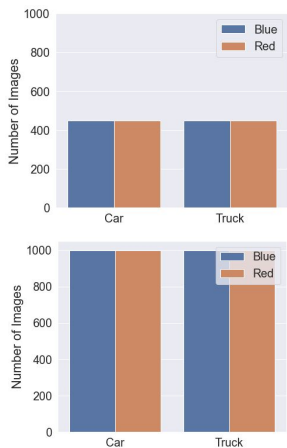
- Distribution shift between train and test.
- How to prevent the model from learning the correlation between Y (Car) and B (Blue).

Bias Mimicking

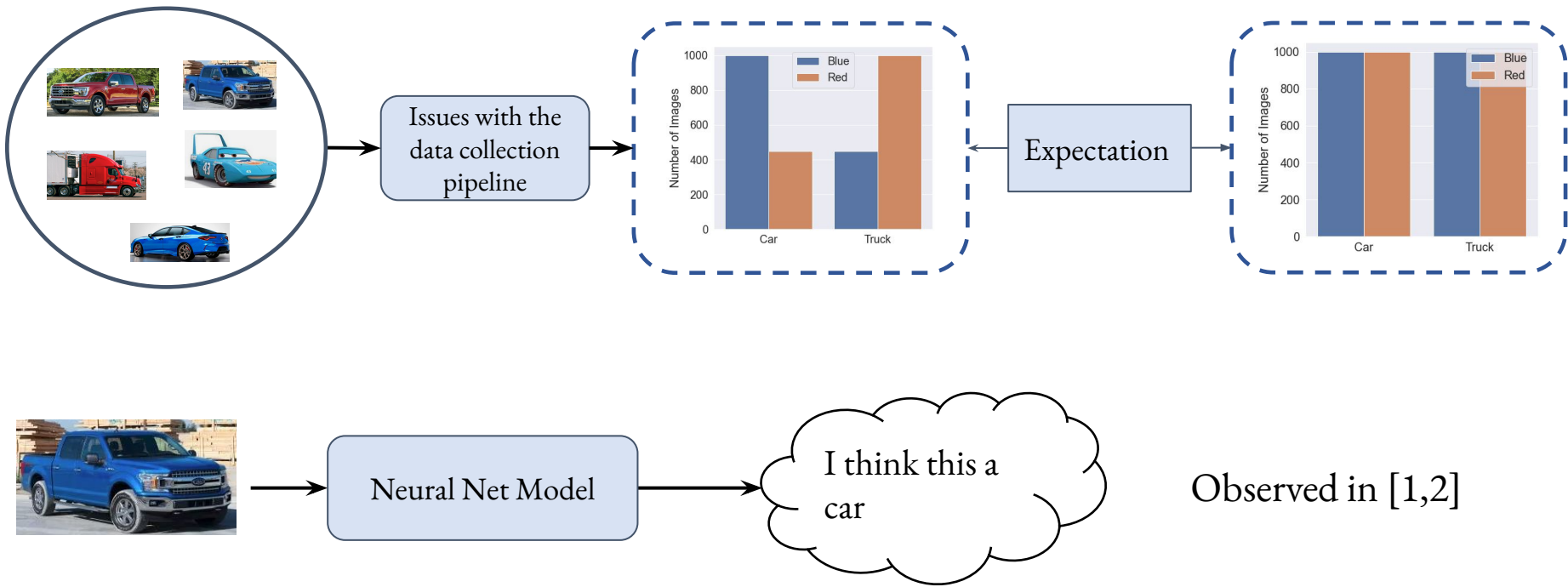


- Guarantees statistical independence.
- We introduce a novel training procedure that:
 - Addresses prior sampling methods shortcomings.
 - Bridges the performance gap between sampling and non sampling methods.
 - Maintains sampling methods low cost and simplicity.

Simple Sampling Methods

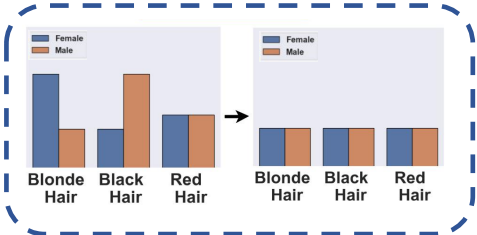


- $P_D(Y|B) = P_D(Y)$
- Simple to implement and cost effective.
- Suffer from shortcomings like overfitting over repeated samples.

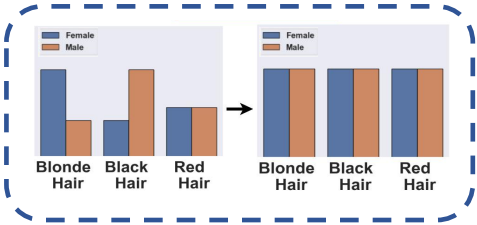


Sampling Methods Overview

- Simple Methods; few lines of code
- Ensure $P_D(Y|B) = P_D(Y)$
- Introduce no additional hyperparameters.
- They are missing from recent work benchmarks



Undersampling

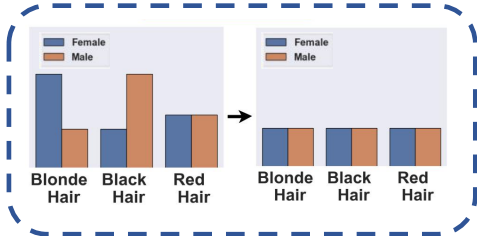


Oversampling

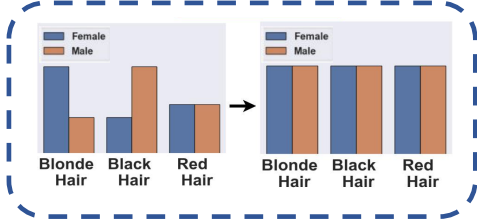
$$E_{x,c,s} \left[\frac{1}{p_D(x \in D_{c,s})} l(x, c; w) \right]$$

Upweighting

Sampling methods shortcomings



Limited distribution exposure per epoch which may compromise predictive performance



Model overfits over repeated samples [1]

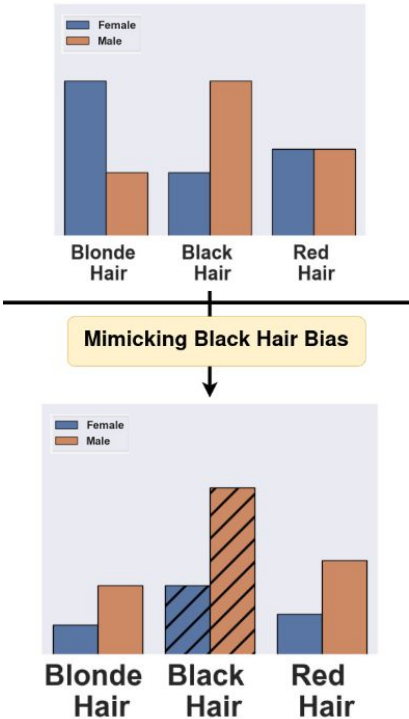
$$E_{x,c,s}[\frac{1}{p_D(x \in D_{c,s})}l(x, c; w)]$$



Instability issues with gradient descent [3]

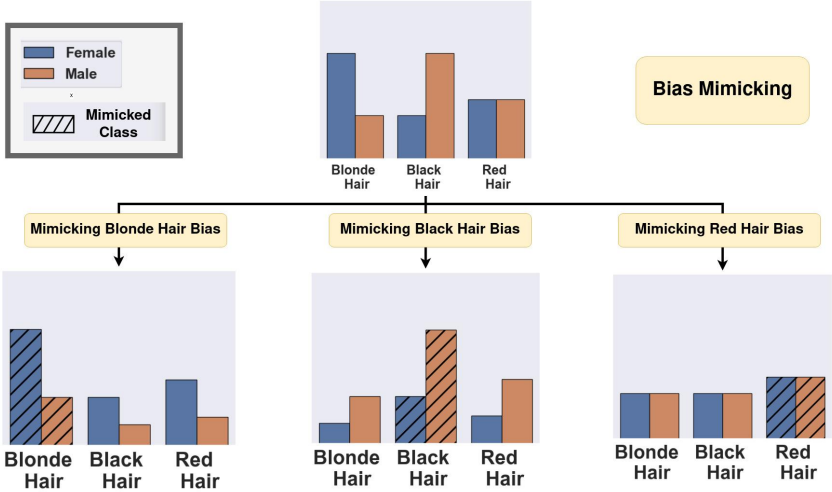
Overview of the Bias Mimicking

- We want to guarantee $P_D(Y|B) = P_D(Y)$
- Bias Mimicking: Given class “c”:
 - Ensure that $P_D(B|Y = c)$ is “mimicked” in each other class
- Turns out: this guarantees statistical independence!



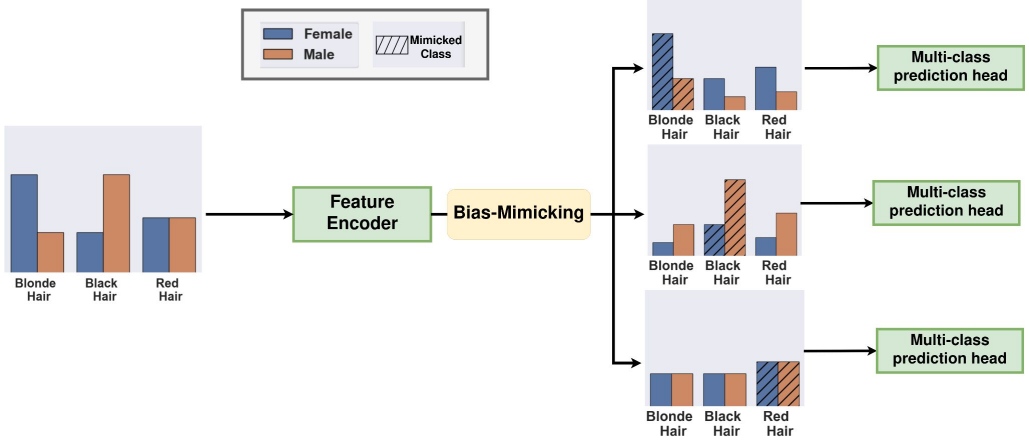
Training with BM

- We produce a distribution for each class.
- Each distribution preserves class c samples.
- Use all the distributions for training.
- Each distribution does not repeat samples.



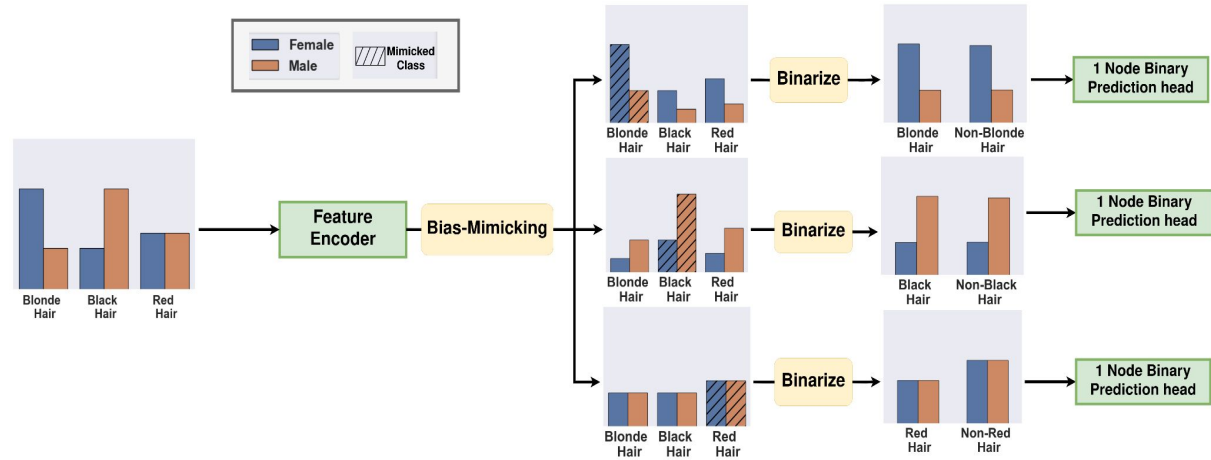
Training with BM

- Dedicate a multi class head for each distribution.
- **Issue:** Too many additional parameters.



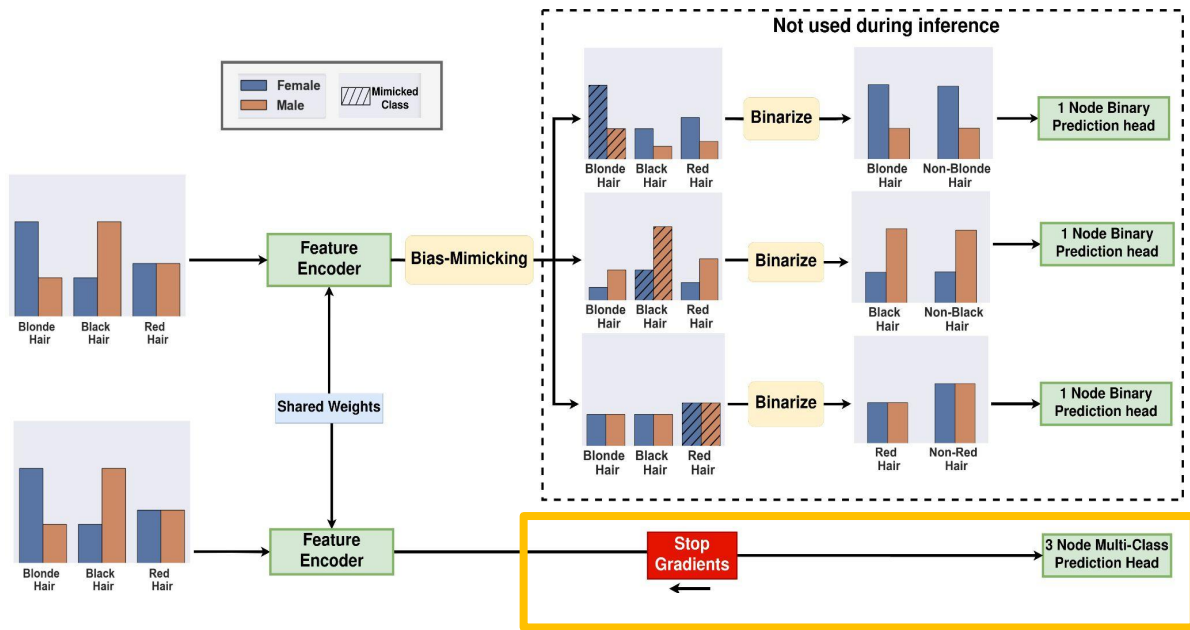
Training with BM

- Dedicate a binary head for each distribution.
- The binary prediction heads combined are equivalent to one multi class prediction head.

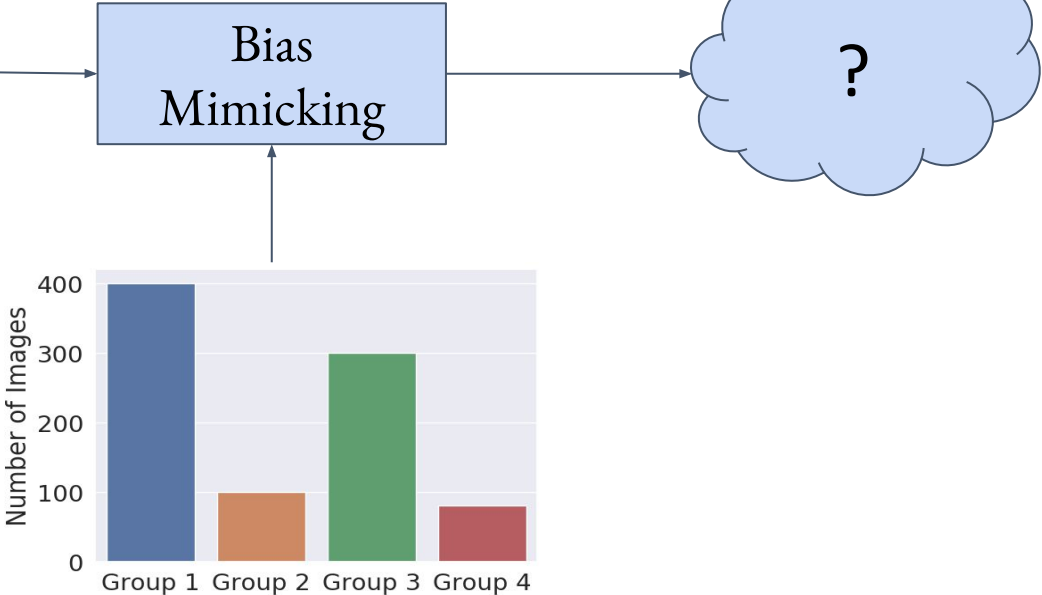
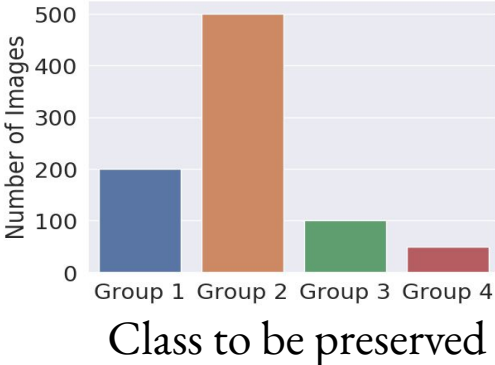


Inference with BM

- Issue:** the scores may not be calibrated with respect to each other.
- Solution:** train a multi-class prediction head on top of the debiased feature encoder.



How to Bias Mimick?



How to Bias Mimick?

- We constrain the solution space such that the solution retains the most number of samples.
- We obtain the set of solutions using a linear program.

c : preserved class
 c' : mimicked class
 s : bias
 l : count

Set of biases

$$\begin{aligned}
 &\max \sum_s l_s^{c'} \\
 \text{s.t.} \quad &l_s^{c'} \leq |D_{c',s}|
 \end{aligned}$$

$$s \in S$$

$$\frac{l_s^{c'}}{\sum_s l_s^{c'}} = P_D(B = s | Y = c) \quad s \in S$$

How to measure performance?

- Compute accuracy per group.
- **Unbiased Accuracy (UA):** Take the mean of accuracies over subgroups. [2]
- **Bias Conflict (BC):** Take the mean accuracies of the under-represented subgroups only. [2]

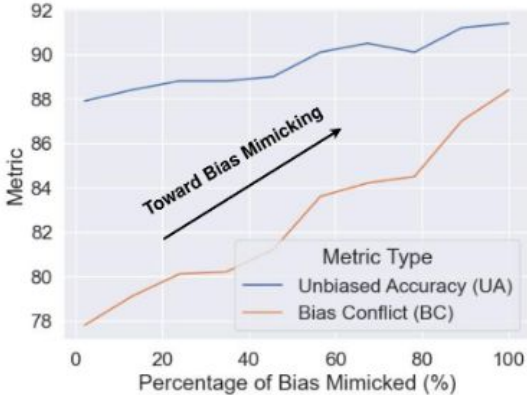
- Sampling methods show strong performance on some baselines.
- They lag behind on other benchmarks.
- Bias Mimicking shows consistent good performance unlike other sampling methods.

		Non-Sampling Methods					Sampling Methods			
		Vanilla	Adv	G-DRO	DI	BC+BB	OS	UW	US	BM
Utk-Face Age	UA	72.8	70.2	74.2	75.5	<u>78.9</u>	76.6	78.8	78.2	<u>79.7</u>
	BC	47.1	44.1	<u>75.9</u>	58.8	71.4	58.1	77.2	69.8	<u>79.1</u>
Utk-Face Race	UA	88.4	86.1	90.8	90.7	<u>91.4</u>	<u>91.3</u>	89.7	90.8	90.8
	BC	80.8	77.1	90.2	<u>90.9</u>	<u>90.6</u>	90.0	89.2	89.3	<u>90.7</u>
CelebA Blonde	UA	82.4	82.4	90.4	<u>90.9</u>	90.4	88.1	<u>91.6</u>	91.1	90.8
	BC	66.3	66.3	<u>89.4</u>	86.1	86.5	80.1	<u>88.3</u>	<u>88.5</u>	87.1
CIFAR-S	UA	88.7	81.8	89.1	<u>92.1</u>	90.9	87.8	86.5	88.2	<u>91.6</u>
	BC	82.8	72.0	88.0	<u>91.9</u>	89.5	82.5	80.0	83.7	<u>91.1</u>
Average	UA	83.0	80.1	86.1	87.3	87.9	85.9	86.6	87.0	88.2
	BC	69.2	64.8	85.8	81.9	84.5	77.6	83.6	82.8	87.0

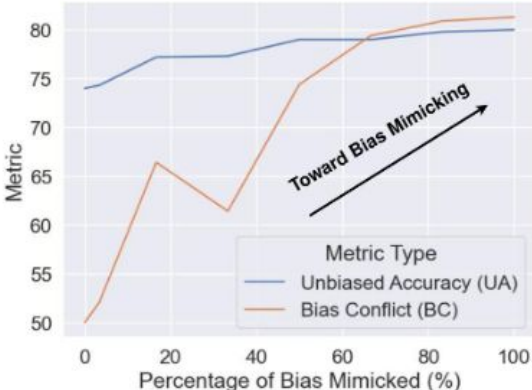
How sensitive is the model to the mimicking condition?

We vary the amount of bias mimicked between a percentage where:

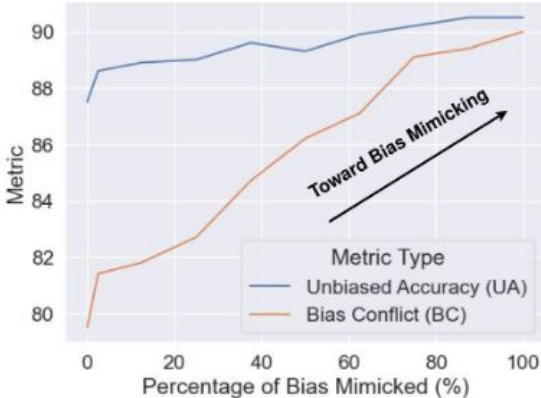
- 0%: distribution remains the same.
- 100%: Complete bias mimicking.



a) CelebA/Blonde



b) Utk-Face/Age



a) Utk-Face/Race

- **We showed** that simple sampling methods can be competitive on some benchmarks when compared to non sampling state-of-the-art approaches.
- **We introduced** a novel resampling method: Bias Mimicking that bridges the performance gap between sampling and nonsampling methods.
- **We conducted** an extensive empirical analysis of Bias Mimicking that details the method's sensitivity to the Mimicking condition. Refer to the paper for more details.

[1] Wang, Z., Qinami, K., Karakozis, I. C., Genova, K., Nair, P., Hata, K., & Russakovsky, O. (2020). Towards fairness in visual recognition: Effective strategies for bias mitigation. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 8919-8928).

[2] Sagawa, S., Koh, P. W., Hashimoto, T. B., & Liang, P. (2019). Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. arXiv preprint arXiv:1911.08731.

[3] An, J., Ying, L., & Zhu, Y. (2020). Why resampling outperforms reweighting for correcting sampling bias with stochastic gradients. arXiv preprint arXiv:2009.13447.

Link to code on github

