

Removing information -

> A lot of interest is recently raised to the themes of fairness and bias in deep learning, against unethical and discriminatory AI.

Many techniques attempts to tackle these issues, but the metric evaluated is typically accuracy, on balanced datasets: are they really removing information?

We propose a method to remove specific information at the **bottleneck** of deep neural networks.

This method relies on the **estimation** of the **information** we desire to maintain **private**, with the employment of an auxiliary classifier.

Then, we **minimize** a differentiable proxy of the **mutual information**.



Information Removal at the Bottleneck in Deep Neural Networks

Enzo Tartaglione LTCI, Télécom Paris, Institut Polytechnique de Paris

The input is forward-propagated through the whole network, including the two heads. Losses are computed, along with the differentiable proxy, in the form

$$\mathcal{I}(v,\hat{v}) = \sum_{i} \sum_{j} p_{\sigma(v)\hat{v}}(i,j) \log\left(\frac{p_{\sigma(v)\hat{v}}(i,j)}{p_{\sigma(v)}(i)p_{\hat{v}}(j)}\right)$$

At back-propagation, the task-related loss is back-propagated through the task-related head and the backbone.

The information removal loss is backpropagated through the IR head only. The mutual information term is backpropagated to the backbone, passing through the IR head. In the latter, no gradient will be accumulated.





Target	Method	Prediction accuracy of G (trained task) [%](↑)	Gender prediction accuracy of \mathcal{H} (information to remove) [%](\downarrow)
Blond hair	Baseline	95.34±0.07	84.32±2.76
	RUBi	$95.29 {\pm} 0.14$	88.55 ± 1.22
	Rebias	95.59±0.11	88.50 ± 3.78
	LearnedMixin	$90.01 {\pm} 2.66$	$74.09{\pm}2.66$
	IRENE ($\gamma = 0.1$)	$95.37 {\pm} 0.10$	55.47 ± 8.18
	IRENE ($\gamma = 0.5$)	$95.28 {\pm} 0.09$	$53.64{\pm}10.69$
	IRENE $(\gamma = 1)$	$95.24 {\pm} 0.29$	53.58±10.71
Heavy makeup	Baseline	90.58±0.14	92.89±0.36
	RUBi	$90.40 {\pm} 0.08$	95.17 ± 1.11
	Rebias	90.28 ± 0.34	$93.78 {\pm} 2.55$
	LearnedMixin	$84.88 {\pm} 3.28$	68.09 ± 10.55
	IRENE ($\gamma = 0.1$)	$90.32 {\pm} 0.97$	65.13 ± 11.08
	IRENE ($\gamma = 0.5$)	$85.66 {\pm} 2.80$	56.45 ± 9.46
	IRENE $(\gamma = 1)$	83.31±3.41	51.98±9.56
Eyeglasses	Baseline	99.67±0.02	69.51±4.33
	RUBi	$99.64{\pm}0.01$	59.21 ± 5.22
	Rebias	$99.65 {\pm} 0.01$	76.61 ± 8.21
	LearnedMixin	$93.54{\pm}2.94$	61.35 ± 3.67
	IRENE ($\gamma = 0.1$)	$99.68 {\pm} 0.01$	$64.08 {\pm} 1.08$
	IRENE ($\gamma = 0.5$)	99.69±0.02	$61.57{\pm}6.95$
	IRENE $(\gamma = 1)$	$99.68 {\pm} 0.01$	54.66±12.58