

Gradient Based Activations for Accurate Bias Free Learning Vinod K. Kurmi¹, Rishabh Sharma², Yash Vardhan Sharma², Vinay P. Namboodiri³ KU Leuven, Belgium¹ IIT Roorkee, India² University of Bath, UK³

Introduction

- Adversarial methods are very popular to tackle for bias free learning
- We propose to use a biased discriminator to improve the accuracy of adversarial methods while debiasing the feature representation.
- We use gradients of the discriminator and propose a masking scheme for the features, we term as Gradient Based Activation (GBA).
- Use of GBA provides us a masking rule to drop certain features in order to do unbiased training. The effectiveness of this approach is demonstrated and validated





CelebA Dataset

Model	Heavy Makeup				Hair Color			
	Aligned	Conflicting	Mean	Bias Gap	Aligned	Conflicting	Mean	Bias Gap
Vanilla	92.44 ± 0.74	31.46 ± 2.45	61.95 ± 1.28	60.98 ± 2.56	90.58 ± 0.34	57.35 ± 0.21	73.97 ± 0.20	33.23 ± 0.40
LfF (Nam et al. 2020)	83.85 ± 1.68	45.54 ± 4.28	64.69 ± 2.29	38.31 ± 4.60	88.85 ± 1.27	80.24 ± 2.16	84.55 ± 1.25	8.61 ± 2.5
Domain Independent (Wang et al. 2020)	79.88 ± 1.71	43.24 ± 4.33	61.56 ± 2.31	36.64 ± 5.64	90.97±3.71	79.25 ± 3.33	85.11 ± 2.67	7.44 ± 3.21
Group DRO (Sagawa et al. 2020)	79.28 ± 1.20	46.24±3.61	62.76±2.22	33.04 ± 3.22	89.68 ± 0.65	81.41 ± 1.47	85.55 ± 0.88	8.27 ± 2.01
Adversarial	92.07 ± 2.88	33.79 ± 3.81	62.93 ± 2.38	58.28 ± 4.77	93.4 ± 0.91	62.75 ± 3.47	78.08 ± 1.79	30.65 ± 5.59
Adversarial with GBA	81.49 ± 1.91	49.79± 3.15	$\textbf{65.64} \pm \textbf{1.55}$	$\textbf{31.70} \pm \textbf{3.10}$	90.67 ± 1.01	83.28 ± 1.83	$\textbf{86.98} \pm \textbf{1.04}$	$\textbf{7.39} \pm \textbf{2.09}$

CIFAR-I Dataset

Model Name	Model	Aligned	Conflicting	Mean	Bias (GAP)(↓)	
Baseline	N-way Softmax	87.94 ± 0.36	69.39 ± 0.42	78.67 ± 0.27	18.55 ± 0.55	
LfF(Nam et al. 2020)	N-way Softmax	87.01 ± 0.63	56.87 ± 0.72	71.93 ± 0.47	30.14 ± 0.95	
Domain Independent(Wang et al. 2020)	N-way Classifier per Domain	88.39 ± 0.15	78.14 ± 0.13	83.26 ± 0.01	10.25 ± 0.20	
Adversarial	Adversarial Gradient Reversal		75.74 ± 0.29	80.63 ± 0.35	9.78 ± 0.71	
Adversarial with GBA	Proposed	88.81 ± 0.19	$\textbf{79.46} \pm \textbf{0.22}$	$\textbf{84.41} \pm \textbf{0.14}$	$\textbf{9.35} \pm \textbf{0.29}$	

Proposed Approach

Using the Biased Discriminator

- When the discriminator is correct: Mask the features used by the discriminator When the discriminator is incorrect
- Emphasize the features used by the discriminator

$$\begin{split} g_i^d &= \frac{\partial \hat{y}}{\partial f_i} \\ &= \begin{cases} 0, & \text{if } (g_i^d.ind^d) > 0. \\ 1 & \text{if } (g_i^d.ind^d) \leq 0 \end{cases} \end{split}$$

$$f_i^{cis} = f_i * a_i^a$$

$$\hat{y}_i^c = C(f_i^{cls}, \theta_c);$$

Results

36TH MAI CONFERENCE ON ARTIFICIAL INTELLIGENCE A VIRTUAL CONFERENCE BRUARY 22 - MARCH 1, 2022

+ve **Predicted Class** θ_d i - - - - $= y_d$

Feature Map

R

Revgrad

$$\mathcal{L}_c = \frac{1}{N} \sum_{x_i \in \mathcal{D}} \mathcal{L}(\hat{y}_i^c, y_i^c) \quad \mathcal{L}_d = \frac{1}{N} \sum_{x_i \in \mathcal{D}} \mathcal{L}(\hat{y}_i^d, y_i^d)$$

Conclusion

Predicted Domain

- Identify a key drawback of using adversarial methods to debias models
- Developed a novel masking scheme to address the same
- We gave extensive experimental evaluation to validate our findings
- Provide insights to the learning of a classifier



