

On Mitigating Bias Using Causal Generative Models

Ajay Gopi

Linwei Wang

Abstract

While deep learning models have demonstrated superior performance across computer vision tasks, the higher performance comes at the cost of poor generality towards out-of-distribution and / or underrepresented minority samples in a dataset. Causality-inspired research in deep learning presents a lucrative avenue to address these problems by modeling cause-effect relationships in data and providing tools for interventional strategies to answer counterfactual queries. This research aims to evaluate the performance of existing causality-inspired models in datasets that are known to exhibit spurious correlations, with the goal of understanding the ability and limitations of these models to address these biases and to improve interpretability and robust modeling across diverse samples within a given dataset.

1 Introduction

Correlation does not imply causation. In deep learning terminology, this means that the correlations made by a deep learning model do not impute an outcome y . Traditional deep learning models are known [2] to make such correlations often addressed as spurious with respect to a task to get better train time accuracies, often resulting in poor generalization in real world scenarios.

Consider the scenario, where we need to classify between waterbirds and landbirds, due to the nature of the data collected, most waterbirds are often found in a lake or water background, while the landbirds are often found in forest or land backgrounds. This results in unintended bias during the data collection process. Given a particular type of bird, we have certain task relevant features for eg, shape of the beak and task irrelevant features eg, background of the subject. Causality-Inspired deep generative models aims to model these task relevant features as causal features and promise generating counterfactuals that are bias conflicting to debias models for downstream computer vision tasks such as classification or generations. This research aims to evaluate the generative capabilities of causal deep-generative models, an emerging but understudied area, particularly in biased settings to debias a model.

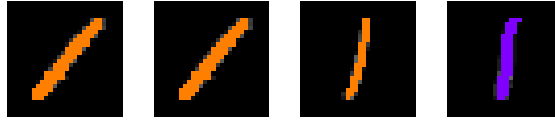
2 Literature Review

2.1 Benchmark Datasets for Debiasing

Several datasets, both synthetic and real world datasets are currently used in debiasing research, and are covered in subsequent sections.

Coloured MNIST (CMNIST), a dataset with task relevant feature i.e shape of the digit, vs task irrelevant features i.e color of the digits is used due to it's simplicity in synthetically generating any combination of bias aligned and conflicting samples as shown in Fig 1a. Similarly, WaterBirds is another synthetically generated dataset by combining bird segmentations from CUB dataset with Places365 as backgrounds to generate bias aligned and conflicting groups as shown in Fig 1b.

CelebA is a real-world dataset with celebrity faces, in which task relevant feature is hair color vs task irrelevant feature gender or vice versa as shown in Fig 1c.



(a) Dataset: CMNIST
Task-relevant: Digit Shape
Task-irrelevant: Digit Color



(b) Dataset: WaterBirds
Task-relevant: Bird Features
Task-irrelevant: BackGround



(c) Dataset: CelebA
Task-relevant: Gender
Task-irrelevant: Hair Color

Figure 1: Examples of biased datasets with task-relevant and task-irrelevant features.

2.2 Existing Debiasing Strategies

Although there have been numerous studies in debiasing groupDRO by Sagawa et al [6] promises to debias a classifier however they have varying definition on how long to train the subgroups to debias

a model between datasets using the same technique, we mainly investigate approaches that rely on generative-augmentation of bias conflicting minority groups. Jaeju et al [1] and Goel et al [3] in their work demonstrated that minority groups augmentation resulted in significant improvement in test time accuracies for bias-conflicting groups. BiasSwap by Kim et al [4] propose an unsupervised bias feature augmentation for model robustness. However, none of the previous strategies use xAI lens to model deep generative models, with Counterfactual Generative Networks by Sauer et al [7] an exception, where they report accuracies on correlated mnist dataset but not on real-world datasets.

2.3 Pearl’s Ladder Of Causation

Pearl’s ladder of causation involves 3 rungs. The first level association, dealing with identifying associations in passively observed data. The second level intervention, commonly referred to by the do-operator, involves actively inferring what happens if I do action A? The third level, Counterfactual generation, involves generating plausible outcomes by modelling what would have happened if I did action B instead of action A?

2.4 Causality In Deep Learning

In the last few years, there have been several works [5] integrating causal inference principles in deep learning to have model interpretability and steerability. Consequently, there have been two notable lines of work in this emerging field, Causal Representation Learning (CRL) and Causal CounterFactual Generations (CCG) where former involves counterfactuals in the representation space, the later involves counterfactual generation of images. The current work studies the impact of these causal generative models, specifically 3 main model families, Causal VAE’s, Causal GAN’s and Causal Diffusion Models for counterfactual augmentation of bias conflicting minority groups.

References

- [1] Jaeju An, Taejune Kim, Donggeun Ko, Sangyup Lee, and Simon S Woo. A²: Adaptive augmentation for effectively mitigating dataset bias. In *Proceedings of the Asian Conference on Computer Vision*, pages 4077–4092, 2022.
- [2] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization, 2020.
- [3] Karan Goel, Albert Gu, Yixuan Li, and Christopher Ré. Model patching: Closing the subgroup performance gap with data augmentation, 2020.
- [4] Eungyeup Kim, Jihyeon Lee, and Jaegul Choo. Biaswap: Removing dataset bias with bias-tailored swapping augmentation, 2021.

- [5] Aneesh Komanduri, Xintao Wu, Yongkai Wu, and Feng Chen. From identifiable causal representations to controllable counterfactual generation: A survey on causal generative modeling, 2024.
- [6] Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization, 2020.
- [7] Axel Sauer and Andreas Geiger. Counterfactual generative networks, 2021.