# Evaluating and Mitigating Bias in Image Classifiers: A Causal Perspective Using Counterfactuals

**Saloni Dash**[1], **Vineeth N Balasubramanian**[2], **Amit Sharma**[1]

t-sadash@microsoft.com, vineethnb@iith.ac.in

amshar@microsoft.com

1 - Microsoft Research, Bangalore, Karnataka, India

2 - Indian Institute of Technology Hyderabad, Telangana, India

## Abstract

Counterfactual examples for an input—perturbations that change specific features but not others—have been shown to be useful for evaluating bias of machine learning models, e.g., against specific demographic groups. However, generating counterfactual examples for images is non-trivial due to the underlying causal structure on the various features of an image. To be meaningful, generated perturbations need to satisfy constraints implied by the causal model. We present a method for generating counterfactuals by incorporating a structural causal model (SCM) in an improved variant of Adversarially Learned Inference (ALI), that generates counterfactuals in accordance with the causal relationships between attributes of an image. Based on the generated counterfactuals, we show how to explain a pre-trained machine learning classifier, evaluate its bias, and mitigate the bias using a counterfactual regularizer. On the Morpho-MNIST dataset, our method generates counterfactuals comparable in quality to prior work on SCM-based counterfactuals (DeepSCM), while on the more complex CelebA dataset our method outperforms DeepSCM in generating high-quality valid counterfactuals. Moreover, generated counterfactuals are indistinguishable from reconstructed images in a human evaluation experiment and we subsequently use them to evaluate the fairness of a standard classifier trained on CelebA data. We show that the classifier is biased w.r.t. skin and hair color, and how counterfactual regularization can remove those biases.

## 1 Introduction

A growing number of studies have uncovered biases in image classifiers, particularly against marginalized demographic groups [Buolamwini and Gebru, 2018, Hendricks et al., 2018, Zhao et al., 2017, Bolukbasi et al., 2016]. To avoid such biases, it is important to explain a classifier's predictions and evaluate its fairness. Given any pre-trained machine learning (ML) classifier, counterfactual reasoning is an important way to explain the classifier's decisions and to evaluate its fairness. Counterfactual reasoning involves simulating an alternative input with some specific changes compared to the original input. For example, to evaluate fairness of a classifier with respect to a sensitive attribute like *skin color*, we can ask how the classifier's output will change if a face that was originally *dark-skinned* is made *light-skinned* while keeping everything else constant. Since the only change to the input is the sensitive attribute, counterfactual reasoning is considered more robust than comparing available faces with different skin colors (association) or comparing simulated inputs with *light skin* or *dark skin* (intervention) since these comparisons may include changes in addition to *skin color*.

However, generating counterfactual (CF) examples for images is non-trivial. For instance, consider the simple task of changing a person's hair color in an image of their face. While Generative Adversarial Networks (GANs) can generate new realistic faces with any hair color [Kocaoglu et al., 2018], they are unable to generate the precise changes needed for a CF, i.e. changing hair color without changing hair style or other features of the face. Other explanation techniques based on perturbations such as occluding pixels [Zeiler and Fergus, 2014] also do not support counterfactual reasoning based on high-level concepts.

There have been recent efforts on using GANs to generate counterfactuals using an added inference step (encoder). Given a pre-trained GAN model, Denton et al. [2019] trained an encoder to match the input of a generated image. However, the latents thus encoded do not directly correspond to the given attributes of an image, and it is difficult to change a specific known attribute to generate a counterfactual. To change an attribute, Joo and Kärkkäinen [Joo and Kärkkäinen, 2020] used the FaderNetwork architecture that inputs attributes of an image separately to the generator. However, their method does not account for causal relationships between attributes. Besides, while both these works use generated images to evaluate biases in a classifier, they do not provide any method to mitigate the found biases. We present a method for generating counterfactuals that is based on their formal causal definition, and present a novel counterfactual-based regularizer to mitigate biases in a given classifier.

Formally, a valid counterfactual example for an image is defined with respect to a Structural Causal Model (SCM) over its attributes. An SCM encodes the domain knowledge about how attributes affect each other in the form of a graph with attributes as the nodes and accompanying functional equations connecting the nodes. Generating a counterfactual, therefore, requires modeling both the underlying SCM for the attributes as well as the generative process that uses the attributes to model the resulting image.

In this paper, we present *ImageCFGen*, a method that combines knowledge from a causal graph and uses an inference mechanism in a GAN-like framework to generate counterfactual images. We first evaluate *ImageCFGen* on the Morpho-MNIST dataset to show that it generates counterfactual images comparable to a prior SCM-based CF generation method (DeepSCM) [Pawlowski et al., 2020]. Moreover, we show that our method is capable of generating high-quality valid counterfactuals for complex datasets like CelebA in comparison to DeepSCM.

Specifically, on the CelebA dataset, *ImageCFGen* can generate CFs for facial attributes like *Black Hair*, *Pale Skin* etc. Based on these counterfactual images, we show that an image classifier for predicting attractiveness on the CelebA dataset exhibits bias with respect to *Pale Skin*, but not with respect to attributes like *Wearing Necklace*. We hence propose and demonstrate a bias mitigation algorithm which uses the counterfactuals to remove the classifier's bias with respect to sensitive attributes like *Pale Skin*. In summary, our contributions include:

- *ImageCFGen*, a method that uses inference in a GAN-like framework to generate counterfactuals based on attributes learned from a known causal graph.
- Theoretical justification that under certain assumptions, CFs generated by *ImageCFGen* satisfy the definition of a counterfactual as in Pearl [Pearl, 2009].
- Detailed experiments on Morpho-MNIST and CelebA datasets that demonstrate the validity of CFs generated by *ImageCFGen* in comparison to DeepSCM [Pawlowski et al., 2020].
- Evaluating fairness of an image classifier and explaining its decisions using counterfactual reasoning.
- A regularization technique using CFs to mitigate bias w.r.t. sensitive attributes in any image classifier.

# 2 Related Work

Our work bridges the gap between generating counterfactuals and evaluating fairness of image classifiers.

## 2.1 Counterfactual Generation

Given the objective to generate a data point $X$ (e.g., an image) based on a set of attributes $A$, Pearl's ladder of causation [Pearl, 2019] describes three types of distributions that can be used: *association* $P(X|A = a)$, *intervention* $P(X|\operatorname{do}(A = a))$ and *counterfactual* $P(X_{A=a}|A = a', X = x')$. In the associational distribution $P(X|A = a)$, $X$ is conditioned on a specific attribute value $a$ in the observed data. For e.g., conditional GAN-based methods [Mirza and Osindero, 2014] implement association between attributes and the image. In the interventional distribution $P(X|\operatorname{do}(A = a))$, $A$ is changed to the value $a$ irrespective of its observed associations with other attributes. Methods like CausalGAN [Kocaoglu et al., 2018] learn an interventional distribution of images after changing specific attributes, and then generate new images that are outside the observed distribution (e.g., women with moustaches on the CelebA faces dataset [Kocaoglu et al., 2018]).

The counterfactual distribution $P(X_{A=a}|A = a', X = x')$ denotes the distribution of images related to a specific image $x'$ and attribute $a'$, if the attribute of the same image is changed to a different value $a$. In general, generating counterfactuals is a harder problem than generating interventional data since we need to ensure that everything except the changed attribute remains the same between an original image and its counterfactual. Pawlowski et al. [2020] recently proposed a method for generating image counterfactuals using a conditional Variational Autoencoder (VAE) architecture. While VAEs allow control over latent variables, GANs have been more successful over recent years in generating high-quality images. Thus, we posit that a GAN-based method is more ideally suited to the task of CF generation in complex image datasets, especially when the generated images need to be realistic to be used for downstream applications such as fairness evaluation. We test this hypothesis in Section 5.

Independent of our goal, there is a second interpretation of a "counterfactual" example w.r.t. a ML classifier [Wachter et al., 2017], referring to the smallest change in an input that changes the prediction of the ML classifier. [Liu et al., 2019] use a standard GAN with a distance-based loss to generate CF images close to the original image. However, the generated CFs do not consider the underlying causal structure – terming such images as CFs can be arguable from a causal perspective. Besides, these CFs are not perceptually interpretable – ideally, a counterfactual image should be able to change only the desired attribute while keeping the other attributes constant, which we focus on in this work.

## 2.2 Fairness of Image Classifiers

Due to growing concerns on bias against specific demographics in image classifiers [Buolamwini and Gebru, 2018], methods have been proposed to inspect what features are considered important by a classifier [Sundararajan et al., 2017], constrain classifiers to give importance to the correct or unbiased features [Ross et al., 2017], or enhance fairness by generating realistic images from underrepresented groups [McDuff et al., 2019]. Explanations, to study the fairness of a trained model, can also be constructed by perturbing parts of an image that change the classifier's output, e.g., by occluding an image region [Zeiler and Fergus, 2014] or by changing the smallest parts of an image that change the classifier's output [Fong and Vedaldi, 2017, Luss et al., 2019, Goyal et al., 2019]. Such perturbation-based methods, however, are not designed to capture high-level concepts and do

not enable study of fairness of classifiers w.r.t. human-understandable concepts (e.g. gender or race).

Existing work closest to our efforts include two adversarially-trained generative models to generate counterfactuals for a given image. [Denton et al., 2019] changed attributes for a given image by linear interpolations of latent variables using a standard progressive GAN [Karras et al., 2018]. Similarly, [Joo and Kärkkäinen, 2020] used a Fader network architecture [Lample et al., 2017] to change attributes. However, both these works ignore the causal structure associated with attributes of an image. In analyzing bias against an attribute, it is important to model the downstream changes *caused* by changing that attribute [Kusner et al., 2017]. For instance, for a chest MRI classifier, age of a person may affect the relative size of their organs [Pawlowski et al., 2020]; it will not be realistic to analyze the effect of age on the classifier's prediction without learning the causal relationship from age to organ size. Hence, in this work, we present a different architecture that can model causal relationships between attributes and provide valid counterfactuals w.r.t. an assumed structural causal model. In addition, using these counterfactuals, we present a simple regularization technique that can be used to decrease bias in any given classifier.

# 3   SCM-Based Counterfactuals

Let $\mathbf{X} = \mathbf{x} \in \mathcal{X}$ denote the image we want to generate the counterfactual for, and let $\mathbf{A} = \mathbf{a} = \{a_i\}_{i=1}^n \in \mathcal{A}$ be its corresponding attributes. In the case of human faces, attributes can be binary variables ($\in \{0, 1\}$) like *Smiling*, *Brown Hair*; or in the case of MNIST digits, continuous attributes ($\in \mathbb{R}$) like *thickness*, *intensity*, etc. A continuous attribute is scaled so that it lies in the range of $[0, 1]$. We have a training set $\mathcal{D}$ containing $(\mathbf{x}, \mathbf{a})$ (image, attribute) tuples. Given an image $(\mathbf{x}, \mathbf{a})$, the goal is to generate a counterfactual image with the attributes changed to $\mathbf{a}_c$.

## 3.1   SCM over Attributes

We assume an SCM for the true data-generating process that defines the relationship among the attributes $\mathbf{a}$, and from the attribute to the image $\mathbf{x}$. For instance, with two binary attributes (*Young*, *Gray Hair*) for an image $\mathbf{x}$ of a face, the true causal graph can be assumed to be *Young* $\rightarrow$ *Gray Hair* $\rightarrow$ $\mathbf{x}$. We separate out the graph into two parts: relationships amongst the attributes ($\mathcal{M}_a$), and relationships from the attributes to the image ($\mathcal{M}_x$). We call $\mathcal{M}_a$ as the *Attribute-SCM* and model $\mathcal{M}_x$ as a generative model, given the attributes.

The Attribute-SCM ($\mathcal{M}_a$) consists of a causal graph structure and associated structural assignments. We assume that the causal graph structure is known. Given the graph structure, Attribute-SCM learns the structural assignments between attributes. E.g., given *Young* $\rightarrow$ *Gray Hair*, it learns the function $g$ such that *Gray Hair* $= g(Young, \epsilon)$, where $\epsilon$ denotes independent random noise. We use the well-known Maximum Likelihood Estimation procedure in Bayesian Networks [Ding, 2010] to estimate these functions for the attributes, but other methods such as Normalizing Flows [Rezende and Mohamed, 2015, Pawlowski et al., 2020] can also be used.

Note that counterfactual estimation requires knowledge of both the *true* causal graph and the *true* structural equations; two SCMs may entail exactly the same observational and interventional distributions, but can differ in their counterfactual values [Karimi et al., 2020]. In most applications, however, it is impractical to know the true structural equations for each edge of a causal graph. Therefore, here we make a simplifying empirical assumption: while we assume a known causal graph for attributes, we estimate the structural equations from observed data. That is, we assume that the data is generated from a subset of all possible SCMs (e.g., linear SCMs with Gaussian noise) such that the structural equations can be uniquely estimated from data. Details on Attribute-SCM

are in Appendix A.

## 3.2   Image Generation from Attribute-SCM

For modeling the second part of the SCM from attributes to the image ($\mathcal{M}_x$), we build a generative model that contains an encoder-generator architecture. Given an Attribute-SCM (either provided by the user or learned partially, as in Sec 3.1), the proposed method *ImageCFGen* has three steps to generate a counterfactual for an image $(\mathbf{x}, \mathbf{a})$ such that the attributes are changed to $\mathbf{a}'$ (architecture in Figure 1).

- An encoder $E : (\mathcal{X}, \mathcal{A}) \rightarrow \mathbf{Z}$ infers the latent vector $\mathbf{z}$ from $\mathbf{x}$ and $\mathbf{a}$, i.e. $\mathbf{z} = E(\mathbf{x}, \mathbf{a})$ where $\mathbf{Z} = \mathbf{z} \in \mathbb{R}^m$.
- The Attribute-SCM intervenes on the desired subset of attributes that are changed from $\mathbf{a}$ to $\mathbf{a}'$, resulting in output $\mathbf{a}_c$. Specifically, let $\mathbf{a}_k \subseteq \mathbf{a}$ be the subset of attributes that changed between the inputs $\mathbf{a}$ and $\mathbf{a}'$. For every $a_i \in \mathbf{a}_k$, set its value to $a_i'$, then change the value of its descendants in the SCM graph by plugging in the updated values in the structural equations (see Lines 5-6 in Algorithm 1, Appendix A).
- Generator $G : (\mathbf{Z}, \mathcal{A}) \rightarrow \mathcal{X}$ takes as input $(\mathbf{z}, \mathbf{a}_c)$ and generates a counterfactual $\mathbf{x}_c$, where $\mathbf{z} \in \mathbf{Z} \subseteq \mathbb{R}^m$.
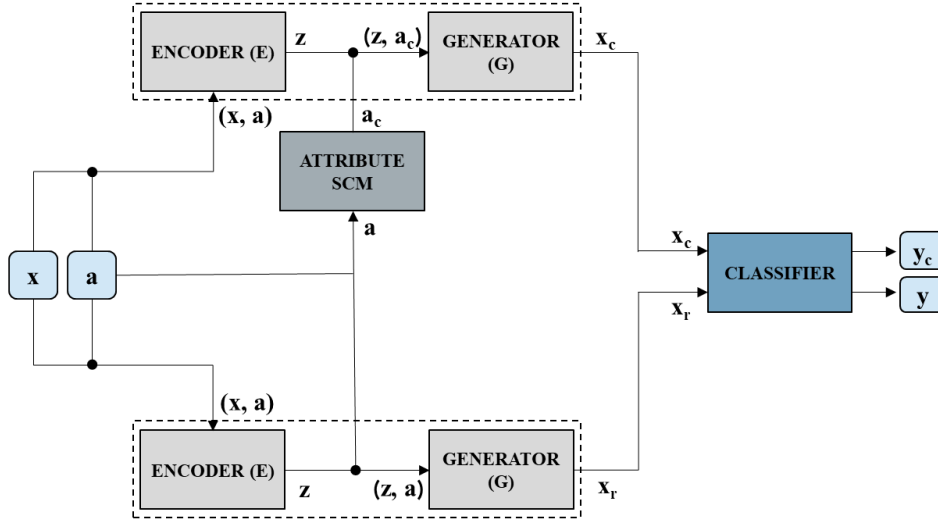


Figure 1: **Counterfactual Generation using *ImageCFGen*.** The top half of the figure shows the CF generation procedure, and the bottom half of the figure shows the reconstruction procedure. Finally, the reconstructed image $x_r$ and the counterfactual image $x_c$ are used for a downstream task like fairness evaluation of a classifier.

The above method for CF generation can be written as:

$$\mathbf{x}_c = G(E(\mathbf{x}, \mathbf{a}), \mathbf{a}_c) \tag{1}$$

The complete algorithm is shown in Algorithm 1 in Appendix A. For our experiments, we use a novel improved variant of Adversarially Learned Inference [Dumoulin et al., 2016] to train the encoder and generator. However, *ImageCFGen* can be extended to any encoder-generator (decoder) architecture.

## 3.3 Correspondence to Counterfactual Theory

The above architecture maps directly to the three steps for generating a valid counterfactual, for $(\mathbf{x}, \mathbf{a})$ as in [Pearl, 2009]:

- **Abduction**: Infer latent $\mathbf{z}$ given the input $(\mathbf{x}, \mathbf{a})$ using the encoder.
- **Action**: Let $\mathbf{a}_k \subseteq \mathbf{a}$ be the set of $k$ attributes that one wants to intervene on. Set attribute $a_i \to a_i' \; \forall a_i \in \mathbf{a}_k$, where $\mathbf{a}_k' = \{a_i'\}_{i=1}^k$.
- **Prediction**: Modify all the descendants of $\mathbf{a}_k$ according to the SCM equations learned by Attribute-SCM. This outputs $\mathbf{a}_c$, the intervened attributes. Use $\mathbf{z}$ from the encoder and $\mathbf{a}_c$ from the Attribute-SCM and input it to the generator to obtain the counterfactual $x_c$.

The proof that Equation 1 corresponds to generating a valid counterfactual is in Appendix B.

## 3.4 Implementing the Encoder and Generator

Many studies have reported generating high-quality, realistic images using GANs [Wang et al., 2018, Karras et al., 2018, 2019]. However, vanilla GANs lack an inference mechanism where the input $\mathbf{x}$ can be mapped to its latent space representation $\mathbf{z}$. We hence use Adversarially Learned Inference (ALI) [Dumoulin et al., 2016], which integrates the inference mechanism of variational methods like VAEs in a GAN-like framework, thereby leveraging the generative capacity of GANs as well as providing an inference mechanism. We generate the images using a conditional variant of ALI where the model is conditioned on the attributes $\mathbf{a}$ while generating an image.

An ALI-based method, however, has two limitations: (1) generation capabilities are limited when compared to state-of-the-art [Karras et al., 2019]; and (2) reconstructions are not always faithful reproductions of the original image [Dumoulin et al., 2016]. Image reconstructions are important in the counterfactual generation process because they indicate how good the inferred latent space variable $z$ is, which is used in the abduction step of generating counterfactuals. We address both these issues by using a style-based generator for better generation and a cyclic cost minimization algorithm for improved reconstructions. We refer to our ALI model as Cyclic Style ALI (CSALI) and describe each of these components below.

**Adversarially Learned Inference.** ALI uses an encoder $E$ and a generator $G$ in an adversarial framework, where the encoder learns to approximate the latent space distribution from the input $\mathbf{x}$ and attributes $\mathbf{a}$, and the generator learns to generate realistic images from the latent space distribution and attributes $\mathbf{a}$. The discriminator $D$ is optimized to differentiate between pairs of tuples containing {the real image, the corresponding approximated latent space variable, attributes} from joint samples of {the generated image, the input latent variable, attributes}. The Encoder and Generator are optimized to fool the discriminator. Unlike [Dumoulin et al., 2016] which uses an embedding network, we directly pass the attributes to the Generator, Encoder and Discriminator since we found that it helped in conditioning the model on the attributes better. The conditional ALI hence optimizes:

$$
\min_{G,E} \max_D V(D, G, E) = \mathbb{E}_{q(\mathbf{x})p(\mathbf{a})}[\log(D(\mathbf{x}, E(\mathbf{x}, \mathbf{a}), \mathbf{a}))]
$$
$$
+ \mathbb{E}_{p(\mathbf{z})p(\mathbf{a})}[\log(1 - D(G(\mathbf{z}, \mathbf{a}), \mathbf{z}, \mathbf{a}))]
\tag{2}
$$

where $G$ is the generator, $E$ is the encoder and $D$ is the discriminator. $q(\mathbf{x})$ is the distribution for the images, $p(z) \sim \mathcal{N}(0, 1)$ and $p(a)$ is the distribution of the attributes. Image reconstructions are defined as:

$$
\mathbf{x}_r = G(E(\mathbf{x}, \mathbf{a}), \mathbf{a})
\tag{3}
$$

**Style-Based Generator.** Style-based generator architectures (StyleGANs) implicitly learn separations of high-level attributes of images like hair color, pose, etc [Karras et al., 2019, 2020], and

generate images that are indistinguishable from real images [Zhou et al., 2019]. To improve generation, we replace the ALI generator with the style-based generator architecture in [Karras et al., 2019]. Details of the architecture are provided in the Appendix D.

**Cyclic Cost Minimization.** To improve image reconstructions (which in turn indicate the goodness of the learned latents $\mathbf{z}$), we employ the cyclic cost minimization algorithm in [Dogan and Keles, 2020] after training the style-based ALI model. The generator is fixed, and the encoder is fine-tuned to minimize a reconstruction loss computed using: (i) error in the image space $\mathcal{L}_{\mathbf{x}} = \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \|\mathbf{x} - G(E(\mathbf{x}, \mathbf{a}))\|_2$; and (ii) error in the latent space $\mathcal{L}_{\mathbf{z}} = \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \|\mathbf{z} - E(\mathbf{x}, \mathbf{a})\|$, where $G(E(\mathbf{x}, \mathbf{a}))$ is the reconstructed image $\mathbf{x}_r$ according to Eqn 3 and $E(\mathbf{x}, \mathbf{a})$ is the encoder's output $\mathbf{z}_r$, which is expected to capture the image's latent space distribution. We fine-tune the encoder using the above reconstruction loss *post-hoc* after obtaining a good generator in order to explicitly improve image reconstructions.

# 4    Applications of Generated CFs

We now show how the counterfactuals generated using *ImageCFGen* can be used to evaluate fairness of, as well as explain a given image classifier. We will also present a method to mitigate any fairness biases in the classifier. Suppose we are given a pre-trained image classifier $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y}$, such that $\hat{f}(\mathbf{x}) = \hat{y}$, where $\mathbf{x} \in \mathcal{X}$ refers to the images and $\hat{y} \in \mathcal{Y}$ refers to the classifier's discrete outcome. Let $\mathbf{a} \in \mathcal{A}$ be the corresponding image attributes, and let $\mathbf{a}_S \in \mathcal{A}_S \subseteq \mathcal{A}$ be the set of sensitive attributes we want to evaluate the classifier on.

## 4.1    Evaluating Fairness of a Classifier

We can use the generated CFs to estimate biases in a given classifier that predicts some discrete outcome $\hat{y} = y$ (like *Attractive*). In an ideal scenario, the latent variable $\mathbf{z}$ for the real image and its reconstructed image would match exactly. However, experiments using ALI demonstrate that the reconstructed images are not perfect reproductions of the real image [Dumoulin et al., 2016, Donahue and Simonyan, 2019, Dandi et al., 2020]. Therefore, for objective comparison, we compare classification labels for reconstructed images ($\mathbf{x}_r$ from Eqn 3) and counterfactual images ($\mathbf{x}_c$ from Eqn 1), since the reconstructed images share the exact same latent $\mathbf{z}$ as the CF image (and hence the CF will be valid). We hence refer to the reconstructed images (which share the latent $\mathbf{z}$ with the CF) as *base* images for the rest of the paper.

We characterize a classifier as *biased* w.r.t. an attribute if: (a) it changes its classification label for the CF image (obtained by changing that attribute); and (b) if it changes the label to one class from another class more often than vice versa (for CFs across test images obtained by changing the considered attribute). (To illustrate the second condition, if setting hair color as blonde makes test images consistently be classified as attractive more often than otherwise, this indicates bias.) We capture these intuitions as a formula for the degree of bias in a binary classifier w.r.t. a considered attribute:

$$\text{bias} = p(y_r \neq y_c)(p(y_r = 0, y_c = 1 | y_r \neq y_c) \\ -p(y_r = 1, y_c = 0 | y_r \neq y_c)) \tag{4}$$

where $y_r$ is the classification label for the reconstructed image, and $y_c$ is the classification label for the CF image. Using Bayes Theorem, Eqn 4 reduces to:

$$\text{bias} = p(y_r = 0, y_c = 1) - p(y_r = 1, y_c = 0) \tag{5}$$

The bias defined above ranges from -1 to 1. It is 0 in the ideal scenario when the probability of CF label changing from 0 to 1 and vice-versa is the same (=0.5). The bias is 1 in the extreme case when the CF label always changes to 1, indicating that the classifier is biased *towards* the counterfactual change. Similarly if the CF label always changes to 0, the bias will be -1, indicating that the classifier is biased *against* the counterfactual change. In Appendix C, we show that a classifier with zero bias in the above metric is fair w.r.t. the formal definition of counterfactual fairness [Kusner et al., 2017].

## 4.2 Explaining a Classifier

We can also use the CFs from *ImageCFGen* to generate explanations for a classifier. For any input $\mathbf{x}$, a local *counterfactual importance score* for an attribute $\mathbf{a}_i$ states how the classifier's prediction changes upon changing the value of $\mathbf{a}_i$. Assuming $\mathbf{a}_i$ can take binary values $a' = 1$ and $a = 0$, the local *CF importance score* of $\mathbf{a}_i$ is given by:

$$
\mathbb{E}_Y[Y_{\mathbf{a}_i \leftarrow a'}|\mathbf{x}, \mathbf{a}] - \mathbb{E}_Y[Y_{\mathbf{a}_i \leftarrow a}|\mathbf{x}, \mathbf{a}]
$$
$$
= y_{\mathbf{a}_i \leftarrow a'}|\mathbf{x}, \mathbf{a} - y_{\mathbf{a}_i \leftarrow a}|\mathbf{x}, \mathbf{a}
\tag{6}
$$

where $Y$ is the random variable for the classifier's output, $y$ is a value for the classifier's output, and the above equality is for a deterministic classifier. For a given $(\mathbf{x}, \mathbf{a})$, the score for each attribute (feature) can be ranked to understand the relative importance of features. To find global feature importance, we average the above score over all inputs.

## 4.3 Bias Mitigation for a Classifier

Finally, in addition to evaluating a classifier for bias, CFs generated using *ImageCFGen* can be used to remove bias from a classifier w.r.t. a sensitive attribute. Here we propose a *counterfactual* regularizer to ensure that an image and its counterfactual over the sensitive attribute obtain the same prediction from the image classifier. For an image $\mathbf{x}$, let $logits(\mathbf{x})$ be the output of the classifier $\hat{f}$ before the sigmoid activation layer. To enforce fairness, we can finetune the classifier by adding a regularizer that the logits of the image and its counterfactual should be the same, i.e.

$$
\text{BCE}(y_{true}, \hat{f}(\mathbf{x})) + \lambda \text{MSE}(logits(\mathbf{x}_r), logits(\mathbf{x}_c))
\tag{7}
$$

where BCE is the binary cross-entropy loss, $y_{true}$ is the ground truth label for the real image $\mathbf{x}$, $\lambda$ is a regularizing hyperparameter, MSE is the mean-squared error loss, and $\mathbf{x}_r$ and $\mathbf{x}_c$ are defined in Eqns 3 and 1 respectively.

# 5 Experiments and Results

Considering the limited availability of datasets with known causal graphs, we study *ImageCFGen* on the Morpho-MNIST dataset (a simple dataset for validating our approach), and on the CelebA dataset (which provides an important context for studying bias and fairness in image classifiers). Specifically, we study the following:

- **Validity of *ImageCFGen* CFs.** We use the Morpho-MNIST dataset which adds causal structure on the MNIST images, to compare counterfactuals from *ImageCFGen* to the Deep-SCM method [Pawlowski et al., 2020]. We show that CFs from *ImageCFGen* are comparable to those from DeepSCM, thus validating our approach.

- **Quality of *ImageCFGen* CFs.** On the more complex CelebA dataset, we evaluate the quality of *ImageCFGen* CFs by quantifying the generation and reconstruction, using established benchmark metrics. We find that using the proposed CSALI architecture offers significant advantages over the standard ALI model. We also contrast the quality and validity of the generated CFs with those of DeepSCM, and find that *ImageCFGen* outperforms DeepSCM.
- **Fairness Evalution and Explanation of a ML Classifier Using *ImageCFGen* CFs.** We show using *ImageCFGen* that a standard pre-trained classifier on CelebA that predicts whether a face is attractive or not, has bias w.r.t. attribute *Pale Skin* across all three hair colors (*Black Hair*, *Blond Hair*, *Brown Hair*). We also explain the classifier's predictions using CFs.
- **Bias Mitigation of a ML Classifier Using *ImageCFGen* CFs.** Finally, we show how our proposed method can be used to decrease detected bias in the classifier for the attributes mentioned above.

**Baselines and Performance Metrics**. We compare to DeepSCM's [Pawlowski et al., 2020] results on Morpho-MNIST and CelebA. We present both quantitative and qualitative performance of our method for these datasets. While we follow the metrics of [Pawlowski et al., 2020] for Morpho-MNIST, for CelebA we report quantitative scores like Fréchet Inception Distance [Heusel et al., 2017] (FID) and Mean Squared Error (MSE) to compare generation and reconstruction quality with the base ALI method. For measuring quality of generated counterfactuals, we report human evaluation scores, in addition to qualitative results. For bias evaluation, we compare *ImageCFGen* to affine image transformations like horizontal flip and brightness that are commonly used data augmentation techniques.

**Datasets**. *Morpho-MNIST* [Castro et al., 2019] is a publicly available dataset based on MNIST [LeCun et al., 1998] with interpretable attributes like *thickness*, *intensity*, etc. It was extended by [Pawlowski et al., 2020] to introduce morphological transformations with a known causal graph. The attributes are *thickness (t)* and *intensity (i)*, where *thickness $\rightarrow$ intensity* ($\rightarrow$ indicating causal effect). We extend this dataset by introducing an independent morphological attribute—*slant (s)* from the original Morpho-MNIST dataset and digit *label (l)* as an attribute. The causal graph for the dataset is given in Fig 2a.
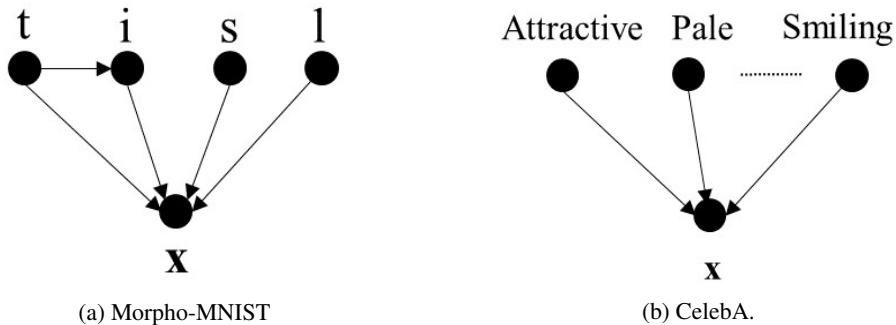


(a) Morpho-MNIST                    (b) CelebA.

Figure 2: **Causal Graphs for Morpho-MNIST and CelebA.** Attributes for Morpho-MNIST are *thickness* t, *intensity* i, *slant* s and *label* l; for CelebA are *Pale*, *Black Hair*, etc. In both graphs, attributes cause the image x.

*CelebA [Liu et al., 2015]* is a dataset of 200k celebrity images annotated with 40 attributes like *Black Hair*, *Wearing Hat*, *Smiling* etc. We train an image classifier on the dataset that predicts the attribute *Attractive* as done in [Sattigeri et al., 2019, Ehrlich et al., 2016]. While explaining the classifier's decisions, we generate CFs for all attributes excluding *Male*, *Young* and *Blurry*. For

fairness evaluations, we focus on generating CFs for the attributes *Black Hair*, *Blond Hair*, *Brown Hair*, *Pale* and *Bangs*. Similar to [Denton et al., 2019], we do not generate CFs for *Male* because of inconsistent social perceptions surrounding gender, thereby making it difficult to define a causal graph not influenced by biases. Therefore, all attributes we consider have a clear causal structure (Fig 2b shows the causal graph). Additionally, our method can also be utilized in the setting where the attributes are connected in a complex causal graph structure, unlike [Denton et al., 2019, Joo and Kärkkäinen, 2020]. We show this by conducting a similar fairness and explanation analysis for a *Attractive* classifier in Appendix O, where *Young* affects other visible attributes like *Gray hair*.

Details of implementation, architecture (including ALI) and training are provided in Appendix D.

## 5.1  Validity of *ImageCFGen* CFs on Morpho-MNIST

We generate CFs using *ImageCFGen* on images from the Morpho-MNIST dataset by intervening on all four attributes - *thickness*, *intensity*, *slant* and *label* and observe how the image changes with these attributes. Fig 3 demonstrates CFs for a single image with label 0. Along the first column vertically, the label is changed from 0 to $\{1, 4, 6, 9\}$ while the thickness, intensity and slant are kept constant. Then, as we proceed to the right in each row, the attributes of thickness, intensity and slant are changed sequentially but the label is kept constant. Visually, the generated counterfactuals change appropriately according to the intervened attributes. For instance, according to the causal graph in Fig 2a, changing the digit label should not change the digit thickness intensity and slant. That is precisely observed in the first column of Fig 3. Whereas, changing the thickness should also change the intensity of the image which is observed in the third and fourth columns of Fig 3. Results of latent space interpolations and digit reconstructions are provided in Appendix E and F. We find that the encoder learns meaningful latent space representations, and the reconstructions are faithful reproductions of Morpho-MNIST digits.
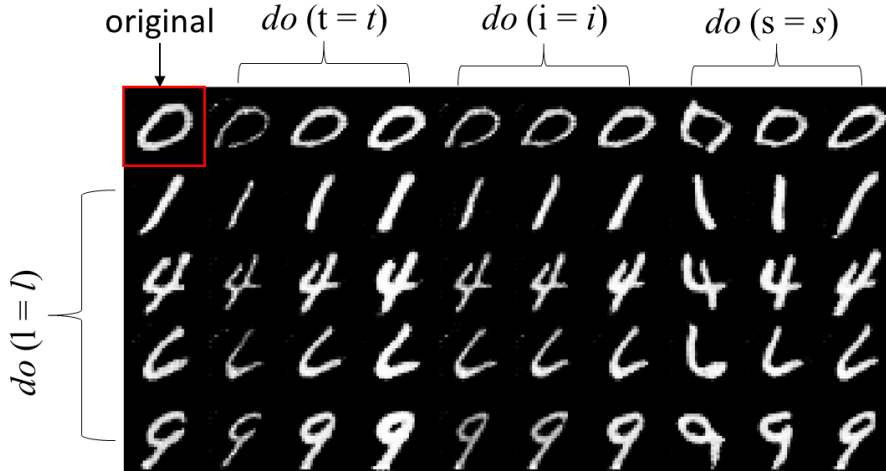


Figure 3: **Morpho-MNIST Counterfactuals.** Top-left cell shows a real image sampled from the test set. Vertically, rows correspond to interventions on the label, $do(l = 1, 4, 6, 9)$. Moving horizontally, columns correspond to interventions on thickness: $do$ ($t$=1, 3, 5), intensity: $do$ ($i$ = 68, 120, 224), and slant: $do$ ($s$ = -0.7, 0, 1) respectively.
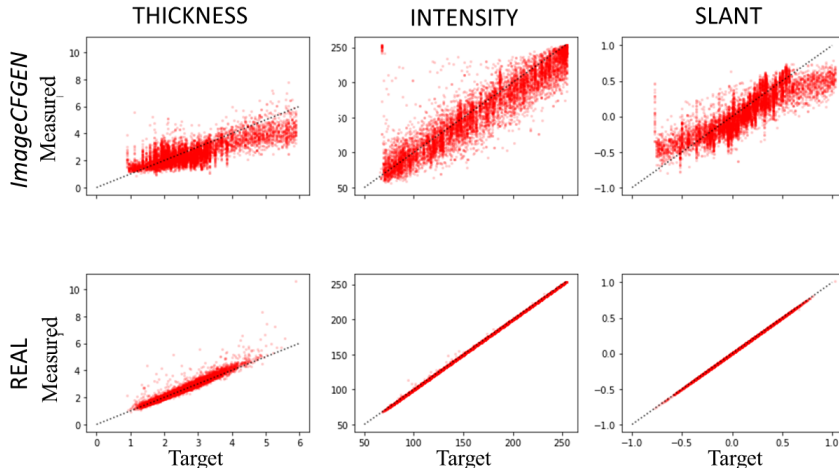
Figure 4: **Morpho-MNIST CFs.** For each attribute, Target ($x$ axis) is the desired value for the CF and Measured ($y$ axis) is the attribute value obtained from the generated counterfactual using a formula provided in the Morpho-MNIST dataset [Pawlowski et al., 2020]. We compare between ground-truth CFs and CSALI-generated CFs. In an ideal scenario (real samples), points should lie along the $y = x$ line.

To quantify these observations, we randomly sample hundred values for *thickness*, *intensity* and *slant* attributes and generate corresponding CFs for each test image. Fig 4 plots the target ground-truth attribute vs the measured attribute in the CF image (measured using a formula from Morpho-MNIST [Pawlowski et al., 2020]), and shows that all attributes are on an average clustered along the reference line of y = x with some variability. We quantify this variability in Table 1 using median absolute error, by comparing the CFs generated using *ImageCFGen* vs those using the DeepSCM method [Pawlowski et al., 2020] (we choose Median Absolute Error to avoid skewed measurements due to outliers). *ImageCFGen* and DeepSCM are comparable on attributes of *thickness* and *intensity*, showing the validity of our CFs. We could not compare *slant* since [Pawlowski et al., 2020] do not use *slant* in their studies.

|  | $do$ (thickness) | $do$ (intensity) |
|---|---|---|
| *ImageCFGen* | $0.408 \pm 0.004$ | $\mathbf{10.740 \pm 0.161}$ |
| DeepSCM | $\mathbf{0.253 \pm 0.002}$ | $12.519 \pm 0.105$ |

Table 1: **Median Abs. Error: *ImageCFGen* and DeepSCM.** Lower is better.

## 5.2 Quality of *ImageCFGen* CFs on CelebA

We now evaluate *ImageCFGen* on CelebA by comparing the quality of the generated CF images against the vanilla ALI model and the DeepSCM model [Pawlowski et al., 2020].
**Generation Quality.** We show real images and their corresponding reconstructions using ALI and CSALI in rows (I), (II) and (IV) of Fig 13 in Appendix K. While the reconstructions of both ALI and CSALI are not perfect, those of CSALI are significantly better than ALI. Moreover, they capture

11

majority of the high-level features of the real images like *Hair Color*, *Smiling*, *Pale*, etc. We quantify this in Table 2 using Fréchet Inception Distance [Heusel et al., 2017] (FID) score for generated images, Mean Squared Error (MSE) between the real images $\mathbf{x}$ and reconstructed images $\mathbf{x}_r$ and Mean Absolute Error (MAE) between the real latent variable $\mathbf{z}$ and the encoder's approximation of the latent variable $\mathbf{z}_r$. We randomly sample 10k generated images to calculate FID scores, and randomly sample 10k real images for which we generate reconstructions and report the MSE and MAE metrics. We report these metrics on the baseline model of ALI as well. We observe a significant improvement in generation quality after using a style-based generator and significant improvements in reconstruction with the cyclic cost minimization algorithm (refer to Appendix K for ablation study). We report MSE on images and MAE on latent space variables since the cyclic cost minimization algorithm [Dogan and Keles, 2020] uses these metrics to improve reconstructions.

|  | FID | MSE $(x, x_r)$ | MAE $(z, z_r)$ |
|---|---|---|---|
| ALI | 67.133 | 0.177 | 1.938 |
| CSALI | **21.269** | **0.103** | **0.940** |

Table 2: **FID, MSE and MAE scores for ALI and Cyclic Style ALI (CSALI).** Lower is better.

**Counterfactual Quality.** We contrast the quality and validity of the CFs from *ImageCFGen* with the CFs from DeepSCM [Pawlowski et al., 2020] (refer to Appendix H for our implementation of DeepSCM on CelebA). To generate the counterfactual images, we intervene on the attributes of *Black Hair*, *Blond Hair*, *Brown Hair*, *Pale* and *Bangs*. We observe in Figure 5 that the CFs from *ImageCFGen* are qualitatively better than those from DeepSCM. *ImageCFGen* CFs successfully change the hair color and skin color, in addition to adding bangs to the face (refer to Figure 12 in Appendix J for more *ImageCFGen* CFs and Appendix I for more comparisons with DeepSCM). In contrast, the CFs from DeepSCM only partially change the hair color and the skin color in columns (a) through (f) and fail to add bangs in column (g).

We also perform a human evaluation experiment of the generated counterfactuals in Appendix L, which showed that the distribution of counterfactuals is identical to the distribution of their corresponding base images.

## 5.3 Bias Evaluation & Explaining a Classifier

We train a classifier on the CelebA dataset to predict attractiveness of an image w.r.t. an attribute (architecture and training details in Appendix M). We then use the generated CFs to identify biases in the classifier. We sample 10k points from the test set and generate seven CFs for each of them as shown in Fig 5 for different attributes. We only consider those images for which the corresponding attribute was absent in order to avoid redundancies. For instance, we filter out CF (c) of the second sample from Fig 5 since blond hair was already present in the base image. We then provide the generated CFs along with the base (reconstructed) image to the attractive classifier and compare their labels. As a baseline comparison, we also pass images with affine transformations like horizontal flip and increasing brightness to the classifier. We treat the classifier's outputs as the probability of being labeled as attractive. Fig. 6 shows these probabilities for the base image, affine transformations, and the CFs. If the classifier is fair w.r.t. these attributes, all points should be clustered along the y=x line.

For the CF plots, for *do* (black hair =1, pale = 1) and *do* (brown = 1, pale = 1) almost all points are above the reference line, suggesting that the probability of being classified as attractive increases after applying these interventions. Not only does the probability increase, for *do* (black hair = 1,
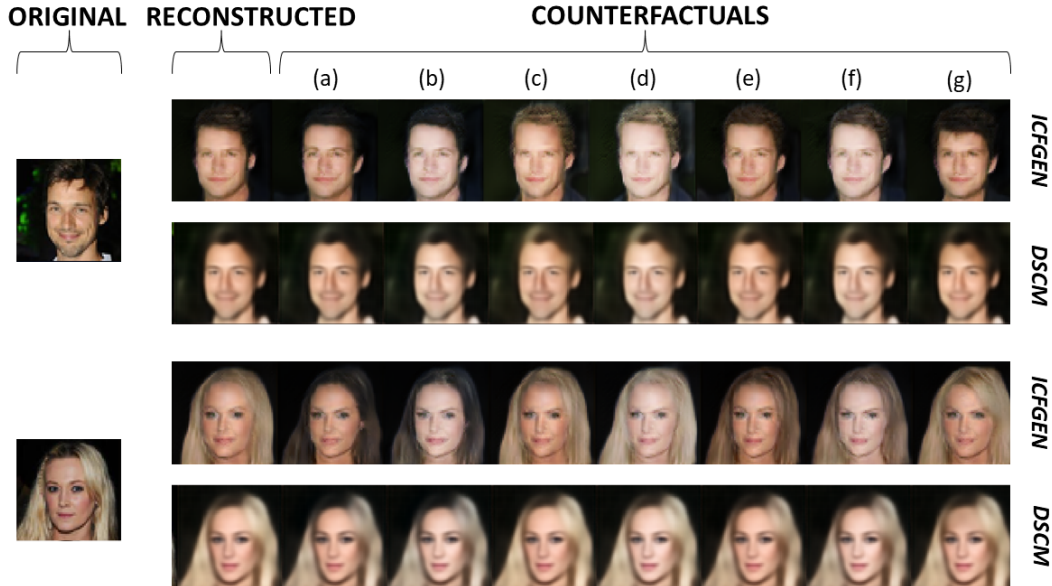
Figure 5: *ImageCFGen* **and DeepSCM Counterfactuals.** (a) denotes $do$ (black hair = 1) and (b) denotes $do$ (black hair = 1, pale =1). Similarly (c) denotes $do$ (blond hair = 1); (d) denotes $do$ (blond hair = 1, pale = 1); (e) denotes $do$ (brown hair = 1); (hf denotes $do$ (brown hair = 1, pale = 1); and (g) denotes $do$ (bangs = 1).

pale = 1), 18% of the CF labels are flipped from the base ones, and *94%* of these labels are changed from not attractive to attractive. In case of $do$ (brown hair = 1, pale = 1), 19% of the CF labels are flipped and *94%* of the labels are flipped from not attractive to attractive. For the CF $do$ (blond hair = 1, pale = 1), 16% of the labels are flipped and 74% of the labels are flipped from attractive to not attractive. In comparison, the affine transformations are unable to provide a clear picture of bias. For horizontal flip, the points are equally distributed on both sides of the reference line y = x. In the case of brightness, there is more variation.

In Table 3, we quantify these observations using Eqn 4. Our metric for bias measurement gives an overall estimate of the bias in classifier, and provides an interpretable and uniform scale to compare biases among different classifiers. The reported bias values reflect the observations from Fig 6. We observe that the CF of $do$ (brown hair = 1, pale = 1) has the highest positive bias amongst all CFs and affine transformations, i.e. the classifier is biased towards labeling these CFs as more attractive in contrast to the base image. Using CFs, we are able to detect other significant biases towards setting skin color as $pale = 1$ for all hair colors (black, blond and brown). In contrast, using the baseline transformations, we are unable to detect skin color bias in the classifier since the calculated bias values are negligible.

We also use CFs to explain the classifier's decisions for predicting attractiveness of a face. Fig 7 shows the top 4 positive and top 4 negative influences for classifying a face as attractive. We can see that *Pale* skin is the top attribute that contributes to the classifier predicting a face as attractive, while *Bald* is the top attribute that contributes to the classifier prediction of not attractive.
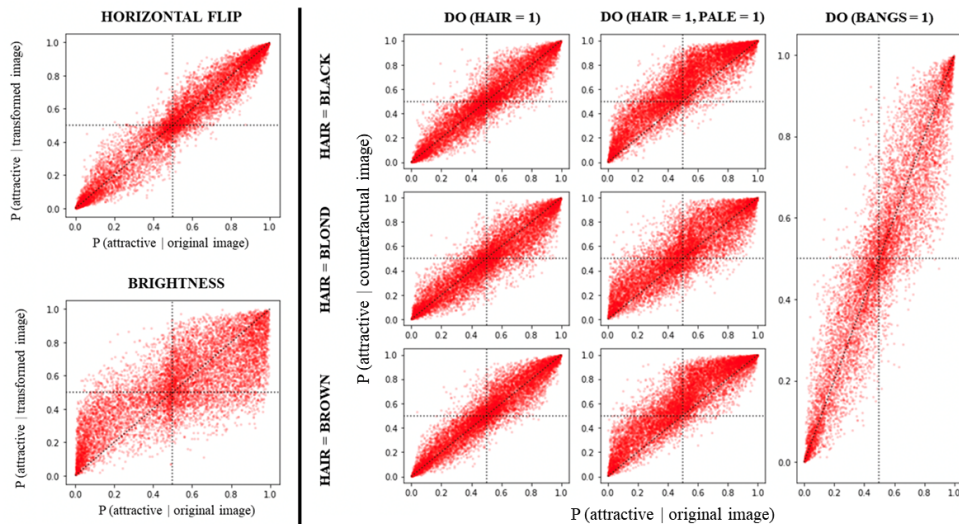
13

Figure 6: **Fairness Analysis.** Affine transformations *(left)*, CFs *(right)*. Each point in a scatter plot is a pair of a base image and its corresponding CF image. In the ideal case, all points should lie along the y = x line. To analyze the figures, divide the scatter plot into four quadrants formed by lines x = 0.5 and y = 0.5. Any point in the top left quadrant signifies that the attractive label was changed from 0 to 1 and vice-versa for the bottom right quadrant.
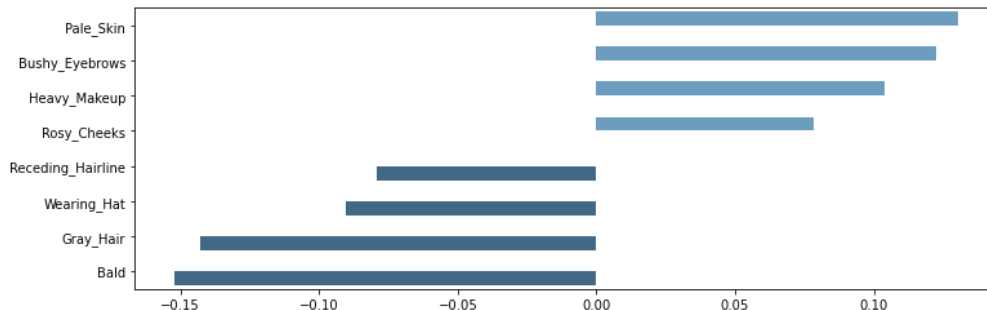


Figure 7: **Explaining a Classifier.** Attribute ranking of top 4 positive and top 4 negative influential attributes.

## 5.4  Bias Mitigation for a Classifier

Finally, using Eqn. 7, we employ the generated CFs to remove the identified biases in the *Attractive* classifier on *do* (black = 1, pale = 1), *do* (blond = 1, pale = 1) and *do* (brown = 1, pale = 1) . The CF-regularized classifier reports a bias score of 0.032 for black hair and pale (against 0.159 for the original classifier) and 0.012 for brown hair and pale (against 0.154 for the original classifier). Also, the reduced biases are no longer significant, without reducing the accuracy (82.3% versus 82.6%). Details are in Appendix N.

Additionally, our method can also be utilized in the setting where the attributes are connected

|  | $p(a_r \neq a_c)$ | $p(0 \rightarrow 1)$ | bias |
|---|---|---|---|
| horizontal_flip | 0.073 | 0.436 | -0.009 |
| brightness | 0.192 | 0.498 | -0.001 |
| black_h | 0.103 | 0.586 | 0.018 |
| black_h, pale | 0.180 | 0.937 | **0.158** |
| blond_h | 0.115 | 0.413 | - 0.02 |
| blond_h, pale | 0.155 | 0.738 | **0.073** |
| brown_h | 0.099 | 0.704 | 0.041 |
| brown_h, pale | 0.186 | 0.942 | **0.164** |
| bangs | 0.106 | 0.526 | 0.005 |

Table 3: **Bias Estimation.** Bias values above a threshold of 5% are considered significant.

in a complex causal graph structure, unlike [Denton et al., 2019, Joo and Kärkkäinen, 2020]. We conduct a similar fairness and explanation analysis for a *Attractive* classifier in Appendix O, where *Young* affects other visible attributes like *Gray hair*.

# 6 Conclusion

We propose a framework *ImageCFGen* for generating counterfactuals based on an underlying SCM, utilizing the generative capacity of GANs. We demonstrate how the counterfactuals can be used to evaluate and mitigate a classifier's biases and explain its decisions. That said, we acknowledge two limitations , 1) Our CF generation method relies on accurate knowledge of the causal graph; 2) It uses a statistical model that can have have unknown failure modes in generating meaningful counterfactuals. Therefore, this work should be considered as a prototype early work in generating counterfactuals, and is not suitable for deployment.

# References

M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Is-ard, et al. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pages 265–283, 2016.

T. Bolukbasi, K.-W. Chang, J. Y. Zou, V. Saligrama, and A. T. Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in neural information processing systems*, pages 4349–4357, 2016.

J. Buolamwini and T. Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91, 2018.

D. C. Castro, J. Tan, B. Kainz, E. Konukoglu, and B. Glocker. Morpho-mnist: Quantitative as-sessment and diagnostics for representation learning. *Journal of Machine Learning Research*, 20 (178):1–29, 2019.

F. Chollet et al. Keras. https://keras.io, 2015.

Y. Dandi, H. Bharadhwaj, A. Kumar, and P. Rai. Generalized adversarially learned inference. *arXiv preprint arXiv:2006.08089*, 2020.

E. Denton, B. Hutchinson, M. Mitchell, and T. Gebru. Detecting bias with generative counterfactual face attribute augmentation. *arXiv preprint arXiv:1906.06439*, 2019.

J. Ding. Probabilistic inferences in bayesian networks. *InA. Rebai (ed.), Bayesian Network*, pages 39–53, 2010.

Y. Dogan and H. Y. Keles. Semi-supervised image attribute editing using generative adversarial networks. *Neurocomputing*, 401:338–352, 2020.

J. Donahue and K. Simonyan. Large scale adversarial representation learning. *arXiv preprint arXiv:1907.02544*, 2019.

V. Dumoulin, I. Belghazi, B. Poole, O. Mastropietro, A. Lamb, M. Arjovsky, and A. Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016.

M. Ehrlich, T. J. Shields, T. Almaev, and M. R. Amer. Facial attributes classification using multi-task representation learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 47–55, 2016.

R. C. Fong and A. Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3429–3437, 2017.

Y. Goyal, Z. Wu, J. Ernst, D. Batra, D. Parikh, and S. Lee. Counterfactual visual explanations. In *ICML*, 2019.

L. A. Hendricks, K. Burns, K. Saenko, T. Darrell, and A. Rohrbach. Women also snowboard: Overcoming bias in captioning models. In *European Conference on Computer Vision*, pages 793–811. Springer, 2018.

M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *arXiv preprint arXiv:1706.08500*, 2017.

J. Joo and K. Kärkkäinen. Gender slopes: Counterfactual fairness for computer vision models by attribute manipulation. *arXiv preprint arXiv:2005.10430*, 2020.

A.-H. Karimi, J. von Kügelgen, B. Schölkopf, and I. Valera. Algorithmic recourse under imperfect causal knowledge: a probabilistic approach. *Advances in Neural Information Processing Systems*, 33, 2020.

T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *International Conference on Learning Representations*, 2018.

T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4401–4410, 2019.

T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. Analyzing and improving the image quality of stylegan, 2020.

D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

M. Kocaoglu, C. Snyder, A. G. Dimakis, and S. Vishwanath. Causalgan: Learning causal implicit generative models with adversarial training. In *International Conference on Learning Representations*, 2018.

M. J. Kusner, J. Loftus, C. Russell, and R. Silva. Counterfactual fairness. In *Advances in neural information processing systems*, pages 4066–4076, 2017.

G. Lample, N. Zeghidour, N. Usunier, A. Bordes, L. Denoyer, and M. Ranzato. Fader networks: Manipulating images by sliding attributes. In *Advances in neural information processing systems*, pages 5967–5976, 2017.

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

S. Liu, B. Kailkhura, D. Loveland, and Y. Han. Generative counterfactual introspection for explainable deep learning. In *2019 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pages 1–5. IEEE, 2019.

Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.

R. Luss, P.-Y. Chen, A. Dhurandhar, P. Sattigeri, Y. Zhang, K. Shanmugam, and C.-C. Tu. Generating contrastive explanations with monotonic attribute functions. *arXiv preprint arXiv:1905.12698*, 2019.

D. McDuff, S. Ma, Y. Song, and A. Kapoor. Characterizing bias in classifiers using generative models. In *Advances in Neural Information Processing Systems*, pages 5403–5414, 2019.

M. Mirza and S. Osindero. Conditional generative adversarial nets. arxiv 2014. *arXiv preprint arXiv:1411.1784*, 2014.

N. Pawlowski, D. C. Castro, and B. Glocker. Deep structural causal models for tractable counterfactual inference. *arXiv preprint arXiv:2006.06485*, 2020.

J. Pearl. *Causality*. Cambridge university press, 2009.

J. Pearl. Bayesian networks. 2011.

J. Pearl. The seven tools of causal inference, with reflections on machine learning. *Communications of the ACM*, 62(3):54–60, 2019.

D. J. Rezende and S. Mohamed. Variational inference with normalizing flows. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning-Volume 37*, pages 1530–1538, 2015.

A. S. Ross, M. C. Hughes, and F. Doshi-Velez. Right for the right reasons: training differentiable models by constraining their explanations. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 2662–2670, 2017.

P. Sattigeri, S. C. Hoffman, V. Chenthamarakshan, and K. R. Varshney. Fairness gan: Generating datasets with fairness properties using a generative adversarial network. *IBM Journal of Research and Development*, 63(4/5):3–1, 2019.

M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3319–3328, 2017.

S. Wachter, B. Mittelstadt, and C. Russell. Counterfactual explanations without opening the black box: Automated decisions and the gdpr. *Harv. JL & Tech.*, 31:841, 2017.

T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018.

M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014.

J. Zhao, T. Wang, M. Yatskar, V. Ordonez, and K.-W. Chang. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2979–2989, 2017.

S. Zhou, M. L. Gordon, R. Krishna, A. Narcomey, L. Fei-Fei, and M. S. Bernstein. Hype: A benchmark for human eye perceptual evaluation of generative models. *arXiv preprint arXiv:1904.01121*, 2019.

In this appendix, we discuss the following details, which could not be included in the main paper owing to space constraints.

# A  *ImageCFGen* Algorithm

As described in Section 3.1, our method involves two components: learning the Attribute-SCM ($\mathcal{M}_a$) that models relationships between attributes, and learning the generator ($\mathcal{M}_x$) that produces counterfactual images given the attributes.

## A.1  Learning the Attribute-SCM

We assume that the causal graph structure is provided by the user, but the causal functions relating the attributes are not available. For each attribute $a_i$, the graph entails a set of attributes $P_{a_i}$ that cause the attribute (its *parents* in the graph structure) that should be included while estimating the value of the attribute. It leads to the following generating equation for each attribute,

$$a_i = g_i(P_{a_i}, \epsilon_i) \tag{8}$$

where $P_{a_i}$ are the parents of attribute $a_i$ and $\epsilon_i$ is independent noise. The goal is to learn the unknown function $g_i$ for each attribute. Thus, given $n$ attributes, we obtain a set of $n$ equations, each of which can be independently estimated using Maximum Likelihood Estimation from Bayesian networks [Ding, 2010].

## A.2  Generating Counterfactuals

Algorithm 1 below shows how *ImageCFGen* is implemented, in extension to the discussion in Sec 3.2 of the main paper. Here we assume that the Attribute-SCM, the Encoder $E$, and Generator $G$ are pre-learnt and provided as input to the algorithm. $(\mathbf{x}^*, \mathbf{a}^*)$ denotes the input value of the image and the attributes. Note that the Attribute SCM operations (Steps 1-3) are based on the formal procedure to generate counterfactuals from [Pearl, 2009].

In Step 1 (Abduction), the encoder uses the input data $(\mathbf{x}^*, \mathbf{a}^*)$ to create the latent vector $\mathbf{z}$.

In Step 2 (Action), given the indices $K$ of sensitive attributes, the goal is to remove the dependence of these sensitive attributes on any other attribute and then change their value to the desired value. To do so, the causal graph of Attribute-SCM is modified to remove all incoming arrows to $\mathbf{a}_K$. That is, the structural equations for each $a_i \in \mathbf{a}_K$ changes from $a_i = g_i(P_{a_i}, \epsilon_i)$ to $a_i = \epsilon'_i$, yielding a modified Attribute-SCM $\mathcal{M}'_a$. Then each $a_i \in \mathbf{a}_K$ is set to the desired modified value $a_i^{**}$ (Line 5; we do this without changing its parents or other variables, hence it is called an *intervention*).

In Step 3 (Prediction), we propagate the change in sensitive attributes to all attributes caused by them. That is, values of all descendants of sensitive attributes $\mathbf{a}_K$ in the causal graph of Attribute-SCM $\mathcal{M}'_a$ are changed, based on the modified value of the sensitive attributes from Step 2. We do so level-by-level: first changing the values of children of $\mathbf{a}_K$ based on the modified value of $\mathbf{a}_K$, then the values of children of children of $\mathbf{a}_K$, and so on. Note that line 13 ensures that an attribute's value is changed only when all its parents' values have been modified. The corresponding equation for each descendant $a_j$ is given in Line 14,

$$a_j^{**} \rightarrow g_j(P_{a_j^{**}}, \hat{\epsilon}_j) \tag{9}$$

where $P_{a_j^{**}}$ refers to the modified values of the parents of $a_j$, $g_j$ is the pre-learnt Attribute-SCM function, and $\hat{\epsilon}_j$ is the estimated error term from the original attribute data (i.e., it is the same error

---

**Algorithm 1:** *ImageCFGen*

---

**Input:** Input data-attribute pair $(\mathbf{x}^*, \mathbf{a}^*)$; Indices of sensitive attributes to modify $K$; Modified attribute values $\mathbf{a}_K^{**}$, Attribute SCM $\mathcal{M}_a = \{g_i(P_{a_i^*}, \hat{\epsilon}_i)\} \forall a_i^* \in \mathbf{a}^*, P_{a_i^*} \subseteq \mathbf{a}^*$ where $P_{a_i^*}$ denotes parents of $a_i^*$ and $\hat{\epsilon}_i$ is estimated from data, Encoder $E$, Generator $G$

**Output:** $\mathbf{x}_c$

**1  Step 1: Abduction**

**2**  $\mathbf{z} = E(\mathbf{x}^*, \mathbf{a}^*)$

**3  Step 2: Action**

**4**  Remove all incoming arrows to $\mathbf{a}_K$ in $\mathcal{M}_a$ to yield $\mathcal{M}_a'$

**5**  $a_i \to a_i^{**} \quad \forall i \in K, a_i^{**} \in \mathbf{a}_K^{**}$ /* Intervene on sensitive attributes and set
    to modified attribute values                                                */

**6  Step 3: Prediction**

**7**  $D = desc(\mathbf{a}_K)$ /* descendants of $\mathbf{a}_K$                              */

**8**  $currset = children(\mathbf{a}_K)$

**9**  $changed[\mathbf{a}^*] = False$

**10 while** *currset is not empty* **do**

**11**      $cs = \emptyset$;

**12**      **foreach** $a_j \in currset$ **do**

            /* Proceed only if the values of all parents of $a_j$ are already
            changed                                                    */

**13**          **if** *not $any_l(P_{a_j}^{(l)} \in D$ & $changed[P_{a_j}^{(l)}] == False)$* **then**

**14**              $a_j^{**} \to g_j(P_{a_j^{**}}, \hat{\epsilon}_j)$

**15**              $cs.append(children(a_j))$

**16**              $changed[a_j] = True$

**17**          **end**

**18**      **end**

**19**      $currset = cs$ /* If there exist children of the current nodes, repeat
        while loop                                                */

**20 end**

**21**  $\mathbf{a}_c = \mathbf{a}_K^{**} \cup \mathbf{a}_{-K}^{**}$ /* $\mathbf{a}_{-K}^{**}$ denotes the final values for all non-sensitive
    attributes                                                     */

**22**  $\mathbf{x}_c = G(\mathbf{z}, \mathbf{a}_c)$

---

value that satisfies $a_j^* \to g_j(P_{a_j^*}, \hat{\epsilon}_j)$). Finally, the updated values of all attributes are used to create $\mathbf{a}_c$, the counterfactual attribute vector. This counterfactual attribute vector is then combined with latent $z$ from Step 1 and provided as input to the generator G, to obtain the counterfactual image $\mathbf{x}_c$ (Line 22).

# B  Counterfactual Generation Proof

**Proposition 1.** *Assume that the true SCM $\mathcal{M}$ belongs to a class of SCMs where the structural equations can be identified uniquely given the causal graph. Further, assume that the Encoder $E$, Generator $G$ and Attribute SCM are optimal and thus correspond to the true structural equations of the SCM $\mathcal{M}$. Then Equation 1 generates a valid counterfactual image for any given input $(\mathbf{x}, \mathbf{a})$ and the requested modified attributes $\mathbf{a}_k'$.*

*Proof.* Let $\mathcal{M} = \{\mathcal{M}_x, \mathcal{M}_a\}$ be the true SCM that generates the data $(\mathbf{x}, \mathbf{a})$. Let $\mathbf{a}_k \subset \mathbf{a}$ be the attributes you want to intervene on. Let $\mathbf{a}_{-k} = \mathbf{a} \setminus \mathbf{a}_k$ be the remaining attributes. The corresponding

equations for $\mathcal{M}$ are:

$$a_i := g_i(p_{a_i}, \epsilon_i), \quad \forall i = 1..n \tag{10}$$
$$\mathbf{x} := g(\mathbf{a}_k, \mathbf{a}_{-k}, \epsilon)$$

where $\epsilon$ and $\epsilon_i$ refer to independent noise, g and $g_i$ refer to structural assignments of the SCMs $\mathcal{M}_x$ and $\mathcal{M}_a$ respectively while $p_{a_i}$ refers to parent attributes of $a_i$. Given an input $(\mathbf{x}, a_k, a_{-k})$, we generate a counterfactual using $\mathcal{M}$ and show that it is equivalent to Equation 1.

**Abduction:** Infer $\epsilon$ for $\mathbf{x}$ and $\epsilon_i$ for all $a_i$ from Equation 10 using the following equations:

$$\hat{\epsilon}_i := g_i^{-1}(a_i, p_{a_i}), \quad \forall i = 1..n \tag{11}$$
$$\hat{\epsilon} := g^{-1}(\mathbf{x}, \mathbf{a}_k, \mathbf{a}_{-k})$$

**Action:** Intervene on all attributes $\mathbf{a}_k$ by setting them to the requested modified attributes $\mathbf{a}'_k$.

$$a_i \to a'_i \quad \forall a_i \in \mathbf{a}_k \quad \text{and} \quad \forall a'_i \in \mathbf{a}'_k \tag{12}$$

**Prediction:** The following equation then changes the values of all descendants of $\mathbf{a}_k$.

$$desc(a_i) \to g_i(p_{desc(a_i)}, \hat{\epsilon}_i) \quad \forall a_i \in \mathbf{a}_k \tag{13}$$

where $desc(a_i)$ are descendants of $a_i$ and $p_{desc(a_i)}$ are parents of the descendants of $a_i$, $\forall a_i \in \mathbf{a}_k$. Let $\mathbf{a}_c = \mathbf{a}'_k \cup \mathbf{a}'_{-k}$ where $\mathbf{a}'_{-k}$ are (possibly) modified values of the other attributes according to Equation 13. Therefore, the counterfactual of $\mathbf{x}$, $\mathbf{x}_c$ can be generated as:

$$\mathbf{x}_c := g(\mathbf{a}_c, \hat{\epsilon}) \tag{14}$$

We now show that Equation 1 produces the same $\mathbf{x}_c$. By the assumption in the theorem statement, the Attribute-SCM corresponds to the structural assignments $\{g_i, g_i^{-1}\}$, $\forall i = 1..n$ of SCM $\mathcal{M}_a$ while the Generator $G$ learns the structural assignment $g$ and the Encoder $E$ learns $g^{-1}$ of the SCM $\mathcal{M}_x$. Hence, the Attribute-SCM, Generator and Encoder learn the combined SCM $\mathcal{M}$.

When the SCM assignments learned by the Attribute-SCM are optimal, i.e. Attribute-SCM = $\mathcal{M}_a$ then:

$$\hat{a}_c = a_c$$

Similarly, under optimal Generator, $G = g$ and $E = g^{-1}$:

$$\mathbf{x}_c = g(\mathbf{a}_c, g^{-1}(\mathbf{x}, \mathbf{a}_k, \mathbf{a}_{-k})) \tag{15}$$
$$= G(\mathbf{a}_c, E(\mathbf{x}, \mathbf{a})) \quad (\text{as} \quad \mathbf{a} = \mathbf{a}_k \cup \mathbf{a}_{-k})$$

which is the same as Equation 1. □

# C  Counterfactual Fairness Proof

**Definition 1.** *Counterfactual Fairness from [Kusner et al., 2017]. Let $\mathcal{A}$ be the set of attributes, comprising of sensitive attributes $\mathcal{A}_S \subseteq \mathcal{A}$ and other non-sensitive attributes $\mathcal{A}_N$. The classifier $\hat{f}$ is counterfactually fair if under any context $\mathbf{X} = \mathbf{x}$ and $\mathbf{A} = \mathbf{a}$, changing the value of the sensitive features to $\mathbf{A}_S \leftarrow a'_s$ counterfactually does not change the classifier's output distribution $Y$.*

$$P(Y_{A_S \leftarrow a_s} = y | \mathbf{X} = \mathbf{x}, \mathbf{A}_S = \mathbf{a}_s, \mathbf{A}_N = \mathbf{a}_N)$$
$$= P(Y_{A_S \leftarrow a'_s} = y | \mathbf{X} = \mathbf{x}, \mathbf{A}_S = \mathbf{a}_s, \mathbf{A}_N = \mathbf{a}_N) \tag{16}$$

*for all $y$ and for any value $\mathbf{a}'_s$ attainable by $\mathbf{A}_S$.*

**Proposition 2.** *Under the assumptions of Proposition 1 for the encoder E, generator G, and Attribute SCM $\mathcal{M}_a$, a classifier $\hat{f}(\mathbf{X}) : \mathcal{X} \to \mathcal{Y}$ that satisfies zero bias according to Equation 4 is counterfactually fair with respect to $\mathcal{M}$.*

*Proof.* To evaluate fairness, we need to reason about the output $Y = y$ of the classifier $\hat{f}$. Therefore, we add a functional equation to the SCM $\mathcal{M}$ from Equation 10.

$$
\begin{aligned}
a_i &:= g_i(p_{a_i}, \epsilon_i), \quad \forall i = 1..n \\
\mathbf{x} &:= g(\mathbf{a}_k, \mathbf{a}_{-k}, \epsilon) \\
y &\leftarrow \hat{f}(\mathbf{x})
\end{aligned}
\tag{17}
$$

where $\epsilon$ and $\epsilon_i$ are independent errors as defined in Equation 10, and $P_{\mathbf{a}_i}$ refers to the parent attributes of an attribute $a_i$ as per the Attribute-SCM. The SCM equation for $y$ does not include an noise term since the ML classifier $\hat{f}$ is a deterministic function of $\mathbf{X}$. In the above equations, the attributes $\mathbf{a}$ are separated into two subsets: $\mathbf{a}_k$ are the attributes specified to be changed in a counterfactual and $\mathbf{a}_{-k}$ refers to all other attributes.

Based on this SCM, we now generate a counterfactual $y_{a_k \leftarrow a'_k}$ for an arbitrary new value $a'_k$. Using the **Prediction** step, the counterfactual output label $y_c$ for an input $(\mathbf{x}, \mathbf{a})$ is given by: $y_c = \hat{f}(\mathbf{X} = \mathbf{x}_c)$. From Theorem 1, under optimality of the encoder $E$, generator $G$ and learned functions $g_i$, we know that $\mathbf{x}_c$ generated by the following equation is a valid counterfactual for an input $(\mathbf{x}, \mathbf{a})$,

$$
\begin{aligned}
\mathbf{x}_c &= G(E(\mathbf{x}, \mathbf{a}), \mathbf{a}_c) \\
&= X_{\mathbf{A}_k \leftarrow a'_k} | (\mathbf{X} = \mathbf{x}, \mathbf{A}_k = \mathbf{a}_k, \mathbf{A}_{-k} = \mathbf{a}_{-k})
\end{aligned}
\tag{18}
$$

where $\mathbf{a}_c$ represents the modified values of the attributes under the action $\mathbf{A}_k \leftarrow a'_k$. Therefore, $y_{a_k \leftarrow a'_k} | (\mathbf{X} = \mathbf{x}, \mathbf{A}_k = \mathbf{a}_k, \mathbf{A}_{-k} = \mathbf{a}_{-k})$ is given by $y_c = \hat{f}(\mathbf{x}_c)$.

Using the above result, we now show that the bias term from Equation 4 and 5 reduces to the counterfactual fairness definition from Definition 1,

$$
\begin{aligned}
&P(y_r = 0, y_c = 1) - P(y_r = 1, y_c = 0) \\
&= [P(y_r = 0, y_c = 1) + P(y_r = 1, y_c = 1)] - \\
&\quad [P(y_r = 1, y_c = 0) + (P(y_r = 1, y_c = 1)] \\
&= P(y_c = 1) - P(y_r = 1) \\
&= [P(Y_{\mathbf{A}_k \leftarrow a'_k} = 1 | \mathbf{X} = \mathbf{x}, \mathbf{A}_k = \mathbf{a}_k, \mathbf{A}_{-k} = \mathbf{a}_{-k}) - \\
&\quad P(Y_{\mathbf{A}_k \leftarrow a_k} = 1 | \mathbf{X} = \mathbf{x}, \mathbf{A}_k = \mathbf{a}_k, \mathbf{A}_{-k} = \mathbf{a}_{-k})] \\
&= [P(Y_{\mathbf{A}_S \leftarrow a'_s} = 1 | \mathbf{X} = \mathbf{x}, \mathbf{A}_S = \mathbf{a}_s, \mathbf{A}_N = \mathbf{a}_n) - \\
&\quad P(Y_{\mathbf{A}_S \leftarrow a_s} = 1 | \mathbf{X} = \mathbf{x}, \mathbf{A}_S = \mathbf{a}_s, \mathbf{A}_N = \mathbf{a}_n)]
\end{aligned}
\tag{19}
$$

where the second equality is since $y_r, y_c \in \{0, 1\}$, the third equality is since the reconstructed $y_r$ is the output prediction when $\mathbf{A} = \mathbf{a}$, and the last equality is $\mathbf{A}_k$ being replaced by the sensitive attributes $\mathbf{A}_S$ and $\mathbf{A}_{-k}$ being replaced by $\mathbf{A}_N$. We can prove a similar result for $y_c = 0$. Hence, when bias term is zero, the ML model $\hat{f}$ satisfies counterfactual fairness (Definition 1). $\square$

# D    Architecture Details

Here we provide the architecture details for the base Adversarially Learned Inference (ALI) model trained on Morpho-MNIST and CelebA datasets. For Cyclic Style ALI, we replace the ALI generator with the style-based generator architecture in [Karras et al., 2019]. We however do not use

progressive growing or other regularization strategies suggested in [Karras et al., 2019, 2020] while training our model. For details on the style-based generation architecture, please refer [Karras et al., 2019].

Overall, the architectures and hyperparameters are similar to the ones used by [Dumoulin et al., 2016], with minor variations. Instead of using the Embedding network from the original paper, the attributes are directly passed on to the Encoder, Generator and Discriminator. We found that this enabled better conditional generation in our experiments. All experiments were implemented using Keras 2.3.0 [Chollet et al., 2015] with Tensorflow 1.14.0 [Abadi et al., 2016]. All models were trained using a Nvidia Tesla P100 GPU.

## D.1 Morpho-MNIST

Tables 4, 5 and 6 show the Generator, Encoder and Discriminator architectures respectively for generating Morpho-MNIST counterfactuals. Conv2D refers to Convolution 2D layers, Conv2DT refers to Transpose Convoloution layers, F refers to number of filters, K refers to kernel width and height, S refers to strides, BN refers to Batch Normalization, D refers to dropout probability and A refers to the activation function. LReLU denotes the Leaky ReLU activation function. We use the Adam optimizer [Kingma and Ba, 2014] with a learning rate of $10^{-4}$, $\beta_1 = 0.5$ and a batch size of 100 for training the model. For the LeakyReLU activations, $\alpha = 0.1$. The model converges in approximately 30k iterations. All weights are initialized using the Keras truncated normal initializer with $mean = 0.0$ and $stddev = 0.01$. All biases are initialized with zeros.

| Layer | F | K | S | BN | D | A |
|-------|-----|-------|-------|----|-----|---------|
| Conv2DT | 256 | (4,4) | (1,1) | Y | 0.0 | LReLU |
| Conv2DT | 128 | (4,4) | (2,2) | Y | 0.0 | LReLU |
| Conv2DT | 64 | (4,4) | (1,1) | Y | 0.0 | LReLU |
| Conv2DT | 32 | (4,4) | (2,2) | Y | 0.0 | LReLU |
| Conv2DT | 32 | (1,1) | (1,1) | Y | 0.0 | LReLU |
| Conv2D | 1 | (1,1) | (1,1) | Y | 0.0 | Sigmoid |

Table 4: **Architecture for Morpho-MNIST Generator**

| Layer | F | K | S | BN | D | A |
|-------|-----|-------|-------|----|-----|--------|
| Conv2D | 32 | (5,5) | (1,1) | Y | 0.0 | LReLU |
| Conv2D | 64 | (4,4) | (2,2) | Y | 0.0 | LReLU |
| Conv2D | 128 | (4,4) | (1,1) | Y | 0.0 | LReLU |
| Conv2D | 256 | (4,4) | (2,2) | Y | 0.0 | LReLU |
| Conv2D | 512 | (3,3) | (1,1) | Y | 0.0 | LReLU |
| Conv2D | 512 | (1,1) | (1,1) | Y | 0.0 | Linear |

Table 5: **Architecture for Morpho-MNIST Encoder**

## D.2 CelebA

Tables 7, 8 and 9 show the Generator, Encoder and Discriminator architectures respectively for generating CelebA counterfactuals.

| Layer | F | K | S | BN | D | A |
|---|---|---|---|---|---|---|
| $D_z$ | | | | | | |
| Conv2D | 512 | (1,1) | (1,1) | N | 0.2 | LReLU |
| Conv2D | 512 | (1,1) | (1,1) | N | 0.5 | LReLU |
| $D_x$ | | | | | | |
| Conv2D | 32 | (5,5) | (1,1) | N | 0.2 | LReLU |
| Conv2D | 64 | (4,4) | (2,2) | Y | 0.2 | LReLU |
| Conv2D | 128 | (4,4) | (1,1) | Y | 0.5 | LReLU |
| Conv2D | 256 | (4,4) | (2,2) | Y | 0.5 | LReLU |
| Conv2D | 512 | (3,3) | (1,1) | Y | 0.5 | LReLU |
| $D_{x\_z}$ | | | | | | |
| Conv2D | 1024 | (1,1) | (1,1) | N | 0.2 | LReLU |
| Conv2D | 1024 | (1,1) | (1,1) | N | 0.2 | LReLU |
| Conv2D | 1 | (1,1) | (1,1) | N | 0.5 | Sigmoid |

Table 6: **Architecture for Morpho-MNIST Discriminator.** $D_x$, $D_z$ and $D_{xz}$ are the discriminator components to process the image $x$, the latent variable $z$ and the output of $D_x$ and $D_z$ concatenated, respectively.

| Layer | F | K | S | BN | D | A |
|---|---|---|---|---|---|---|
| Conv2DT | 512 | (4,4) | (1,1) | Y | 0.0 | LReLU |
| Conv2DT | 256 | (7,7) | (2,2) | Y | 0.0 | LReLU |
| Conv2DT | 256 | (5,5) | (2,2) | Y | 0.0 | LReLU |
| Conv2DT | 128 | (7,7) | (2,2) | Y | 0.0 | LReLU |
| Conv2DT | 64 | (2,2) | (1,1) | Y | 0.0 | LReLU |
| Conv2D | 3 | (1,1) | (1,1) | Y | 0.0 | Sigmoid |

Table 7: **Architecture for CelebA Generator**

| Layer | F | K | S | BN | D | A |
|---|---|---|---|---|---|---|
| Conv2D | 64 | (2,2) | (1,1) | Y | 0.0 | LReLU |
| Conv2D | 128 | (7,7) | (2,2) | Y | 0.0 | LReLU |
| Conv2D | 256 | (5,5) | (2,2) | Y | 0.0 | LReLU |
| Conv2D | 256 | (7,7) | (2,2) | Y | 0.0 | LReLU |
| Conv2D | 512 | (4,4) | (1,1) | Y | 0.0 | LReLU |
| Conv2D | 512 | (1,1) | (1,1) | Y | 0.0 | Linear |

Table 8: **Architecture for CelebA Encoder**

# E   Morpho-MNIST Latent Space Interpolations

We also plot latent space interpolations between pairs of images sampled from the test set. Figure 8 shows that the model has learned meaningful latent space representations and the transitional images look realistic as well.

| Layer | F | K | S | BN | D | A |
|---|---|---|---|---|---|---|
| $D_z$ | | | | | | |
| Conv2D | 1024 | (1,1) | (1,1) | N | 0.2 | LReLU |
| Conv2D | 1024 | (1,1) | (1,1) | N | 0.2 | LReLU |
| $D_x$ | | | | | | |
| Conv2D | 64 | (2,2) | (1,1) | N | 0.2 | LReLU |
| Conv2D | 128 | (7,7) | (2,2) | Y | 0.2 | LReLU |
| Conv2D | 256 | (5,5) | (2,2) | Y | 0.2 | LReLU |
| Conv2D | 256 | (7,7) | (2,2) | Y | 0.2 | LReLU |
| Conv2D | 512 | (4,4) | (1,1) | Y | 0.2 | LReLU |
| $D_{x\_z}$ | | | | | | |
| Conv2D | 2048 | (1,1) | (1,1) | N | 0.2 | LReLU |
| Conv2D | 2048 | (1,1) | (1,1) | N | 0.2 | LReLU |
| Conv2D | 1 | (1,1) | (1,1) | N | 0.2 | Sigmoid |

Table 9: **Architecture for CelebA Discriminator.** $D_x$, $D_z$ and $D_{xz}$ are the discriminator components to process the image $x$, the latent variable $z$ and the output of $D_x$ and $D_z$ concatenated, respectively.
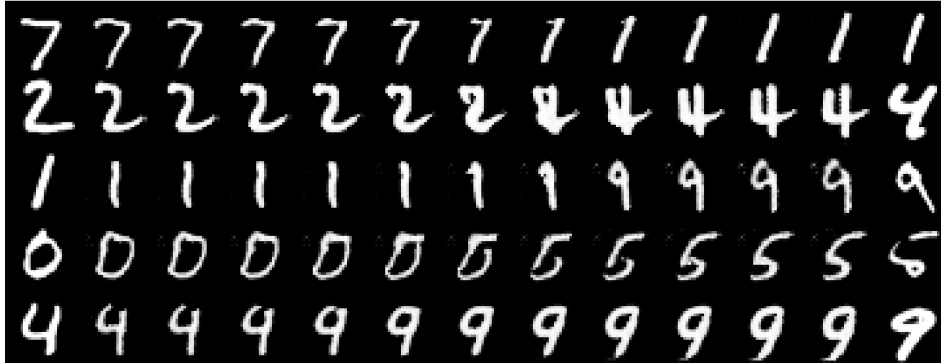


Figure 8: **Morpho-MNIST Interpolations.** The columns on the extreme left and right denote real samples from the test set and the columns in between denote generated images for the linearly interpolated latent space representation **z**.

# F   Morpho-MNIST Reconstructions

We qualitatively evaluate the inference model (Encoder) by sampling images along with their attributes from the test set and passing them through the encoder to obtain their latent space representations. These representations are passed to the generator which outputs reconstructions of the original image. The reconstructions are showed in Figure 9. Overall, reconstructions for Morpho-MNIST are faithful reproductions of the real images.
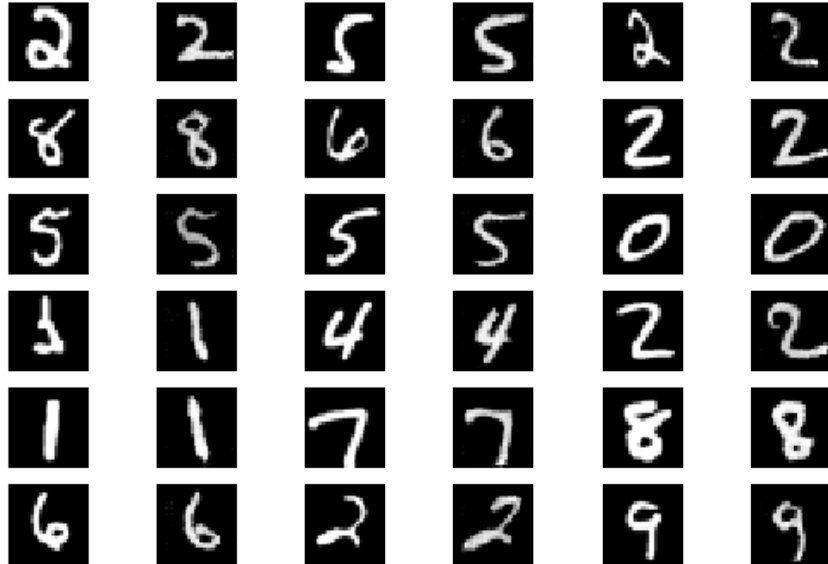
Figure 9: **Morpho-MNIST Reconstructions.** Odd columns denote real images sampled from the test set, and even columns denote reconstructions for those real images.

# G  Morpho-MNIST: Evaluating Label CFs

As shown in Figure 3, we empirically evaluated the counterfactual generation on Morpho-MNIST by generating CFs that change the digit label for an image. To check whether the generated counterfactual image corresponds to the desired digit label, we use the output of a digit classifier trained on Morpho-MNIST images. Here we provide details of this pre-trained digit classifier. The classifier architecture is shown in Table 10. The classifier converges in approximately 1k iterations with a validation accuracy of 98.30%. We then use the classifier to predict the labels of the counterfac-

| Layer | F | K | S | BN | D | A |
|---|---|---|---|---|---|---|
| Conv2D | 32 | (3,3) | (1,1) | N | 0.0 | ReLU |
| Conv2D | 64 | (3,3) | (1,1) | N | 0.0 | ReLU |
| MaxPool2D | - | (2,2) | - | N | 0.0 | - |
| Dense | 256 | - | - | N | 0.5 | ReLU |
| Dense | 128 | - | - | N | 0.5 | ReLU |
| Dense | 10 | - | - | N | 0.5 | Softmax |

Table 10: **Morpho-MNIST Label Classifier**

tual images and compare them to the labels that they were supposed to be changed to (target label). Overall, 97.30% of the predicted labels match the target label. Since the classifier is not perfect, the difference between the CF image's (predicted) label and the target label may be due to an error in the classifier's prediction or due to an error in CF generation, but this is minimal. We show some images for which predicted labels do not match the target labels in Figure 10. Most of these images are digits that can be confused as another digit; for instance, the first image in row 2 is a counterfactual

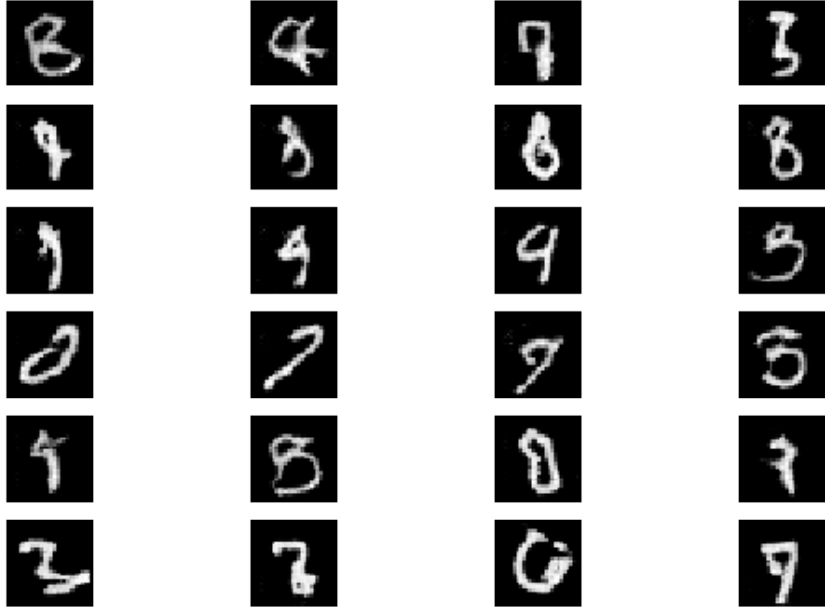image with a target label of 9, but was predicted by the digit classifier as a 1.



Figure 10: **Misclassified Counterfactuals.** Counterfactual images for Morpho-MNIST on which target label does not match predicted label.

# H DeepSCM Implementation on CelebA

Since DeepSCM was not implemented on CelebA, we reproduce their method based on their available code and choose an appropriate VAE architecture for CelebA data. Tables 11 and 12 show the Decoder and Encoder architectures respectively for generating CelebA counterfactuals. Specifically, we use Equation A.4 from the DeepSCM paper [Pawlowski et al., 2020] to model the image $\mathbf{x}$. We use a conditional VAE with a fixed variance Gaussian decoder to output a *bias* to further reparametrize a Gaussian distribution using a location-scale transform (as stated in [Pawlowski et al., 2020]). In our implementation of DeepSCM, we do not need to model the attributes $\mathbf{a}$ using Normalizing Flows, since the attributes independently influence the images in the underlying SCM in case of CelebA (as explained in Section 5). We set the latent dimension to 256, use the Adam optimizer [Kingma and Ba, 2014] with a learning rate of 0.0005 to train the conditional VAE with a batch size of 128. Similar to Pawlowski et al. [2020], we adopt a constant variance assumption and set $\log \sigma^2 = -5$. We also preprocess the images by scaling them between [0,1]; the preprocessing flow in Equation A.4 in [Pawlowski et al., 2020] is hence altered by normalizing its outputs between [0,1).

| Layer   | F    | K     | S     | BN | D    | A       |
|---------|------|-------|-------|----|------|---------|
| Dense   | 4096 | -     | -     | N  | 0.0  | Linear  |
| Conv2DT | 256  | (3,3) | (2,2) | Y  | 0.25 | LReLU   |
| Conv2DT | 128  | (3,3) | (2,2) | Y  | 0.25 | LReLU   |
| Conv2DT | 64   | (3,3) | (2,2) | Y  | 0.25 | LReLU   |
| Conv2DT | 32   | (3,3) | (2,2) | Y  | 0.25 | LReLU   |
| Conv2D  | 3    | (3,3) | (1,1) | Y  | 0.0  | Sigmoid |

Table 11: **Architecture for CelebA DeepSCM Decoder**

| Layer  | F   | K     | S     | BN | D    | A     |
|--------|-----|-------|-------|----|------|-------|
| Conv2D | 32  | (3,3) | (2,2) | Y  | 0.25 | LReLU |
| Conv2D | 64  | (3,3) | (2,2) | Y  | 0.25 | LReLU |
| Conv2D | 128 | (3,3) | (2,2) | Y  | 0.25 | LReLU |
| Conv2D | 256 | (3,3) | (2,2) | Y  | 0.25 | LReLU |
| Conv2D | 256 | (1,1) | (1,1) | N  | 0.0  | LReLU |

Table 12: **Architecture for CelebA DeepSCM Encoder**

# I  Additional DeepSCM vs *ImageCFGen* CelebA Counterfactuals

In addition to Figure 5 of the main paper, we provide additional counterfactual images obtained using *ImageCFGen* and DeepSCM in Figure 11. As observed in Section 5.2, the counterfactuals generated using *ImageCFGen* are more true to the corresponding intervention than DeepSCM.

# J  Additional *ImageCFGen* CelebA Counterfactuals

Figure 12 shows more CF examples obtained using *ImageCFGen*. Note how the CF image is different w.r.t. the base (reconstructed) image only in terms of the intervened attribute. Consequently, if the attribute in the base image is already present, the CF image is exactly the same as the original image. For instance, (Ib) and (Ig) in Fig 12 are exactly the same since (Ia) already has brown hair, hence intervening on brown hair has no effect on the hair color.

# K  Ablation Study of Style-Based Generator and Cyclic Cost Minimization

We plot the reconstructed images of real images, produced by ALI, StyleALI (ALI with style based generator) and Cyclic Style ALI (Style ALI with Cyclic Cost Minimization) in Figure 13. We observe that using a style-based generator significantly improves the quality of the generated images and applying the cyclic cost minimization algorithm on top of it improves reconstruction of the real image.
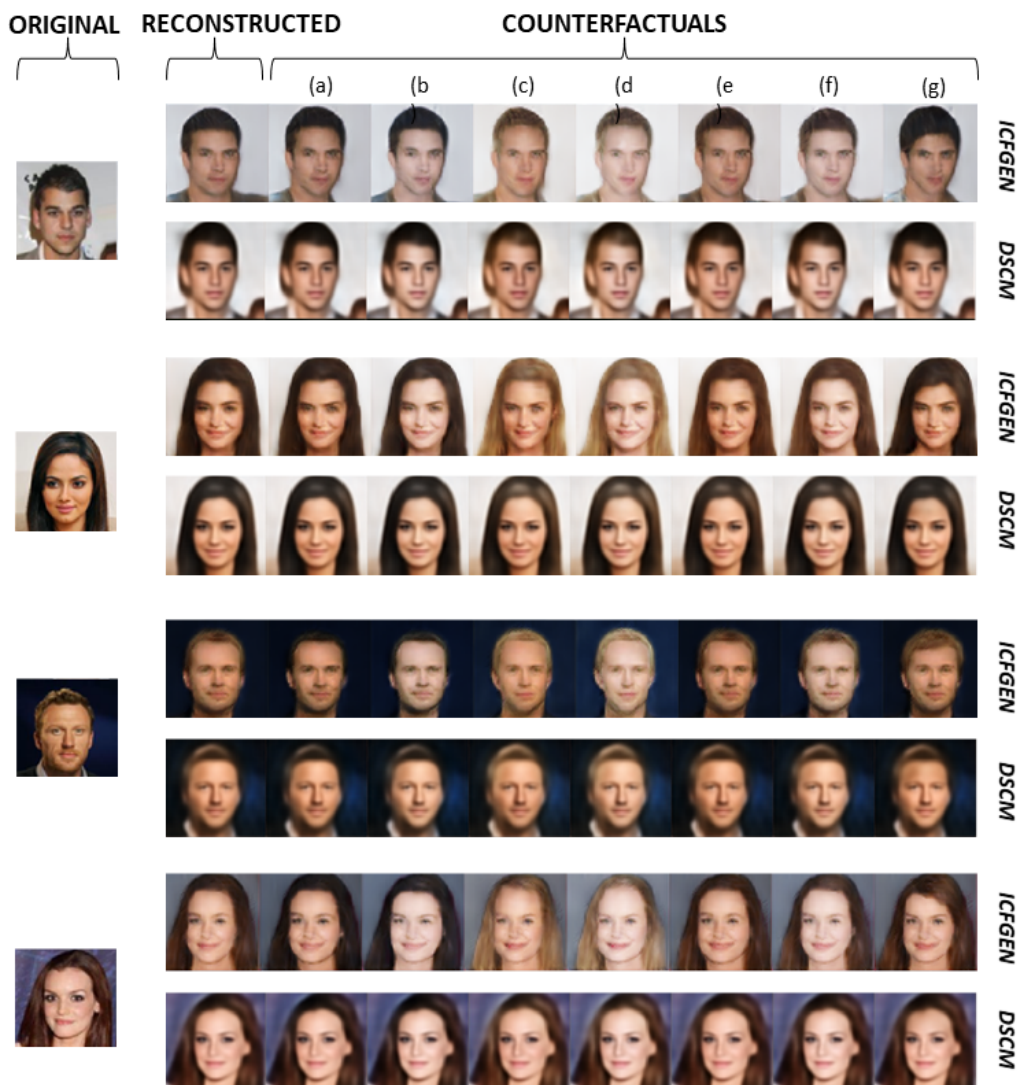
Figure 11: *ImageCFGen* and DeepSCM Counterfactuals. (a) denotes *do* (black hair = 1) and (b) denotes *do* (black hair = 1, pale =1). Similarly (c) denotes *do* (blond hair = 1); (d) denotes *do* (blond hair = 1, pale = 1); (e) denotes *do* (brown hair = 1); (hf denotes *do* (brown hair = 1, pale = 1); and (g) denotes *do* (bangs = 1).

# L    Human Evaluation of Counterfactuals

To quantitatively evaluate the counterfactuals, we asked human evaluators to pick the "edited" (counterfactual) image from 10 randomly sampled pairs of reconstructed and counterfactual images to human evaluators and asked them to pick the *edited* (counterfactual) image. Overall, we got 66 responses resulting in an average score of 5.15 correct answers out of 10 with a standard deviation
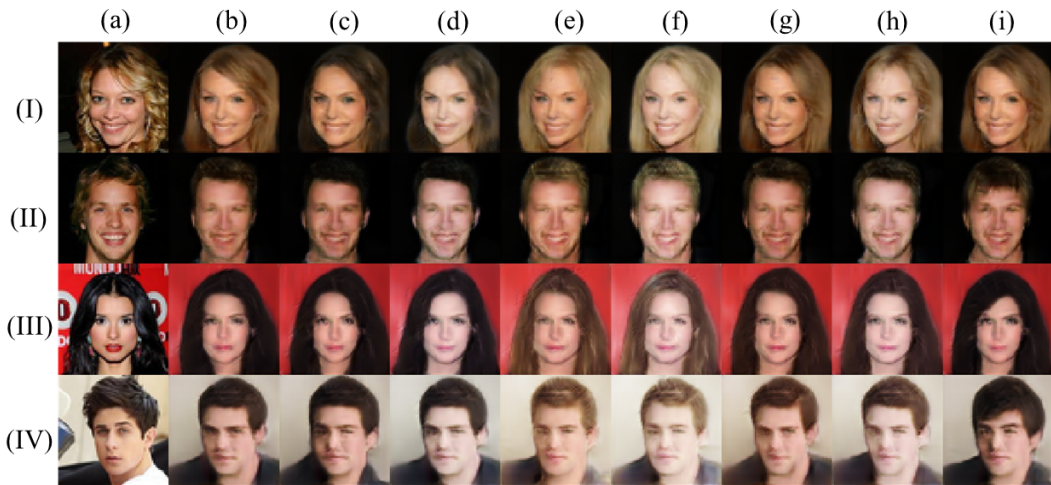
Figure 12: **CelebA Counterfactuals.** Column (a) across all rows I, II, III, IV represent the real image and (b) represents the reconstructed image. (c) denotes $do$ (black hair = 1) and (d) denotes $do$ (black hair = 1, pale =1). Similarly (e) denotes $do$ (blond hair = 1); (f) denotes $do$ (blond hair = 1, pale = 1); (g) denotes $do$ (brown hair = 1); (h) denotes $do$ (brown hair = 1, pale = 1); and (i) denotes $do$ (bangs = 1).
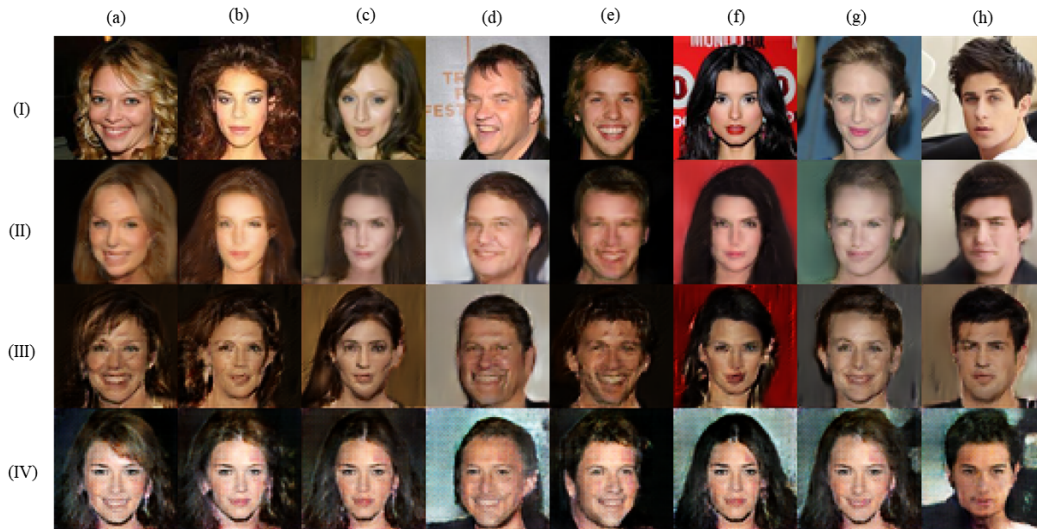


Figure 13: **Ablation Study.** Row (I) consists of real images, (II) consists of reconstructions generated by Cyclic Style ALI (CSALI), (III) Style ALI reconstructions and (IV) has ALI reconstructions

of 1.64, indicating that the generated counterfactual distribution is perceptually indistinguishable from the reconstructed one. Figure 14 shows sample questions on the form circulated for our human evaluation studies. For each question, the user chose one of two images that seems edited to the human eye. If the counterfactual can fool human perception, it indicates better performance of
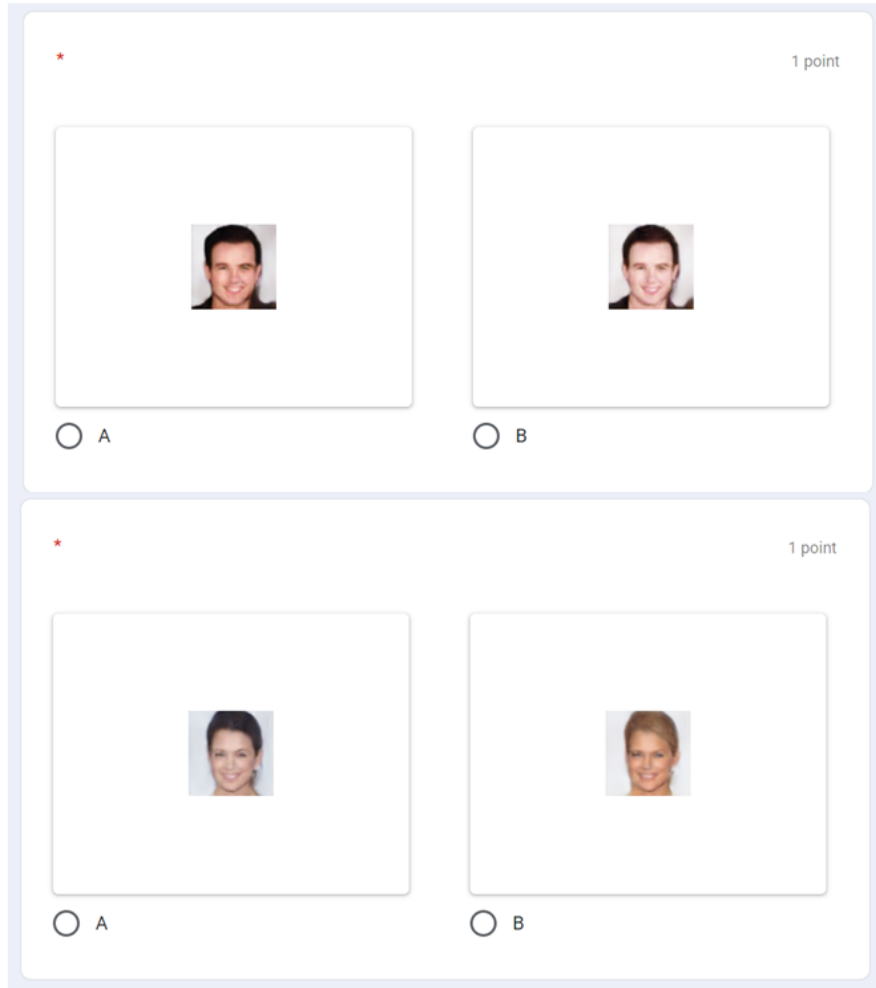
counterfactual generation.'



Figure 14: **Edited or Not?** Sample questions on the form circulated for human evaluation studies

## M *Attractive* Classifier Architecture

We describe the architecture and training details for the *attractiveness* classifier whose fairness was evaluated in Figure 6. The same classifier's output was explained in Figure 7. *Attractive* is a binary attribute associated with every image in CelebA dataset. The architecture for the binary classifier is shown in Table 13. We use the Adam optimizer with a learning rate of $10^{-4}$, $\beta_1 = 0.5$ and a batch size of 128 for training the model. For the LeakyReLU activations, $\alpha = 0.02$. The model converges in approximately 20k iterations. All weights are initialized using the Keras truncated normal initializer with $mean = 0.0$ and $stddev = 0.01$. All biases are initialized with zeros.

| Layer | F | K | S | BN | D | A |
|---|---|---|---|---|---|---|
| Conv2D | 64 | (2,2) | (1,1) | Y | 0.2 | LReLU |
| Conv2D | 128 | (7,7) | (2,2) | Y | 0.2 | LReLU |
| Conv2D | 256 | (5,5) | (2,2) | Y | 0.2 | LReLU |
| Conv2D | 256 | (7,7) | (2,2) | Y | 0.2 | LReLU |
| Conv2D | 512 | (4,4) | (1,1) | Y | 0.5 | LReLU |
| Conv2D | 512 | (1,1) | (1,1) | Y | 0.5 | Linear |
| Dense | 128 | - | - | N | 0.5 | ReLU |
| Dense | 128 | - | - | N | 0.5 | ReLU |
| Dense | 1 | - | - | N | 0.5 | Sigmoid |

Table 13: **Architecture for the *Attractiveness* Classifier**

# N    Bias Mitigation Details

We now present more details of our bias mitigation results in continuation to the discussion in Sec 5.4 of the main paper. To pick the optimal model while finetuning the classifier with the proposed bias mitigation regularizer, we set an accuracy threshold of 80% and pick the model with the lowest bias. We use $\lambda = 1.0$ and use an Adam optimizer with a learning rate of $10^{-4}$, $\beta_1 = 0.5$ and a batch size of 128 for training the model.

|  | Fair Classifier | Biased Classifier |
|---|---|---|
| black_h, pale | 0.032 | 0.159 |
| blond_h, pale | -0.041 | 0.077 |
| brown_h, pale | 0.012 | 0.154 |

Table 14: **Bias Values after Bias Mitigation.** Lower bias is better. Absolute bias values less than 5% are not considered significant.

# O    Complex Attribute SCM

Our method can also be utilized in the setting where the attributes are connected in a complex causal graph structure, unlike [Denton et al., 2019, Joo and Kärkkäinen, 2020]. We now conduct a fairness analysis w.r.t age for the *Attractive* classifier, assuming that *Young* affects other visible attributes like *Gray hair*.

*Defining the Causal Graph.* Fig 15 shows a potential causal graph for the attribute *Young*. While we ignored the attribute *Young* in the previous experiments, we now propose a SCM that describes the relationship between *Young* and a given image, mediated by a few facial attributes. Among the two demographic attributes, we choose *Young* over *Male* since *Young* has some descendants that can be agreed upon in general, whereas using *Male* would have forced us to explicitly model varying societal conventions surrounding this attribute. For e.g., the SCM proposed by [Kocaoglu et al., 2018] describes a causal relationship from *gender* to *smiling* and *narrow eyes*, which is a problematic construct. For the attribute *Young*, we specifically choose *Gray Hair*, *Eyeglasses* and *Receding Hairline* since it is reasonable to assume that older people are more likely to have these attributes, as compared to younger ones.
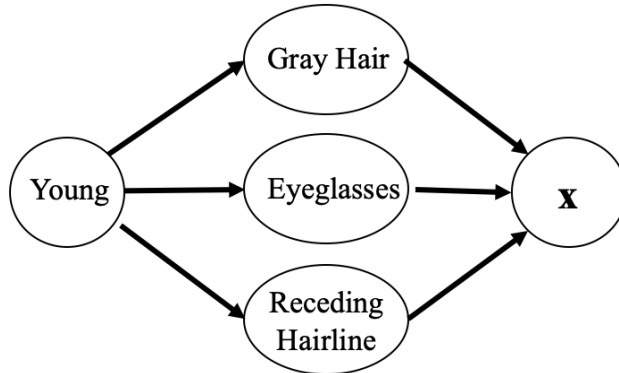
Figure 15: *Young* **Attribute-SCM.** SCM that describes the relationship between the attribute *Young* and the generated image.

*Learning the Attribute SCM.* To learn the model parameters for the SCM shown in Figure 15, we estimate conditional probabilities for each edge between *Young* and the other attributes using maximum likelihood over the training data, as done in Bayesian networks containing only binary-valued nodes [Pearl, 2011]. This SCM in Figure 15, combined with the SCM from Fig 2b that connects other facial attributes to the image, provides us the *augmented* Attribute SCM that we use for the downstream task of CF generation. Once the Attribute SCM is learned, the rest of the steps in Algorithm 1 remain the same as before. That is, based on this augmented Attribute SCM, the counterfactuals are generated according to the Eqn 1. For example, to generate a CF changing *Young* from 1 to 0, the Prediction step involves changing the values of gray hair, receding hairline and eyeglasses based on the modified value of *Young*, according to the learned parameters (conditional probabilities) of the SCM.
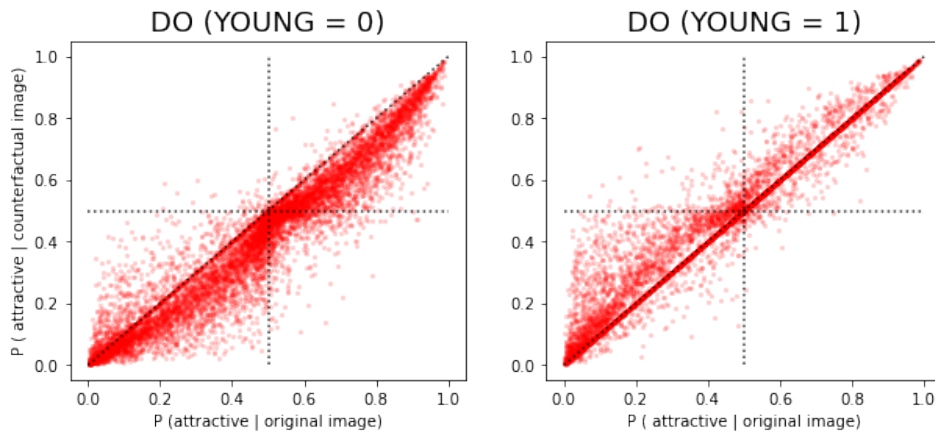


Figure 16: **Fairness Analysis for Complex SCM.** Counterfactual images that change *Young* from 1 to 0 (left), have a lower prediction score from the *attractiveness* classsifier, while counterfactual images that change *Young* from 0 to 1 (right), have a slightly higher prediction score.

*Fairness Analysis.* We conduct a fairness analysis similar as above, using the above attractiveness

classifier. We generate counterfactuals for *Young = 1* and *Young = 0* according to the causal graph in Fig 15. The analysis is showed in Fig 16; we observe that the classifier is evidently biased against *Young = 0* and biased towards *Young = 1*, when predicting *attractive=1*. We quantify this bias using Eqn 5. Using counterfactuals that change *Young* from 1 to 0, we get a negative bias of -0.136 and for changing *Young* from 0 to 1, we get a positive bias of 0.087 which are both substantial biases assuming a threshold of 5%. Therefore, given the causal graph from Figure 15, our method is able to generate counterfactuals for complex high-level features such as age, and use them to detect any potential biases in a machine learning classifier.