

Unifying Image Counterfactuals and Feature Attributions with Latent-Space Adversarial Attacks

Jeremy Goldwasser* Giles Hooker†

Abstract

Counterfactuals are a popular framework for interpreting machine learning predictions. These *what if* explanations are notoriously challenging to create for computer vision models: standard gradient-based methods are prone to produce adversarial examples, in which imperceptible modifications to image pixels provoke large changes in predictions. We introduce a new, easy-to-implement framework for counterfactual images that can flexibly adapt to contemporary advances in generative modeling. Our method, *Counterfactual Attacks*, resembles an adversarial attack on the representation of the image along a low-dimensional manifold. In addition, given an auxiliary dataset of image descriptors, we show how to accompany counterfactuals with feature attribution that quantify the changes between the original and counterfactual images. These importance scores can be aggregated into global counterfactual explanations that highlight the overall features driving model predictions. While this unification is possible for any counterfactual method, it has particular computational efficiency for ours. We demonstrate the efficacy of our approach with the MNIST and CelebA datasets.

1 Introduction

In broad terms, a counterfactual statement is a *what if*. It elucidates causal relationships by highlighting that if some factor had been different, then the outcome would have changed. This intuitive form of explanation was extended to machine learning (ML) by Wachter et al. (2018). Like SHAP and LIME, counterfactuals interpret a model’s prediction on an individual sample (Lundberg and Lee, 2017; Ribeiro et al., 2016). Their local explanations modify the input in a concise, realistic, actionable way that provokes a large change in the model output (Verma et al., 2024). Ideally, these alterations illuminate the salient reason(s) behind the prediction. Such interpretations may empower humans to learn from ML models, build trust in them, or identify their failure modes.

Most prior works produce counterfactuals by solving a regularized regression problem. The loss encourages the model’s output on the counterfactual to be near a pre-specified target, i.e. flipping the prediction. The regularization term penalizes examples that are too far from the original input. Conventional notions of proximity include low overall deviation (Wachter et al., 2018), sparse interventions (Dandl et al., 2020), and plausible achievability (Asemota and Hooker, 2024).

While tuning the regularization hyperparameter may be challenging or arbitrary, in general this is an effective approach for explaining tabular models. In fact, the counterfactual changes may even be interpreted as attribution scores (Mothilal et al., 2021). Features whose values changed a lot for the counterfactual are deemed more important.

On image data, however, regularized optimization faces unique challenges. Neural networks are prone to adversarial attacks on computer vision tasks (Goodfellow et al., 2015). Using gradient ascent methods like FGSM, potentially imperceptible changes can yield wildly different outputs. Therefore it may be possible to construct counterfactuals with no substantive alteration.

Making any sort of substantive edit to an image is challenging in its own right. Images are very high dimensional: 256x256 RGB images, for example, have roughly 200,000 pixels. Even without the presence of adversarial examples, blindly editing pixels will likely not correspond to any semantically meaningful change.

*Department of Statistics, University of California, Berkeley

†Department of Statistics and Data Science, University of Pennsylvania

Nor can these changes be induced by clever regularization. Sparsity penalties, for example, would render it impossible to learn counterfactuals with simple transformations in global features like lighting.

To address this, we propose new methodology to produce highly realistic image counterfactuals. Our approach, Counterfactual Attacks, runs gradient ascent on a latent representation of the image, in essence generating an adversarial example on the data manifold. Unlike previous methods, ours enables use of modern generative models like StyleGAN3 entirely off-the-shelf, and has no regularization hyperparameters to tune. Each counterfactual is the endpoint of a smooth trajectory starting with the original image; the intermediate images along this path may be visualized as well.

Furthermore, our method can be unified with feature attributions when auxiliary labels are present. That is, the content of the counterfactual is captured in a set of importance scores which explain what changed from the original image. This unification allows counterfactuals to be analyzed in an automated fashion, for example to attain global importance scores without manually annotating hundreds of images. This concise summary provides novel insights into the salient factors driving model behavior. We show how to produce important scores not only for our algorithm, but for any counterfactual methodology. To our knowledge, this is the first work quantifying image counterfactuals with feature importance scores.¹

2 Background

Formally, let x be a sample input to model f . The goal is to generate a similar, feasible input x' whose output $f(x')$ is near the desired target y' .

2.1 Standard Counterfactual Approaches

By definition, counterfactuals must balance multiple objectives. A loss objective such as squared error ensures $f(x')$ is close to y' . In addition, one or more proximity objectives ensures x' is close to x . For example, Wachter et al. (2018) proposed using a weighted L1 norm to measure the size of the change $x' - x$ that produces a desired change $f(x')$. Subsequently, Dandl et al. (2020) added a L_0 penalty for sparsity, as well as the distance to the nearest input for feasibility.

Some works simultaneously optimize the objectives via genetic algorithms, e.g. Dandl et al. (2020); Mothilal et al. (2020); Asemota and Hooker (2024). More common, however, is the original framework for ML counterfactuals, introduced by Wachter et al. (2018). In this setting, loss and regularizer objectives, \mathcal{L} and \mathcal{R} , are combined into a hyperparameter-weighted sum. The expression is solved with gradient-based optimization to produce a counterfactual:

$$x' = \operatorname{argmin}_{\tilde{x}} \mathcal{L}(f(\tilde{x}), y') + \lambda \mathcal{R}(\tilde{x}, x). \quad (1)$$

Another seminal paper, DiCE, proposes generating a set of counterfactuals, rather than just one (Mothilal et al., 2020). The rationale for doing so is the so-called Rashomon Effect – a given outcome may have multiple explanations, each valid in its own right but at odds with one another. In this case, different modifications to an input may cause the prediction to change in a desired direction.

2.2 Image Counterfactuals

A number of works have extended the notion of counterfactuals to images, leveraging powerful generative models. Several use Generative Adversarial Networks, or GANs (Goodfellow et al., 2014). Singla et al. (2020) trained a Conditional GAN that encodes an image, along with a desired amount to change its prediction on some model. The generator reconstructs the image, training to both fool the discriminator and shift the prediction accordingly. In contrast, the approach by Liu et al. (2019) ignores the predictor in question, training an unconditional DCGAN off-the-shelf (Radford et al., 2016). To generate counterfactuals, it solves an optimization problem akin to Equation (1) in the latent space. They use cross-entropy loss for classification error and a pixel-wise L1-norm for the penalty.

¹Code for this work is available at <https://github.com/jeremy-goldwasser/counterfactual-attacks>.

In recent years, diffusion models have been shown to outperform GANs for image generation (Dhariwal and Nichol, 2021). Analogous to GANs, counterfactual works propose training conditional diffusion models and targeted modifications from unconditional models. For the former, Augustin et al. (2022) trains a class-conditional diffusion model, sampling through the reverse process with altered classes to generate counterfactuals. They introduce various regularization techniques to control the stability of generation. For the case of diffusion models, Jeanneret et al. (2024) simply train a unconditional Denoising Diffusion Probabilistic Model (DDPM) on the data, then sample a counterfactual using guided diffusion (Dhariwal and Nichol, 2021). In sampling through the reverse path, guided diffusion moves in the direction of the negative gradient of a loss function. Here, the loss is the hyperparameter-weighted sum of classification and proximity terms, again akin to Equation (1).

While these methods are capable of generating high-quality counterfactuals, they suffer a number of drawbacks. Firstly, none of them guarantee the counterfactual will have the desired score. Rather, they either minimize regularized losses or direct sampling towards y' . In addition, they all have implementation challenges. GANs, for example, are notoriously difficult to train (Mescheder et al., 2018), and the methodology of Augustin et al. (2022) is quite complicated. Liu et al. (2019) and Jeanneret et al. (2024) circumvent these issues somewhat by leveraging off-the-shelf models; however, both require tuning the regularization hyperparameter, which lacks reasonable heuristics in this context. Moreover, defining the perceptive loss is unintuitive. Pixel-wise loss, suggested by Liu et al. (2019), restricts the space of possible counterfactuals, as in the aforementioned lighting example.

A separate line of work, albeit further from counterfactuals, takes a different approach. It discovers latent image features in generative models like age and gender, then uses them to edit to an image. InterFaceGAN (Shen et al., 2020) extracted the latent space representations of face images used to train PGGAN (Karras et al., 2018) and StyleGAN (Karras et al., 2019). With these relatively low-dimensional vectors, they fit linear models to features like age, glasses, and gender. These directions can then be used to edit images according to these features.

The InterFaceGAN method relies on having labels for semantically meaningful image concepts. While datasets like CelebA have such features, it is not true in general. To discover latent features, StyleEx (Lang et al., 2021) merely perturbs latent features one at a time in the StyleSpace of StyleGAN2 (Karras et al., 2020). This space appears after a few layers of transformations to the initial latent space in StyleGAN models. It has been shown to be highly disentangled, making StyleGAN a natural choice for image manipulation (Wu et al., 2021).

2.3 Dimensionality Reduction

Our methodology rests upon the observation that neural networks are capable of representing data in a low-dimensional manifold. A number of architectures are capable of providing such encodings. For example, data compression with autoencoders dates back decades (Lecun, 1987; Ballard, 1987; Bourlard and Kamp, 1988; Hinton and Zemel, 1993; Hinton and Salakhutdinov, 2006). VAEs (Kingma and Welling, 2014) greatly improved upon vanilla autoencoders by encoding data into a smooth, probabilistic latent space. This enables the decoder to be more effectively used as a generative model. While autoencoders generally struggle to generate natural images, recent years have seen a number of significant advancements, e.g. Adversarial Autoencoders (Makhzani et al., 2016), VQ-VAE (van den Oord et al., 2017), and NVAE (Vahdat and Kautz, 2021).

GANs also perform dimensionality reduction, and are more tailored for images. The latent code input to the generator has far lower dimension than the input. Similarly, the default implementation of the StyleGAN models produces a StyleSpace with dimension 512.

The latent space of diffusion models like DDPMs tends to have the same dimensionality as the images, due to their forward and reverse noising process (Ho et al., 2020). Therefore diffusion generally cannot be conceived as a dimensionality reduction technique. Latent Diffusion Models (Rombach et al., 2022) and Diffusion Autoencoders (Preechakul et al., 2022) are notable exceptions, and may be cast into our framework.

Having trained an autoencoder, the low-dimensional representation for an image is obtained simply by running it through the encoder. Similarly, for LDMs and DAEs, the encoded image is then passed through the forward noising process. The task of ‘‘GAN inversion,’’ however, may be more challenging (Xia et al., 2023). Some GAN architectures have an encoder that maps to the generator, in which case the

embedding may be obtained directly. Such an encoder may be trained jointly with the GAN itself, a technique first proposed independently by Donahue et al. (2017) and Dumoulin et al. (2017). Encoders may also be trained retrospectively, freezing the pre-trained weights of the generator and discriminator (Zhu et al., 2016; Richardson et al., 2021). Alternatively, encoder-free methods directly optimize the latent code for each image. For example, Abdal et al. (2019) learns a StyleSpace vector whose reconstruction minimizes a loss with pixel-wise and perceptual terms; the latter term seeks to reproduce the internal activations of a pre-trained VGG model (Johnson et al., 2016).

3 Counterfactual Attacks

Our approach first trains any generative model that yields a low-dimensional representation of an image. This may be done entirely off-the-shelf, enabling easy use of state-of-the-art models. Any model will suffice so long as it has a low-dimensional latent space that can accurately encode and manipulate an image.

Formally, let \mathcal{E} be an encoding procedure: the encoder of a VAE, for example, or a GAN inversion strategy like latent optimization. Recall the notation in Section 2, wherein a counterfactual on input x to model f aims to satisfy $f(x') = y'$. Encode x with the latent vector $\mathcal{E}(x) = z \in \mathbb{R}^d$. The generator \mathcal{G} maps z back to image space, such that

$$x^{\text{Recon}} = \mathcal{G}(\mathcal{E}(x)) \approx x.$$

Given a model predicting $f(x) = y$, the counterfactual is an image x' whose prediction is y' . Our approach modifies the latent representation z to z' in such a way that the reconstruction $x' = \mathcal{G}(z)$ satisfies $f(x') = y'$. To do so, it runs gradient updates on the latent score itself until the score is attained. Suppose without loss of generality that f outputs a scalar, and that $y' > f(x)$. With SGD and learning rate η , the update is

$$z \leftarrow z + \eta \nabla_z f(\mathcal{G}(z)). \tag{2}$$

Algorithm 1 Counterfactual Attacks

Require: Input x , predictor f , target $y' > f(x)$, encoder \mathcal{E} , generator \mathcal{G} , learning rate η

Ensure: Counterfactual x' such that $f(x') \geq y'$

- 1: $z \leftarrow \mathcal{E}(x)$
 - 2: **while** $f(\mathcal{G}(z)) \leq y'$ **do**
 - 3: $z \leftarrow z + \eta \nabla_z f(\mathcal{G}(z))$ {Latent gradient ascent}
 - 4: **end while**
 - 5: $x' \leftarrow \mathcal{G}(z)$ {Generate counterfactual image}
 - 6: **return** x'
-

Several adjustments may be made to the basic approach presented in Algorithm 1. Optimizers beyond SGD may be used, though the gradient computation is the same. If $y' < f(x)$, then the goal is to shrink the output of f , so gradient *descent* is run rather than ascent. Further consider the case in which f outputs a vector, as in multiclass classification. In this setting, modify Equation (2) to take only the coordinate of f corresponding to the class of interest. y' can be set to have a high value, i.e. reclassifying the input. Finally, it may be prudent to stop searching if no counterfactual is attained after a large number of iterations. This may be for computational convenience, or to exclude unrecognizably distant counterfactuals.

The natural image manifold of GANs allows for smooth interpolation between latent variables (Zhu et al., 2016). Our method is iterative, so intermediate steps may be visualized as well. This permits a smooth visual shift from original to counterfactual image.

4 Feature Attributions for Image Counterfactuals

We next introduce a framework to quantify image counterfactuals with feature attributions. These scores describe the important distinctions between the original image and its counterfactual. This approach applies for any counterfactual method, but has strong computational efficiency for our algorithm.

4.1 General Schema

In many image datasets, auxiliary feature labels are present. This is especially common for face datasets, tagging primarily categorical attributes like whether the person is smiling, has glasses, wears lipstick, etc. We use these labels to quantify the content of counterfactual explanations with feature importance scores.

Formally, let \mathcal{A} be the set of attributes for which labels are present, such as the aforementioned examples. Then, the images themselves X are accompanied by attribute labels $Y_a \forall a \in \mathcal{A}$. Our approach first trains a machine learning model $g_a : \mathcal{X} \rightarrow \mathcal{Y}_a$ to predict each attribute. These attribute predictors may input images in pixel space or a lower-dimensional representation.

We next consider the counterfactual explanation x' for image x . This counterfactual may be obtained by *any* method, not merely our Counterfactual Attack algorithm. Define each attribute’s score as the change in prediction, normalized by the range if the feature is numeric. That is, the importance scores are

$$\phi_a(x, x') = \frac{g_a(x') - g_a(x)}{\text{diam}(\mathcal{Y}_a)}. \quad (3)$$

The diameter of \mathcal{Y}_a is the breadth of possible values it can take. For categorical data, this is just 1. For numeric features, finite limits may be known outright, for example a restricted set of pose angles. If not, it may be taken as the difference between the highest and lowest values of $g_a(x)$ observed in a background dataset.

As defined, the scores in Equation (3) only apply to 1-dimensional outputs of $g_a(\cdot)$. This presents challenges for categorical attributes with more than two levels. To handle multi-class attribute models, one could analyze the levels in isolation, or aggregate them into a single attribute score, e.g. with the norm. An alternate approach would instead fit multiple binary classifiers.

The importance scores in Equation (3) may be aggregated to produce global scores. A simple formula would take the sample average; or, one could normalize each image’s scores by the sum across all attributes. Either way, it is imperative to account for the fact that different counterfactuals move the model’s predictions in different directions. In the binary setting, images predicted in the positive class are moved to the negative class, and vice versa. A straightforward approach would use absolute values:

$$\psi_a = \sum_{i=1}^n |\phi_a(x_i, x'_i)|. \quad (4)$$

This approach loses a notion of directionality. One cannot discern from the global scores whether the presence of an attribute tends to move the prediction in the positive or negative direction. To accomplish this, one can weight the scores $\phi_a(x_i, x'_i)$ by the direction of the counterfactual.

$$\psi_a = \sum_{i=1}^n \phi_a(x_i, x'_i) s_i, \text{ where } s_i = \begin{cases} 1, & \text{if } f(x_i) < 0.5 \\ -1, & \text{otherwise} \end{cases} \quad (5)$$

4.2 Considerations for Counterfactual Attacks

To ensure the importance scores are reliable, each attribute model $g_a(x)$ must generalize well. This can only be achieved by neural networks for sophisticated image prediction tasks. Of course, training $|\mathcal{A}|$ neural networks may come at a very high computational cost. Fortunately, when Counterfactual Attacks is used to generate counterfactuals, we can circumvent this computational burden almost entirely.

Shen et al. (2020) showed that fitting linear classifiers in StyleGAN’s latent space revealed meaningful feature directions. Having trained StyleGAN or a similar model for counterfactual explanations, it can be easily repurposed to fit these cheap attribute models. Moreover, far fewer samples may be necessary to fit them because the latent space is relatively low-dimensional.

To do so, embed each image in the latent space with $z_i = \mathcal{E}(x_i)$. Then, fit generalized linear models h_a predicting each y_a from z . The image-to-attribute scores are $g_a(x) = h_a(\mathcal{E}(x))$. The standard choices of GLM are logistic regression for categorical data, linear regression for numeric data, and Poisson regression for count data. GLM coefficients may also be regularized with Lasso and/or Ridge penalties.

5 Experiments

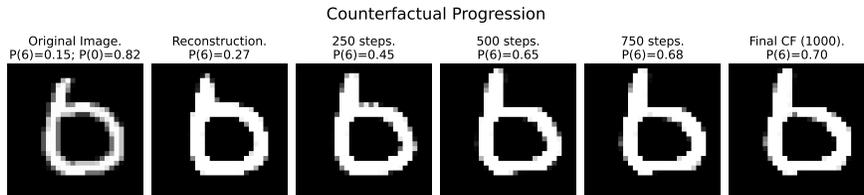
We demonstrate the utility of Counterfactual Attacks on the MNIST and CelebA datasets. We did not benchmark it against a wide range of other methods, because its aim is not to produce the highest-quality counterfactuals; rather, it is to be the easiest to implement, a more subjective notion. Moreover, our methodology to depict expansions with feature importance scores is the first of its kind.

However, Appendix C demonstrates the effectiveness of our feature attribution strategy on counterfactuals generated with InterFaceGAN (Shen et al., 2020). We also compared its counterfactuals to those from Counterfactual Attacks themselves. This is somewhat of an apples-to-oranges comparison, as InterFaceGAN uses a different classifier; nevertheless, we see that Counterfactual Attacks produces better counterfactual images on certain classes.

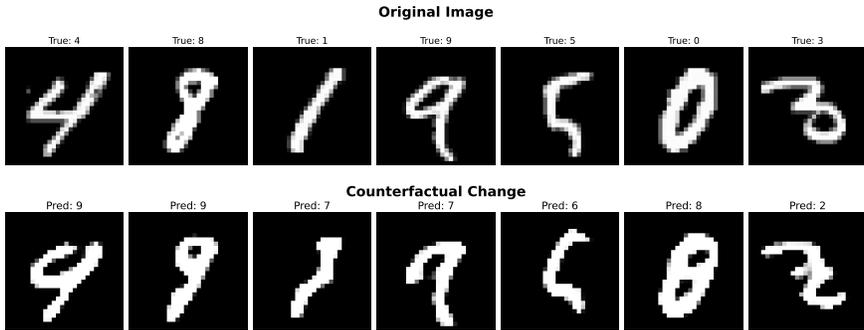
5.1 MNIST

To show its efficacy in simple settings, we ran Counterfactual Attacks on the MNIST dataset. These small grayscale images have dimensionality 28x28. We trained a straightforward neural network with 2 convolutional layers followed by 2 linear layers. Training stopped after only 1 epoch to ensure there were a reasonable number of misclassifications for failure analysis. The top-1 accuracy for both train and test sets was 97%.

For the manifold representation, we trained a variational autoencoder with a 64-dimensional latent space. Also a simple architecture, its encoder had 3 convolutional layers and 1 linear layer; the decoder reversed this process. To ensure its reconstructions were sparse, we added an L1 penalty to the loss. We also penalized the output of the Sobel operator in order to encourage smooth edges.



(a) Counterfactual path towards correct classification.



(b) Transforming various digits.

Figure 1: Running Counterfactual Attacks on MNIST dataset.

Figure 1 illustrates the capabilities of Counterfactual Attacks. Subfigure 1a examines a failure mode in which the original digit 6 was misclassified by the network as a 0. We visualize the intermediate steps, as the model’s prediction is gradually nudged towards the correct class. These images capture how the digit’s structure is altered, with each step increasing the probability of predicting 6 by elongating the protrusion on the top left. This progression highlights how the method not only corrects misclassifications but also provides a transparent sequence that sheds light on the salient transformations. Appendix A displays a wide array of similar failure analyses.



Figure 2: Counterfactuals for smiling classifier.

Subfigure 1b showcases a different use case: transforming correctly-classified digits to a new number. For instance, a 4 is systematically turned into a 9 by curling and extending its left prong; a slender 8 is turned into a 9 by closing its bottom gap; and a 0 is turned into an 8 by squeezing its center and adding a connective bar. These results demonstrate our method’s flexibility in generating realistic digit morphs along the data manifold, while remaining coherent at each iteration.

5.2 CelebA

We ran our methods on CelebA, a dataset of over 200,000 celebrity face images (Liu et al., 2015). The purpose of doing so was two-fold: To evaluate Counterfactual Attacks on a more complex dataset, as well as to demonstrate the capability of the feature attribution framework.

CelebA contains 40 binary attribute labels, describing the physical features, accessories, and expressions of each image. We trained separate 5-layer CNNs on four of the attributes: Young, Attractive, Male, and Smiling. The first three are more subjective than descriptive attributes like baldness. As a result, they pose useful targets for our importance-scoring methodology, to elucidate what facial features drive predictions.

For the generative model, we trained StyleGAN3 on CelebA (Karras et al., 2021). Works like InterFaceGAN (Shen et al., 2020) and StyleEx (Lang et al., 2021) demonstrated the capacity of the StyleGAN models for semantic face editing. StyleGAN3 improved upon its predecessors by removing the capacity for aliasing, which had tied details to absolute image coordinates.

After training, we represented each image in StyleSpace via the projection method of Abdal et al. (2019). Using `sklearn`, we fit logistic regression models with L2 penalties on 10,000 projected images. Counterfactuals were scored on the objective attributes, removing unhelpful or redundant features. Appendix B.1 provides a deeper account of our preprocessing and modeling choices.

Figure 2 shows various faces transformed by the smiling classifier. The top row shows the original image, followed by the altered counterfactual. In some cases, the counterfactual transformation is obvious as in the first three images. The latter three are much more subtle, yet nevertheless distinctive. Upon close inspection, it is clear that the 4th and 6th counterfactuals no longer smile, but the 5th does. Because smiling is a straightforward concept, we did not accompany the counterfactuals with importance scores.

Figure 3 displays counterfactuals with top-5 importance scores on the attractiveness, age, and gender classifiers. In all cases, Counterfactual Attacks produces realistical counterfactuals. Furthermore, the counterfactuals are all accurately described by the importance scores. Examining the top-scoring features that accompany each counterfactual, the older woman has whiter hairs and wears less makeup; the attractiveness classifier puts large emphasis on weight; and the transformation to male is in part defined by the lack of makeup, lipstick, and wavy hair. Appendix B.2 displays more examples of counterfactuals - both labeled as in Figure 3, as well as for the Smiling classifier.

Figure 4 showcases global counterfactual explanations for these three classifiers. It averages 100 local counterfactuals according to Equation (5). The top-5 highlighted features are highly intuitive. The model characterizes youth, attractiveness, and femininity in celebrity photos with makeup, lipstick, and arched

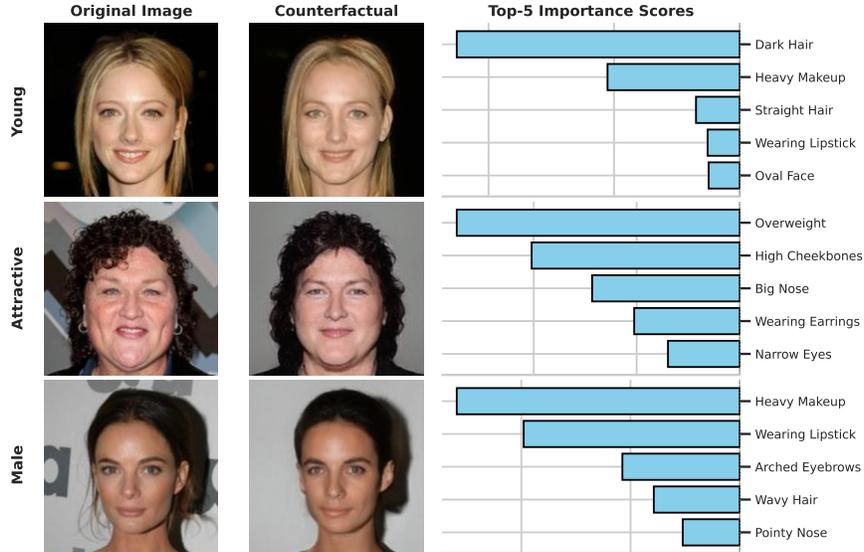


Figure 3: Counterfactual images with accompanying importance scores. Each row presents an individual counterfactual on a separate CNN classifier. All scores are negative, indicating the removal of their features.

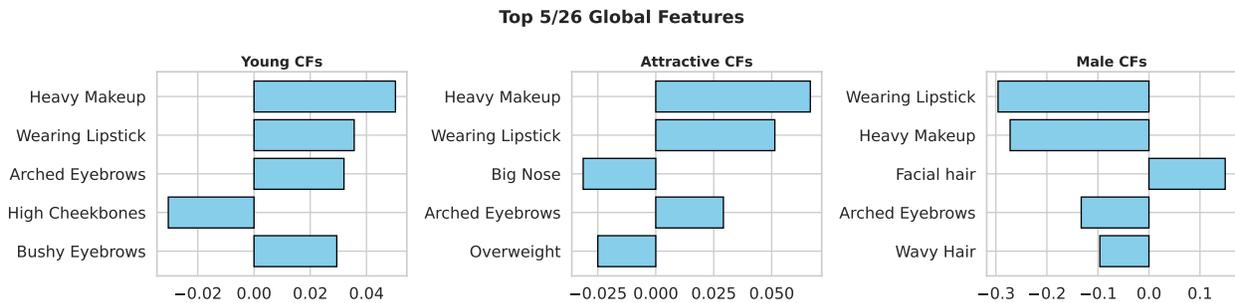


Figure 4: Global counterfactual explanations for three CelebA classifiers. The direction indicates whether the feature is added or removed.

eyebrows. The model also associates young celebrities with bushy eyebrows and lower cheekbones. Indeed, eyebrows thin with age, and structural features like cheekbones may become more pronounced. Attractiveness is further characterized by being less overweight and having a slender nose, and masculinity with facial hair.

While powerful, our counterfactual techniques have their limitations, discussed in depth in Appendix B.3. Firstly, the counterfactuals are only as good as the generative model they utilize. In our experiments with StyleGAN3, not every image projected neatly onto the latent space. This motivates the use of more sophisticated architectures, such as those with built-in encoders. Secondly, the labels must describe all possible transformations in order to be as useful as possible. While generally strong for youth and attractiveness, the provided CelebA attributes are limited for the gender classifier, as they do not include hair length. Wavy hair is the nearest proxy, hence its presence in Figures 3 and 4. Figures 7 and 8 visualize these two failure modes.

6 Discussion

Both counterfactuals and adversarial attacks make small changes to an image such that the prediction changes (Goodfellow et al., 2015; Freiesleben, 2022). We created a method that leverages this similarity, inspired by existing work in adversarial ML. Prior works on counterfactuals, and many adversarial methods, minimize a regularized loss function (Szegedy et al., 2014). Ours is the first counterfactual method, however, to mimic the approach of crafting examples via iterative gradient ascent steps (Kurakin et al., 2017). The primary difference

is that our “attacks” are performed within a low-dimensional latent space. This ensures that gradient steps walk along the natural image manifold – presuming that the latent space is a good representation thereof.

Our algorithm is more straightforward to implement than other methods for image counterfactuals. State-of-the-art generative models may be trained off-the-shelf, with no adjustment. Our experiments on CelebA used StyleGAN3, but other dimensionality reduction models may be used. Future work could explore the performance of Counterfactual Attacks across image autoencoders, GANs, LDMs, and DAEs. Moreover, because of the simplicity of gradient ascent, it is not necessary to tune regularization hyperparameters, as in prior work. Our approach may also satisfy the Rashomon Effect with a diverse set of counterfactuals, owing to the random noise inserted in the generation process.

Another novel contribution of our work is its unification of image counterfactuals with feature importance scores. These scores may serve as a useful objective complement or substitute to visual inspection. They can also be aggregated for global attributions with Equation (5) (van der Linden et al., 2019). The feature rankings from these global scores can be verified with the retrospective procedure from Goldwasser and Hooker (2025). In general, global image counterfactuals are an underexplored area of research. The only prior work we are aware of is Sobieski and Biecek (2024), which finds directions for global transformations in DAEs.

Global attributions facilitate an efficient, objective analysis of a model’s overall behavior by removing the need to manually inspect individual counterfactuals. This capability is particularly valuable for identifying and understanding patterns behind systematic errors, such as recurring types of misclassifications. By summarizing the underlying reasons for these failures, global attributions can generate actionable insights that guide model debugging and improvement. For instance, if the model consistently misclassifies images of people wearing glasses, this insight could prompt the collection of more training data featuring glasses. Similarly, if poor lighting conditions lead to degraded performance, data augmentation techniques could be employed to simulate such conditions during training, thereby improving the model’s robustness.

In general, counterfactual explanations provide more actionable insights than perturbation- or gradient-based methods like LIME or GradCAM (Ribeiro et al., 2016; Selvaraju et al., 2017). Rather than merely associating a score to pixels or patches, methods like Counterfactual Attacks provide examples of how to change a model’s prediction. Using the paper’s examples, a person could learn how to write zeros that looks less like sixes, or to present themselves more youthfully. This paper’s works go a step further by also providing feature importance scores. These highlight the importance of factors that visual inspection might otherwise overlook.

The importance scores we introduce describe the contents of the counterfactual post-hoc. Alternatively, one could investigate using them to generate the counterfactuals themselves. A simple way to do this would be to perform edits using the InterFaceGAN strategy until the desired prediction is reached. This may require composing edits on more than one feature. A more sophisticated approach would modify Counterfactual Attacks by projecting each gradient step onto the subspace defined by the attribute vectors. That way, the resulting counterfactual could be entirely accounted for by the labeled features.

References

- Abdal, R., Qin, Y., and Wonka, P. (2019). Image2stylegan: How to embed images into the stylegan latent space? In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
- Asemota, A. and Hooker, G. (2024). Longitudinal counterfactuals: Constraints and opportunities. *KDD '24 Workshop on Human-Interpretable AI*.
- Augustin, M., Boreiko, V., Croce, F., and Hein, M. (2022). Diffusion visual counterfactual explanations. In *Proceedings of the 36th Conference on Neural Information Processing Systems (NeurIPS 2022)*.
- Ballard, D. H. (1987). Modular learning in neural networks. In *Proceedings of the Sixth National Conference on Artificial Intelligence, Volume 1*, pages 279–284. AAAI.
- Bourlard, H. and Kamp, Y. (1988). Auto-association by multilayer perceptrons and singular value decomposition. *Biological Cybernetics*, 59:291–294.
- Dandl, S., Molnar, C., Binder, M., and Bischl, B. (2020). Multi-objective counterfactual explanations. In Bäck, T., Preuss, M., Deutz, A., Wang, H., Doerr, C., Emmerich, M., and Trautmann, H., editors, *Parallel Problem Solving from Nature – PPSN XVI*, pages 448–469, Cham. Springer International Publishing.
- Dhariwal, P. and Nichol, A. (2021). Diffusion models beat gans on image synthesis. In *Proceedings of the 35th Conference on Neural Information Processing Systems (NeurIPS 2021)*. OpenAI.
- Donahue, J., Krähenbühl, P., and Darrell, T. (2017). Adversarial feature learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Dumoulin, V., Belghazi, I., Poole, B., Mastropietro, O., Lamb, A., Arjovsky, M., and Courville, A. (2017). Adversarially learned inference. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Freiesleben, T. (2022). The intriguing relation between counterfactual explanations and adversarial examples. *Minds & Machines*, 32:77–109.
- Goldwasser, J. and Hooker, G. (2025). Statistical significance of feature importance rankings.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. In *Proceedings of the 28th International Conference on Neural Information Processing Systems (NIPS)*, pages 2672–2680. Volume 2.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. (2015). Explaining and harnessing adversarial examples.
- Hinton, G. E. and Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507.
- Hinton, G. E. and Zemel, R. S. (1993). Autoencoders, minimum description length and helmholtz free energy. In *Advances in Neural Information Processing Systems*.
- Ho, J., Jain, A., and Abbeel, P. (2020). Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Jeanneret, G., Simon, L., and Jurie, F. (2024). Diffusion models for counterfactual explanations. *Computer Vision and Image Understanding*, 249:104207.
- Johnson, J., Alahi, A., and Fei-Fei, L. (2016). Perceptual losses for real-time style transfer and super-resolution. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Karras, T., Aila, T., Laine, S., and Lehtinen, J. (2018). Progressive growing of gans for improved quality, stability, and variation. In *International Conference on Learning Representations (ICLR)*.

- Karras, T., Aittala, M., Laine, S., Härkönen, E., Hellsten, J., Lehtinen, J., and Aila, T. (2021). Alias-free generative adversarial networks. In *Proc. NeurIPS*.
- Karras, T., Laine, S., and Aila, T. (2019). A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., and Aila, T. (2020). Analyzing and improving the image quality of StyleGAN. In *Proc. CVPR*.
- Kingma, D. P. and Ba, L. J. (2015). Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, page 13.
- Kingma, D. P. and Welling, M. (2014). Auto-encoding variational bayes. In *Proceedings of the 2nd International Conference on Learning Representations (ICLR 2014)*.
- Kurakin, A., Goodfellow, I., and Bengio, S. (2017). Adversarial examples in the physical world.
- Lang, O., Gandelsman, Y., Yarom, M., Wald, Y., Elidan, G., Hassidim, A., Freeman, W. T., Isola, P., Globerson, A., Irani, M., and Mosseri, I. (2021). Explaining in style: Training a gan to explain a classifier in stylespace. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Lecun, Y. (1987). Modèles connexionnistes de l'apprentissage. *Intellectica Revue de l'Association pour la Recherche Cognitive*, 2(1).
- Liu, S., Kailkhura, B., Loveland, D., and Han, Y. (2019). Generative counterfactual introspection for explainable deep learning.
- Liu, Z., Luo, P., Wang, X., and Tang, X. (2015). Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.
- Lundberg, S. M. and Lee, S. (2017). A unified approach to interpreting model predictions. In Guyon, I., von Luxburg, U., Bengio, S., Wallach, H. M., Fergus, R., Vishwanathan, S. V. N., and Garnett, R., editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 4765–4774.
- Makhzani, A., Shlens, J., Jaitly, N., Goodfellow, I., and Frey, B. (2016). Adversarial autoencoders.
- Mescheder, L., Geiger, A., and Nowozin, S. (2018). Which training methods for gans do actually converge? In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, Stockholm, Sweden. PMLR.
- Mothilal, R. K., Mahajan, D., Tan, C., and Sharma, A. (2021). Towards unifying feature attribution and counterfactual explanations: Different means to the same end. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 652–663. ACM.
- Mothilal, R. K., Sharma, A., and Tan, C. (2020). Explaining machine learning classifiers through diverse counterfactual explanations. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 607–617, New York, NY, USA. Association for Computing Machinery.
- Preechakul, K., Chatthee, N., Wizadwongsa, S., and Suwajanakorn, S. (2022). Diffusion autoencoders: Toward a meaningful and decodable representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Radford, A., Metz, L., and Chintala, S. (2016). Unsupervised representation learning with deep convolutional generative adversarial networks.
- Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). "why should I trust you?": Explaining the predictions of any classifier. In Krishnapuram, B., Shah, M., Smola, A. J., Aggarwal, C. C., Shen, D., and Rastogi, R., editors, *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pages 1135–1144. ACM.

- Richardson, E., Alaluf, Y., Patashnik, O., Nitzan, Y., Azar, Y., Shapiro, S., and Cohen-Or, D. (2021). Encoding in style: a stylegan encoder for image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. (2017). Grad-cam: Visual explanations from deep networks via gradient-based localization. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 618–626.
- Shen, Y., Gu, J., Tang, X., and Zhou, B. (2020). Interpreting the latent space of gans for semantic face editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Singla, S., Pollack, B., Chen, J., and Batmanghelich, K. (2020). Explanation by progressive exaggeration. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Sobieski, B. and Biecek, P. (2024). Global counterfactual directions. In *Computer Vision – ECCV 2024*, pages 72–90.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I. J., and Fergus, R. (2014). Intriguing properties of neural networks. In Bengio, Y. and LeCun, Y., editors, *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*.
- Vahdat, A. and Kautz, J. (2021). Nvae: A deep hierarchical variational autoencoder. In *Advances in Neural Information Processing Systems*.
- van den Oord, A., Vinyals, O., and Kavukcuoglu, K. (2017). Neural discrete representation learning. In *Neural Information Processing Systems*.
- van der Linden, I., Haned, H., and Kanoulas, E. (2019). Global aggregations of local explanations for black box models.
- Verma, S., Boonsanong, V., Hoang, M., Hines, K., Dickerson, J., and Shah, C. (2024). Counterfactual explanations and algorithmic recourses for machine learning: A review. *ACM Computing Surveys*, 56(12):1–42. Published: 03 October 2024.
- Wachter, S., Mittelstadt, B., and Russell, C. (2018). Counterfactual explanations without opening the black box: Automated decisions and the gdpr. *Harvard Journal of Law & Technology*, 31(2). Spring 2018.
- Wu, Z., Lischinski, D., and Shechtman, E. (2021). Stylespace analysis: Disentangled controls for stylegan image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Xia, W., Zhang, Y., Yang, Y., Xue, J.-H., Zhou, B., and Yang, M.-H. (2023). Gan inversion: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(3):3121–3138.
- Zhu, J.-Y., Krahenbühl, P., Shechtman, E., and Efros, A. A. (2016). Generative visual manipulation on the natural image manifold. In *European Conference on Computer Vision (ECCV)*.

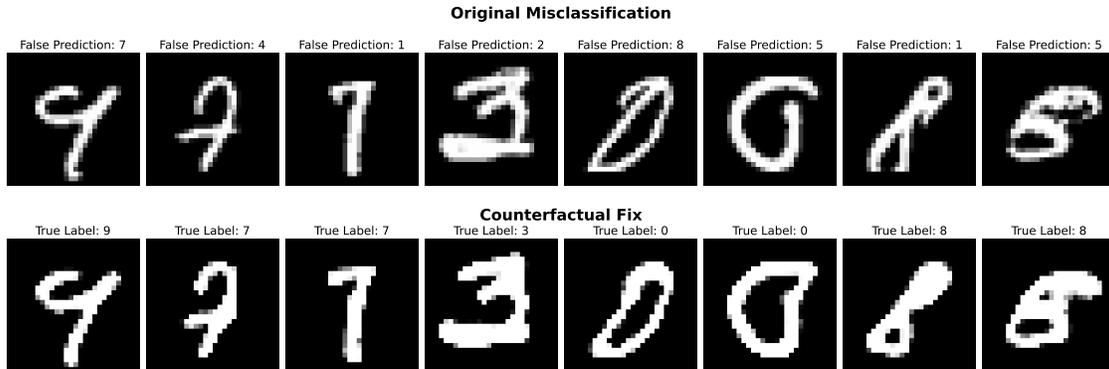


Figure 5: Counterfactuals that alter misclassified inputs to correct model predictions.

A MNIST

Figure 5 displays a number of misclassified inputs, and their accompanying counterfactuals that steer them to the correct classification. We used the Adam optimizer (Kingma and Ba, 2015), and a threshold of 0.5 for each image’s counterfactual. The counterfactual approach was the same as in section 5.1. The counterfactuals all make intuitive sense in explaining the model failures. For example, the 7 had been misclassified as a 1 because its tip dipped down; the 0 had been misclassified as an 8 because its top circled over; and the 8 had been misclassified as a 1 because its base was not connected.

B CelebA

B.1 Experimental Setup

To train StyleGAN3, we used CelebA centered images of dimensionality 256x256. Training took 3 days on 4 A4000 GPUs, fine-tuning from weights pre-trained on the FlickrFaces (FFHQ) dataset. We trained the model to be translation and rotation equivariant, with a batch size of 32.

The latent attribute predictors were fit using the default `sklearn` implementation of logistic regression. We combined several features into individual predictors:

- **Facial hair** encompassed the results of “sideburns,” “goatee,” “5 o’clock shadow,” “mustache,” and “no beard.” Naturally, the “no beard” feature being present was marked as the absence of a beard.
- **Dark hair** encompassed “black hair,” “brown hair,” and the reverse of “blond hair.”
- **Overweight** encompassed “chubby” and “double chin.”

We also ignored attributes regarding expressions, image characteristics and more abstract features. This encompassed “attractive,” “blurry,” “young,” “male,” “mouth slightly open,” and “smiling.” Doing so was necessary to hone in on concrete facial features.

Next, we trained neural networks for attractiveness, smiling, youth, and gender. The networks were all 5-layer CNNs, the first 3 of which were convolutional. All activation functions were ReLU, and dropout was applied before the final linear layer. We trained for 5 epochs on the full 200k-image dataset with an 80/20 train/test split. The resulting test accuracies were 80% for attractive, 91% for smiling, 90% for gender, and 87% for youth. In contrast, CelebA’s class imbalances were roughly 50% for attractive and smiling, 60% women, and 75% young.

To generate counterfactuals, we ran Counterfactual Attacks with the Adam optimizer, initializing the learning rate to 0.01. We targeted a prediction of 0.75 if the original image was in predicted the negative class, i.e. if $f(x) < 0.5$. Otherwise, positive predictions were transformed to a value of 0.25.

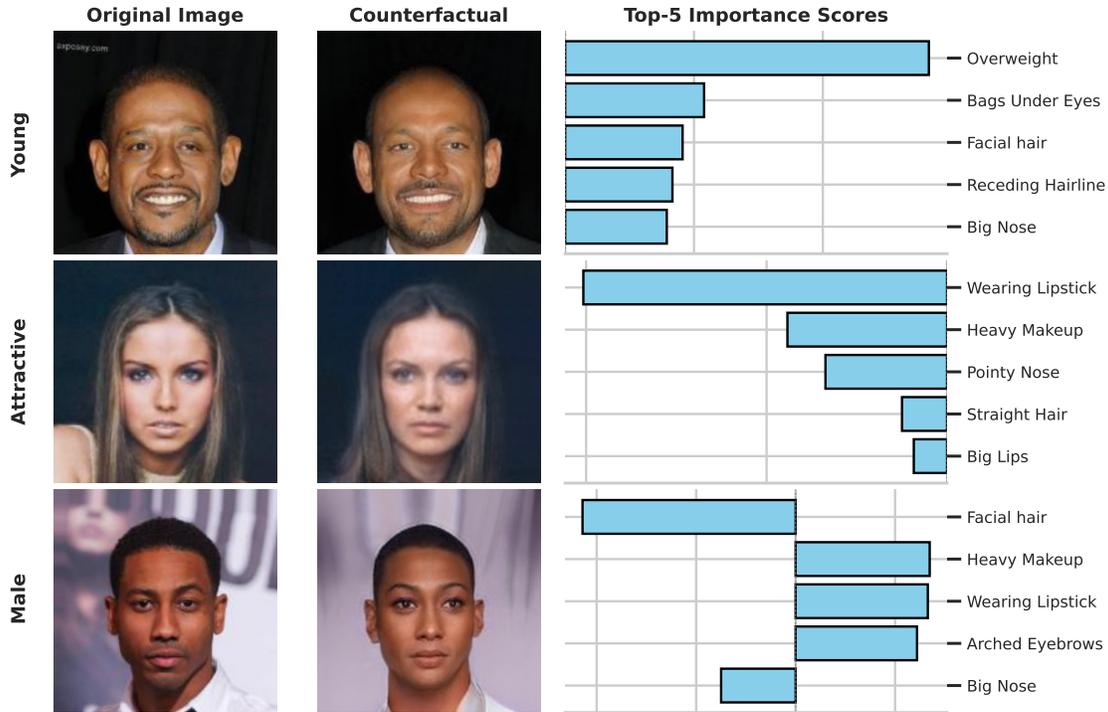


Figure 6: More examples of labeled counterfactuals.

B.2 Further Results

Figure 6 displays three more counterfactuals with importance scores. From top to bottom, the older man’s face is more rounded, and he clearly has a receding hairline. The woman scored as unattractive no longer wears lipstick or makeup, and her nose is less slim. The transformed woman has lost stubble, wears lipstick and makeup, and styles arched eyebrows.

B.3 Limitations

Here, we showcase two potential failure modes of our method. Figure 7 displays three images whose StyleSpace projections do not represent the original image particularly well. Respectively, the reconstructions alter the profile, eyes, and hair of the three faces. This could be remedied with a vast number of methods that better represent images in a low-dimensional space. Nevertheless, we present this to highlight that Counterfactual Attacks - or any counterfactual method, for that matter - will not perform well on images poorly represented in latent space.

More specific to our contributions, Figure 8 shows a counterfactual on the perceived gender classifier. The realistic counterfactual retains most facial characteristics, but notably truncates the person’s hair below the ears. However, this is not reflected in the shown feature importance scores, because hair length is not included in the attribute set. This demonstrates that the importance scores are only useful if they cover all reasonable explanations for the model’s prediction to flip.

C Experimentation with InterFaceGAN

We showcased the effectiveness of our feature attribution method to describe counterfactuals in general, not limited to Counterfactual Attacks. In particular, we applied it to counterfactual images generated with InterFaceGAN (Shen et al., 2020). Here, we describe InterFaceGAN in more detail than in Section 2, and show how it can be repurposed to produce counterfactuals.

Poor Reconstructions

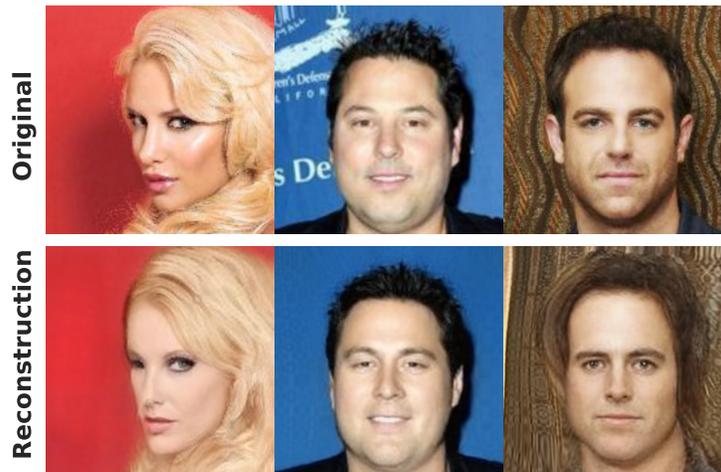


Figure 7: Poor reconstructions of various CelebA images.

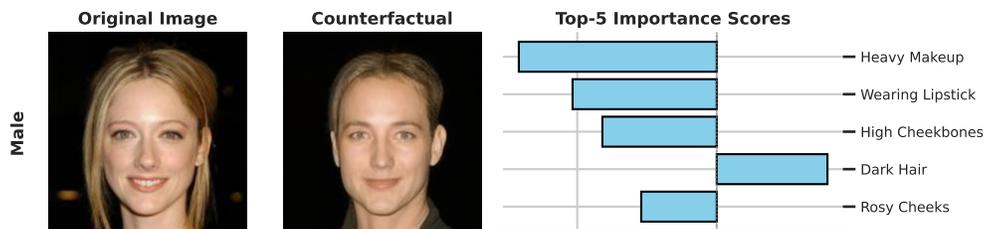


Figure 8: Gender counterfactual does not label the shortened hair length.

InterFaceGAN operates in the StyleSpace of a StyleGAN model. This space represents an image x with an encoding z . To edit a face image with regards to some feature y , this method fits a linear classifier f in StyleSpace. Assume without loss of generality that the label of interest y is binary; further assume it is predicted with logistic regression, outputting $\mathbb{P}(y = 1|z)$. In this context, $\hat{\beta}$ and \hat{c} are fit such that for logistic function σ ,

$$f(z) = \sigma(\langle z, \hat{\beta} \rangle + \hat{c}). \quad (6)$$

Presuming the model in Equation (6) is sufficiently accurate, $\hat{\beta}$ represents the *direction* of the feature in StyleSpace. InterFaceGAN modifies the presence of the feature in the original image by moving its latent representation in that direction. For some value of $\eta \in \mathbb{R}$, InterFaceGAN generates the image

$$\mathcal{G}(z + \eta\hat{\beta}).$$

When η is positive, the predicted probability that the feature is present increases. This aims to turn the feature *on* if it is not present. Conversely, the feature may be removed from an image by moving in the negative direction of $\hat{\beta}$.

How much should the image be moved along this axis? While InterFaceGAN does not employ the language of counterfactuals, it can easily be extended to this framework. Let p be the predicted probability for image z under the logistic regression model, and let p' be the desired prediction of its counterfactual. We rearrange Equation (6) to identify the value η' that yields this probability:

$$\begin{aligned} p' &= \sigma(\langle z + \eta'\hat{\beta}, \hat{\beta} \rangle + \hat{c}) \\ &= \sigma(\langle z, \hat{\beta} \rangle + \hat{c} + \eta'\|\hat{\beta}\|_2^2) \\ &= \sigma(\sigma^{-1}(p) + \eta'\|\hat{\beta}\|_2^2) \end{aligned}$$

Defining the inverse logistic function $\sigma^{-1}(y) = \log(\frac{y}{1-y})$,

$$\sigma^{-1}(p') = \sigma^{-1}(p) + \eta'\|\hat{\beta}\|_2^2.$$

Solving for η' reveals its value,

$$\eta' = \frac{\sigma^{-1}(p') - \sigma^{-1}(p)}{\|\hat{\beta}\|_2^2}. \quad (7)$$

On CelebA, we fit three latent logistic regression models for the Attractiveness, Youth, and Male attributes. Their Encoding a number of images, we identified values of η' that would flip these models' predictions. As with the neural network classifiers, our counterfactuals changed positive predictions to 0.25, and negative predictions to 0.75. Each case produced a counterfactual latent vector, $z' = z + \eta'\hat{\beta}$. Their corresponding images were simply regenerated with $x' = \mathcal{G}(z')$.

We employed our feature importance methodology to quantify the contents of these counterfactuals. To recap, our approach fits models g_a for each attribute, and takes the change in prediction as the score. In this context, we used the same attribute predictors as for our Counterfactual Attacks experiments (Appendix B.1). These were the logistic regression models that input images in latent space. (For our InterFaceGAN experiments, the latent space in question was the StyleSpace of StyleGAN3.)

Figure 9 displays labeled InterFaceGAN counterfactuals with the same images as in Figures 3 and 6. Upon visual inspection, many of these counterfactuals seem to successfully flip the attribute label in question. Furthermore, the importance scores unanimously provide accurate description of their contents. For example, facial hair is clearly removed in the second and sixth counterfactuals. On all of the others, the absence or presence of lipstick is apparent. This indicates that our proposed method is capable of describing counterfactuals in general, and is not restricted to Counterfactual Attacks.

It is worth noting that Counterfactual Attacks algorithm seems to generate better counterfactuals for Youth and Attractiveness than InterFaceGAN. This may be explained by the discrepancies in the classifiers they interpret. For these features, the neural networks have roughly 3% higher test accuracy. In fact, the attractiveness classifier gives roughly the same prediction for the original and counterfactual images in the third row.

The poorer performance of logistic regression could be due to Youth and Attractiveness being nonlinear features in StyleSpace. In fact, the non-linear direction learned by the Youth classifier seems to conflate older age with masculinity. Strikingly, the woman’s hair is truncated; more subtly, the man loses facial hair and adds a hint of lipstick. Operating in image space, the neural networks have more flexible modeling capacity, and thus can avoid these linearity issues.

Finally, we aggregate local scores to produce global importance scores for InterFaceGAN counterfactuals. This follows an identical approach as for Counterfactual Attacks, in Figure 4. On the same 100 images, we compute each latent-space counterfactual with InterFaceGAN, then take an average with Equation (5). Figure 10 visualizes the results. Intuitively, the model identifies lipstick, makeup, and arched eyebrows as characteristically feminine, and facial hair as masculine.

These features are all important for the youth classifier as well, per its unfortunate conflation of age and gender. They also make the top-5 for the attractiveness classifier. While lipstick, makeup, and arched eyebrows are all sensible, the inclusion of facial hair suggests that classifier has learned to conflate attractiveness and femininity. This could explain in part why the classifier does not perform as well as the neural network. It also associates smaller noses with attractiveness.

InterFaceGAN counterfactuals

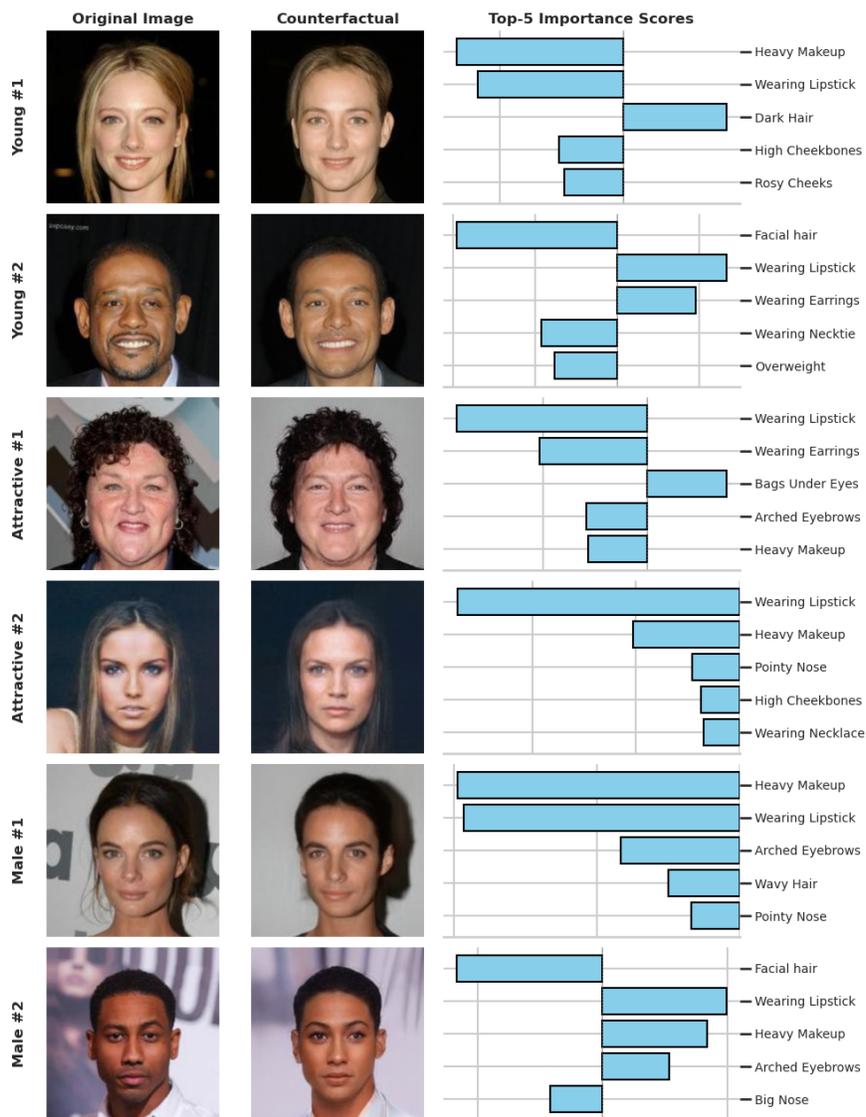


Figure 9: InterFaceGAN counterfactuals with accompanying importance scores from our methodology in Section 4. Each row presents an individual counterfactual on a separate logistic regression classifier, fit in StyleSpace.

Top 5/26 Global Features, InterFaceGAN Counterfactuals

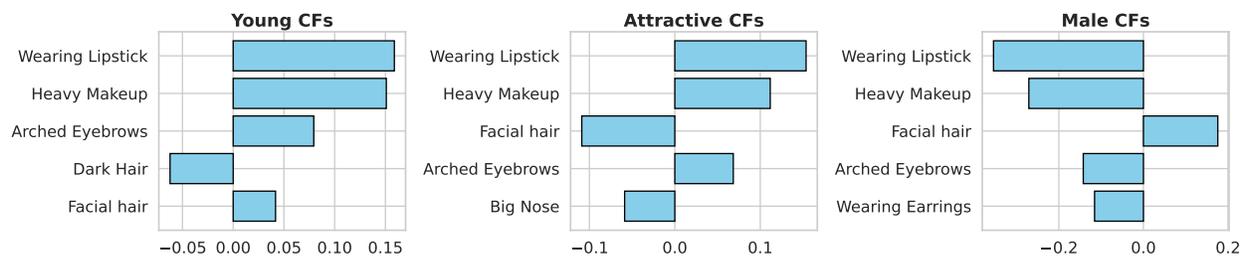


Figure 10: Top-5 global importance scores for InterFaceGAN counterfactual explanations. Three classifiers are logistic regression models fit in StyleSpace. The direction of each score indicates whether the feature is added or removed.