

Counterfactual Training: Teaching Models Plausible and Actionable Explanations

Author information scrubbed for double-blind reviewing

No Institute Given

Abstract. We propose a novel model objective that leverages counterfactual explanations to increase the explanatory capacity of trained models. Counterfactual explanations have emerged as a popular tool to explain predictions made by opaque machine learning models: they inform how factual inputs would need to change in order for a model to produce some desired output. Much existing research has focused on generating counterfactuals that are not only valid but also deemed acceptable with respect to the stakeholder requirements and plausible with respect to the underlying data. Recent work has shown that under this premise, the task of learning plausible explanations—ones that respect the conditional distribution of the target class—is effectively reassigned from the model itself to the post-hoc counterfactual explainer. Building on that work, we propose a novel objective function that employs counterfactuals (ad-hoc) during the training phase to minimize the divergence between the learned representations and the explanations. Through extensive experiments, we demonstrate that our proposed method facilitates training models that deliver inherently plausible explanations while maintaining high predictive performance.

Keywords: Counterfactual Explanations · Explainable AI · Representation Learning

1 Introduction

Today’s prominence of artificial intelligence (AI) has largely been driven by advances in **representation learning**: instead of relying on features and rules that are carefully hand-crafted by humans, modern machine learning (ML) models are tasked with learning these representations from scratch, guided by narrow objectives such as predictive accuracy [14]. Modern advances in computing have made it possible to provide such models with ever greater degrees of freedom to achieve that task, which has often led them to outperform traditionally more parsimonious models. Unfortunately, in doing so they also learn increasingly complex and highly sensitive representations that we can no longer easily interpret.

This trend towards complexity for the sake of performance has come under serious scrutiny in recent years. At the very cusp of the deep learning revolution, [34] showed that artificial neural networks (ANN) are sensitive to adversarial

examples: counterfactuals of model inputs that yield vastly different model predictions despite being “imperceptible” in that they are semantically indifferent from their factual counterparts. Despite partially effective mitigation strategies such as **adversarial training** [15], truly robust deep learning (DL) remains unattainable even for models that are considered shallow by today’s standards [23].

Part of the problem is that high degrees of freedom provide room for many solutions that are locally optimal with respect to narrow objectives [37]¹. Based purely on predictive performance, these solutions may seem to provide compelling explanations for the data, when in fact they are based on purely associative, semantically meaningless patterns. This poses two related challenges: firstly, it makes these models inherently opaque, since humans cannot simply interpret what type of explanation the complex learned representations correspond to; secondly, even if we could resolve the first challenge, it is not obvious how to mitigate models from learning representations that correspond to meaningless and implausible explanations.

The first challenge has attracted an abundance of research on **explainable AI** (XAI) which aims to develop tools to derive explanations from complex model representations. This can mitigate a scenario in which we deploy opaque models and blindly rely on their predictions. On countless occasions, this scenario has already occurred in practice and caused real harm to people who were affected adversely and often unfairly by automated decision-making systems (ADMS) involving opaque models [28]. Effective XAI tools can aide us in monitoring models and providing recourse to individuals to turn adverse outcomes (e.g. “loan application rejected”) into positive ones (“application accepted”). [36] propose **counterfactual explanations** as an effective approach to achieve this: they explain how factual inputs need to change in order for some fitted model to produce some desired output, typically involving minimal perturbations.

To our surprise, the second challenge has not yet attracted any consolidated research effort. Specifically, there has been no concerted effort towards improving model **explainability**, which we define here as the degree to which learned representations correspond to explanations that are interpretable and deemed **plausible** by humans (see Definition 1). Instead, the choice has typically been to improve the capacity of XAI tools to identify the subset explanations that are both plausible and valid for any given model, independent of whether the learned representations are also compatible with implausible explanations [3]. Fortunately, recent findings indicate that explainability can arise as byproduct of regularization techniques aimed at other objectives such as robustness, generalization and generative capacity [33].

Building on these findings, we introduce **counterfactual training**: a novel regularization technique geared explicitly towards aligning model representations with plausible explanations. Our contributions are as follows:

¹ For clarity: we follow standard ML convention in using “degrees of freedom” to refer to the number of parameters estimated from data.

- We discuss existing related work on improving models and consolidate it through the lens of counterfactual explanations (Section 2).
- We present our proposed methodological framework that leverages faithful counterfactual explanations during the training phase of models to achieve the explainability objective (Section 3).
- Through extensive experiments we demonstrate the counterfactual training improve model explainability while maintaining high predictive performance. We run ablation studies and grid searches to understand how the underlying model components and hyperparameters affect outcomes. (Section 4).

Despite limitations of our approach discussed in Section 5, we conclude that counterfactual training provides a practical framework for researchers and practitioners interested in making opaque models more trustworthy Section 6. We also believe that this work serves as an opportunity for XAI researchers to reevaluate the premise of improving XAI tools without improving models.

2 Related Literature

To the best of our knowledge, our proposed framework for counterfactual training represents the first attempt to use counterfactual explanations during training to improve model explainability. In high-level terms, we define model explainability as the extent to which valid explanations derived for an opaque model are also deemed plausible with respect to the underlying data and stakeholder requirements. To make this more concrete, we follow [4] in tieing the concept of explainability to the quality of counterfactual explanations that we can generate for a given model. The authors show that counterfactual explanations—understood here as minimal input perturbations that yield some desired model prediction—are generally more meaningful if the underlying model is more robust to adversarial examples. We can make intuitive sense of this finding when looking at adversarial training (AT) through the lens of representation learning with high degrees of freedom: by inducing models to “unlearn” representations that are susceptible to worst-case counterfactuals (i.e. adversarial examples), AT effectively removes some implausible explanations from the solution space.

2.1 Adversarial Examples are Counterfactual Explanations

This interpretation of the link between explainability through counterfactuals on one side, and robustness to adversarial examples on the other, is backed by empirical evidence. [32] demonstrate that using counterfactual images during classifier training improves model robustness. Similarly, [1] argue that counterfactuals represent potentially useful training data in machine learning, especially in supervised settings where inputs may be reasonably mapped to multiple outputs. They, too, demonstrate the augmenting the training data of image classifiers can improve generalization. [35] propose an approach using counterfactuals in training that does not rely on data augmentation: they argue that counterfactual pairs

typically already exist in training datasets. Specifically, their approach relies on, firstly, identifying similar input samples with different annotations and, secondly, ensuring that the gradient of the classifier aligns with the vector between pairs of counterfactual inputs using the cosine distance as a loss function.

In the natural language processing (NLP) domain, counterfactuals have similarly been used to improve models through data augmentation: [38], propose *Polyjuice*, a general-purpose counterfactual generator for language models. They demonstrate empirically that augmenting training data through *Polyjuice* counterfactuals improves robustness in a number of NLP tasks. [5] also use *Polyjuice* to augment NLP datasets through diverse counterfactuals and show that classifier robustness improves up to 20%. Finally, [26] introduce Counterfactual Adversarial Training (CAT), which also aims at improving generalization and robustness of language models. Specifically, they propose to proceed as follows: firstly, they identify training samples that are subject to high predictive uncertainty; secondly, they generate counterfactual explanations for those samples; and, finally, they fine-tune the given language model on the augmented dataset that includes the generated counterfactuals.

There have also been several attempts at formalizing the relationship between counterfactual explanations (CE) and adversarial examples (AE). Pointing to clear similarities in how CE and AE are generated, [13] makes the case for jointly studying the opaqueness and robustness problem in representation learning. Formally, AE can be seen as the subset of CE, for which misclassification is achieved [13]. Similarly, [29] show that CE and AE are equivalent under certain conditions and derive theoretical upper bounds on the distances between them.

Two recent works are closely related to ours in that they use counterfactuals during training with the explicit goal of affecting certain properties of post-hoc counterfactual explanations. Firstly, [31] propose a way to train models that are guaranteed to provide recourse for individuals to move from an adverse outcome to some positive target class with high probability. The approach proposed by [31] builds on adversarial training, where in this context susceptibility to targeted adversarial examples for the positive class is explicitly induced. The proposed method allows for imposing a set of actionability constraints ex-ante: for example, users can specify that certain features (e.g. *age*, *gender*, ...) are immutable. Secondly, [18] are the first to propose an end-to-end training pipeline that includes counterfactual explanations as part of the training procedure. In particular, they propose a specific network architecture that includes a predictor and CE generator network, where the parameters of the CE generator network are learnable. Counterfactuals are generated during each training iteration and fed back to the predictor network. In contrast to [18], we impose no restrictions on the neural network architecture at all.

2.2 Beyond Robustness

Improving the adversarial robustness of models is not the only path towards aligning representations with plausible explanations. In a work closely related to this one, [3] show that explainability can be improved through model averaging

and refined model objectives. The authors propose a way to generate counterfactuals that are maximally **faithful** to the model in that they are consistent with what the model has learned about the underlying data. Formally, they rely on tools from energy-based modelling to minimize the divergence between the distribution of counterfactuals and the conditional posterior over inputs learned by the model. Their proposed counterfactual explainer, *ECCCo*, yields plausible explanations if and only if the underlying model has learned representations that align with them. They find that both deep ensembles [24] and joint energy-based models (JEMs) [16] tend to do well in this regard.

Once again it helps to look at these findings through the lens of representation learning with high degrees of freedom. Deep ensembles are approximate Bayesian model averages, which are most called for when models are underspecified by the available data [37]. Averaging across solutions mitigates the aforementioned risk of relying on a single locally optimal representations that corresponds to semantically meaningless explanations for the data. Previous work by [33] similarly found that generating plausible (“interpretable”) counterfactual explanations is almost trivial for deep ensembles that have also undergone adversarial training. The case for JEMs is even clearer: they involve a hybrid objective that induces both high predictive performance and generative capacity [16]. This is closely related to the idea of aligning models with plausible explanations and has inspired our proposed counterfactual training objective, as we explain in Section 3.

3 Counterfactual Training

Counterfactual training combines ideas from adversarial training, energy-based modelling and counterfactuals explanations with the explicit objective of aligning representations with plausible explanations that comply with user requirements. In the context of CE, plausibility has broadly been defined as the degree to which counterfactuals comply with the underlying data generating process [30,17,3]. Plausibility is a necessary but insufficient condition for using CE to provide algorithmic recourse (AR) to individuals affected by opaque models in practice. This is because for recourse recommendations to be **actionable**, they need to not only result in plausible counterfactuals but also be attainable. A plausible CE for a rejected 20-year-old loan applicant, for example, might reveal that their application would have been accepted, if only they were 20 years older. Ignoring all other features, this complies with the definition of plausibility if 40-year-old individuals are in fact more credit-worthy on average than young adults. But of course this CE does not qualify for providing actionable recourse to the applicant since *age* is not a mutable feature. For our intents and purposes, counterfactual training aims at improving model explainability by aligning models with counterfactuals that meet both desiderata, plausibility and actionability. Formally, we define explainability as follows:

Definition 1 (Model Explainability).

Let $M_\theta : \mathcal{X} \mapsto \mathcal{Y}$ denote a supervised classification model that maps from the D -dimensional input space \mathcal{X} to representations $\phi(\mathbf{x}; \theta)$ and finally to the

K -dimensional output space \mathcal{Y} . Assume that for any given input-output pair $\{\mathbf{x}, \mathbf{y}\}_i$ there exists a counterfactual $\mathbf{x}' = \mathbf{x} + \Delta : \mathbf{M}_\theta(\mathbf{x}') = \mathbf{y}^+ \neq \mathbf{y} = \mathbf{M}_\theta(\mathbf{x})$ where $\arg \max_y \mathbf{y}^+ = y^+$ and y^+ denotes the index of the target class.

We say that \mathbf{M}_θ is **explainable** to the extent that faithfully generated counterfactuals are plausible (i.e. consistent with the data) and actionable. Formally, we define these properties as follows:

1. (Plausibility) $\int^A p(\mathbf{x}'|\mathbf{y}^+) d\mathbf{x} \rightarrow 1$ where A is some small region around \mathbf{x}' .
2. (Actionability) Permutations Δ are subject to actionability constraints.

We consider counterfactuals as faithful to the extent that they are consistent with what the model has learned about the input data. Let $p_\theta(\mathbf{x}|\mathbf{y}^+)$ denote the conditional posterior over inputs, then formally:

3. (Faithfulness) $\int^A p_\theta(\mathbf{x}'|\mathbf{y}^+) d\mathbf{x} \rightarrow 1$ where A is defined as above.

The definitions of faithfulness and plausibility in Definition 1 are the same as in [3], with adapted notation. Actionability constraints in Definition 1 vary and depend on the context in which \mathbf{M}_θ is deployed. In this work, we focus on domain and mutability constraints for individual features x_d for $d = 1, \dots, D$. We limit ourselves to classification tasks for reasons discussed in Section 5.

3.1 Our Proposed Objective

Let \mathbf{x}'_t for $t = 0, \dots, T$ denote a counterfactual explanation generated through gradient descent over T iterations as initially proposed by [36]. For our purposes, we let T vary and consider the counterfactual search as converged as soon as the predicted probability for the target class has reached a pre-determined threshold, $\tau: \mathcal{S}(\mathbf{M}_\theta(\mathbf{x}'))[y^+] \geq \tau^2$.

To train models with high explainability as defined in Definition 1, we propose to leverage counterfactuals in the following objective,

$$\min_{\theta} \text{yloss}(\mathbf{M}_\theta(\mathbf{x}), \mathbf{y}) + \lambda_{\text{div}} \text{div}(\mathbf{x}, \mathbf{x}'_T, y; \theta) + \lambda_{\text{adv}} \text{advloss}(\mathbf{M}_\theta(\mathbf{x}'_{\leq T}), \mathbf{y}) \quad (1)$$

where $\text{yloss}(\cdot)$ is any conventional classification loss that induces discriminative performance (e.g. cross-entropy). The two additional components in Equation 1 are explained in more detail below. For now, they can be sufficiently described as inducing explainability directly and indirectly by penalizing: 1) the contrastive divergence, $\text{div}(\cdot)$, between counterfactuals \mathbf{x}'_T and observed samples \mathbf{x} and, 2) the adversarial loss, $\text{advloss}(\cdot)$, with respect to nascent counterfactuals $\mathbf{x}'_{\leq T}$. The tradeoff between the different components can be governed by adjusting the strengths of the penalties λ_{div} and λ_{adv} .

² For detailed background information on gradient-based counterfactual search and convergence see Section B.1.

Directly Inducing Explainability through Contrastive Divergence [16] observe that any classifier can be re-interpreted as a joint energy-based model (JEM) that learns to discriminate output classes conditional on inputs and generate inputs. They show that JEMs can be trained to perform well at both tasks by directly maximizing the joint log-likelihood factorized as $\log p_\theta(\mathbf{x}, \mathbf{y}) = \log p_\theta(\mathbf{y}|\mathbf{x}) + \log p_\theta(\mathbf{x})$. The first factor can be optimized using conventional cross-entropy as in Equation 1. To optimize $\log p_\theta(\mathbf{x})$ [16] minimize the contrastive divergence between samples drawn from $p_\theta(\mathbf{x})$ and training observations, i.e. samples from $p(\mathbf{x})$.

A key empirical finding in [3] was that JEMs tend to do well with respect to the plausibility objective in Definition 1. If we consider samples drawn from $p_\theta(\mathbf{x})$ as counterfactuals, this is an expected finding, because the JEM objective effectively minimizes the divergence between the conditional posterior and $p(\mathbf{x}|\mathbf{y}^+)$. To generate samples, [16] rely on Stochastic Gradient Langevin Dynamics (SGLD) using an uninformative prior for initialization. This is where we depart from their methodology: instead of generating samples through SGLD, we propose using counterfactual explainers to generate counterfactuals for observed training samples. Specifically, we have

$$\text{div}(\mathbf{x}, \mathbf{x}'_T, y; \theta) = \mathcal{E}_\theta(\mathbf{x}, y) - \mathcal{E}_\theta(\mathbf{x}'_T, y) \quad (2)$$

where $\mathcal{E}_\theta(\cdot)$ denotes the energy function. In particular, we set $\mathcal{E}_\theta(\mathbf{x}, \mathbf{y}) = -\mathbf{M}_\theta(\mathbf{x})[y^+]$ where y^+ denotes the index of the target class. We generate samples \mathbf{x}'_T by first randomly sampling the target class $y^+ \sim p(y)$ and then generating a counterfactual explanation for that target over T iterations using a gradient-based counterfactual generator. This is similar to how conditional sampling is used to draw from $p_\theta(\mathbf{x})$ in [16].

Intuitively, the gradient of Equation 2 decreases the energy of observed training samples (positive samples) while at same time increasing the energy of counterfactuals (negative samples) [12]. As the generated counterfactuals get more plausible (Definition 1) over the course of training, these two opposing effects gradually balance each out [25].

The departure from SGLD allows us to tap into the vast repertoire of explainers that have been proposed in the literature to meet different desiderata. Typically, these methods facilitate the imposition of domain and mutability constraints, for example. In principle, any existing approach for generating counterfactual explanations is viable, so long as it does not violate the faithfulness condition. Like JEMs [27], counterfactual training can be considered as a form of contrastive representation learning.

Indirectly Inducing Explainability through Adversarial Robustness
Based on our analysis in Section 2, counterfactuals \mathbf{x}' can be repurposed as additional training samples [26,5] or adversarial examples [13,29]. This leaves some flexibility with respect to the exact choice for $\text{advloss}(\cdot)$ in Equation 1. An intuitive functional form to use, though likely not the only reasonable choice, is inspired by adversarial training:

$$\begin{aligned} \text{advloss}(\mathbf{M}_\theta(\mathbf{x}'_{t \leq T}), \mathbf{y}; \varepsilon) &= \text{yloss}(\mathbf{M}_\theta(\mathbf{x}'_{t_\varepsilon}), \mathbf{y}) \\ t_\varepsilon &= \max_t \{t : \|\Delta_t\|_\infty < \varepsilon\} \end{aligned} \quad (3)$$

Under this choice, we consider nascent counterfactuals $\mathbf{x}'_{t \leq T}$ as adversarial examples as long as the magnitude of the perturbation to any individual feature is at most ε . This is closely aligned with [34], who define an adversarial attack as an “imperceptible non-random perturbation”. Thus, we choose to work with a different distinction between CE and AE than [13], who considers misclassification as the key distinguishing feature of AE. One of the key observations in this work is that we can leverage counterfactual explanations during training and get adversarial examples, essentially for free.

3.2 Encoding Actionability Constraints

Many existing counterfactual explainers support domain and mutability constraints out-of-the-box. In fact, both types of constraints can be implemented for any counterfactual explainer that relies on gradient descent in the feature space for optimization [2]. In this context, domain constraints can be imposed by simply projecting counterfactuals back to the specified domain, if the previous gradient step resulted in updated feature values that were out-of-domain. Mutability constraints can similarly be enforced by setting partial derivatives to zero to ensure that features are only mutated in the allowed direction, if at all.

Since actionability constraints are binding at test time, we should also impose them when generating \mathbf{x}' during each training iteration to align model representations with user requirements. Through their effect on \mathbf{x}' , both types of constraints influence model outcomes through Equation 2. Here it is crucial that we avoid penalizing implausibility that arises due to mutability constraints. For any mutability-constrained feature d this can be achieved by enforcing $\mathbf{x}[d] - \mathbf{x}'[d] := 0$ whenever perturbing $\mathbf{x}'[d]$ in the direction of $\mathbf{x}[d]$ would violate mutability constraints. Specifically, we set $\mathbf{x}[d] := \mathbf{x}'[d]$ if

1. Feature d is strictly immutable in practice.
2. We have $\mathbf{x}[d] > \mathbf{x}'[d]$ but feature d can only be decreased in practice.
3. We have $\mathbf{x}[d] < \mathbf{x}'[d]$ but feature d can only be increased in practice.

From a Bayesian perspective, setting $\mathbf{x}[d] := \mathbf{x}'[d]$ can be understood as assuming a point mass prior for $p(\mathbf{x})$ with respect to feature d . Intuitively, we think of this simply in terms ignoring implausibility costs with respect to immutable features, which effectively forces the model to instead seek plausibility with respect to the remaining features. This in turn results in lower overall sensitivity to immutable features, which we demonstrate empirically for different classifiers in Section 4. Under certain conditions, this results holds theoretically[For the proof, see the supplementary appendix.]:

Proposition 1 (Protecting Immutable Features).

Let $f_\theta(\mathbf{x}) = \mathcal{S}(\mathbf{M}_\theta(\mathbf{x})) = \mathcal{S}(\Theta\mathbf{x})$ denote a linear classifier with softmax activation \mathcal{S} (i.e. multinomial logistic regression) where $y \in \{1, \dots, K\} = \mathcal{K}$ and $\mathbf{x} \in \mathbb{R}^D$. If we assume multivariate Gaussian class densities with common diagonal covariance matrix $\Sigma_k = \Sigma$ for all $k \in \mathcal{K}$, then protecting an immutable feature from the contrastive divergence penalty (Equation 2) will result in lower classifier sensitivity to that feature relative to the remaining features, provided that at least one of those is mutable and discriminative.

It is worth highlighting that Proposition 1 assumes independence of features. This raises a valid concern about the effect of protecting immutable features in the presence of proxy features that remain unprotected. We discuss this limitation in Section 5.

3.3 Illustration

To better convey the intuition underlying our proposed method, we illustrate different model outcomes in Example 1.

Example 1 (Prediction of Consumer Credit Default).

Suppose we are interested in predicting the likelihood that loan applicants default on their credit. We have access to historical data on previous loan takers comprised of a binary outcome variable ($y \in \{1 = \text{default}, 2 = \text{no default}\}$) two input features: 1) the subjects' *age*, which we define as immutable, and 2) the subjects' existing level of *debt*, which we define as mutable.

We have simulated this scenario using synthetic data with independent features and Gaussian class-conditional densities in Figure 1. The four panels in Figure 1 show the outcomes for different training procedures using the same model architecture each time (a linear classifier). In each case, we show the linear decision boundary (green) and the training data colored according to their ground-truth label: orange points belong to the target class, $y^+ = 2$, blue points belong to the non-target class, $y^- = 1$. Stars indicate counterfactuals in the target class generated at test time using generic gradient descent until convergence.

In panel (a), we have trained our model conventionally, and we do not impose mutability constraints at test time. The generated counterfactuals are all valid, but not plausible: they are clearly distinguishable from the ground-truth data. In panel (b), we have trained our model with counterfactual training, once again not imposing mutability constraints at test time. We observe that the counterfactuals are clearly plausible, therefore meeting the first objective of Definition 1.

In panel (c), we have used conventional training again, this time imposing the mutability constraint on *age* at test time. Counterfactuals are valid but involve some substantial reductions in *debt* for some individuals, in particular very young applicants. By comparison, counterfactual paths are shorter on average in panel (d), where we have used counterfactual training and protected immutable features as described in Section 3.2. In particular, we observe that due to the classifier's lower sensitivity to *age*, recourse recommendations with respect to *debt* are much more homogenous, in that they do not disproportionately punish

younger individuals. The counterfactuals are also plausible with respect to the mutable feature. Thus, we consider the model in panel (d) as the most explainable according to Definition 1.

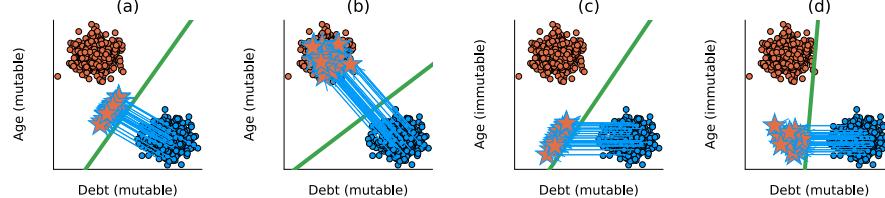


Fig. 1: Visual illustration of how counterfactual training improves explainability. See Example 1 for details.

4 Experiments

In this section, we present experiments that we have conducted in order to answer the following research questions:

Corollary 1 (Plausibility).

Does our proposed counterfactual training objective (Equation 1) induce models to learn plausible explanations?

Corollary 2 (Actionability).

Does our proposed counterfactual training objective (Equation 1) yield more favorable algorithmic recourse outcomes in the presence of actionability constraints?

Beyond this, we are also interested in understanding how robust our answers to RQ 1 and RQ 2 are:

Corollary 3 (Hyperparameters).

What are the effects of different hyperparameter choices with respect to Equation 1?

4.1 Experimental Setup

4.2 Experimental Results

5 Discussion

1. Limited to classification models.
2. Proxy attributes of immutable features.
3. Increased training time.
4. Training instabilities
5. Fairness and caveats (aware it's not a classical approach in this context, but there is a clear link).

6 Conclusion

A Notation

- y^+ : The target class and also the index of the target class.
- y^- : The non-target class and also the index of non-the target class.
- \mathbf{y}^+ : The one-hot encoded output vector for the target class.
- θ : Model parameters (unspecified).
- Θ : Matrix of parameters.

B Technical Details of Our Approach

B.1 Generating Counterfactuals through Gradient Descent

In this section, we provide some background on gradient-based counterfactual generators (Section B.1) and discuss how we define convergence in this context (Section B.1).

Background Gradient-based counterfactual search was originally proposed by [36]. It generally solves the following unconstrained objective,

$$\min_{\mathbf{z}' \in \mathcal{Z}^L} \{ \text{yloss}(\mathbf{M}_\theta(g(\mathbf{z}')), \mathbf{y}^+) + \lambda \text{cost}(g(\mathbf{z}')) \}$$

where $g : \mathcal{Z} \mapsto \mathcal{X}$ is an invertible function that maps from the L -dimensional counterfactual state space to the feature space and $\text{cost}(\cdot)$ denotes one or more penalties that are used to induce certain properties of the counterfactual outcome. As above, \mathbf{y}^+ denotes the target output and $\mathbf{M}_\theta(\mathbf{x})$ returns the logit predictions of the underlying classifier for $\mathbf{x} = g(\mathbf{z})$.

For all generators used in this work we use standard logit crossentropy loss for $\text{yloss}(\cdot)$. All generators also penalize the distance (ℓ_1 -norm) of counterfactuals from their original factual state. For *Generic* and *ECCo*, we have $\mathcal{Z} := \mathcal{X}$ and $g(\mathbf{z}) = g(\mathbf{z})^{-1} = \mathbf{z}$, that is counterfactual are searched directly in the feature space. Conversely, *REVISE* traverses the latent space of a variational autoencoder (VAE) fitted to the training data, where $g(\cdot)$ corresponds to the decoder [22]. In addition to the distance penalty, *ECCo* uses an additional penalty component that regularizes the energy associated with the counterfactual, \mathbf{x}' [3].

Convergence An important consideration when generating counterfactual explanations using gradient-based methods is how to define convergence. Two common choices are to 1) perform gradient descent over a fixed number of iterations T , or 2) conclude the search as soon as the predicted probability for the target class has reached a pre-determined threshold, τ : $\mathcal{S}(\mathbf{M}_\theta(\mathbf{x}'))[y^+] \geq \tau$. We prefer the latter for our purposes, because it explicitly defines convergence in terms of the black-box model, $\mathbf{M}(\mathbf{x})$.

Defining convergence in this way allows for a more intuitive interpretation of the resulting counterfactual outcomes than with fixed T . Specifically, it allows us to think of counterfactuals as explaining ‘high-confidence’ predictions by the model for the target class y^+ . Depending on the context and application, different choices of τ can be considered as representing ‘high-confidence’ predictions.

B.2 Protecting Mutability Constraints with Linear Classifiers

In Section 3.2 we explain that to avoid penalizing implausibility that arises due to mutability constraints, we impose a point mass prior on $p(\mathbf{x})$ for the corresponding feature. We argue in Section 3.2 that this approach induces models to be less sensitive to immutable features and demonstrate this empirically in Section 4. Below we derive the analytical results in Proposition 1.

Proof. Let d_{mtbl} and d_{immtbl} denote some mutable and immutable feature, respectively. Suppose that $\mu_{y^-, d_{\text{immtbl}}} < \mu_{y^+, d_{\text{immtbl}}}$ and $\mu_{y^-, d_{\text{mtbl}}} > \mu_{y^+, d_{\text{mtbl}}}$, where $\mu_{k,d}$ denotes the conditional sample mean of feature d in class k . In words, we assume that the immutable feature tends to take lower values for samples in the non-target class y^- than in the target class y^+ . We assume the opposite to hold for the mutable feature.

Assuming multivariate Gaussian class densities with common diagonal covariance matrix $\Sigma_k = \Sigma$ for all $k \in \mathcal{K}$, we have for the log likelihood ratio between any two classes $k, m \in \mathcal{K}$ [19]:

$$\log \frac{p(k|\mathbf{x})}{p(m|\mathbf{x})} = \mathbf{x}^\top \Sigma^{-1} (\mu_k - \mu_m) + \text{const} \quad (4)$$

By independence of x_1, \dots, x_D , the full log-likelihood ratio decomposes into:

$$\log \frac{p(k|\mathbf{x})}{p(m|\mathbf{x})} = \sum_{d=1}^D \frac{\mu_{k,d} - \mu_{m,d}}{\sigma_d^2} x_d + \text{const} \quad (5)$$

By the properties of our classifier (*multinomial logistic regression*), we have:

$$\log \frac{p(k|\mathbf{x})}{p(m|\mathbf{x})} = \sum_{d=1}^D (\theta_{k,d} - \theta_{m,d}) x_d + \text{const} \quad (6)$$

where $\theta_{k,d} = \Theta[k, d]$ denotes the coefficient on feature d for class k .

Based on Equation 5 and Equation 6 we can identify that $(\mu_{k,d} - \mu_{m,d}) \propto (\theta_{k,d} - \theta_{m,d})$ under the assumptions we made above. Hence, we have that $(\theta_{y^-, d_{\text{immtbl}}} - \theta_{y^+, d_{\text{immtbl}}}) < 0$ and $(\theta_{y^-, d_{\text{mtbl}}} - \theta_{y^+, d_{\text{mtbl}}}) > 0$

Let \mathbf{x}' denote some randomly chosen individual from class y^- and let $y^+ \sim p(y)$ denote the randomly chosen target class. Then the partial derivative of the contrastive divergence penalty Equation 2 with respect to coefficient $\theta_{y^+, d}$ is equal to

$$\frac{\partial}{\partial \theta_{y^+, d}} (\text{div}(\mathbf{x}, \mathbf{x}', \mathbf{y}; \theta)) = \frac{\partial}{\partial \theta_{y^+, d}} ((-\mathbf{M}_\theta(\mathbf{x})[y^+]) - (-\mathbf{M}_\theta(\mathbf{x}')[y^+])) = x'_d - x_d \quad (7)$$

and equal to zero everywhere else.

Since $(\mu_{y^-, d_{\text{immtbl}}} < \mu_{y^+, d_{\text{immtbl}}})$ we are more likely to have $(x'_{d_{\text{immtbl}}} - x_{d_{\text{immtbl}}}) < 0$ than vice versa at initialization. Similarly, we are more likely to have $(x'_{d_{\text{mtbl}}} - x_{d_{\text{mtbl}}}) > 0$ since $(\mu_{y^-, d_{\text{mtbl}}} > \mu_{y^+, d_{\text{mtbl}}})$.

This implies that if we do not protect feature d_{immtbl} , the contrastive divergence penalty will decrease $\theta_{y^-, d_{\text{immtbl}}}$ thereby exacerbating the existing effect $(\theta_{y^-, d_{\text{immtbl}}} - \theta_{y^+, d_{\text{immtbl}}}) < 0$. In words, not protecting the immutable feature would have the undesirable effect of making the classifier more sensitive to this feature, in that it would be more likely to predict class y^- as opposed to y^+ for lower values of d_{immtbl} .

By the same rationale, the contrastive divergence penalty can generally be expected to increase $\theta_{y^-, d_{\text{mtbl}}}$ exacerbating $(\theta_{y^-, d_{\text{mtbl}}} - \theta_{y^+, d_{\text{mtbl}}}) > 0$. In words, this has the effect of making the classifier more sensitive to the mutable feature, in that it would be more likely to predict class y^- as opposed to y^+ for higher values of d_{mtbl} .

Thus, our proposed approach of protecting feature d_{immtbl} has the net affect of decreasing the classifier’s sensitivity to the immutable feature relative to the mutable feature (i.e. no change in sensitivity for d_{immtbl} relative to increased sensitivity for d_{mtbl}).

B.3 Domain Constraints

We apply domain constraints on counterfactuals during training and evaluation. There are at least two good reasons for doing so. Firstly, within the context of explainability and algorithmic recourse, real-world attributes are often domain constrained: the *age* feature, for example, is lower bounded by zero and upper bounded by the maximum human lifespan. Secondly, domain constraints help mitigate training instabilities commonly associated with energy-based modelling [16,3].

For our image datasets, features are pixel values and hence the domain is constrained by the lower and upper bound of values that pixels can take depending on how they are scaled (in our case $[-1, 1]$). For all other features d in our synthetic and tabular datasets, we automatically infer domain constraints $[x_d^{\text{LB}}, x_d^{\text{UB}}]$ as follows,

$$\begin{aligned} x_d^{\text{LB}} &= \arg \min_{x_d} \{\mu_d - n_{\sigma_d} \sigma_d, \arg \min_{x_d} x_d\} \\ x_d^{\text{UB}} &= \arg \max_{x_d} \{\mu_d + n_{\sigma_d} \sigma_d, \arg \max_{x_d} x_d\} \end{aligned} \quad (8)$$

where μ_d and σ_d denote the sample mean and standard deviation of feature d . We set $n_{\sigma_d} = 3$ across the board but higher values and hence wider bounds may be appropriate depending on the application.

B.4 Training Details

In this section, we describe the training procedure in detail. While the details laid out here are not crucial for understanding our proposed approach, they are of importance to anyone looking to implement counterfactual training.

C Detailed Results

C.1 Qualitative Findings for Image Data

Note

Figure A1 shows much more plausible (faithful) counterfactuals for a model with CT than the model with conventional training (Figure A2). In fact, this is not even using *ECCo+* and still showing better results than the best results we achieved in our AAAI paper for JEM ensembles.

C.2 Grid Searches

Generator Parameters The hyperparameter grid with varying generator parameters during training is shown in Note 1. The corresponding evaluation grid used for these experiments is shown in Note 2.

Note 1: Training Phase

- Generator Parameters:
 - Decision Threshold: 0.75, 0.9, 0.95
 - λ_{energy} : 0.1, 0.5, 5.0, 10.0, 20.0
 - Maximum Iterations: 5, 25, 50
- Generator: `ecco`, `generic`, `revise`
- Model: `mlp`
- Training Parameters:
 - Objective: `full`, `vanilla`

Note 2: Evaluation Phase

- Generator Parameters:
 - λ_{energy} : 0.1, 0.5, 1.0, 5.0, 10.0

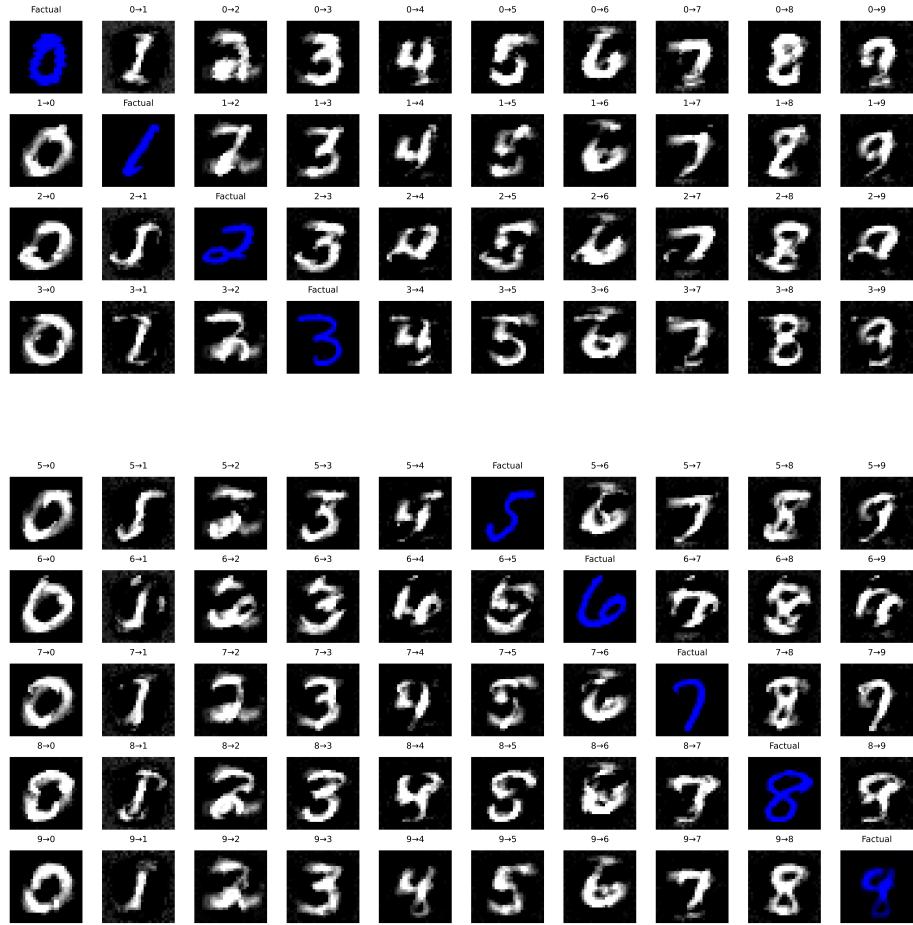


Fig. A1: Counterfactual images for *MLP* with counterfactual training. The underlying generator, *ECCo*, aims to generate counterfactuals that are faithful to the model [3].

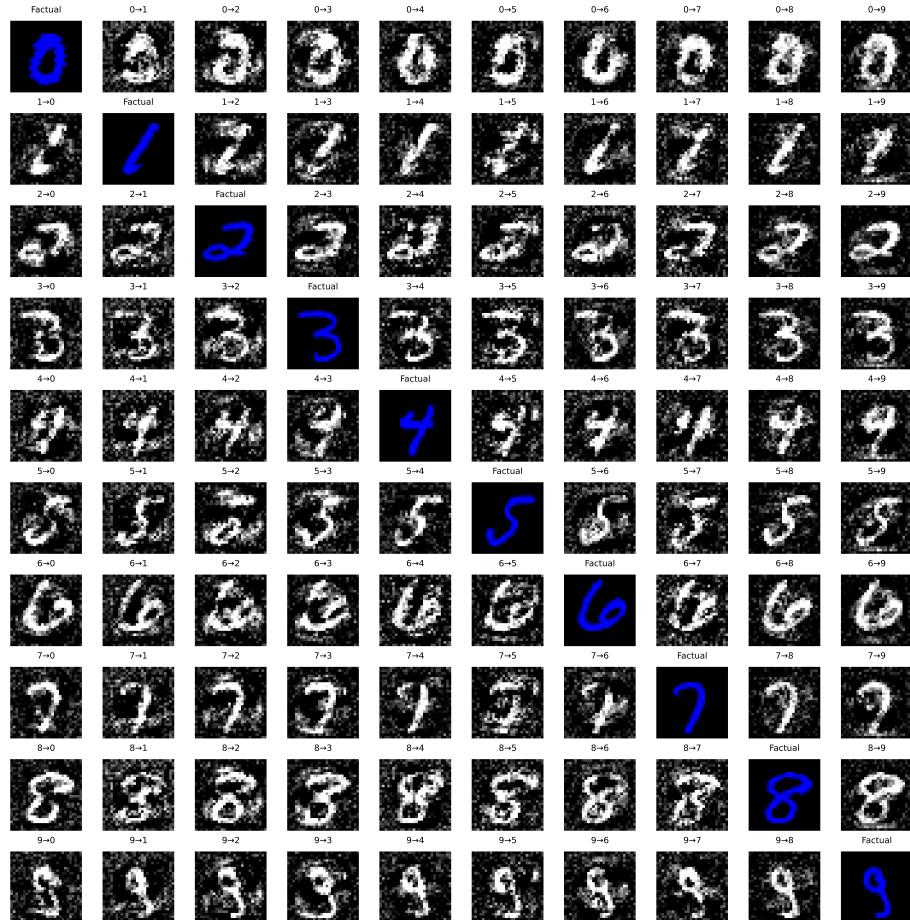
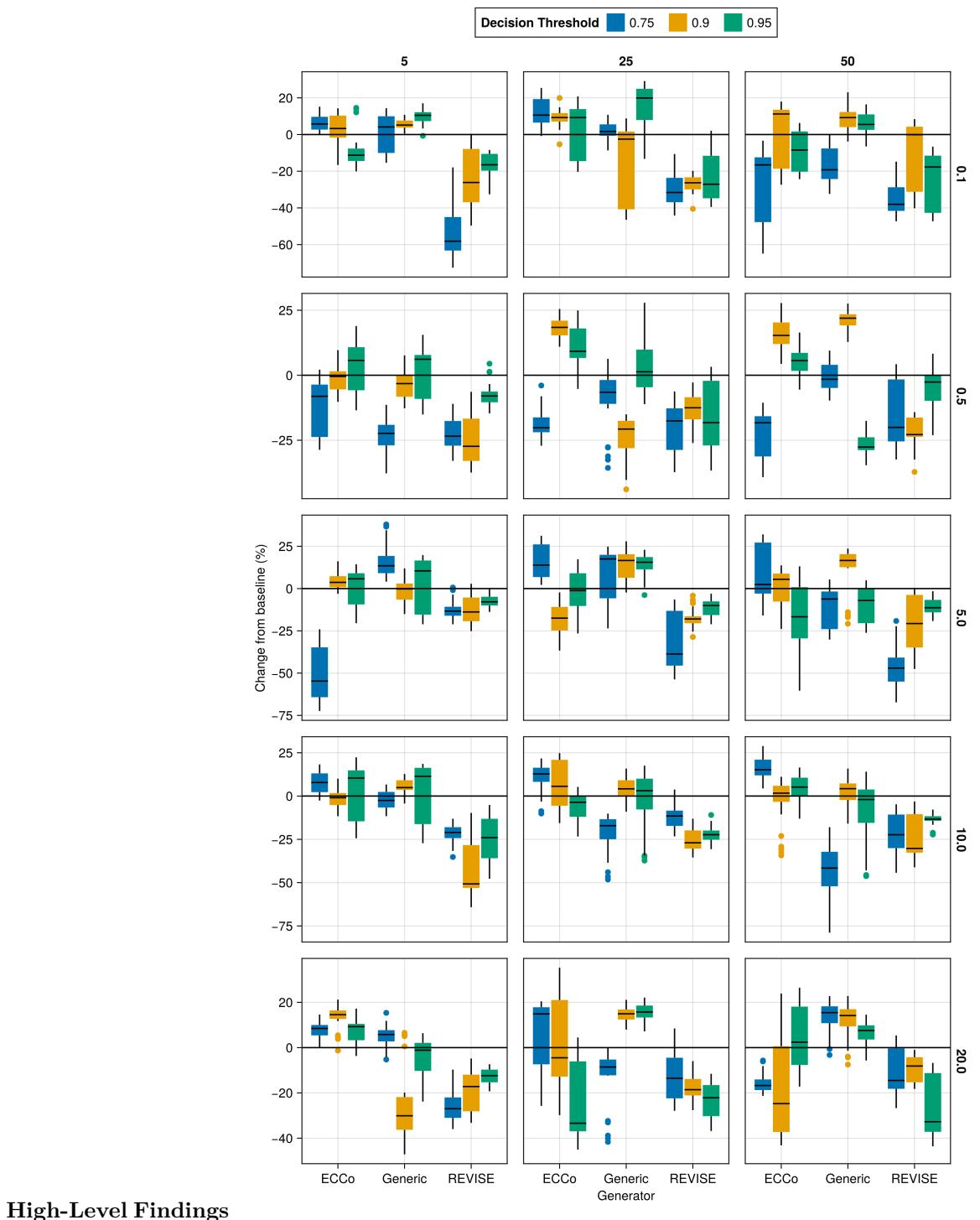
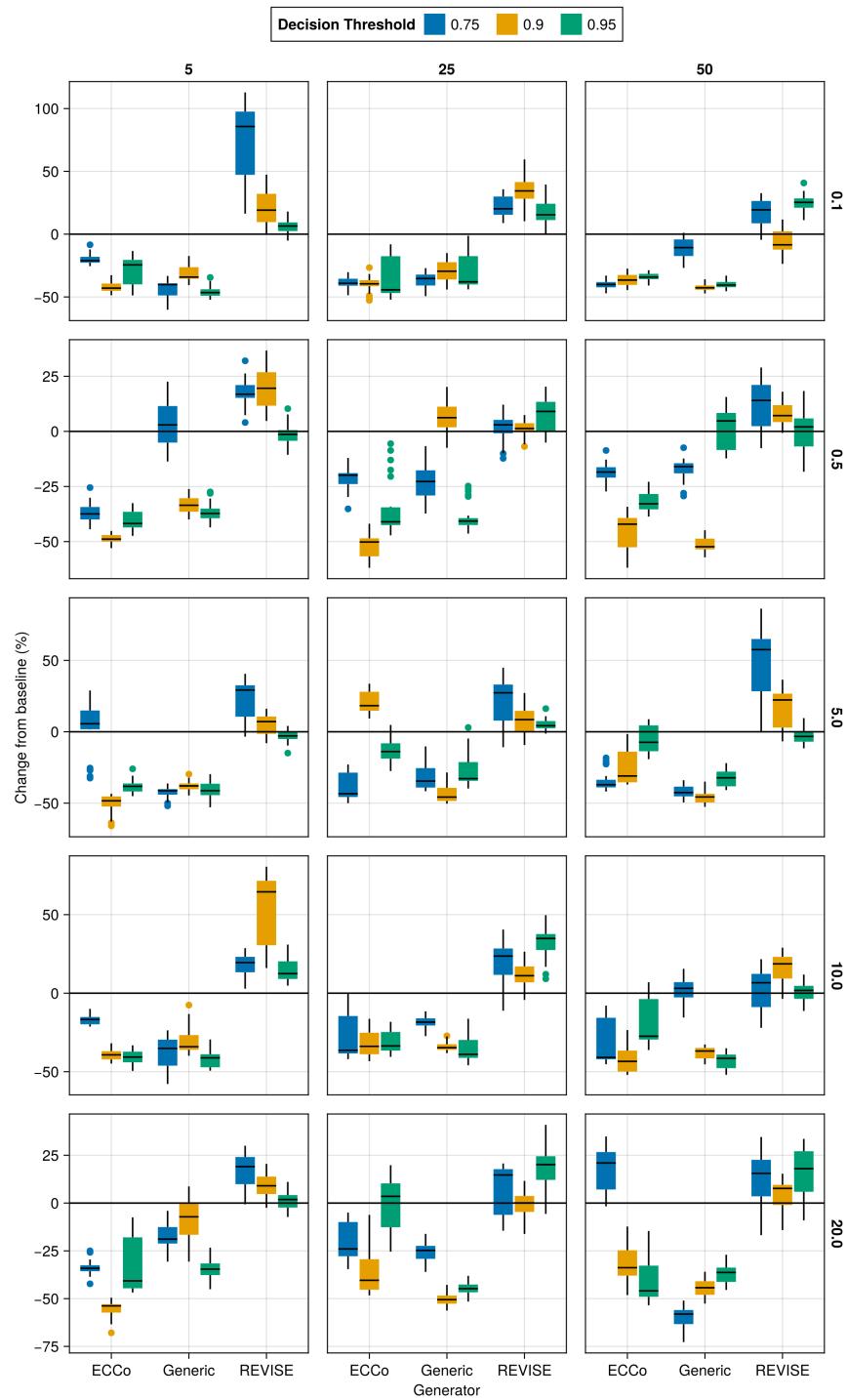
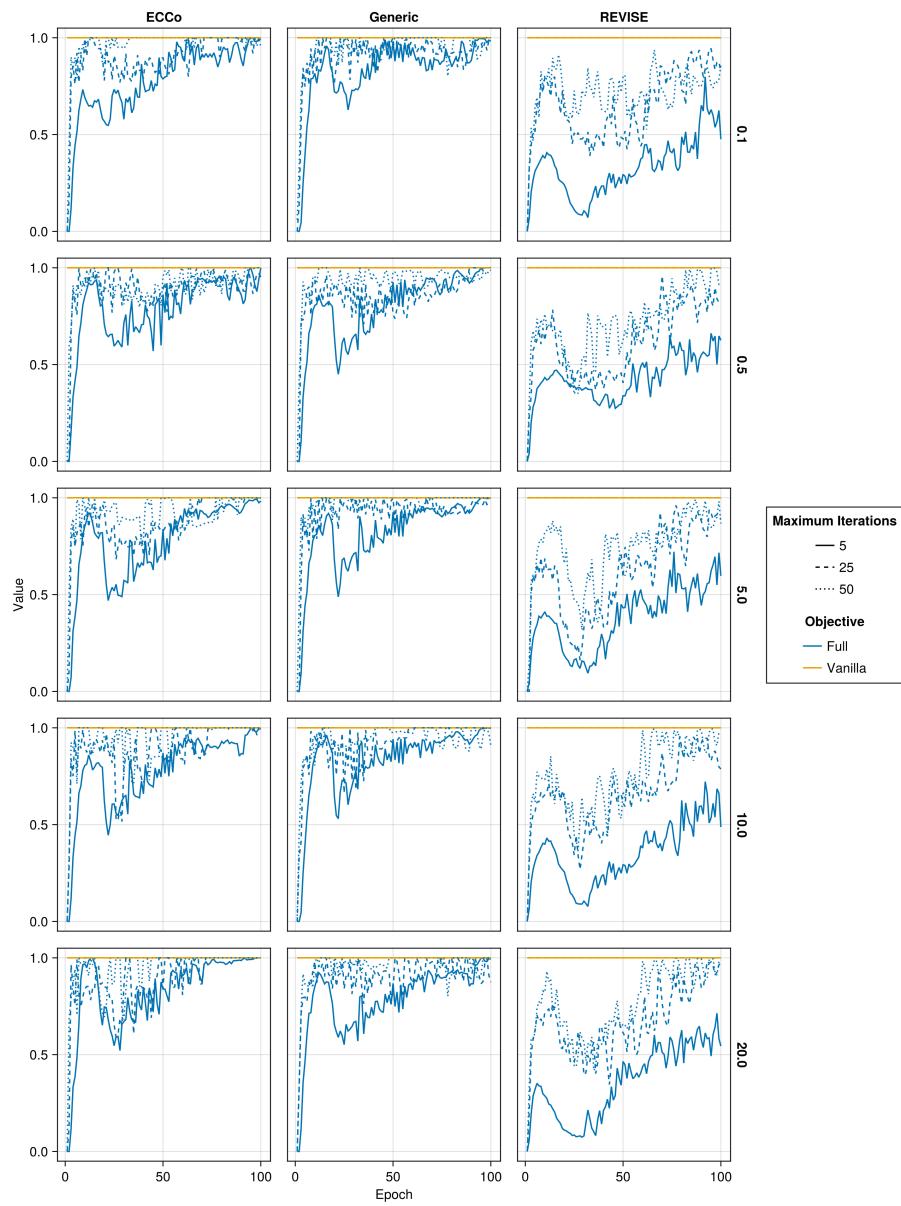
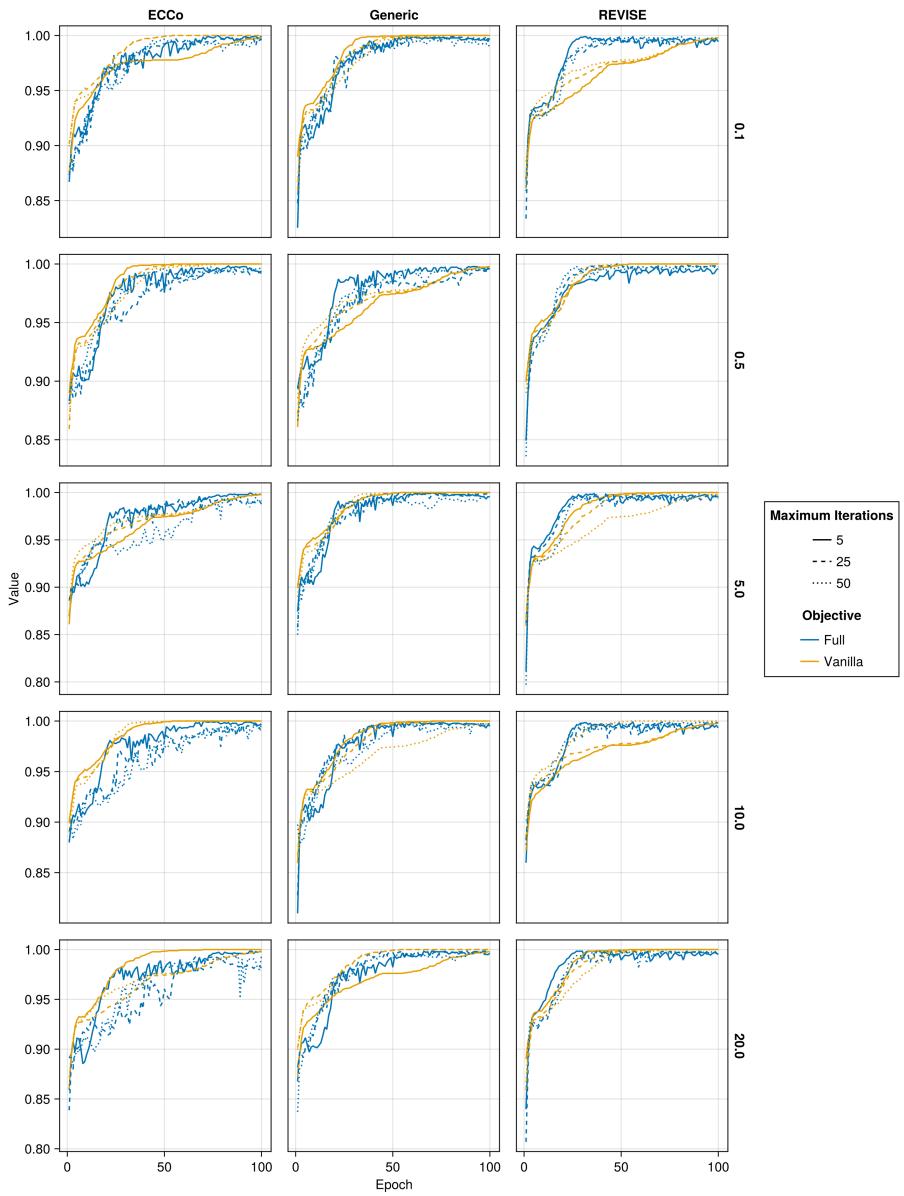


Fig. A2: Counterfactual images for *MLP* with conventional training. The underlying generator, *ECCo*, aims to generate counterfactuals that are faithful to the model [3].









Detailed Findings

Linearly Separable

Table A1: Results for Linearly Separable data by energy penalty.

Decision Threshold	λ_{energy}	Maximum Iterations	Generator	Value	Std
0.75	0.1	5	<i>ECCo</i>	-15.9	16.7
0.75	0.1	5	Generic	41.9	8.13
0.75	0.1	5	<i>REVISE</i>	-64.4	3.58
0.75	0.1	25	ECCo	28.9	2.28
0.75	0.1	25	<i>Generic</i>	-11.3	11.5
0.75	0.1	25	<i>REVISE</i>	-55.3	4.18
0.75	0.1	50	ECCo	10.9	7.47
0.75	0.1	50	<i>Generic</i>	-32.3	16.6
0.75	0.1	50	<i>REVISE</i>	-44.1	4.24
0.9	0.1	5	<i>ECCo</i>	-13.7	13.7
0.9	0.1	5	Generic	22.3	4.87
0.9	0.1	5	<i>REVISE</i>	-27.8	6.26
0.9	0.1	25	<i>ECCo</i>	8.84	7.13
0.9	0.1	25	Generic	30.7	5.77
0.9	0.1	25	<i>REVISE</i>	-72.7	34.9
0.9	0.1	50	ECCo	17.1	7.98
0.9	0.1	50	<i>Generic</i>	3.82	8.86
0.9	0.1	50	<i>REVISE</i>	-45.8	2.79
0.95	0.1	5	ECCo	35.7	7.56
0.95	0.1	5	<i>Generic</i>	-13.3	5.73
0.95	0.1	5	<i>REVISE</i>	-54.9	5.43
0.95	0.1	25	ECCo	37	5.31
0.95	0.1	25	<i>Generic</i>	29	8.6
0.95	0.1	25	<i>REVISE</i>	-40.5	41
0.95	0.1	50	<i>ECCo</i>	13.2	4.89
0.95	0.1	50	Generic	35.5	4.27
0.95	0.1	50	<i>REVISE</i>	-60.4	45.8
0.75	0.5	5	ECCo	-0.94	17.6
0.75	0.5	5	<i>Generic</i>	-11.7	17.4
0.75	0.5	5	<i>REVISE</i>	-37.7	3.92
0.75	0.5	25	<i>ECCo</i>	14.7	6.32
0.75	0.5	25	Generic	32.3	3.49
0.75	0.5	25	<i>REVISE</i>	-39.9	7.21
0.75	0.5	50	ECCo	28.3	5.05
0.75	0.5	50	<i>Generic</i>	-9.72	9.19
0.75	0.5	50	<i>REVISE</i>	-30	5.22
0.9	0.5	5	<i>ECCo</i>	-7.38	13.6
0.9	0.5	5	Generic	-0.0952	17.1
0.9	0.5	5	<i>REVISE</i>	-50.6	4.77
0.9	0.5	25	ECCo	39.1	8.3

Continuing table below.

Decision Threshold	λ_{energy}	Maximum Iterations	Generator	Value	Std
0.9	0.5	25	<i>Generic</i>	38.6	6.03
0.9	0.5	25	<i>REVISE</i>	-56.9	4.28
0.9	0.5	50	<i>ECCo</i>	-6.53	10.9
0.9	0.5	50	Generic	13.6	12.8
0.9	0.5	50	<i>REVISE</i>	-24.1	12.5
0.95	0.5	5	<i>ECCo</i>	0.596	1.11
0.95	0.5	5	Generic	4.57	5.88
0.95	0.5	5	<i>REVISE</i>	-34.3	6.09
0.95	0.5	25	<i>ECCo</i>	16.5	3.53
0.95	0.5	25	Generic	19.1	4.43
0.95	0.5	25	<i>REVISE</i>	-37.2	39.2
0.95	0.5	50	<i>ECCo</i>	-31	14.2
0.95	0.5	50	Generic	24.8	4.78
0.95	0.5	50	<i>REVISE</i>	-1.28	2.7
0.75	5	5	ECCo	37	9.21
0.75	5	5	<i>Generic</i>	0.977	21.1
0.75	5	5	<i>REVISE</i>	-55.2	4.75
0.75	5	25	ECCo	21.8	5.97
0.75	5	25	<i>Generic</i>	17.8	5.33
0.75	5	25	<i>REVISE</i>	-68.1	8.69
0.75	5	50	<i>ECCo</i>	0.864	9.51
0.75	5	50	Generic	8.56	13.6
0.75	5	50	<i>REVISE</i>	-71.8	3.18
0.9	5	5	<i>ECCo</i>	-3.59	20.3
0.9	5	5	Generic	-1.18	16.7
0.9	5	5	<i>REVISE</i>	-58.6	3.55
0.9	5	25	<i>ECCo</i>	-1.73	9.07
0.9	5	25	Generic	21.9	7.48
0.9	5	25	<i>REVISE</i>	-44.4	4.23
0.9	5	50	ECCo	30.8	10.7
0.9	5	50	<i>Generic</i>	16.3	11.1
0.9	5	50	<i>REVISE</i>	-53.3	6.73
0.95	5	5	ECCo	-5.29	24.5
0.95	5	5	<i>Generic</i>	-5.71	17.4
0.95	5	5	<i>REVISE</i>	-36.5	4.01
0.95	5	25	ECCo	26.9	4.08
0.95	5	25	<i>Generic</i>	19.5	5.49
0.95	5	25	<i>REVISE</i>	-25	5
0.95	5	50	<i>ECCo</i>	4.6	7.46
0.95	5	50	Generic	16.6	4.3
0.95	5	50	<i>REVISE</i>	-50.4	46.4
0.75	10	5	ECCo	2.47	8.57

Continuing table below.

Decision Threshold λ_{energy}	Maximum Iterations	Generator	Value	Std
0.75	10	<i>Generic</i>	-7.23	9.13
0.75	10	<i>REVISE</i>	-67.2	5.75
0.75	10	ECCo	29.9	15.8
0.75	25	<i>Generic</i>	4.47	8.22
0.75	25	<i>REVISE</i>	-52	5.51
0.75	50	<i>ECCo</i>	6.23	6.11
0.75	10	Generic	15.2	10.2
0.75	50	<i>REVISE</i>	-17	5.46
0.9	5	<i>ECCo</i>	-6.93	14.7
0.9	10	Generic	6.88	5.11
0.9	5	<i>REVISE</i>	-45.6	47.8
0.9	25	<i>ECCo</i>	6.2	6.75
0.9	10	Generic	34.7	9.74
0.9	25	<i>REVISE</i>	-55.7	6.04
0.9	50	<i>ECCo</i>	30.1	14.2
0.9	10	Generic	35	7.11
0.9	50	<i>REVISE</i>	-48.8	3.3
0.95	10	ECCo	35.7	3.83
0.95	5	<i>Generic</i>	2.85	7.9
0.95	5	<i>REVISE</i>	-34.9	4.56
0.95	10	ECCo	33	6.45
0.95	25	<i>Generic</i>	24.8	4.59
0.95	25	<i>REVISE</i>	-55.7	44
0.95	10	ECCo	8.07	7.68
0.95	50	<i>Generic</i>	-7.11	9.2
0.95	50	<i>REVISE</i>	-26.1	18.6
0.75	20	<i>ECCo</i>	-10.8	20.8
0.75	20	Generic	12.6	17.6
0.75	5	<i>REVISE</i>	-21.6	3.34
0.75	20	ECCo	17.2	4.25
0.75	25	<i>Generic</i>	-10.3	9.59
0.75	25	<i>REVISE</i>	-59	3.66
0.75	20	ECCo	15.6	8.38
0.75	50	<i>Generic</i>	3.92	6.2
0.75	50	<i>REVISE</i>	-70.9	5.16
0.9	20	ECCo	26.3	6.62
0.9	5	<i>Generic</i>	2.35	7.04
0.9	5	<i>REVISE</i>	-64.8	4.31
0.9	25	<i>ECCo</i>	21.1	11.6
0.9	20	Generic	39	5.98
0.9	25	<i>REVISE</i>	-50.7	13.8
0.9	20	ECCo	15.1	7.46

Continuing table below.

Decision Threshold λ_{energy}	Maximum Iterations	Generator	Value	Std
0.9	20	<i>Generic</i>	8.53	6.95
0.9	20	<i>REVISE</i>	-36	6.14
0.95	20	<i>ECCo</i>	-23.8	17.8
0.95	20	5	Generic	7.23 4.74
0.95	20	<i>REVISE</i>	-34	4.11
0.95	20	25	ECCo	19.3 14.6
0.95	20	<i>Generic</i>	-1.35	7.71
0.95	20	<i>REVISE</i>	-52.8	5.11
0.95	20	<i>ECCo</i>	18.7	3.84
0.95	20	50	Generic	37.3 7.78
0.95	20	<i>REVISE</i>	-28.3	15

D Computation Details

D.1 Hardware

We performed our experiments on a high-performance cluster. Details about the cluster will be disclosed upon publication to avoid revealing information that might interfere with the double-blind review process. Since our experiments involve highly parallel tasks and rather small models by today’s standard, we have relied on distributed computing across multiple central processing units (CPU). Graphical processing units (GPU) were not required. For larger grid searches we used up to but typically less than 100 CPUs.

D.2 Software

All computations were performed in the Julia Programming Language [6]. We have developed a package for counterfactual training that leverages and extends the functionality provided by several existing packages, most notably CounterfactualExplanations.jl [2] and the Flux.jl library for deep learning [20,21]. For data-wrangling and presentation-ready tables we relied on DataFrames.jl [7] and PrettyTables.jl [9], respectively. For plots and visualizations we used both Plots.jl [10] and Makie.jl [11], in particular AlgebraOfGraphics.jl. To distribute computational tasks across multiple processors, we have relied on MPI.jl [8].

References

1. Abbasnejad, E., Teney, D., Parvaneh, A., Shi, J., van den Hengel, A.: Counterfactual vision and language learning. In: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 10041–10051 (2020). <https://doi.org/10.1109/CVPR42600.2020.01006>
2. Altmeyer, P., van Deursen, A., et al.: Explaining black-box models through counterfactuals. In: Proceedings of the JuliaCon Conferences. vol. 1, p. 130 (2023)

3. Altmeyer, P., Farmanbar, M., van Deursen, A., Liem, C.C.: Faithful model explanations through energy-constrained conformal counterfactuals. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 38, pp. 10829–10837 (2024)
4. Augustin, M., Meinke, A., Hein, M.: Adversarial robustness on in-and out-distribution improves explainability. In: European Conference on Computer Vision. pp. 228–245. Springer (2020)
5. Balashankar, A., Wang, X., Qin, Y., Packer, B., Thain, N., Chi, E., Chen, J., Beutel, A.: Improving classifier robustness through active generative counterfactual data augmentation. In: Findings of the Association for Computational Linguistics: EMNLP 2023. pp. 127–139 (2023)
6. Bezanson, J., Edelman, A., Karpinski, S., Shah, V.B.: Julia: A fresh approach to numerical computing. SIAM review **59**(1), 65–98 (2017), <https://doi.org/10.1137/141000671>
7. Bouchet-Valat, M., Kamiński, B.: Dataframes.jl: Flexible and fast tabular data in julia. Journal of Statistical Software **107**(4), 1–32 (2023). <https://doi.org/10.18637/jss.v107.i04>, <https://www.jstatsoft.org/index.php/jss/article/view/v107i04>
8. Byrne, S., Wilcox, L.C., Churavy, V.: Mpi.jl: Julia bindings for the message passing interface. Proceedings of the JuliaCon Conferences **1**(1), 68 (2021). <https://doi.org/10.21105/jcon.00068>, <https://doi.org/10.21105/jcon.00068>
9. Chagas, R.A.J., Baumgold, B., Hertz, G., Ranocha, H., Wells, M., Boyer, N., Ritchie, N., Christensen, Z.P., gustaphe, pdeffebach, Pasquier, B., Kawczynski, C., Pinyol, D., Sambrani, D., Vanguelov, D., Hanson, E., Butterworth, I., Peters, J., Kraak, J., Storopoli, J., TagBot, J., Amin, N., Christ, S., Lienart, T., Christensen, T., Flatberg, T., Gagnon, Y.L., Mengali, A., jariji, jondeuce: ronisbr/prettytables.jl: v2.4.0 (Sep 2024). <https://doi.org/10.5281/zenodo.13835553>, <https://doi.org/10.5281/zenodo.13835553>
10. Christ, S., Schwabeneder, D., Rackauckas, C., Borregaard, M.K., Breloff, T.: Plots.jl – a user extendable plotting api for the julia programming language (2023). <https://doi.org/https://doi.org/10.5334/jors.431>, <https://openresearchsoftware.metajnl.com/articles/10.5334/jors.431/>
11. Danisch, S., Krumbiegel, J.: Makie.jl: Flexible high-performance data visualization for Julia. Journal of Open Source Software **6**(65), 3349 (2021). <https://doi.org/10.21105/joss.03349>, <https://doi.org/10.21105/joss.03349>
12. Du, Y., Mordatch, I.: Implicit generation and generalization in energy-based models (2020)
13. Freiesleben, T.: The intriguing relation between counterfactual explanations and adversarial examples. Minds and Machines **32**(1), 77–109 (2022)
14. Goodfellow, I., Bengio, Y., Courville, A.: Deep Learning. MIT Press (2016)
15. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples (2014)
16. Grathwohl, W., Wang, K.C., Jacobsen, J.H., Duvenaud, D., Norouzi, M., Swersky, K.: Your classifier is secretly an energy based model and you should treat it like one. In: International Conference on Learning Representations (2020)
17. Guidotti, R.: Counterfactual explanations and how to find them: literature review and benchmarking. Data Mining and Knowledge Discovery pp. 1–55 (2022)
18. Guo, H., Nguyen, T.H., Yadav, A.: Counternet: End-to-end training of prediction aware counterfactual explanations. In: Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining. p. 577–589. KDD ’23, Association for Computing Machinery, New York, NY, USA (2023). <https://doi.org/10.1145/3580305.3599290>, <https://doi.org/10.1145/3580305.3599290>

19. Hastie, T., Tibshirani, R., Friedman, J.: *The Elements of Statistical Learning*. Springer New York (2009). <https://doi.org/10.1007/978-0-387-84858-7>
20. Innes, M., Saba, E., Fischer, K., Gandhi, D., Rudiloso, M.C., Joy, N.M., Karmali, T., Pal, A., Shah, V.: *Fashionable modelling with flux* (2018)
21. Innes, M.: *Flux: Elegant machine learning with Julia*. *Journal of Open Source Software* **3**(25), 602 (2018). <https://doi.org/10.21105/joss.00602>
22. Joshi, S., Koyejo, O., Vijitbenjaronk, W., Kim, B., Ghosh, J.: Towards realistic individual recourse and actionable explanations in black-box decision making systems (2019)
23. Kolter, Z.: Keynote Addresses: SaTML 2023 . In: 2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML). pp. xvi–xvi. IEEE Computer Society, Los Alamitos, CA, USA (Feb 2023). <https://doi.org/10.1109/SaTML54575.2023.00009>
24. Lakshminarayanan, B., Pritzel, A., Blundell, C.: Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems* **30** (2017)
25. Lippe, P.: UvA Deep Learning Tutorials. <https://uvadlc-notebooks.readthedocs.io/en/latest/> (2024)
26. Luu, H.L., Inoue, N.: Counterfactual adversarial training for improving robustness of pre-trained language models. In: Proceedings of the 37th Pacific Asia Conference on Language, Information and Computation. pp. 881–888 (2023)
27. Murphy, K.P.: *Probabilistic Machine Learning: An Introduction*. MIT Press (2022)
28. O’Neil, C.: *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. Crown (2016)
29. Pawelczyk, M., Agarwal, C., Joshi, S., Upadhyay, S., Lakkaraju, H.: Exploring counterfactual explanations through the lens of adversarial examples: A theoretical and empirical analysis. In: Camps-Valls, G., Ruiz, F.J.R., Valera, I. (eds.) *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics. Proceedings of Machine Learning Research*, vol. 151, pp. 4574–4594. PMLR (28–30 Mar 2022), <https://proceedings.mlr.press/v151/pawelczyk22a.html>
30. Poiiadzi, R., Sokol, K., Santos-Rodriguez, R., De Bie, T., Flach, P.: FACE: Feasible and actionable counterfactual explanations. In: *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*. pp. 344–350 (2020)
31. Ross, A., Lakkaraju, H., Bastani, O.: Learning models for actionable recourse. In: *Proceedings of the 35th International Conference on Neural Information Processing Systems. NIPS ’21*, Curran Associates Inc., Red Hook, NY, USA (2024)
32. Sauer, A., Geiger, A.: Counterfactual generative networks (2021), <https://arxiv.org/abs/2101.06046>
33. Schut, L., Key, O., Mc Grath, R., Costabello, L., Sacaleanu, B., Gal, Y., et al.: Generating Interpretable Counterfactual Explanations By Implicit Minimisation of Epistemic and Aleatoric Uncertainties. In: *International Conference on Artificial Intelligence and Statistics*. pp. 1756–1764. PMLR (2021)
34. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks (2013)
35. Teney, D., Abbasnedjad, E., van den Hengel, A.: Learning what makes a difference from counterfactual examples and gradient supervision. In: *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X 16*. pp. 580–599. Springer (2020)
36. Wachter, S., Mittelstadt, B., Russell, C.: Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harv. JL & Tech.* **31**, 841 (2017). <https://doi.org/10.2139/ssrn.3063289>

37. Wilson, A.G.: The case for bayesian deep learning (2020)
38. Wu, T., Ribeiro, M.T., Heer, J., Weld, D.: Polyjuice: Generating counterfactuals for explaining, evaluating, and improving models. In: Zong, C., Xia, F., Li, W., Navigli, R. (eds.) Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 6707–6723. Association for Computational Linguistics, Online (Aug 2021). <https://doi.org/10.18653/v1/2021.acl-long.523>, <https://aclanthology.org/2021.acl-long.523>