
COUNTERFACTUAL TRAINING: TEACHING MODELS PLAUSIBLE AND ACTIONABLE EXPLANATIONS*

January 22, 2026

ABSTRACT

We propose a novel training regime termed counterfactual training that leverages counterfactual explanations to increase the explanatory capacity of models. Counterfactual explanations have emerged as a popular post-hoc explanation method for opaque machine learning models: they inform how factual inputs would need to change in order for a model to produce some desired output. To be useful in real-world decision-making systems, counterfactuals should be plausible with respect to the underlying data and actionable with respect to the feature mutability constraints. Much existing research has therefore focused on developing post-hoc methods to generate counterfactuals that meet these desiderata. In this work, we instead hold models directly accountable for the desired end goal: counterfactual training employs counterfactuals during the training phase to minimize the divergence between learned representations and plausible, actionable explanations. We demonstrate empirically and theoretically that our proposed method facilitates training models that deliver inherently desirable counterfactual explanations and additionally exhibit improved adversarial robustness.

Keywords Counterfactual Training • Counterfactual Explanations • Algorithmic Recourse • Explainable AI • Representation Learning

*This work has been accepted for publication at the IEEE Conference on Secure and Trustworthy Machine Learning (SaTML). The final version will be available on IEEE Xplore.

1 Introduction

Today’s prominence of artificial intelligence (AI) has largely been driven by the success of representation learning with high degrees of freedom: instead of relying on features and rules hand-crafted by humans, modern machine learning (ML) models are tasked with learning highly complex representations directly from the data, guided by narrow objectives such as predictive accuracy (Goodfellow, Bengio, and Courville 2016). These models tend to be so complex that humans cannot easily interpret their decision logic.

Counterfactual explanations (CE) have become a key part of the broader explainable AI (XAI) toolkit (Molnar 2022) that can be applied to make sense of this complexity. They prescribe minimal changes for factual inputs that, if implemented, would prompt some fitted model to produce an alternative, more desirable output (Wachter, Mittelstadt, and Russell 2017). This is useful and necessary to not only understand how opaque models make their predictions, but also to provide algorithmic recourse to individuals subjected to them: a retail bank, for example, could use CE to provide meaningful feedback to unsuccessful loan applicants that were rejected based on an opaque automated decision-making (ADM) system (Figure 1).

For such feedback to be meaningful, counterfactual explanations need to fulfill certain desiderata (Verma et al. 2022; Karimi et al. 2021)—they should be faithful to the model (Altmeyer et al. 2024), plausible (Joshi et al. 2019), and actionable (Ustun, Spangher, and Liu 2019). Plausibility is typically understood as counterfactuals being *in-domain*: unsuccessful loan applicants that implement the provided recourse should end up with credit profiles that are genuinely similar to that of individuals who have successfully repaid their loans in the past. Actionable explanations further comply with practical constraints: a young, unsuccessful loan applicant cannot increase their age in an instant.

Existing state-of-the-art (SOTA) approaches in the field have largely focused on designing model-agnostic CE methods that identify subsets of counterfactuals, which comply with specific desiderata. This is problematic because the narrow focus on any specific desideratum can adversely affect others: it is possible, for example, to generate plausible counterfactuals for models that are also highly vulnerable to implausible, possibly adversarial counterfactuals (Altmeyer et al. 2024). Indeed, existing approaches generally fail to guarantee that the representations learned by a model are compatible with truly meaningful explanations.

In this work, we propose an approach to bridge this gap, embracing the paradigm that models—as opposed to explanation methods—should be held accountable for explanations that are plausible and actionable. While previous work has shown that at least plausibility can be indirectly achieved through existing techniques aimed at models’ generative capacity, generalization and robustness (Altmeyer et al. 2024; Augustin, Meinke, and Hein 2020; Schut et al. 2021), we directly incorporate both plausibility and actionability in the training objective of models to improve their overall explanatory capacity.

Specifically, we introduce **counterfactual training (CT)**: a novel training regime that leverages counterfactual explanations on-the-fly to ensure that differentiable models learn plausible and actionable explanations for the underlying data, while at the same time being more robust to adversarial examples (AE). Figure 1 illustrates the outcomes of CT compared to a conventionally trained model. First, in panel (a), faithful and valid counterfactuals end up near the decision boundary forming a clearly distinguishable cluster in the target class (orange). In panel (b), CT is applied to the same underlying linear classifier architecture resulting in much more plausible counterfactuals. In panel (c), the classifier is again trained conventionally and we have introduced a mutability constraint on the *age* feature at test time—counterfactuals are valid but the classifier is roughly equally sensitive to both features. By contrast, the decision boundary in panel (d) is tilted, making the model trained with CT relatively less sensitive to the immutable *age* feature. To achieve these outcomes, CT draws inspiration from the literature on contrastive and robust learning: we contrast faithful CEs with ground-truth data while protecting immutable features, and capitalize on methodological links between CE and AE by penalizing the model’s adversarial loss on interim (*nascent*) counterfactuals. To the best of our knowledge, CT represents the first venture in this direction with promising empirical and theoretical results.

The remainder of this manuscript is structured as follows. Section 2 presents related work, focusing on the links to contrastive and robust learning. Then follow our two principal contributions. In Section 3, we introduce our methodological framework and show theoretically that it can be employed to respect global actionability constraints. In our experiments (Section 4), we find that thanks to counterfactual training, (1) the implausibility of CEs decreases by up to 90%; (2) the cost of reaching valid counterfactuals with protected features decreases by 19% on average; and (3) models’ adversarial robustness improves across the board. Finally, we discuss open challenges in Section 5 and conclude in Section 6.

2 Related Literature

To make the desiderata for CT more concrete, we follow previous work, tying the explanatory capacity of models to the quality of CEs that can be generated for them (Altmeyer et al. 2024; Augustin, Meinke, and Hein 2020).

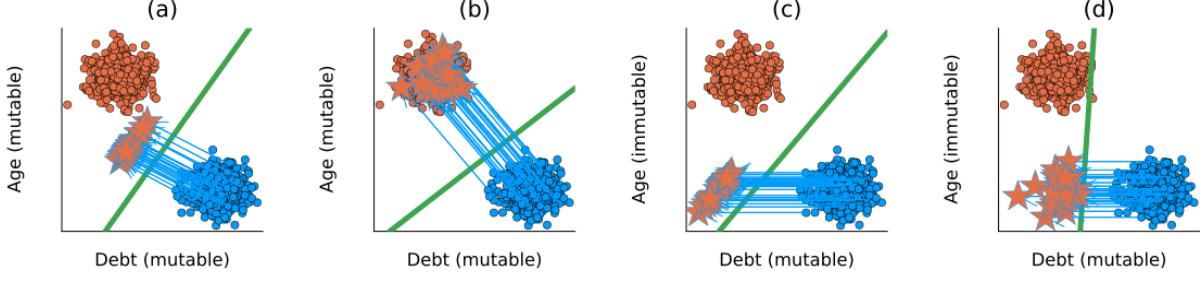


Figure 1 Counterfactual explanations (stars) for linear classifiers trained under different regimes on synthetic data: (a) conventional training, all mutable; (b) CT, all mutable; (c) conventional, *age* immutable; (d) CT, *age* immutable. The linear decision boundary is shown in green along with training data colored according to ground-truth labels: y^- = "loan withheld" (blue) and y^+ = "loan provided" (orange). Class and feature annotations (*debt* and *age*) are for illustrative purposes.

2.1 Explanatory Capacity and Contrastive Learning

A closely related work shows that model averaging and, in particular, contrastive model objectives can produce models that have a higher explanatory capacity, and hence ones that are more trustworthy (Altmeyer et al. 2024). The authors propose a way to generate counterfactuals that are maximally faithful in that they are consistent with what models have learned about the underlying data. Formally, they rely on tools from energy-based modelling (Teh et al. 2003) to minimize the contrastive divergence between the distribution of counterfactuals and the conditional posterior over inputs learned by a model. Their algorithm, *ECCCo*, yields plausible counterfactual explanations if and only if the underlying model has learned representations that align with them. The authors find that both deep ensembles (Lakshminarayanan, Pritzel, and Blundell 2017) and joint energy-based models (JEMs) (Grathwohl et al. 2020), a form of contrastive learning, do well in this regard.

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 particularly effective when models are underspecified by the available data (Wilson 2020). Averaging across solutions mitigates the risk of overrelying on a single locally optimal representation that corresponds to semantically meaningless explanations. Likewise, it has been shown that generating plausible ("interpretable") CEs is almost trivial for deep ensembles that have undergone adversarial training (Schut et al. 2021). The case for JEMs is even clearer: they optimize a hybrid objective that induces both high predictive performance and strong generative capacity (Grathwohl et al. 2020), resembling the idea of aligning models with plausible explanations. This was an inspiration for CT.

2.2 Explanatory Capacity and Robust Learning

Prior work has shown that counterfactual explanations tend to be more meaningful ("explainable") if the underlying model is more robust to adversarial examples (Augustin, Meinke, and Hein 2020). Once again, we can make intuitive sense of this finding if we look at adversarial training (AT) through the lens of representation learning with high degrees of freedom: highly complex and flexible models may learn representations that make them sensitive to implausible or even adversarial examples (Szegedy et al. 2014). Thus, by inducing models to "unlearn" susceptibility to such examples, adversarial training can effectively remove implausible explanations from the solution space.

This interpretation of the link between explanatory capacity through counterfactuals on the one side, and robustness to adversarial examples on the other is backed by empirical evidence. Firstly, prior work has shown that using counterfactual images during classifier training improves model robustness (Sauer and Geiger 2021). Similarly, related work has shown that counterfactuals represent potentially useful training data in machine learning tasks, especially in supervised settings where inputs may be reasonably mapped to multiple outputs (Abbasnejad et al. 2020). The authors show that augmenting the training data of (image) classifiers can improve generalization performance. Finally, another related work has demonstrated that counterfactual pairs tend to exist in training data (Teney, Abbasnejad, and Hengel 2020). Hence, the proposed approach aims to identify similar inputs with different annotations and ensure that the gradient of the classifier aligns with the vector between such pairs of inputs using a cosine distance loss function.

CEs have also been used to improve models in the natural language processing domain. A well-known paper in this domain has proposed *Polyjuice* (Wu et al. 2021), a general-purpose CE generator for language models. The authors demonstrate that the augmentation of training data with *Polyjuice* improves robustness in a number of tasks. Related work has introduced the *Counterfactual Adversarial Training* (CAT) framework (Luu and Inoue 2023), which aims to

improve generalization and robustness of language models by generating counterfactuals for training samples that are subject to high predictive uncertainty.

There have also been several attempts at formalizing the relationship between counterfactual explanations and adversarial examples. Pointing to clear similarities in how CEs and AEs are generated, prior work makes the case for jointly studying the opaqueness and robustness problems in representation learning (Freiesleben 2022). Formally, the authors show that AEs can be seen as the subset of CEs for which misclassification is achieved (Freiesleben 2022). Similarly, others have shown that CEs and AEs are equivalent under certain conditions (Pawelczyk et al. 2022).

Two other works are closely related to ours in that they use counterfactuals during training with the explicit goal of affecting certain properties of the post-hoc counterfactual explanations. The first closely related work has proposed a way to train models that guarantee recourse to a positive target class with high probability (Ross, Lakkaraju, and Bastani 2024). The approach builds on adversarial training by explicitly inducing susceptibility to targeted AEs for the positive class. Additionally, the method allows for imposing a set of actionability constraints ex-ante. For example, users can specify that certain features are immutable. A second closely related work has introduced the first end-to-end training pipeline that includes CEs as part of the training procedure (Guo, Nguyen, and Yadav 2023); the *CounterNet* network architecture includes a predictor and a CE generator, where the parameters of the CE generator are learnable. Counterfactuals are generated during each training iteration and fed back to the predictor. In contrast, we impose no restrictions on the artificial neural network architecture at all.

3 Counterfactual Training

This section introduces the counterfactual training framework, applying ideas from contrastive and robust learning to counterfactual explanations. CT produces models whose learned representations align with plausible explanations that comply with user-defined actionability constraints.

Counterfactual explanations are typically generated by solving variations of the following optimization problem,

$$\min_{\mathbf{x}' \in \mathcal{X}^D} \{y\text{loss}(\mathbf{M}_\theta(\mathbf{x}'), \mathbf{y}^+) + \lambda \text{reg}(\mathbf{x}')\} \quad (1)$$

where $\mathbf{M}_\theta : \mathcal{X} \mapsto \mathcal{Y}$ denotes a classifier, \mathbf{x}' denotes the counterfactual with D features and $\mathbf{y}^+ \in \mathcal{Y}$ denotes some target class. The $y\text{loss}(\cdot)$ function quantifies the discrepancy between current model predictions for \mathbf{x}' and the target class (a conventional choice is cross-entropy). Finally, we use $\text{reg}(\cdot)$ to denote any form of regularization used to induce certain properties on the counterfactual. The seminal CE paper, (Wachter, Mittelstadt, and Russell 2017), proposes regularizing the distance between counterfactuals and their original factual values to ensure that individuals seeking recourse through CE face minimal costs in terms of feature changes. Different variations of Equation 1 have been proposed in the literature to address many desiderata including the ones discussed above (faithfulness, plausibility and actionability). Much like in the seminal work (Wachter, Mittelstadt, and Russell 2017), most of these approaches rely on gradient descent to optimize Equation 1, and this holds true for all approaches tested in this work. We introduce them briefly in Section 4.1, but refer the reader to the supplementary appendix for details. In the following, we describe how counterfactuals are generated and used in CT.

3.1 Proposed Training Objective

The goal of CT is to improve the explanatory capacity of models by aligning the learned representations with faithful explanations that are plausible and actionable. For simplicity, we refer to models with high explanatory capacity as **explainable** in this manuscript. We define explainability as follows:

Definition 3.1 (Model Explainability). Let $\mathbf{M}_\theta : \mathcal{X} \mapsto \mathcal{Y}$ 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} . Let \mathbf{x}'_0 denote a factual input and assume that for any given input-output pair $\{\mathbf{x}'_0, \mathbf{y}\}_i$ there exists a counterfactual $\mathbf{x}' = \mathbf{x}'_0 + \Delta : \mathbf{M}_\theta(\mathbf{x}') = \mathbf{y}^+ \neq \mathbf{y} = \mathbf{M}_\theta(\mathbf{x})$, where $\arg \max_y \mathbf{y}^+ = y^+$ is the index of the target class.

We say that \mathbf{M}_θ has an **explanatory capacity** to the extent that faithfully generated, valid counterfactuals are also plausible and actionable. We define these properties as:

- (Faithfulness) $P(\mathbf{x}' \in \mathcal{X}_\theta | \mathbf{y}^+) = 1 - \delta$, where δ is some small value, and $\mathcal{X}_\theta | \mathbf{y}^+$ is the conditional posterior distribution over inputs (adapted from (Altmeyer et al. 2024), Def. 4.1).
- (Plausibility) $P(\mathbf{x}' \in \mathcal{X} | \mathbf{y}^+) = 1 - \delta$, where δ is some small value, and $\mathcal{X} | \mathbf{y}^+$ is the conditional distribution of inputs in the target class (adapted from (Altmeyer et al. 2024), Def. 2.1).
- (Actionability) Perturbations Δ may be subject to some actionability constraints.

Intuitively, plausible counterfactuals are consistent with the data, and faithful counterfactuals are consistent with what the model has learned about the input data. Actionability constraints in Definition 3.1 depend on the context in which \mathbf{M}_θ is deployed (e.g., specified by end-users or model owners). We consider two types of actionability constraints: on the domain of features and on their mutability. The former naturally arise in automated decision-making systems whenever a feature can only take a specific range of values. For example, *age* is lower bounded by zero and upper bounded by the maximum human lifespan. Specifying such domain constraints can also help address training instabilities commonly associated with energy-based modelling (Grathwohl et al. 2020). The latter arise when a feature cannot be freely modified. Continuing the example, *age* of a person can only increase, but it may even be considered as an immutable feature: waiting many years for an improved outcome is hardly feasible for individuals affected by algorithmic decisions. We choose to only consider domain and mutability constraints for individual features x_d for $d = 1, \dots, D$. Of course, this is a simplification since feature values may correlate, e.g., higher *age* may be associated with higher *level of completed education*. We address this challenge in Section 5, where we also explain why we restrict this work to classification settings.

Let \mathbf{x}'_t for $t = 0, \dots, T$ denote a counterfactual generated through gradient descent over T iterations as originally proposed (Wachter, Mittelstadt, and Russell 2017). CT adopts gradient-based CE search in training to generate on-the-fly model explanations \mathbf{x}' for the training samples. We use the term *nascent* to denote interim counterfactuals \mathbf{x}'_{AE} that have not yet converged. As we explain below, these nascent counterfactuals can be stored and repurposed as adversarial examples. Conversely, we consider counterfactuals \mathbf{x}'_{CE} as *mature* explanations if they have converged within the T iterations by reaching a pre-specified threshold, τ , for the predicted probability of the target class: $\mathcal{S}(\mathbf{M}_\theta(\mathbf{x}'))[y^+] \geq \tau$, where \mathcal{S} is the softmax function.

Formally, we propose the following counterfactual training objective to train explainable (as in Definition 3.1) models,

$$\begin{aligned} & \min_{\theta} yloss(\mathbf{M}_\theta(\mathbf{x}), \mathbf{y}) + \lambda_{div} \text{div}(\mathbf{x}^+, \mathbf{x}'_{CE}, y^+; \theta) \\ & + \lambda_{adv} \text{advloss}(\mathbf{M}_\theta(\mathbf{x}'_{AE}), \mathbf{y}_{AE}) + \lambda_{reg} \text{ridge}(\mathbf{x}^+, \mathbf{x}'_{CE}, y; \theta) \end{aligned} \quad (2)$$

where $yloss(\cdot)$ is any classification loss that induces discriminative performance (e.g., cross-entropy). The second and third terms are explained in detail in the following subsections. For now, they can be summarized as inducing explainability directly and indirectly by penalizing (1) the contrastive divergence, $\text{div}(\cdot)$, between mature counterfactuals \mathbf{x}'_{CE} and observed samples $\mathbf{x}^+ \in \mathcal{X}^+ = \{\mathbf{x} : y = y^+\}$ in the target class y^+ , and (2) the adversarial loss, $\text{advloss}(\cdot)$, wrt. nascent counterfactuals \mathbf{x}'_{AE} and their corresponding labels \mathbf{y}_{AE} . Finally, $\text{ridge}(\cdot)$ denotes a Ridge penalty (squared ℓ_2 -norm) that regularizes the magnitude of the energy terms involved in the contrastive divergence, $\text{div}(\cdot)$, term (Du and Mordatch 2020):

$$\frac{1}{n_{CE}} \sum_{i=1}^{n_{CE}} (\mathcal{E}_\theta(\mathbf{x}^+, y^+)^2 + \mathcal{E}_\theta(\mathbf{x}'_{CE}, y^+)^2) \quad (3)$$

The trade-offs between these components are adjusted through penalties λ_{div} , λ_{adv} , and λ_{reg} .

The full counterfactual training regime is sketched out in Algorithm 1. During each iteration, we do the following steps. Firstly, we randomly draw a subset of $n_{CE} \leq n$ factuals \mathbf{x}'_0 from \mathbf{X} of size n , for which we uniformly draw a target class y^+ (ensuring that it does not coincide with the class currently predicted for \mathbf{x}'_0) and a corresponding training sample from the target class, $\mathbf{x}^+ \sim \mathbf{X}^+ = \{\mathbf{x} \in \mathbf{X} : y = y^+\}$. Secondly, we conduct the counterfactual search by solving (Equation 1) through gradient descent. Thirdly, we sample mini-batches $(\mathbf{x}_i, \mathbf{y}_i)_{i=1}^{n_b}$ from the training dataset $\mathcal{D} = (\mathbf{X}, \mathbf{Y})$ for conventional training and distribute the tuples composed of counterfactuals, their target labels and corresponding training samples, as well as adversarial examples and corresponding labels, $(\mathbf{x}'_{CE_i}, y^+_i, \mathbf{x}'_{AE_i}, \mathbf{y}_{AE_i}, \mathbf{x}^+_i)_{i=1}^{n_{CE}}$, across the mini-batches. Finally, we backpropagate through (Equation 2).

Algorithm 1 Pseudo-Code for Counterfactual Training

Require: Training dataset \mathcal{D} , initialize model M_θ

- 1: **while** not converged **do**
- 2: Sample $x'_0 \sim X$, $y^+ \sim \mathcal{U}(\mathcal{Y})$ and $x^+ \sim X^+$
- 3: **for** $t = 1$ to T **do**
- 4: Backpropagate $\nabla_{x'}$ through equation (1)
- 5: Store x'_{CE}, x'_{AE}, y_{AE}
- 6: **end for**
- 7: Sample mini-batches $(x_i, y_i)_{i=1}^{n_b}$ from dataset \mathcal{D}
- 8: Distribute $(x'_{CE_i}, y^+_{i=1}, x'_{AE_i}, y_{AE_i}, x^+_{i=1})_{i=1}^{n_{CE}}$
- 9: **for** each batch **do**
- 10: Backpropagate ∇_θ through equation (2)
- 11: **end for**
- 12: **end while**
- 13: **return** M_θ

By limiting ourselves to a subset of n_{CE} counterfactuals, we reduce runtimes; this approach has previously been shown to improve efficiency in the context of adversarial training (Kurakin, Goodfellow, and Bengio 2017; Kaufmann et al. 2022). To improve runtimes even more, we choose to first generate counterfactuals and then distribute them across mini-batches to benefit from greater degrees of parallelization during the counterfactual search. Alternatively, it is possible to generate counterfactuals separately for each mini-batch.²

3.2 Directly Inducing Explainability: Contrastive Divergence

As observed in prior related work (Grathwohl et al. 2020), any classifier can be re-interpreted as a joint energy-based model that learns to discriminate output classes conditional on the observed (training) samples from $p(x)$ and the generated samples from $p_\theta(x)$. The authors show that JEMs can be trained to perform well at both tasks by directly maximizing the joint log-likelihood: $\log p_\theta(x, y) = \log p_\theta(y|x) + \log p_\theta(x)$, where the first term can be optimized using cross-entropy as in Equation 2. To optimize $\log p_\theta(x)$, they minimize the contrastive divergence between the observed samples from $p(x)$ and samples generated from $p_\theta(x)$.

To generate samples, the paper introducing JEMs (Grathwohl et al. 2020) suggests relying on Stochastic Gradient Langevin Dynamics (SGLD) with an uninformative prior for initialization but we depart from this methodology: we propose to leverage counterfactual explainers to generate counterfactuals of observed training samples. Specifically, we have:

$$\text{div}(x^+, x'_{CE}, y^+; \theta) = \mathcal{E}_\theta(x^+, y^+) - \mathcal{E}_\theta(x'_{CE}, y^+) \quad (4)$$

where $\mathcal{E}_\theta(\cdot)$ denotes the energy function defined as $\mathcal{E}_\theta(x, y^+) = -M_\theta(x)[y^+]$, with y^+ denoting the index of the randomly drawn target class, $y^+ \sim p(y)$. Conditional on the target class y^+ , x'_{CE} denotes a mature counterfactual for a randomly sampled factual from a non-target class generated with a gradient-based CE generator for up to T iterations. Intuitively, the gradient of Equation 4 decreases the energy of observed training samples (positive samples) while increasing the energy of counterfactuals (negative samples) (Du and Mordatch 2020). As the counterfactuals get more plausible (Definition 3.1) during training, these opposing effects gradually balance each other out (Lippe 2024).

Since the maturity of counterfactuals in terms of a probability threshold is often reached before T , this form of sampling is not only more closely aligned with Definition 3.1., but can also speed up training times compared to SGLD. The departure from SGLD also allows us to tap into the vast repertoire of explainers that have been proposed in the literature to meet different desiderata. For example, many methods support domain and mutability constraints. In principle, any approach for generating CEs is viable, so long as it does not violate the faithfulness condition. Like JEMs (Murphy 2022), counterfactual training can be viewed as a form of contrastive representation learning.

3.3 Indirectly Inducing Explainability: Adversarial Robustness

Based on our analysis in Section 2, counterfactuals x' can be repurposed as additional training samples (Balashankar et al. 2023; Luu and Inoue 2023) or adversarial examples (Freiesleben 2022; Pawelczyk et al. 2022). This leaves some

²During initial prototyping of CT we also tested an implementation that relies on generating counterfactuals and adversarial examples at the batch level with no discernible difference in outcomes, but increased training times.

flexibility with regards to the choice for the advloss(\cdot) term in Equation 2. An intuitive functional form, but likely not the only sensible choice, is inspired by adversarial training:

$$\begin{aligned} \text{advloss}(\mathbf{M}_\theta(\mathbf{x}'_{\text{AE}}), \mathbf{y}; \varepsilon) &= \text{yloss}(\mathbf{M}_\theta(\mathbf{x}'_{t_\varepsilon}), \mathbf{y}) \\ t_\varepsilon &= \max_t \{t : \|\Delta_t\|_\infty < \varepsilon\} \end{aligned} \quad (5)$$

Under this choice, we consider nascent counterfactuals \mathbf{x}'_{AE} as AEs as long as the magnitude of the perturbation at time t (Δ_t) to any single feature is at most ε . The most strongly perturbed counterfactual $\mathbf{x}'_{t_\varepsilon}$ that still satisfies the condition is used as an adversarial example \mathbf{x}'_{AE} . This formalization is closely aligned with seminal work on adversarial machine learning (Szegedy et al. 2014), which defines an adversarial attack as an “imperceptible non-random perturbation”. Thus, we work with a different distinction between CE and AE than the one proposed in prior work (Freiesleben 2022), which considers misclassification as the distinguishing feature of adversarial examples. One of the key observations of our work is that we can leverage CEs during training and get AEs essentially for free to reap the benefits of adversarial training, leading to improved adversarial robustness and plausibility.

3.4 Encoding Actionability Constraints

Many existing counterfactual explainers support domain and mutability constraints. In fact, both types of constraints can be implemented for any explainer that relies on gradient descent in the feature space for optimization (Altmeyer, Deursen, and Liem 2023). 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. Similarly, mutability constraints can be enforced by setting partial derivatives to zero to ensure that features are only perturbed in the allowed direction, if at all.

As actionability constraints are binding at test time, we must also impose them when generating \mathbf{x}' during each training iteration to inform model representations. Through their effect on \mathbf{x}' , both types of constraints influence model outcomes via Equation 4. It is crucial that we avoid penalizing implausibility that arises from 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 defined for d . Specifically, we set $\mathbf{x}^+[d] := \mathbf{x}'[d]$ if:

1. Feature d is strictly immutable in practice.
2. $\mathbf{x}^+[d] > \mathbf{x}'[d]$, but d can only be decreased in practice.
3. $\mathbf{x}^+[d] < \mathbf{x}'[d]$, but 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 , i.e., we can model this as absolute certainty that the value $\mathbf{x}^+[d]$ remains the same as in the neighbor, $\mathbf{x}'[d]$, but it could be equivalently seen as masking changes to feature d . Intuitively, we can think of this as ignoring implausibility costs of immutable features, which effectively forces the model to instead seek plausibility through the remaining features. This can be expected to produce a classifier with relatively lower sensitivity to immutable features, and the higher relative sensitivity to mutable features should make mutability-constrained recourse less costly (see Section 4). Under certain conditions, this result also holds theoretically (for the proof, see the supplementary appendix):

Proposition 3.1 (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} where $y \in \{1, \dots, K\} = \mathcal{K}$, $\mathbf{x} \in \mathbb{R}^D$ and Θ is the matrix of coefficients with $\theta_{k,d} = \Theta[k, d]$ denoting the coefficient on feature d for class k . Assume multivariate Gaussian class densities with a common diagonal covariance matrix $\Sigma_k = \Sigma$ for all $k \in \mathcal{K}$, then protecting an immutable feature from the contrastive divergence penalty will result in lower classifier sensitivity to that feature relative to the remaining features, provided that at least one of those is discriminative and mutable.*

4 Experiments

We start by introducing the experimental setup, including performance metrics, datasets, algorithms, and explain our approach to evaluation in Section 4.1. Then, we address the research questions (RQ). Two questions relating to the principal goals of counterfactual training are presented in Section 4.2:

RQ 4.1. To what extent does the CT objective in Equation 2 induce models to learn plausible explanations?

RQ 4.2. To what extent does CT result in more favorable algorithmic recourse outcomes in the presence of actionability constraints

Next, in Section 4.3 we consider the performance of models trained with CT, focusing on their adversarial robustness but also commenting on the validity of generated CEs.

RQ 4.3. To what extent does CT influence the adversarial robustness of trained models?

Finally, in Section 4.4 we perform an ablation of the CT objective and evaluate its sensitivity to hyperparameters:

RQ 4.4. How does the CT objective depends on its individual components? (*ablation*)

RQ 4.5. What are the effects of hyperparameter selection on counterfactual training?

4.1 Experimental Setup

Our focus is the improvement in explainability (Definition 3.1). Thus, we mainly look at the plausibility and cost of faithfully generated counterfactuals at test time, but several other metrics are covered in the supplementary appendix. To measure the cost, we follow the standard proxy of distances (ℓ_1 -norm) between factuals and counterfactuals. For plausibility, we assess how similar CEs are to observed samples in the target domain, $\mathbf{X}' \subset \mathcal{X}^+$. For the evaluation, we rely on the metric proposed in prior work (Altmeyer et al. 2024) with ℓ_1 -norm for distances,

$$\text{IP}(\mathbf{x}', \mathbf{X}') = \frac{1}{|\mathbf{X}'|} \sum_{\mathbf{x} \in \mathbf{X}'} \text{dist}(\mathbf{x}', \mathbf{x}) \quad (6)$$

and introduce a novel divergence-based adaptation,

$$\text{IP}^*(\mathbf{X}', \mathbf{X}') = \text{MMD}(\mathbf{X}', \mathbf{X}') \quad (7)$$

where \mathbf{X}' denotes a collection of counterfactuals and $\text{MMD}(\cdot)$ is the unbiased estimate of the squared population maximum mean discrepancy (Gretton et al. 2012):

$$\begin{aligned} \text{MMD}(\mathbf{X}', \mathbf{X}') &= \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(x_i, x_j) \\ &\quad + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n k(\tilde{x}_i, \tilde{x}_j) \\ &\quad - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, \tilde{x}_j) \end{aligned} \quad (8)$$

with a kernel function $k(\cdot, \cdot)$. We use a characteristic Gaussian kernel with a constant length-scale parameter of 0.5, which means that the metric in Equation 7 is equal to zero if and only if the two distributions are exactly the same, $\mathbf{X}' = \mathbf{X}^+$.

To assess outcomes with respect to actionability for non-linear models, we look at the costs of (just) valid counterfactuals in terms of their distances from factual starting points with $\tau = 0.5$. While this is an imperfect proxy of sensitivity, we hypothesize that CT can reduce these costs by teaching models to seek plausibility with respect to mutable features, much like we observe in Figure 1 in panel (d) compared to (c). We supplement this analysis with estimates using integrated gradients (IG) (Sundararajan, Taly, and Yan 2017). To evaluate predictive performance, we use standard metrics, such as robust accuracy estimated on adversarially perturbed data using the fast gradient sign method (FGSM) (Goodfellow, Shlens, and Szegedy 2015) and projected gradient descent (PGD) (Madry et al. 2017).

We make use of nine classification datasets common in the CE/AR literature. Four of them are synthetic with two classes and different characteristics: linearly separable Gaussian clusters (*LS*), overlapping clusters (*OL*), concentric circles (*Circ*), and interlocking moons (*Moon*). Next, we have four real-world binary tabular datasets: *Adult* (Census data) (Becker and Kohavi 1996), California housing (*CH*) (Pace and Barry 1997), Default of Credit Card Clients (*Cred*) (Yeh 2016), and Give Me Some Credit (*GMSC*) (Kaggle 2011). Finally, for convenient illustration, we use the 10-class *MNIST* (LeCun 1998).

We run experiments with three gradient-based generators: *Generic* (Wachter, Mittelstadt, and Russell 2017) as a simple baseline; *REVISE* (Joshi et al. 2019) that aims to generate plausible counterfactuals using a surrogate Variational Autoencoder (VAE); and *ECCCo* (Altmeyer et al. 2024), targeting faithfulness. In all cases, we use standard logit cross-entropy loss for $y_{\text{loss}}(\cdot)$ and all generators penalize the distance (ℓ_1 -norm) of counterfactuals from their original factual state. *Generic* and *ECCCo* search for counterfactuals directly in the feature space; *REVISE* traverses the latent space of a variational autoencoder (VAE) fitted to the training data, so its outputs depend on the quality of the surrogate model. In addition to the distance penalty, *ECCCo* uses a penalty that regularizes the energy associated with the counterfactual, x' (Altmeyer et al. 2024). We omit the conformal set size penalty proposed in the original paper, since the authors found that faithfulness primarily depends on the energy penalty, freeing us from one additional hyperparameter.

Our method does not aim to be agnostic to the underlying CE generator and, as explained in Section 3.2, the selection of the CE generator can impact the explainability of models. To evaluate the specific value of counterfactual training, we extensively test the method using the three above-mentioned CE generators, which are characterized by varying complexity and desiderata, and we present the complete results in the supplementary appendix. Indeed, we observe that *ECCCo* outclasses the other two generators as the backbone of CT, generally leading to the highest reduction in implausibility. This is not surprising; the goals of *ECCCo* most closely align with the objectives of CT: maximally faithful explanations should also be the most useful for feedback. Conversely, we cannot expect the model to learn much from counterfactual explanations that largely depend on the quality of the surrogate model that is trained for *REVISE*. Similarly, *Generic* is a very simple baseline that optimizes only for minimal changes of features (measured in the original seminal paper (Wachter, Mittelstadt, and Russell 2017) using median absolute deviation).

Thus, while counterfactual training can be used with any gradient-based CE generator to improve the explainability of the resulting model, in Section 4.2 we mainly discuss its effectiveness with *ECCCo*, the strongest identified generator, allowing us to optimize the quality of the models. This constitutes our treatment method, but we still present the complete results for all generators in the supplementary appendix.

To assess the effects of CT, we investigate the improvements in performance metrics when using it on top of a weak baseline (BL), a naively (conventionally) trained multilayer perceptron (*MLP*), as the control method. As we hold all other things constant, this is the best way to get a clear picture of the improvement in explainability that can be directly attributed to CT. It is also consistent with the evaluation practices in the related literature (Goodfellow, Shlens, and Szegedy 2015; Ross, Lakkaraju, and Bastani 2024; Teney, Abbasnedjad, and Hengel 2020).

We also note that counterfactual training involves multiple objectives but our principal goal is high explainability as in Definition 3.1, while improved robustness is a welcome byproduct. We neither aim to outperform state-of-the-art approaches that target any single one of these objectives, nor do we claim that CT can achieve this. Specifically, we do not aim to beat JEMs with respect to their generative capacity, SOTA robust neural networks with respect to (adversarial) robustness, or (quasi-)Bayesian neural networks with respect to uncertainty quantification. As we have already explained in Section 2, existing literature has shown that all of these objectives tend to correlate (explaining some of our positive findings), but we situate counterfactual training squarely in the context of (counterfactual) explainability and algorithmic recourse, where it tackles an important shortcoming of existing approaches.

In terms of computing resources, all of our experiments were executed on a high-performance cluster. We have relied on distributed computing across multiple central processing units (CPU); for example, the hyperparameter grid searches were carried out on 34 CPUs with 2GB memory each. Graphical processing units (GPU) were *not* used. All computations were performed in the Julia Programming Language (Bezanson et al. 2017); our code base (algorithms and experimental settings) has been open-sourced on GitHub.³ We explain more about the hardware, software, and reproducibility considerations in the supplementary appendix.

4.2 Main Results

Our main results for plausibility and actionability for *MLP* models are summarized in Table 1 that presents counterfactual outcomes grouped by dataset along with standard errors averaged across bootstrap samples. Asterisks (*) are used when the bootstrapped 99%-confidence interval of differences in mean outcomes does *not* include zero, so the observed effects are statistically significant at the 0.01 level. As our experimental procedure is (by virtue of the proposed method) relatively complex, we choose to work at this stringent alpha level to demonstrate the high reliability of counterfactual training.

The first two columns (IP and IP*) show the percentage reduction in implausibility for our two metrics when using CT on top of the weak baseline. As an example, consider the first row for *LS* data: the observed positive values

³<https://github.com/JuliaTrustworthyAI/CounterfactualTraining.jl>

Table 1 Key evaluation metrics for valid counterfactual along with bootstrapped standard errors for all datasets. **Plausibility** (columns 1-2): percentage reduction in implausibility for **IP** and **IP***, respectively; **Cost / Actionability** (column 3): percentage reduction in costs when selected features are protected. Outcomes are aggregated across bootstrap samples (100 rounds) and varying degrees of the energy penalty λ_{egy} used for ECCCo at test time. Asterisks (*) indicate that the bootstrapped 99%-confidence interval of differences in mean outcomes does **not** include zero.

Data	IP (-%)	IP* (-%)	Cost (-%)
LS	$26.26 \pm 0.67^*$	$51.28 \pm 2.01^*$	$16.41 \pm 0.57^*$
Circ	$58.88 \pm 0.37^*$	$93.84 \pm 6.70^*$	$42.99 \pm 0.85^*$
Moon	$19.59 \pm 0.73^*$	8.00 ± 9.44	$5.16 \pm 1.00^*$
OL	-1.93 ± 1.12	-27.70 ± 14.59	$40.86 \pm 2.30^*$
Adult	0.19 ± 1.05	$34.35 \pm 5.61^*$	4.03 ± 4.03
CH	$10.65 \pm 1.47^*$	$63.06 \pm 4.25^*$	$44.23 \pm 1.43^*$
Cred	$10.14 \pm 1.59^*$	$50.35 \pm 12.26^*$	$-18.17 \pm 4.40^*$
GMSC	$10.65 \pm 2.28^*$	$24.75 \pm 4.84^*$	$66.01 \pm 1.41^*$
MNIST	$6.36 \pm 1.70^*$	-70.31 ± 217.60	$-35.11 \pm 6.96^*$
Avg.	15.64	25.29	18.49

indicate that faithful counterfactuals are around 26-51% more plausible for models trained with CT, in line with our observations in panel (b) of Figure 1 compared to panel (a).

The third column shows the results for a scenario when mutability constraints are imposed on the selected features. Again, we are comparing CT to the baseline, so reductions in the positive direction imply that valid counterfactuals are “cheaper” (more actionable) when using CT with feature protection. Relating this back to Figure 1, the third column represents the reduction in distances traveled by counterfactuals in panel (d) compared to panel (c). In the following paragraphs, we summarize the results for all datasets.

Plausibility (RQ 4.1)

CT generally produces substantial and statistically significant improvements in plausibility.

Average reductions in IP range from around 6% for *MNIST* to almost 60% for *Circ*. For the real-world tabular datasets they are around 10% for *CH*, *Cred* and *GMSC*; for *Adult* and *OL* we find no significant impact of CT on IP. The former is subject to a large proportion of categorical features, which inhibits the generation of large numbers of valid counterfactuals during training and may therefore explain this finding.

Reductions in IP* are even more substantial and generally statistically significant, although the average degree of uncertainty is higher than for IP: reductions range from around 25% (*GMSC*) to more than 90% (*Circ*). The only negative findings are for *OL* and *MNIST*, but they are insignificant. A qualitative inspection of the counterfactuals in Figure 2 suggests recognizable digits for the model trained with CT (bottom row), unlike the baseline (top row).



Figure 2 *Plausibility*: BL (top row) vs CT using the *ECCCo* generator (bottom row) counterfactuals for a randomly selected factual from class “0” (in blue). CT produces more plausible counterfactuals than BL.

Actionability (RQ 4.2)

CT tends to improve actionability in the presence of immutable features, but this is not guaranteed if the assumptions in Proposition 3.1 are violated.

For synthetic datasets, we always protect the first feature; for all real-world tabular datasets we could identify and protect an *age* variable; for *MNIST*, we protect the five top and five bottom rows of pixels of the full image. Statistically

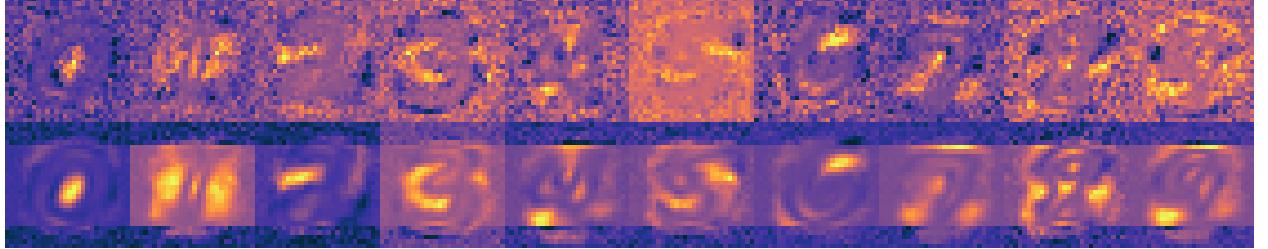


Figure 3 *Actionability*: Sample visual explanations (integrated gradients) for all classes in the *MNIST* dataset. Top and bottom rows of images show the results for BL and CT, respectively. Mutability constraints are imposed on the five top and five bottom rows of pixels. CT is less sensitive to protected features.

significant reductions in costs overwhelmingly point in the positive direction reaching up to around 66% for *GMSC* data. Only in the case of *Cred* and *MNIST*, average costs increase, most likely because any benefits from protecting features are outweighed by an increase in costs required for greater plausibility. With respect to *MNIST* in particular, the weak baseline is susceptible to cheap adversarial attacks that significantly less costly to achieve that plausible counterfactuals. Finally, the findings for *Adult* are insignificant.

To further empirically evaluate the feature protection mechanism of CT beyond linear models covered in Proposition 3.1, we make use of integrated gradients (IG) (Sundararajan, Taly, and Yan 2017). IG calculates the contribution of each input feature towards a specific prediction by approximating the integral of the model output with respect to its input, using a set of samples that linearly interpolate between a test instance and some baseline instance. This process produces a vector of real numbers, one per input feature, which informs about the contribution of each feature to the prediction. The selection of an appropriate baseline is an important design decision (Sundararajan, Taly, and Yan 2017); to remain consistent in our evaluations, we use a baseline drawn at random from the uniform distribution $\mathcal{U}(-1, 1)$ for all datasets, which aligns with standard evaluation practices for IG. As the outputs are not bounded (i.e., they are real numbers), we standardize the integrated gradients across features to allow for a meaningful comparison of the results for different models.

Qualitatively, the class-conditional integrated gradients in Figure 3 suggest that CT has the expected effect even for non-linear models: the model trained with CT (bottom row) is less sensitive (blue) to the five top and five bottom rows of pixels that were protected. Quantitatively, we observe substantial improvements for seven out of nine datasets, and inconclusive results for the remaining two datasets. Table 2 shows the average sensitivity to protected features measured by standardized integrated gradients for CT and BL along with 95% bootstrap confidence intervals: for the synthetic datasets, we observe strong reductions in sensitivity to the protected features for *LS*, *OL* and *OL*, in line with expectations. For the *Moon* dataset, the effect of feature protection is less pronounced but still in the expected direction. We also observe that confidence intervals are in some cases much tighter for models trained with CT: less noisy estimates for integrated gradients likely indicate that the model is more regularized and can be expected to behave more consistently across samples.

For real-world datasets, the sensitivity to the protected *age* variable is reduced by approximately a third for *Adult*, 20% for *CH*, and more than half for protected pixels in *MNIST*, mirroring the qualitative findings in Figure 3. In case of *Cred*, CT fully prevents the model from considering *age* as a factor in classification, with sensitivity reduced to zero. Only for *GMSC*, we observe negative impacts of CT, which we believe is due to any or all of the following: a) data assumptions are violated; b) the impact of other components of the CT objective outweighs expected effects of feature protection; or c) the baseline choice applied consistently to all datasets is not appropriate for *GMSC*.

4.3 Predictive Performance

Adversarial Robustness (RQ 4.3)

Models trained with CT are much more robust to gradient-based adversarial attacks than conventionally-trained (weak) baselines.

Test accuracies on clean and adversarially perturbed test data are shown in Figure 4. The perturbation size, $\varepsilon \in [0, 0.1]$, increases along the horizontal axis, where the case of $\varepsilon = 0$ corresponds to standard test accuracy for non-perturbed data. For synthetic datasets, predictive performance is virtually unaffected by perturbations for all models; those results are therefore omitted from Figure 4 in favor of better illustrations for the real-world data.

Table 2 Median sensitivity to protected features measured by standardized integrated gradients. Square brackets enclose 95% bootstrap confidence intervals.

Dataset	CT	BL
LS	0.21 [0.20, 0.22]	30.69 [12.92, 629.20]
Circ	6.96 [4.88, 20.62]	19.20 [6.48, 193.92]
Moons	0.54 [0.41, 0.68]	0.66 [0.53, 0.92]
Over	0.59 [0.38, 0.79]	24.55 [8.31, 466.26]
Adult	0.48 [0.41, 0.52]	0.74 [0.56, 0.91]
CH	0.04 [0.01, 0.06]	0.05 [0.03, 0.09]
Cred	0.00 [0.00, 0.00]	0.20 [0.18, 0.25]
GMSC	0.71 [0.58, 0.85]	0.16 [0.11, 0.23]
MNIST	0.17 [0.16, 0.17]	0.35 [0.33, 0.37]

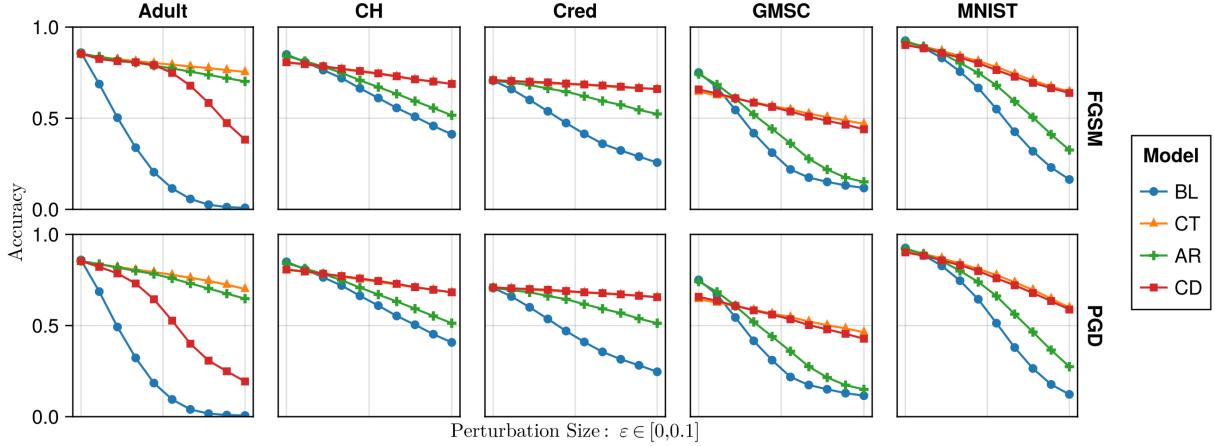


Figure 4 Test accuracies on adversarially perturbed data with varying perturbation sizes for the non-synthetic datasets. Different training objectives are distinguished by color and shape: (1) BL—the weak baseline; (2) CT—the full CT objective; (3) AR—a partial CT objective without contrastive divergence; (4) CD—a partial CT objective without adversarial loss. Top and bottom rows show the results for FGSM and PGD (40 steps at step size $\eta = 0.01$), respectively.

Focusing on the curves for CT and BL in Figure 4 for the moment,⁴ we find that standard test accuracy ($\varepsilon = 0$) is largely unaffected by CT, while robustness against both types of attacks—FGSM (top row) and PGD (bottom row)—is greatly improved: while in some cases robust accuracies for the weak baseline drop to virtually zero (worse than random guessing) for large enough perturbation sizes, accuracies of CT models remain remarkably robust, even though robustness is not the primary objective of counterfactual training. In the only case where standard accuracy on unperturbed test data is substantially reduced for CT (GSMC), we note that robust accuracy decreases particularly fast for the weak baseline as the perturbation size increases. This seems to indicate that the standard accuracy for the weak baseline is inflated by sensitivity to meaningless associations in the data.

We also look at the validity of generated counterfactuals, or the proportion of counterfactuals that attain the target class, as presented in Table 3. We find that in many cases CT leads to substantial reductions in average validity, but this effect does not seem to be strongly influenced by the imposed mutability constraints (columns 1-2 vs columns 3-4). This result does not surprise us: by design, CT shrinks the solution space for valid counterfactual explanations, thus making it “harder” (and yet not “more costly”) to reach validity compared to the baseline model. As further discussed in the supplementary appendix, this should not be seen as a shortcoming of the method for a number of reasons: validity rates can be increased with longer searches; costs of found solutions still generally decrease, as we observe in our

⁴The results for AR and CD are discussed in the context of ablation below.

Table 3 Average validity of counterfactuals for CT vs BL. First two columns correspond to no mutability constraints imposed on the features; last two columns involve mutability constraints imposed on the specified features.

Data	CT mut.	BL mut.	CT constr.	BL constr.
LS	1.0	1.0	1.0	1.0
Circ	1.0	0.51	0.71	0.48
Moon	1.0	1.0	1.0	0.98
OL	0.86	0.98	0.34	0.56
Adult	0.68	0.99	0.7	0.99
CH	1.0	1.0	1.0	1.0
Cred	0.72	1.0	0.74	1.0
GMSC	0.94	1.0	0.97	1.0
MNIST	1.0	1.0	1.0	1.0
Avg.	0.91	0.94	0.83	0.89

experiments; and achieving high validity does not entail that explanations are practical for the recipients (e.g., valid solutions may still be extremely costly) ([Venkatasubramanian and Alfano 2020](#)).

4.4 Ablation and Hyperparameter Settings

In this subsection, we use ablation studies to investigate how the different components of the counterfactual training objective in Equation 2 affect outcomes. Beyond this, we are also interested in understanding how CT depends on various other hyperparameters. To this end, we present the results from extensive grid searches run across all synthetic datasets.

Ablation (RQ 4.4)

All components of the CT objective affect outcomes, even independently, but the full objective achieves the most consistent improvements wrt. our goals.

We ablate the effect of both (1) the contrastive divergence component and (2) the adversarial loss included in the full CT objective in Equation 2. In the following, we refer to the resulting partial objectives as adversarial robustness (AR) and contrastive divergence (CD), respectively. We note that AR corresponds to a form of adversarial training and the CD objective is similar to that of a joint energy-based model. Therefore, the ablation also serves as a comparison of counterfactual training to stronger baselines, although we emphasize again that we do not seek to outperform SOTA methods in the domains of generative or robust machine learning, focusing CT squarely on models with high explainability and actionability in the context of algorithmic recourse.

Firstly, we find that both components play an important role in shaping final outcomes. Both AR and CD can independently improve the plausibility and adversarial robustness of models.

Concerning plausibility, Figure 5 shows the percentage reductions in implausibility for the partial and full objectives compared to the weak baseline. The results for IP and IP* are shown in the top and bottom graphs, respectively, and the datasets are differentiated by color. We find that in the best identified hyperparameter settings, results for the full objective are predominantly affected by the contrastive divergence component, but the inclusion of adversarial loss leads to additional improvements for some datasets (*Adult*, *MNIST*). We penalize contrastive divergence twice as strongly as adversarial loss, which may explain why the former dominates. The outcome for *Adult*, in particular, demonstrates the benefit of including both components: as noted earlier, the large proportion of categorical features in this dataset seems to inhibit the generation of valid counterfactuals, which in turn appears to diminish the effect of the contrastive divergence component.

Looking at AR alone, we find that it produces mixed results for IP, with strong positive results nonetheless dominating overall, reflecting previous findings from the related literature. In particular, for real-world tabular datasets, adversarial robustness seems to substantially benefit plausibility. In these cases, the inclusion of the AR component in the full objective also helps to substantially improve outcomes in relation to the partial CD objective: improvements in plausibility for the *Adult* and *MNIST* datasets are notably higher for full CT. In some cases—most notably *GMSC* and *Cred*—the full CT objective does not outperform the partial objectives, but still achieves the highest levels of adversarial robustness (Figure 4).

Zooming in on adversarial robustness, we find that the full CT objective consistently outperforms the partial objectives, which individually yield improvements. Consistent with the existing literature on JEMs ([Grathwohl et al. 2020](#)), CD

yields substantially more robust models than the weak baseline at varying perturbation sizes (Figure 4). Similarly, AR yields consistent improvements in robustness, as expected. Still, we observe that in cases where either CD or AR show signs of degrading robust accuracy at higher perturbation sizes, the full CT objective maintains robustness. Much like in the context of plausibility, CT benefits from both components, highlighting the effectiveness of our approach to reusing nascent counterfactuals as AEs.

In summary, we find that the full CT objective strikes a balance between both components, thereby leading to the most consistent improvements with respect to plausibility and adversarial robustness.

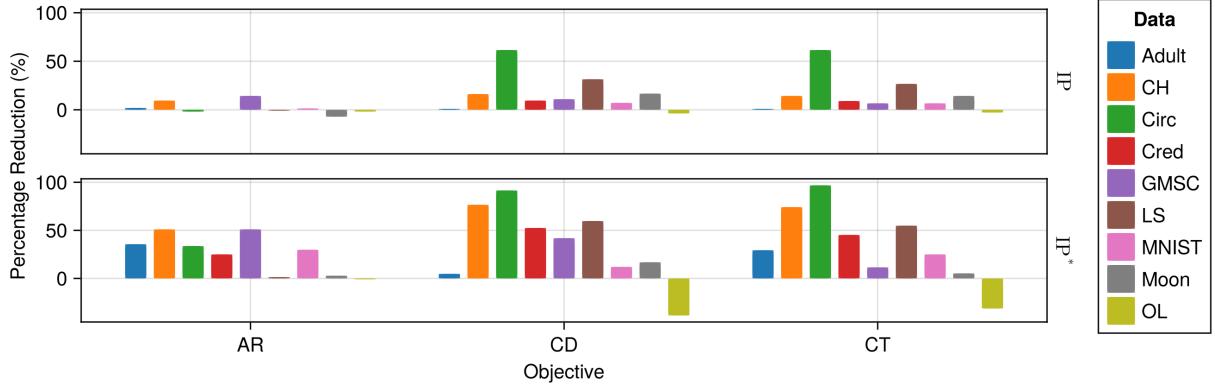


Figure 5 Percentage reductions in implausibility for the partial (AR, CD) and full (CT) objectives compared to the weak baseline. The results for IP and IP* are shown in the top and bottom graphs, respectively, and the datasets are differentiated by color.

Hyperparameter settings (RQ 4.5)

CT is quite sensitive to the choice of a CE generator and its hyperparameters but (1) we observe manageable patterns, and (2) we can usually identify settings that improve either plausibility or actionability, and typically both of them at the same time.

We evaluate the impacts of three types of hyperparameters on CT. In the following, we focus on the highlights and make the full results available in the supplementary appendix.

Firstly, we find that optimal results are generally obtained when using *ECCCo* to generate counterfactuals. Conversely, using a generator that may inhibit faithfulness (*REVISE*), regularly yields smaller improvements in plausibility and is more likely to even increase implausibility. The results of the grid search for *REVISE* also exhibit higher variability than the results for *ECCCo* and *Generic*. As argued above, this finding confirms our intuition that maximally faithful explanations are most suitable for counterfactual training.

Concerning hyperparameters that guide the gradient-based counterfactual search, we find that increasing T , the maximum number of steps, generally yields better outcomes because more CEs can mature. Relatedly, we also find that the effectiveness and stability of CT is positively associated with the total number of counterfactuals generated during each training epoch. The impact of τ , the decision threshold, is more difficult to predict. On “harder” datasets it may be difficult to satisfy high τ for any given sample (i.e., also factuals) and so increasing this threshold does not seem to correlate with better outcomes. In fact, $\tau = 0.5$ generally leads to optimal results as it is associated with high proportions of mature counterfactuals. This is likely because the special case of $\tau = 0.5$ corresponds to equal class probabilities, so a counterfactual is considered mature when the logit for the target class is higher than the logits for all other classes.

Secondly, the strength of the energy regularization, λ_{reg} , is highly impactful and should be set sufficiently high to avoid common problems associated with exploding gradients. The sensitivity with respect to λ_{div} and λ_{adv} is much less evident. While high values of λ_{reg} may increase the variability in outcomes when combined with high values of λ_{div} or λ_{adv} , this effect is not particularly pronounced. These results mirror our observations from the ablation studies and lend further weight to the argument that CT benefits from both components.

Finally, we also observe desired improvements when CT was combined with conventional training and employed only for the final 50% of epochs of the complete training process. Put differently, CT can improve the explainability of models in a post-hoc, fine-tuning manner.

5 Discussion

As our results indicate, counterfactual training achieves its objective of producing models that are more explainable. Nonetheless, these advantages come with certain limitations.

Immutable features may have proxies. We propose a method to modify the sensitivity of a model to certain features, and thus increase the actionability of the generated CEs. However, it requires that model owners define the mutability constraints for (all) features considered by the model. Even if all immutable features are protected, there may exist proxies that are theoretically mutable (and hence should not be protected) but preserve enough information about the principals to hinder these protections. Delineating actionability is a major open challenge in the AR literature (see, e.g., ([Venkatasubramanian and Alfano 2020](#))) impacting the capacity of CT to fulfill its intended goal.

Interventions on features may have implications for fairness. Modifying the sensitivity of a model to certain features may also have implications for the fair and equitable treatment of decision subjects. Model owners could misuse this solution by enforcing explanations based on features that are more difficult to modify by some (group of) decision subjects. For example, consider the *Adult* dataset used in our experiments, where *workclass* or *education* may be more difficult to change for underprivileged groups. When applied irresponsibly, CT could result in an unfairly assigned burden of recourse ([Sharma, Henderson, and Ghosh 2020](#)), threatening the equality of opportunity in the system ([Bell et al. 2024](#)). Nonetheless, these phenomena are not specific to CT.

Plausibility is costly. As noted before, more plausible counterfactuals are inevitably more costly ([Altmeyer et al. 2024](#)). CT improves plausibility and robustness, but this can negatively affect average costs and validity whenever cheap, implausible, and adversarial explanations are removed from the solution space.

CT increases training times. Just like contrastive and robust learning, CT is more resource-intensive than conventional regimes. Three factors mitigate this effect: (1) CT yields itself to parallel execution; (2) it amortizes the cost of CEs for the training samples; and (3) our preliminary findings suggest that it can be used to fine-tune conventionally-trained models.

We also highlight three key directions for future research. Firstly, it is an interesting challenge to extend CT beyond classification settings. Our formulation relies on the distinction between target and non-target classes, requiring the output space to be discrete. Thus, it does not apply to ML tasks where the change in outcome cannot be readily discretized. Classification remains the focus of CE and algorithmic recourse research; other settings have attracted some interest (e.g., regression ([Spooner et al. 2021](#))), but there is little consensus on how to extend the notion of CEs.

Secondly, our analysis covers CE generators with different characteristics, but it is interesting to extend it to more algorithms, including ones that do not rely on computationally costly gradient-based optimization. This should reduce training costs while possibly preserving the benefits of CT.

Finally, we believe that it is possible to considerably improve hyperparameter selection procedures. Our method benefits from the tuning of certain key hyperparameters but we have relied exclusively on grid searches. Future work on CT could benefit from more sophisticated approaches. Notably, CT is iterative, which makes methods such as Bayesian or gradient-based optimization applicable (see, e.g., ([Bischl et al. 2023](#))).

6 Conclusion

State-of-the-art machine learning models are prone to learning complex representations that cannot be interpreted by humans. Existing work on counterfactual explanations has largely focused on designing tools to generate plausible and actionable explanations for any model. In this work, we instead hold models accountable for delivering such explanations. We introduce counterfactual training: a novel training regime that integrates recent advances in contrastive learning, adversarial robustness, and CE to incentivize highly explainable models. Through theoretical results and extensive experiments, we demonstrate that CT satisfies this goal while promoting adversarial robustness of models. Explanations generated from CT-based models are both more plausible (compliant with the underlying data-generating process) and more actionable (compliant with user-specified mutability constraints), and thus meaningful to recipients. In turn, our work highlights the value of simultaneously improving models and their explanations.

Acknowledgment

Some of the authors were partially funded by ICAI AI for Fintech Research, an ING—TU Delft collaboration. Research reported in this work was partially facilitated by computational resources and support of the DelftBlue high-performance computing cluster at TU Delft ([\(DHPC\) 2022](#)).

References

- Abbasnejad, Ehsan, Damien Teney, Amin Parvaneh, Javen Shi, and Anton van den Hengel. 2020. “Counterfactual Vision and Language Learning.” In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 10041–51. <https://doi.org/10.1109/CVPR42600.2020.01006>.
- Altmeyer, Patrick, Arie van Deursen, and Cynthia C. S. Liem. 2023. “Explaining Black-Box Models Through Counterfactuals.” In *Proceedings of the JuliaCon Conferences*, 1:130.
- Altmeyer, Patrick, Mojtaba Farmanbar, Arie van Deursen, and Cynthia C. S. Liem. 2024. “Faithful Model Explanations through Energy-Constrained Conformal Counterfactuals.” In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence*, 38:10829–37. 10. <https://doi.org/10.1609/aaai.v38i10.28956>.
- Augustin, Maximilian, Alexander Meinke, and Matthias Hein. 2020. “Adversarial Robustness on in- and Out-Distribution Improves Explainability.” In *Computer Vision – ECCV 2020*, edited by Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, 228–45. Cham: Springer.
- Balashankar, Ananth, Xuezhi Wang, Yao Qin, Ben Packer, Nithum Thain, Ed Chi, Jilin Chen, and Alex Beutel. 2023. “Improving Classifier Robustness Through Active Generative Counterfactual Data Augmentation.” In *Findings of the Association for Computational Linguistics: EMNLP 2023*, 127–39. ACL. <https://doi.org/10.18653/v1/2023.findings-emnlp.10>.
- Becker, Barry, and Ronny Kohavi. 1996. “Adult.” UCI Machine Learning Repository.
- Bell, Andrew, Joao Fonseca, Carlo Abrate, Francesco Bonchi, and Julia Stoyanovich. 2024. “Fairness in Algorithmic Recourse Through the Lens of Substantive Equality of Opportunity.” <https://arxiv.org/abs/2401.16088>.
- Bezanson, Jeff, Alan Edelman, Stefan Karpinski, and Viral B Shah. 2017. “Julia: A Fresh Approach to Numerical Computing.” *SIAM Review* 59 (1): 65–98. <https://doi.org/10.1137/141000671>.
- Bischl, Bernd, Martin Binder, Michel Lang, Tobias Pielok, Jakob Richter, Stefan Coors, Janek Thomas, et al. 2023. “Hyperparameter optimization: Foundations, algorithms, best practices, and open challenges.” *WIREs Data Mining and Knowledge Discovery* 13 (2): e1484. <https://doi.org/10.1002/widm.1484>.
- Bouchet-Valat, Milan, and Bogumi Kamiski. 2023. “DataFrames.jl: Flexible and Fast Tabular Data in Julia.” *Journal of Statistical Software* 107 (4): 1–32. <https://doi.org/10.18637/jss.v107.i04>.
- Byrne, Simon, Lucas C. Wilcox, and Valentin Churavy. 2021. “MPI.jl: Julia Bindings for the Message Passing Interface.” *Proceedings of the JuliaCon Conferences* 1 (1): 68. <https://doi.org/10.21105/jcon.00068>.
- Chagas, Ronan Arraes Jardim, Ben Baumgold, Glen Hertz, Hendrik Ranocha, Mark Wells, Nathan Boyer, Nicholas Ritchie, et al. 2024. “Ronisbr/PrettyTables.jl: V2.4.0.” Zenodo. <https://doi.org/10.5281/zenodo.1383553>.
- Christ, Simon, Daniel Schwabeneder, Christopher Rackauckas, Michael Krabbe Borregaard, and Thomas Breloff. 2023. “Plots.jl – a User Extendable Plotting API for the Julia Programming Language.” <https://doi.org/https://doi.org/10.5334/jors.431>.
- Danisch, Simon, and Julius Krumbiegel. 2021. “Makie.jl: Flexible High-Performance Data Visualization for Julia.” *Journal of Open Source Software* 6 (65): 3349. <https://doi.org/10.21105/joss.03349>.
- (DHPC), Delft High Performance Computing Centre. 2022. “DelftBlue Supercomputer (Phase 1).” <https://www.tude-lft.nl/dhpc/ark/44463/DelftBluePhase1>.
- Du, Yilun, and Igor Mordatch. 2020. “Implicit Generation and Generalization in Energy-Based Models.” <https://arxiv.org/abs/1903.08689>.
- Freiesleben, Timo. 2022. “The Intriguing Relation Between Counterfactual Explanations and Adversarial Examples.” *Minds and Machines* 32 (1): 77–109.
- Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press.
- Goodfellow, Ian, Jonathon Shlens, and Christian Szegedy. 2015. “Explaining and Harnessing Adversarial Examples.” <https://arxiv.org/abs/1412.6572>.
- Grathwohl, Will, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. 2020. “Your Classifier Is Secretly an Energy Based Model and You Should Treat It Like One.” In *International Conference on Learning Representations*.
- Gretton, Arthur, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. 2012. “A Kernel Two-Sample Test.” *The Journal of Machine Learning Research* 13 (1): 723–73.
- Guo, Hangzhi, Thanh H. Nguyen, and Amulya Yadav. 2023. “CounterNet: End-to-End Training of Prediction Aware Counterfactual Explanations.” In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 577–589. KDD ’23. New York, NY, USA: Association for Computing Machinery. <https://doi.org/10.1145/3580305.3599290>.
- Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2009. *The Elements of Statistical Learning*. Springer New York. <https://doi.org/10.1007/978-0-387-84858-7>.
- Innes, Michael, Elliot Saba, Keno Fischer, Dhairyा Gandhi, Marco Concetto Rudilosso, Neethu Mariya Joy, Tejan Karmali, Avik Pal, and Viral Shah. 2018. “Fashionable Modelling with Flux.” <https://arxiv.org/abs/1811.01457>.

- Innes, Mike. 2018. “Flux: Elegant Machine Learning with Julia.” *Journal of Open Source Software* 3 (25): 602. <https://doi.org/10.21105/joss.00602>.
- Joshi, Shalmali, Oluwasanmi Koyejo, Warut Vigitbenjaronk, Been Kim, and Joydeep Ghosh. 2019. “Towards realistic individual recourse and actionable explanations in black-box decision making systems.” <https://arxiv.org/abs/1907.09615>.
- Kaggle. 2011. “Give Me Some Credit, Improve on the State of the Art in Credit Scoring by Predicting the Probability That Somebody Will Experience Financial Distress in the Next Two Years.” <https://www.kaggle.com/c/GiveMeSomeCredit>; Kaggle. <https://www.kaggle.com/c/GiveMeSomeCredit>.
- Karimi, Amir-Hossein, Gilles Barthe, Bernhard Schölkopf, and Isabel Valera. 2021. “A Survey of Algorithmic Recourse: Definitions, Formulations, Solutions, and Prospects.” <https://arxiv.org/abs/2010.04050>.
- Kaufmann, Maximilian, Yiren Zhao, Ilia Shumailov, Robert Mullins, and Nicolas Papernot. 2022. “Efficient Adversarial Training with Data Pruning.” *arXiv Preprint arXiv:2207.00694*.
- Kurakin, Alexey, Ian Goodfellow, and Samy Bengio. 2017. “Adversarial Machine Learning at Scale.” <https://arxiv.org/abs/1611.01236>.
- Lakshminarayanan, Balaji, Alexander Pritzel, and Charles Blundell. 2017. “Simple and Scalable Predictive Uncertainty Estimation Using Deep Ensembles.” In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 6405–16. NIPS’17. Red Hook, NY, USA: Curran Associates Inc.
- LeCun, Yann. 1998. “The MNIST database of handwritten digits.” <http://yann.lecun.com/exdb/mnist/>.
- Lippe, Phillip. 2024. “UvA Deep Learning Tutorials.” <https://uvadlc-notebooks.readthedocs.io/en/latest/>.
- Luu, Hoai Linh, and Naoya Inoue. 2023. “Counterfactual Adversarial Training for Improving Robustness of Pre-Trained Language Models.” In *Proceedings of the 37th Pacific Asia Conference on Language, Information and Computation*, 881–88. ACL. <https://aclanthology.org/2023.pacific-1.88/>.
- Madry, Aleksander, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2017. “Towards Deep Learning Models Resistant to Adversarial Attacks.” *arXiv Preprint arXiv:1706.06083*.
- Molnar, Christoph. 2022. *Interpretable Machine Learning: A Guide for Making Black Box Models Explainable*. 2nd ed. Christoph Molnar. <https://christophm.github.io/interpretable-ml-book>.
- Murphy, Kevin P. 2022. *Probabilistic Machine Learning: An Introduction*. MIT Press.
- Pace, R Kelley, and Ronald Barry. 1997. “Sparse Spatial Autoregressions.” *Statistics & Probability Letters* 33 (3): 291–97. [https://doi.org/10.1016/s0167-7152\(96\)00140-x](https://doi.org/10.1016/s0167-7152(96)00140-x).
- Pawelczyk, Martin, Chirag Agarwal, Shalmali Joshi, Sohini Upadhyay, and Himabindu Lakkaraju. 2022. “Exploring Counterfactual Explanations Through the Lens of Adversarial Examples: A Theoretical and Empirical Analysis.” In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, edited by Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, 151:4574–94. Proceedings of Machine Learning Research. PMLR. <https://proceedings.mlr.press/v151/pawelczyk22a.html>.
- Ross, Alexis, Himabindu Lakkaraju, and Osbert Bastani. 2024. “Learning Models for Actionable Recourse.” In *Proceedings of the 35th International Conference on Neural Information Processing Systems*. NIPS ’21. Red Hook, NY, USA: Curran Associates Inc.
- Sauer, Axel, and Andreas Geiger. 2021. “Counterfactual Generative Networks.” <https://arxiv.org/abs/2101.06046>.
- Schut, Lisa, Oscar Key, Rory McGrath, Luca Costabello, Bogdan Sacaleanu, Yarin Gal, et al. 2021. “Generating Interpretable Counterfactual Explanations by Implicit Minimisation of Epistemic and Aleatoric Uncertainties.” In *International Conference on Artificial Intelligence and Statistics*, 1756–64. PMLR.
- Sharma, Shubham, Jette Henderson, and Joydeep Ghosh. 2020. “CERTIFAI: A Common Framework to Provide Explanations and Analyse the Fairness and Robustness of Black-Box Models.” In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 166–72. AIES ’20. New York, NY, USA: Association for Computing Machinery. <https://doi.org/10.1145/3375627.3375812>.
- Spooner, Thomas, Danial Dervovic, Jason Long, Jon Shepard, Jiahao Chen, and Daniele Magazzeni. 2021. “Counterfactual Explanations for Arbitrary Regression Models.” <https://arxiv.org/abs/2106.15212>.
- Sturmels, Pascal, Scott Lundberg, and Su-In Lee. 2020. “Visualizing the Impact of Feature Attribution Baselines.” *Distill* 5 (1): e22.
- Sundararajan, Mukund, Ankur Taly, and Qiqi Yan. 2017. “Axiomatic Attribution for Deep Networks.” <https://arxiv.org/abs/1703.01365>.
- Szegedy, Christian, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2014. “Intriguing Properties of Neural Networks.” <https://arxiv.org/abs/1312.6199>.
- Teh, Yee Whye, Max Welling, Simon Osindero, and Geoffrey E. Hinton. 2003. “Energy-Based Models for Sparse Overcomplete Representations.” *J. Mach. Learn. Res.* 4 (null): 1235–60.
- Teney, Damien, Ehsan Abbasnejad, and Anton van den Hengel. 2020. “Learning What Makes a Difference from Counterfactual Examples and Gradient Supervision.” In *Computer Vision - ECCV 2020*, 580–99. Berlin, Heidelberg: Springer-Verlag. https://doi.org/10.1007/978-3-030-58607-2_34.

- Ustun, Berk, Alexander Spangher, and Yang Liu. 2019. “Actionable Recourse in Linear Classification.” In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 10–19. <https://doi.org/10.1145/3287560.3287566>.
- Venkatasubramanian, Suresh, and Mark Alfano. 2020. “The Philosophical Basis of Algorithmic Recourse.” In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 284–93. FAT* ’20. New York, NY, USA: Association for Computing Machinery. <https://doi.org/10.1145/3351095.3372876>.
- Verma, Sahil, Varich Boonsanong, Minh Hoang, Keegan E. Hines, John P. Dickerson, and Chirag Shah. 2022. “Counterfactual Explanations and Algorithmic Recourses for Machine Learning: A Review.” <https://arxiv.org/abs/2010.10596>.
- Wachter, Sandra, Brent Mittelstadt, and Chris Russell. 2017. “Counterfactual Explanations Without Opening the Black Box: Automated Decisions and the GDPR.” *Harv. JL & Tech.* 31: 841. <https://doi.org/10.2139/ssrn.3063289>.
- Wilson, Andrew Gordon. 2020. “The Case for Bayesian Deep Learning.” <https://arxiv.org/abs/2001.10995>.
- Wu, Tongshuang, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2021. “Polyjuice: Generating Counterfactuals for Explaining, Evaluating, and Improving Models.” In *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)*, edited by Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, 6707–23. Online: ACL. <https://doi.org/10.18653/v1/2021.acl-long.523>.
- Yeh, I-Cheng. 2016. “Default of Credit Card Clients.” UCI Machine Learning Repository.

Appendix A Notation

A.1 Variables and Parameters

Below we provide an overview of some notation used frequently throughout the paper:

- \mathcal{Y} : The output domain.
- 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.
- \mathcal{X} : The input domain.
- \mathbf{x} : a single training sample.
- \mathbf{x}' : a counterfactual.
- $t = 1, \dots, T$: Step indicator for counterfactual search iterations.
- \mathbf{x}'_{AE} : a nascent counterfactual, defined as a counterfactual that has not yet converged.
- \mathbf{x}'_{CE} : a mature counterfactual, defined as a counterfactual that has either passed a threshold probability or exhausted all T steps.
- \mathbf{x}^+ : a training sample in the target class (ground-truth).
- \mathbf{y}^+ : The one-hot encoded output vector for the target class.
- θ : Model parameters (unspecified).
- Θ : Matrix of parameters.
- $\mathbf{M}(\cdot)$: linear predictions (logits) of the classifier.

A.2 Formulas

A.2.1 Maximum Mean Discrepancy

Maximum mean discrepancy is defined as follows,

$$\begin{aligned} \text{MMD}(X', \tilde{X}') &= \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(x_i, x_j) \\ &\quad + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n k(\tilde{x}_i, \tilde{x}_j) \\ &\quad - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, \tilde{x}_j) \end{aligned} \tag{9}$$

where $k(\cdot, \cdot)$ is a kernel function (Gretton et al. 2012). We make use of a Gaussian kernel with a constant length-scale parameter of 0.5. In our implementation, Equation 9 is by default applied to the entire subset of the training data for which $y = y^+$.

A.3 Other Conventions

In some place of this appendix, we use the terms *full/Full* (i.e. the full CT objective) and *vanilla/Vanilla* (i.e. vanilla training objective) to refer to models trained with counterfactual training (CT) and the baseline (BL), respectively.

Appendix B Technical Details of Our Approach

B.1 Generating Counterfactuals through Gradient Descent

In this section, we provide some additional background on gradient-based counterfactual generators (Section B.1.1) and discuss how we define convergence in this context (Section B.1.2).

B.1.1 Background

Gradient-based counterfactual search was originally proposed by Wachter, Mittelstadt, and Russell (2017). A more general version of the counterfactual search objective presented in the paper is as follows,

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

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 $\text{reg}(\cdot)$ 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 $y\text{loss}(\cdot)$. All generators also penalize the distance (ℓ_1 -norm) of counterfactuals from their original factual state. For *Generic* and *ECCCo*, we have $\mathcal{Z} := \mathcal{X}$ and $g(\mathbf{z}) = g(\mathbf{z})^{-1} = \mathbf{z}$, that is counterfactuals 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 (Joshi et al. 2019). In addition to the distance penalty, *ECCCo* uses a penalty that regularizes the energy associated with the counterfactual, \mathbf{x}' (Altmeyer et al. 2024). We omit the conformal set size penalty proposed in the original paper, since a) the authors found faithfulness to primarily depend on the energy penalty and hence this alleviates us from one additional hyperparameter.

B.1.2 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 the main paper, 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 that this approach induces models to be relatively less sensitive to immutable features, propose a theoretical result supporting this and provide empirical evidence that strengthens our argument (both in the main paper and additional findings in this appendix). Below we derive the analytical results in Proposition in the main paper.

Proof. Let d_{mbl} and d_{immbl} denote some mutable and immutable feature, respectively. Suppose that $\mu_{y^-, d_{\text{immbl}}} < \mu_{y^+, d_{\text{immbl}}}$ and $\mu_{y^-, d_{\text{mbl}}} > \mu_{y^+, d_{\text{mbl}}}$, 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 $k, m \in \mathcal{K}$ (Hastie, Tibshirani, and Friedman 2009):

$$\log \frac{p(k|\mathbf{x})}{p(m|\mathbf{x})} = \mathbf{x}^\top \Sigma^{-1} (\mu_k - \mu_m) + \text{const} \quad (10)$$

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

$$\log \frac{p(k|\mathbf{x})}{p(m|\mathbf{x})} = \sum_{d=1}^D \frac{\mu_{k,d} - \mu_{m,d}}{\sigma_d^2} x_d + \text{const} \quad (11)$$

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

$$\log \frac{p(k|\mathbf{x})}{p(m|\mathbf{x})} = \sum_{d=1}^D (\theta_{k,d} - \theta_{m,d}) x_d + \text{const} \quad (12)$$

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

Based on Equation 11 and Equation 12 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 $(\theta_{y^-, d_{\text{immbl}}} - \theta_{y^+, d_{\text{immbl}}}) < 0$ and $(\theta_{y^-, d_{\text{mbl}}} - \theta_{y^+, d_{\text{mbl}}}) > 0$.

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 with respect to coefficient $\theta_{y^+, d}$ is equal to

$$\frac{\partial}{\partial \theta_{y^+, d}} (\text{div}(\mathbf{x}^+, \mathbf{x}', \mathbf{y}; \theta)) = \frac{\partial}{\partial \theta_{y^+, d}} ((-\mathbf{M}_\theta(\mathbf{x}^+)[y^+]) - (-\mathbf{M}_\theta(\mathbf{x}') [y^+])) = x'_d - x_d^+ \quad (13)$$

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}). \square

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 (Grathwohl et al. 2020; Altmeyer et al. 2024).

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,

$$\begin{aligned} x_d^{\text{LB}} &= \arg \min_{x_d} \{\mu_d - n_{\sigma_d} \sigma_d, \arg \min_{x_d} x_d\} \\ x_d^{\text{UB}} &= \arg \max_{x_d} \{\mu_d + n_{\sigma_d} \sigma_d, \arg \max_{x_d} x_d\} \end{aligned} \quad (14)$$

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 Hyperparameters

Note 1 presents the default hyperparameters used during training.

Note 1space Training Phase

- Meta Parameters:
 - Generator: `ecco`
 - Model: `mlp`
- Model:
 - Activation: `relu`
 - No. Hidden: 32
 - No. Layers: 1
- Training Parameters:
 - Burnin: 0.0
 - Class Loss: `logitcrossentropy`
 - Convergence: `threshold`
 - Generator Parameters:
 - * Decision Threshold: 0.75
 - * λ_{cst} : 0.001
 - * λ_{egy} : 5.0
 - * Learning Rate: 0.25

- * Maximum Iterations: 30
- * Optimizer: sgd
- * Type: ECCo
- λ_{adv} : 0.25
- λ_{clf} : 1.0
- λ_{div} : 0.5
- λ_{reg} : 0.1
- Learning Rate: 0.001
- No. Counterfactuals: 1000
- No. Epochs: 100
- Objective: full
- Optimizer: adam

B.5 Evaluation Details

B.5.1 Counterfactual Outcomes

For all of our counterfactual evaluations, we proceed as follows: for each dataset we run J bootstrap rounds (“No. Runs”) to account for stochasticity (Note 2); for each bootstrap round, we randomly draw factual and target pairs; then, for each model, we draw samples from the test set (with replacement) for which the model predicts the randomly chosen factual class; finally, we generate multiple counterfactuals (“No. Counterfactuals”) and evaluate the outcomes (Note 2). This is in line with standard practice in the related literature on CE (see e.g. Schut et al. (2021)). For our final results presented in the main paper, we rely on held-out test sets for evaluation. For tuning purposes we rely on training and/or validation sets.

Note 2 presents the default hyperparameters used during evaluation for tuning purposes. For the main results presented in the paper, we use larger evaluations, specifically:

- “No. Runs”: We set the number of bootstrap rounds to $J = 100$ for all datasets.
- “No. Individuals”: In each round we draw 1,250, 500 and 125 samples for synthetic datasets, real-world tabular datasets and *MNIST*, respectively, across five different values for the strength of the energy penalty of *ECCo* at test time, $\lambda_{\text{egy}} \in \{0.1, 0.5, 1.0, 5.0, 10.0\}$.

Note 2space Evaluation Phase

- Convergence: threshold
- Decision Threshold: 0.95
- Maximum Iterations: 50
- No. Individuals: 100
- No. Runs: 5

B.5.2 Predictive Performance

To assess (robust) predictive performance, we evaluate model accuracy on (adversarially perturbed) test data. To generate adversarial examples we use the Fast Gradient Sign Method (FGSM) (Goodfellow, Shlens, and Szegedy 2015). For the main results in the paper, we choose a range of values $\epsilon = [0.0, 0.1]$. In some places of this appendix, you will also find predictive performance evaluations in terms of the F1-score.

Appendix C Details on Main Experiments

C.1 Final Hyperparameters

As discussed the main paper, CT is sensitive to certain hyperparameter choices. We study the effect of many hyperparameters extensively in Section D of this appendix. For the main results, we tune a small set of key hyperparameters (Section E). The final choices for the main results are presented for each data set in Table 4 along with training, test and batch sizes.

C.1.1 Confidence Intervals

Table 5 present the exact confidence intervals (99%) for the difference in mean outcomes on which we base our assessment of statistical significance in the main paper. Grouped by evaluation metrics (Variable) and dataset (Data), the table presents the mean outcomes for CT and BT and finally the lower bound (LB) and upper bound (UB) of the

Table 4 Final hyperparameters used for the main results presented in the main paper. Any hyperparameter not shown here is set to its default value (Note 1).

Data	No. Train	No. Test	Batchsize	Domain	Decision Threshold	No. Counterfactuals	λ_{reg}
LS	3600	600	30	none	0.5	1000	0.01
Circ	3600	600	30	none	0.5	1000	0.5
Moon	3600	600	30	none	0.9	1000	0.25
OL	3600	600	30	none	0.5	1000	0.25
Adult	26049	5010	1000	none	0.75	5000	0.25
CH	16504	3101	1000	none	0.5	5000	0.25
Cred	10617	1923	1000	none	0.5	5000	0.25
GMSC	13371	2474	1000	none	0.5	5000	0.5
MNIST	11000	2000	1000	(-1.0, 1.0)	0.5	5000	0.01

Table 5 Mean outcomes for CT and BL along with bootstrapped confidence intervals (99%) for difference in mean outcomes grouped by dataset and evaluation metric. Column LB and UB show the lower and upper bound of the intervals, respectively, and computed using the percentile method (for significance, interval should not include zero). The underlying counterfactual evaluations are the same as the ones used to produce the main table in the paper.

Variable	Data	CT	BL	LB	UB
Cost	Adult	2.19	2.28	-0.32	0.11
Cost	CH	1.37	2.46	-1.18	-1.0
Cost	Circ	0.7	1.22	-0.55	-0.49
Cost	Cred	2.7	2.29	0.16	0.6
Cost	GMSC	1.03	3.04	-2.37	-1.86
Cost	LS	3.75	4.48	-0.8	-0.67
Cost	MNIST	72.08	53.42	11.15	26.68
Cost	Moon	1.52	1.6	-0.12	-0.05
Cost	OL	1.55	2.62	-1.25	-0.9
IP*	Adult	0.07	0.11	-0.06	-0.02
IP*	CH	0.02	0.06	-0.05	-0.03
IP*	Circ	0.0	0.0	-0.01	-0.0
IP*	Cred	0.03	0.06	-0.05	-0.01
IP*	GMSC	0.05	0.07	-0.02	-0.01
IP*	LS	0.11	0.23	-0.13	-0.11
IP*	MNIST	0.02	0.02	-0.07	0.07
IP*	Moon	0.02	0.02	-0.01	0.0
IP*	OL	0.12	0.09	-0.01	0.05
IP	Adult	15.13	15.16	-0.42	0.39
IP	CH	6.72	7.52	-1.05	-0.6
IP	Circ	0.97	2.36	-1.44	-1.35
IP	Cred	19.79	22.02	-3.17	-1.42
IP	GMSC	7.24	8.1	-1.26	-0.38
IP	LS	2.51	3.4	-0.95	-0.84
IP	MNIST	261.05	278.84	-27.38	-7.51
IP	Moon	1.37	1.71	-0.36	-0.3
IP	OL	4.52	4.44	-0.03	0.19

confidence interval. To compute the intervals, we used the percentile method for bootstrapped confidence intervals: the lower and upper bound represent the $\alpha/2$ - and $(1 - \alpha/2)$ -quantile of the bootstrap distribution, respectively, for $\alpha = 0.01$.

C.1.2 Qualitative Findings for Image Data

Figure 6 shows much more plausible (faithful) counterfactuals for a model with CT than the model with conventional training (Figure 7).

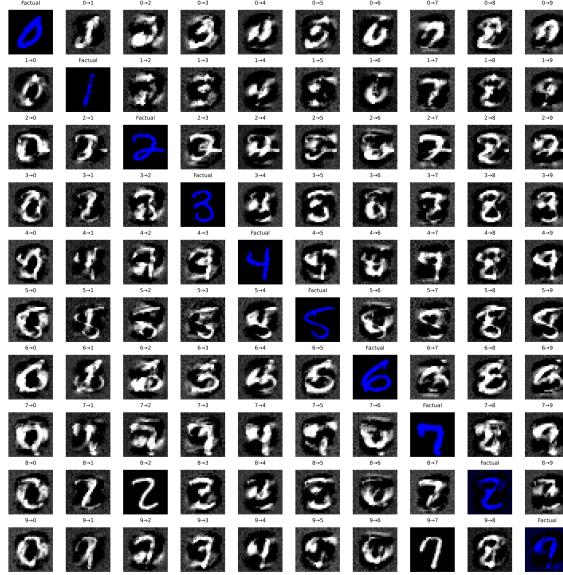


Figure 6 Counterfactual images for *MLP* with counterfactual training. Factual images are shown on the diagonal, with the corresponding counterfactual for each target class (columns) in that same row. The underlying generator, *ECCC0*, aims to generate counterfactuals that are faithful to the model (Altmeyer et al. 2024).

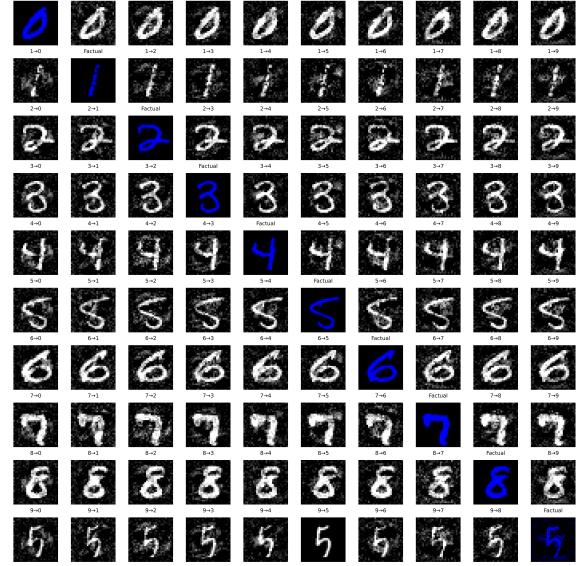


Figure 7 The same setup, factuals, model architecture and generator as in Figure 6, but the model was trained conventionally.

C.1.3 Integrated Gradients

We make use of integrated gradients (IG) proposed by Sundararajan, Taly, and Yan (2017) to empirically evaluate the feature protection mechanism in CT. We choose this approach because it produces theoretically sound results, works well for non-linear models, and remains relatively inexpensive.

IG calculates the contribution of each input feature towards a specific prediction by approximating the integral of the model output with respect to its input, using a set of samples that linearly interpolate between a test instance and some baseline instance (Sundararajan, Taly, and Yan 2017). This process produces a vector of real numbers, one per input feature, which informs about the contribution of each feature to the prediction. For example:

- a large positive value indicates that a feature has strong positive influence on the classification (i.e., increases the score for a class);
- a small negative value indicates that a feature has weak negative influence on the classification (i.e., decreases the score for a class).

To calculate the contributions, IG compares the output to a baseline. The selection of an appropriate baseline is an important design decision — it should produce a “neutral” prediction to avoid capturing effects that cannot be directly attributed to the model (Sundararajan, Taly, and Yan 2017; Sturmels, Lundberg, and Lee 2020). To remain consistent in our evaluations, we use a baseline drawn at random from a uniform distribution, $\mathcal{U}(-1, 1)$, for all datasets. This aligns with standard evaluation practices for IG.

We run IG on models trained on all datasets to compare their sensitivity to features that were protected using CT:

- for synthetic datasets, this is always the first feature
- for real-world tabular datasets, this is always *age*
- for MNIST, this is first five and last five rows of pixels

As IG outputs are not bounded (i.e., they are arbitrary real numbers), it becomes a challenge to meaningfully compare IG outputs of different models — ones that are trained conventionally, and ones that underwent counterfactual training. For our purposes, we observe with reference to our Proposition, that we are interested estimating changes in the relative

contribution of protected features compared to mutable ones. Thus, to meaningfully compare integrated gradients for different models and to accommodate for variable ranges of outputs in absolute terms, we standardize the integrated gradients across features.

Let g_d denote the estimated IG for feature d . Then in the case of 2D synthetic datasets we find that taking the absolute value of the outputs, $|g_d|$, and then dividing them by a $\max(g) - \min(g)$ term allows us to make the most meaningful comparison. In the case of real-world datasets we choose to normalize the values to a $[0, 1]$ range instead. We compare the (average) sensitivity to the features that were protected for CT models. Once again we use bootstrapping (100 rounds, 2500 samples per round) to establish the significance of our results (Figure 8).

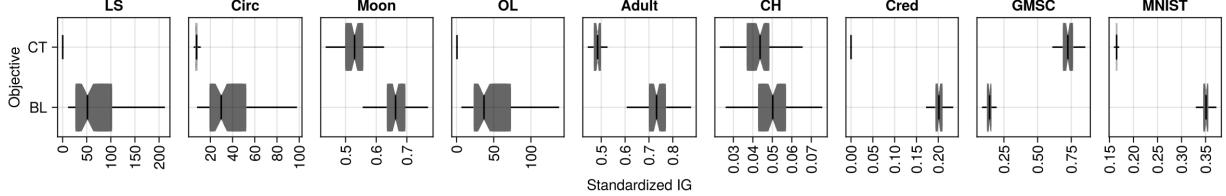


Figure 8 Interquartile ranges of bootstrap outcomes for sensitivity to protected features measured by standardized integrated gradients.

C.1.4 Costs and Validity

In Table 6, we present additional outcomes for common evaluation metrics: Table 6a presents the average reduction in costs of counterfactuals for CT vs. BL with no mutability constraints, i.e. corresponding to the first two columns in the main table of the paper; Table 6b shows the corresponding average validities; finally, Table 6c shows average validities for the case with mutability constraints, i.e. corresponding to the third columns in the main table of the paper.

As noted in the discussion section of the main paper, we observe mixed results here. Average costs in terms of distances from factual values decrease for most datasets, which is positively surprising since improved plausibility requires counterfactuals to travel further into the target domain than minimum distance counterfactuals. It appears that in these cases faithful counterfactuals for the baseline model still end up far away from their initial starting points, but not close enough for samples in the target domain to be plausible. In that sense, CT can be seen to improve both plausibility and costs for faithful CE. In some cases though (*LS*, *CH*, *MNIST*), we do seem to observe the tradeoff between plausibility and costs play out, as we would expect (compare panels (a) and (b) of Figure 1 in the main paper for reference).

Concerning validity, we find that can lead to substantial reductions and only increases average validity compared to the baseline in one case (*Circ*). As noted in the discussion section of the main paper, this result does not surprise us: by design, CT shrinks the solution space for valid counterfactual explanations, thus making it “harder” to reach validity compared to the baseline model. Note that for a number of reasons this should not be seen as problematic:

1. Validity of gradient-based CE is a function on the number of steps and the step size which we both kept fixed during evaluation: simply adjusting $T = 50$ to higher values or choosing a larger step size will lead to higher rates of validity.
2. Even though reaching validity is sometimes “harder” in terms of the necessary number of steps for a given step size, we have already shown that the average distances that counterfactuals need to travel decrease for most datasets. Users care about costs in terms of feature distances, not search iteration steps.
3. From a philosophical perspective on algorithmic recourse, validity in and off itself is not a sufficient desideratum for counterfactuals. In fact, Venkatasubramanian and Alfano (2020) propose introducing an upper bound on costs of the flipset (i.e. the set of valid CE), arguing that valid but highly costly counterfactuals are not useful to individuals in practice. In a similar fashion, it could be argued that there should be an upper bound on the implausibility of counterfactuals in the flipset.

Table 6 Costs and validity.

(a) Reduction in average costs for CT vs. the baseline. Results correspond to the case with no mutability constraints in the main table of the paper.

Data	Cost (-%)
LS	$-27.11 \pm 0.75^*$
Circ	$40.17 \pm 0.85^*$
Moon	$32.54 \pm 1.23^*$
OL	$12.08 \pm 1.58^*$
Adult	-4.59 ± 2.54
CH	$-33.04 \pm 1.96^*$
Cred	$27.43 \pm 1.05^*$
GMSC	$-22.40 \pm 3.64^*$
MNIST	$-40.71 \pm 7.02^*$
Avg.	-1.74

(b) Average validities of counterfactuals for CT and BL. Unconstrained case.

Data	CT	BL
LS	1.0	1.0
Circ	1.0	0.51
Moon	1.0	1.0
OL	0.86	0.98
Adult	0.68	0.99
CH	1.0	1.0
Cred	0.72	1.0
GMSC	0.94	1.0
MNIST	1.0	1.0

(c) Average validities of counterfactuals for CT and BL. Mutability constrained case.

Data	CT	BL
LS	1.0	1.0
Circ	0.71	0.48
Moon	1.0	0.98
OL	0.34	0.56
Adult	0.7	0.99
CH	1.0	1.0
Cred	0.74	1.0
GMSC	0.97	1.0
MNIST	1.0	1.0

Appendix D Grid Searches

To assess the hyperparameter sensitivity of our proposed training regime we ran multiple large grid searches for all of our synthetic datasets. We have grouped these grid searches into multiple categories:

1. **Generator Parameters** (Section D.2): Investigates the effect of changing hyperparameters that affect the counterfactual outcomes during the training phase.
2. **Penalty Strengths** (Section D.3): Investigates the effect of changing the penalty strengths in our proposed training objective.
3. **Other Parameters** (Section D.4): Investigates the effect of changing other training parameters, including the total number of generated counterfactuals in each epoch.

We begin by summarizing the high-level findings in Section D.1.2. For each of the categories, Section D.2 to Section D.4 then present all details including the exact parameter grids, average predictive performance outcomes and key evaluation metrics for the generated counterfactuals.

D.1 Evaluation Details

To measure predictive performance, we compute the accuracy and F1-score for all models on test data (Table 7, Table 8, Table 9). With respect to explanatory performance, we report here our findings for the (im)plausibility and cost of counterfactuals at test time. Since the computation of our proposed divergence-based adaption (IP*) is memory-intensive, we rely on the distance-based metric for the grid searches. For the counterfactual evaluation, we draw factual samples from the training data for the grid searches to avoid data leakage with respect to our final results reported in the body of the paper. Specifically, we want to avoid choosing our default hyperparameters based on results on the test data. Since we are optimizing for explainability, not predictive performance, we still present test accuracy and F1-scores.

D.1.1 Predictive Performance

We find that CT is associated with little to no decrease in average predictive performance for our synthetic datasets: test accuracy and F1-scores decrease by at most ~ 1 percentage point, but generally much less (Table 7, Table 8, Table 9). Variation across hyperparameters is negligible as indicated by small standard deviations for these metrics across the board.

D.1.2 Counterfactual Outcomes

Overall, we find that counterfactual training achieves its key objectives consistently across all hyperparameter settings and also broadly across datasets: plausibility is improved by up to 60 percent (%) for the *Circles* data (e.g. Figure 9), 25-30% for the *Moons* data (e.g. Figure 11) and 10-20% for the *Linearly Separable* data (e.g. Figure 10). At the same time, the average costs of faithful counterfactuals are reduced in many cases by around 20-25% for *Circles* (e.g. Figure 13) and up to 50% for *Moons* (e.g. Figure 15). For the *Linearly Separable* data, costs are generally increased although typically by less than 10% (e.g. Figure 14), which reflects a common tradeoff between costs and plausibility (Altmeyer et al. 2024).

We do observe strong sensitivity to certain hyperparameters, with clear manageable patterns. Concerning generator parameters, we firstly find that using *REVISE* to generate counterfactuals during training typically yields the worst outcomes out of all generators, often leading to a substantial decrease in plausibility. This finding can be attributed to the fact that *REVISE* effectively assigns the task of learning plausible explanations from the model itself to a surrogate VAE. In other words, counterfactuals generated by *REVISE* are less faithful to the model than *ECCCo* and *Generic*, and hence we would expect them to be a less effective and, in fact, potentially detrimental role in our training regime. Secondly, we observe that allowing for a higher number of maximum steps T for the counterfactual search generally yields better outcomes. This is intuitive, because it allows more counterfactuals to reach maturity in any given iteration. Looking in particular at the results for *Linearly Separable*, it seems that higher values for T in combination with higher decision thresholds (τ) yields the best results when using *ECCCo*. But depending on the degree of class separability of the underlying data, a high decision-threshold can also affect results adversely, as evident from the results for the *Overlapping* data (Figure 12): here we find that CT generally fails to achieve its objective because only a tiny proportion of counterfactuals ever reaches maturity.

Regarding penalty strengths, we find that the strength of the energy regularization, λ_{reg} is a key hyperparameter, while sensitivity with respect to λ_{div} and λ_{adv} is much less evident. In particular, we observe that not regularizing energy enough or at all typically leads to poor performance in terms of decreased plausibility and increased costs, in particular for *Circles* (Figure 17), *Linearly Separable* (Figure 18) and *Overlapping* (Figure 20). High values of λ_{reg} can increase the variability in outcomes, in particular when combined with high values for λ_{div} and λ_{adv} , but this effect is less pronounced.

Finally, concerning other hyperparameters we observe that the effectiveness and stability of CT is positively associated with the number of counterfactuals generated during each training epoch, in particular for *Circles* (Figure 25) and *Moons* (Figure 27). We further find that a higher number of training epochs is beneficial as expected, where we tested training models for 50 and 100 epochs. Interestingly, we find that it is not necessary to employ CT during the entire training phase to achieve the desired improvements in explainability: specifically, we have tested training models conventionally during the first half of training before switching to CT after this initial burn-in period.

D.2 Generator Parameters

The hyperparameter grid with varying generator parameters during training is shown in Note 3. The corresponding evaluation grid used for these experiments is shown in Note 4.

Note 3space Training Phase

- Generator Parameters:
 - Decision Threshold: 0.75, 0.9, 0.95
 - λ_{egy} : 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 4space Evaluation Phase

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

D.2.1 Predictive Performance

Predictive performance measures for this grid search are shown in Table 7.

Table 7 Predictive performance measures by dataset and objective averaged across training-phase parameters (Note 3) and evaluation-phase parameters (Note 4).

Dataset	Variable	Objective	Mean	Se
Circ	Accuracy	Full	1.0	0.0
Circ	Accuracy	Vanilla	1.0	0.0
Circ	F1-score	Full	1.0	0.0
Circ	F1-score	Vanilla	1.0	0.0
LS	Accuracy	Full	1.0	0.0
LS	Accuracy	Vanilla	1.0	0.0
LS	F1-score	Full	1.0	0.0
LS	F1-score	Vanilla	1.0	0.0
Moon	Accuracy	Full	1.0	0.0
Moon	Accuracy	Vanilla	1.0	0.0
Moon	F1-score	Full	1.0	0.0
Moon	F1-score	Vanilla	1.0	0.0
OL	Accuracy	Full	0.91	0.0
OL	Accuracy	Vanilla	0.92	0.0
OL	F1-score	Full	0.91	0.0
OL	F1-score	Vanilla	0.92	0.0

D.2.2 Plausibility

The results with respect to the plausibility measure are shown in Figure 9 to Figure 12.

D.2.3 Cost

The results with respect to the cost measure are shown in Figure 13 to Figure 16.

D.3 Penalty Strengths

The hyperparameter grid with varying penalty strengths during training is shown in Note 5. The corresponding evaluation grid used for these experiments is shown in Note 6.

Note 5space Training Phase

- Generator: `ecco`, `generic`, `revise`
- Model: `mlp`
- Training Parameters:
 - λ_{adv} : 0.1, 0.25, 1.0
 - λ_{div} : 0.01, 0.1, 1.0
 - λ_{reg} : 0.0, 0.01, 0.1, 0.25, 0.5
 - Objective: `full`, `vanilla`

Note 6space Evaluation Phase

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

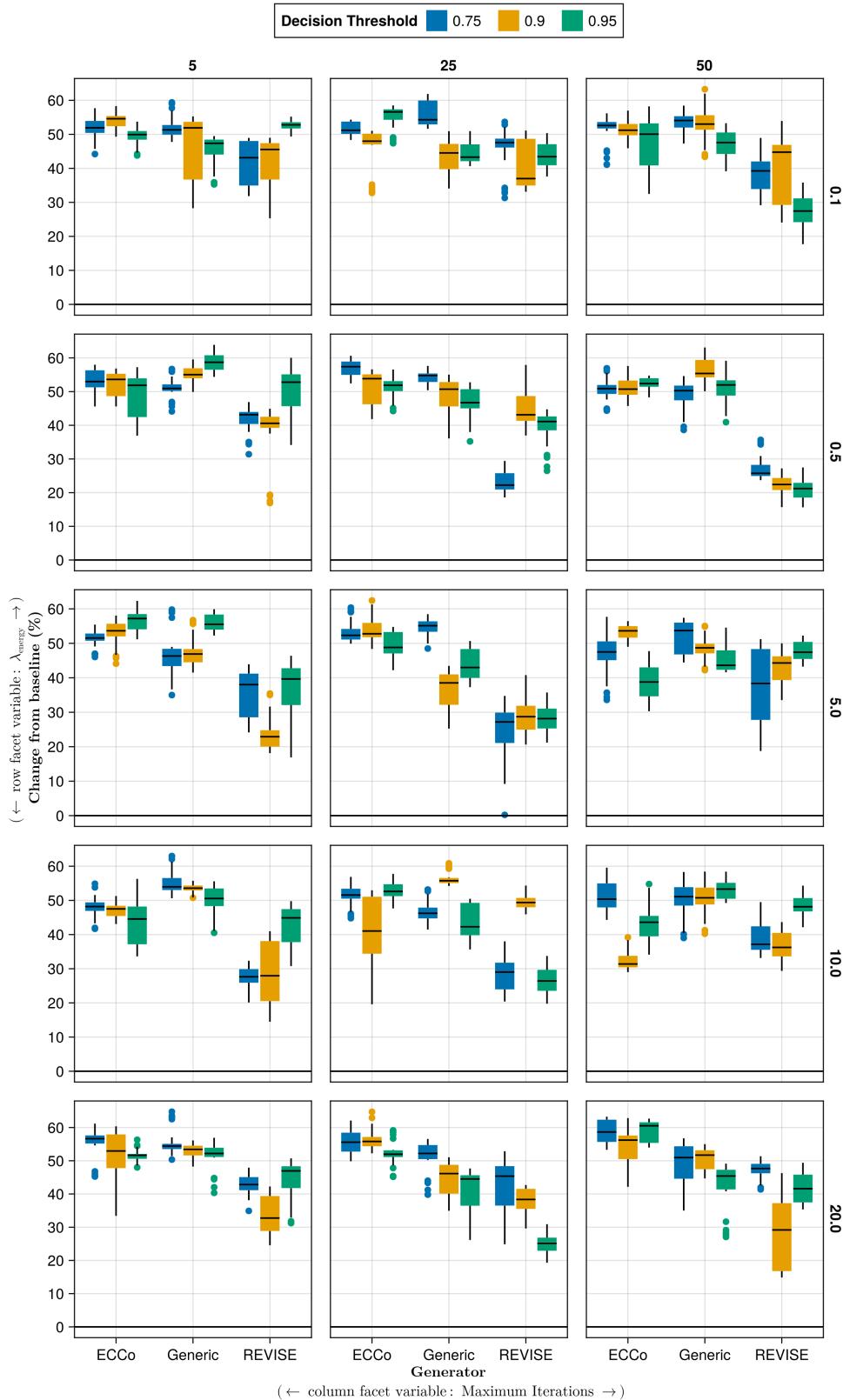


Figure 9 Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Circles.

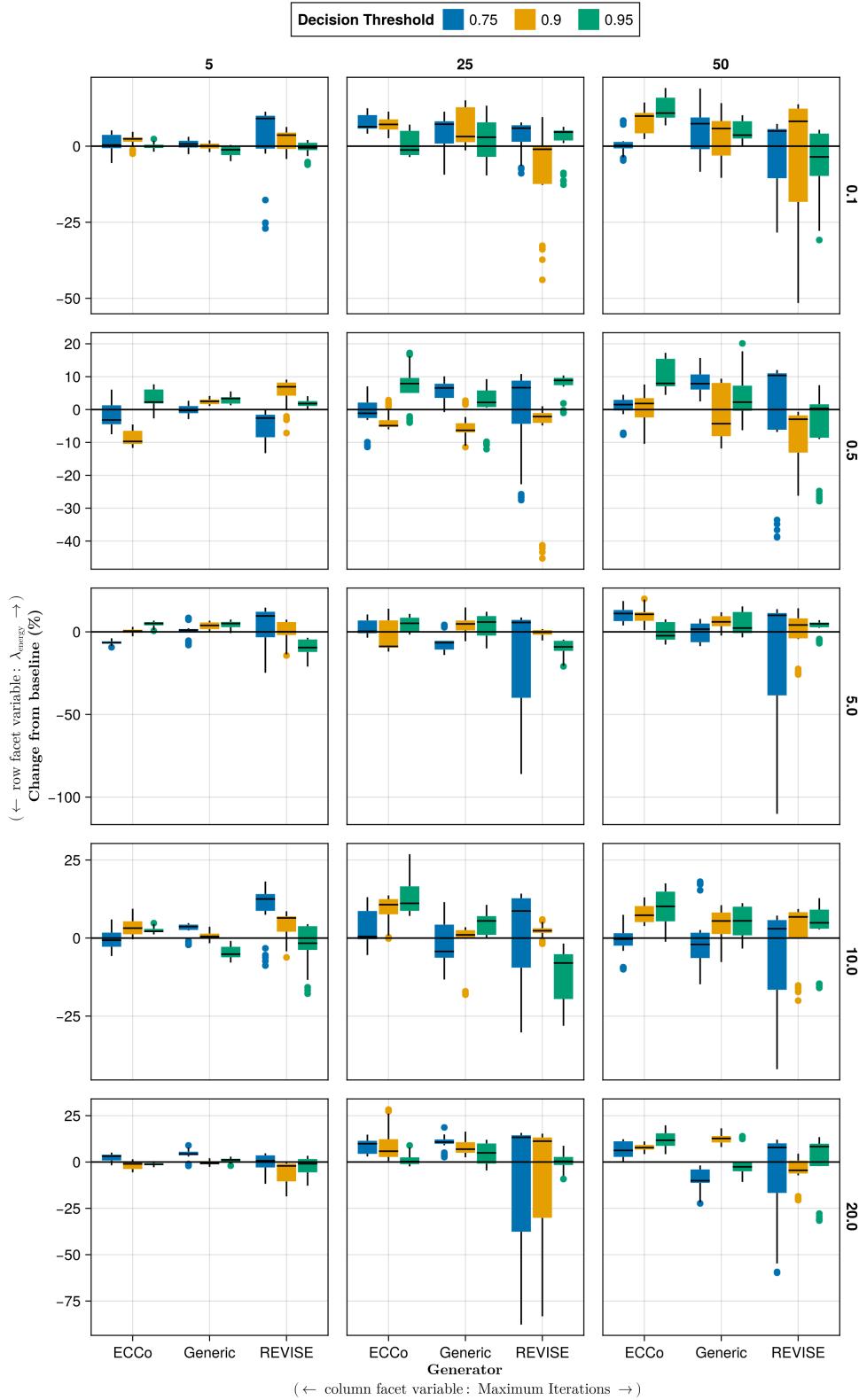


Figure 10 Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Linearly Separable.

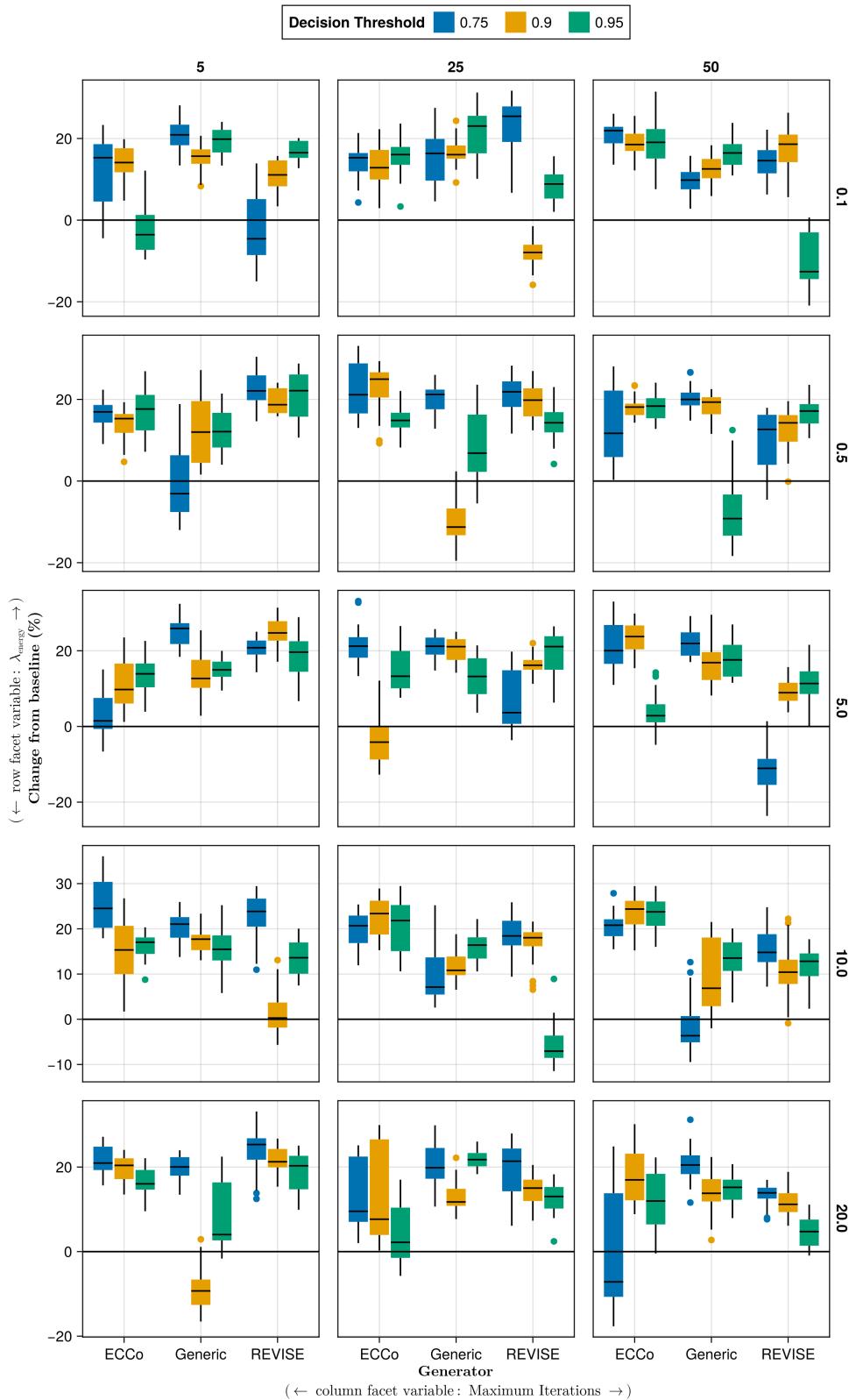


Figure 11 Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Moons.

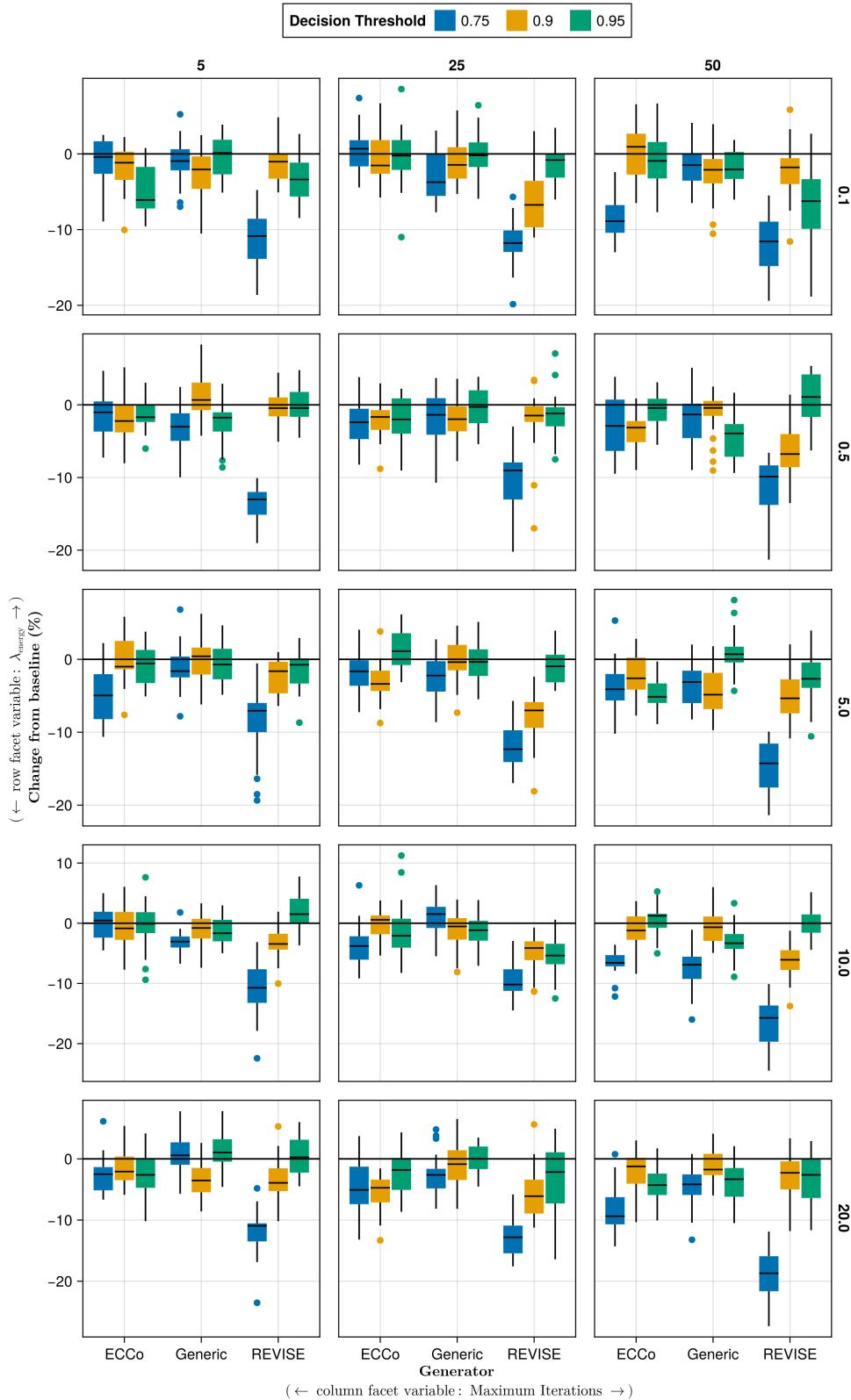


Figure 12 Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

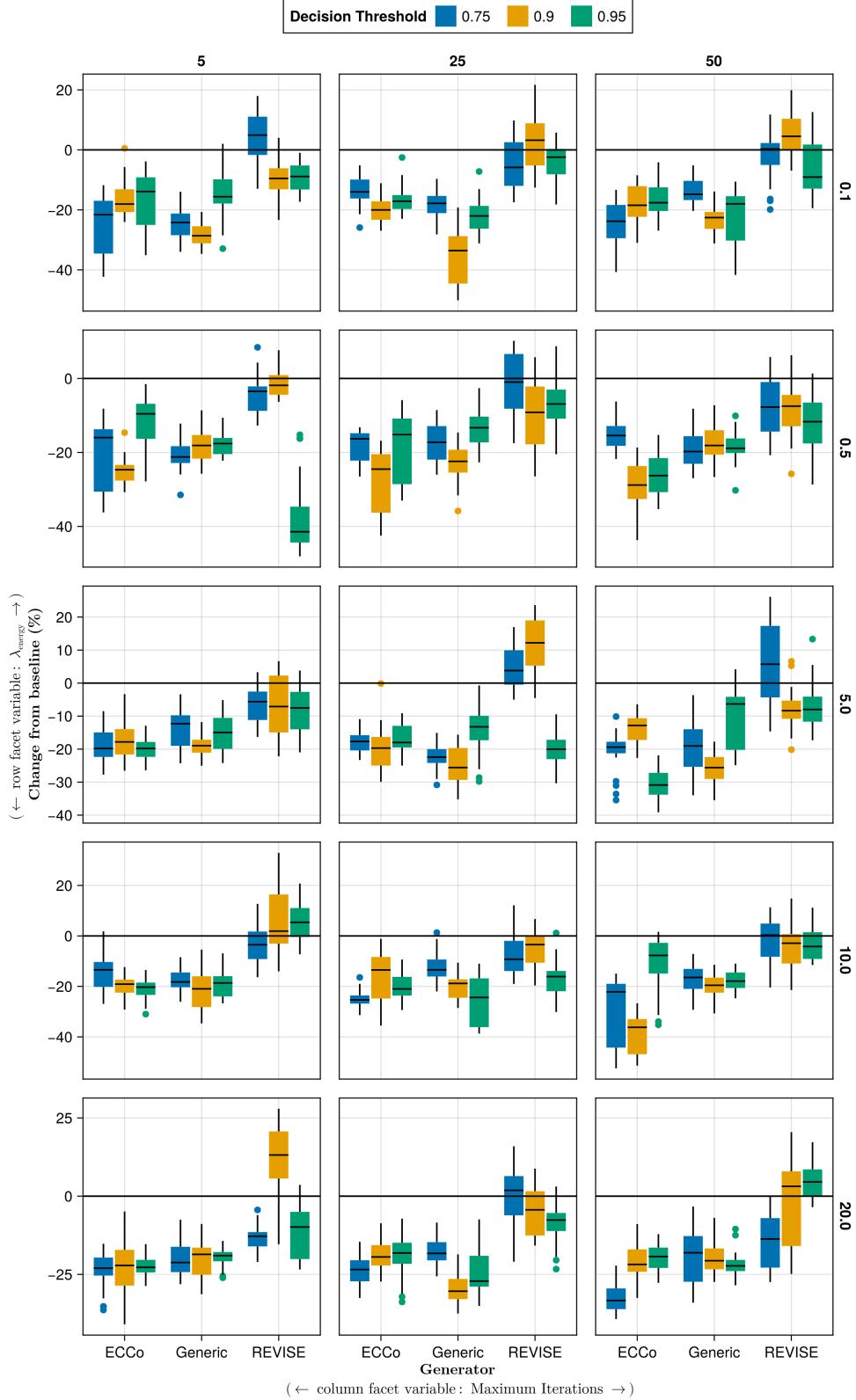


Figure 13 Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Circles.

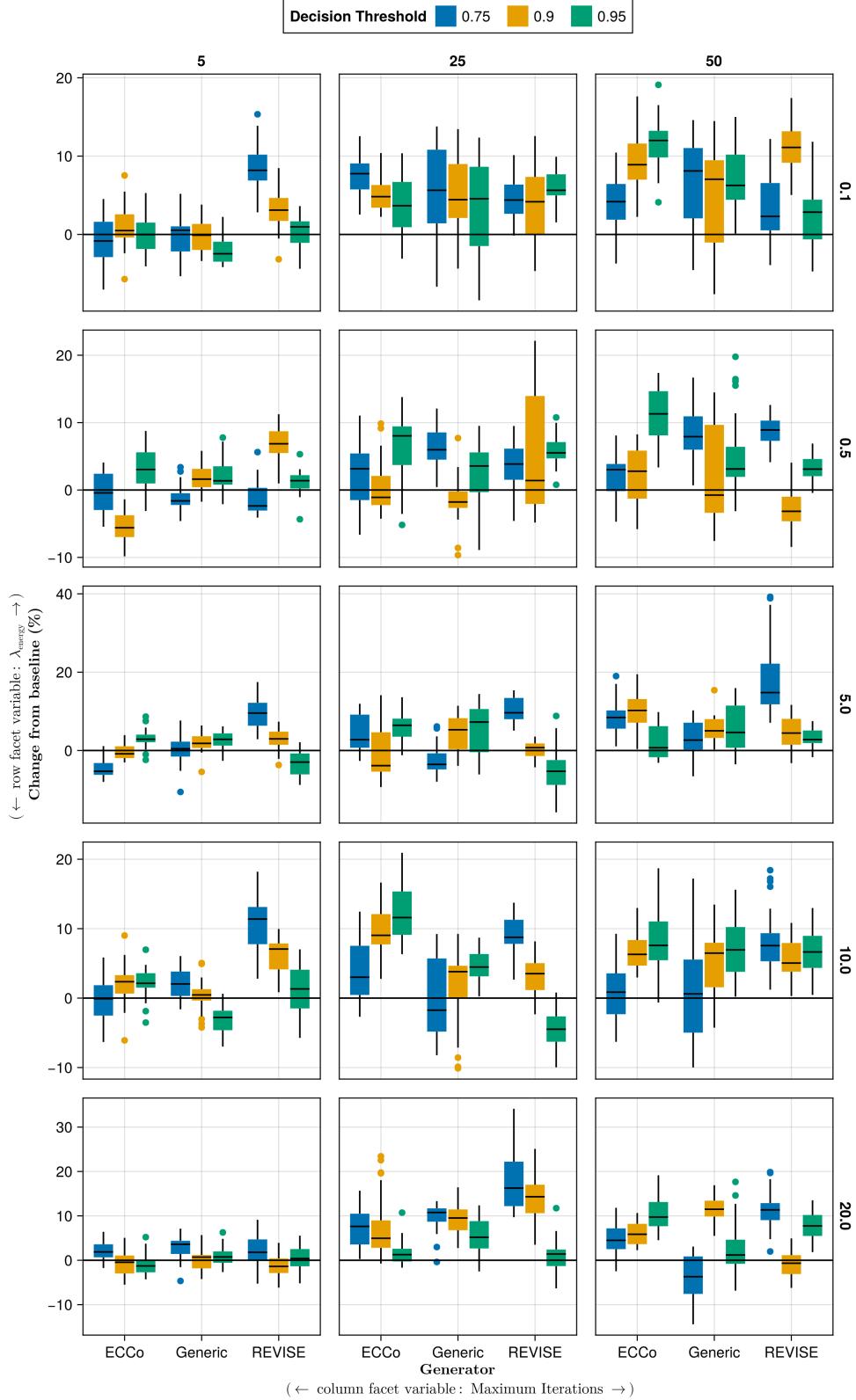


Figure 14 Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Linearly Separable.

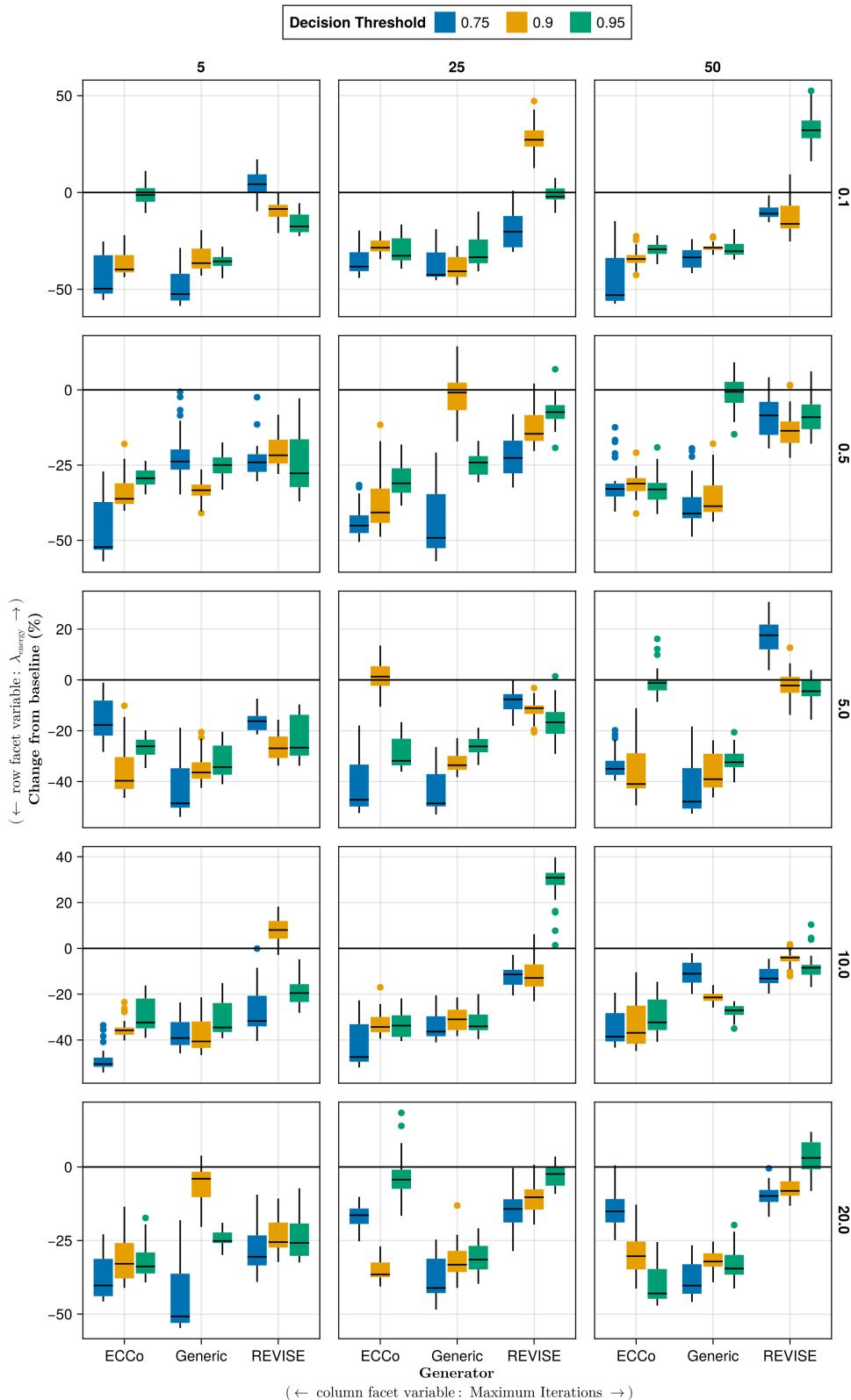


Figure 15 Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Moons.

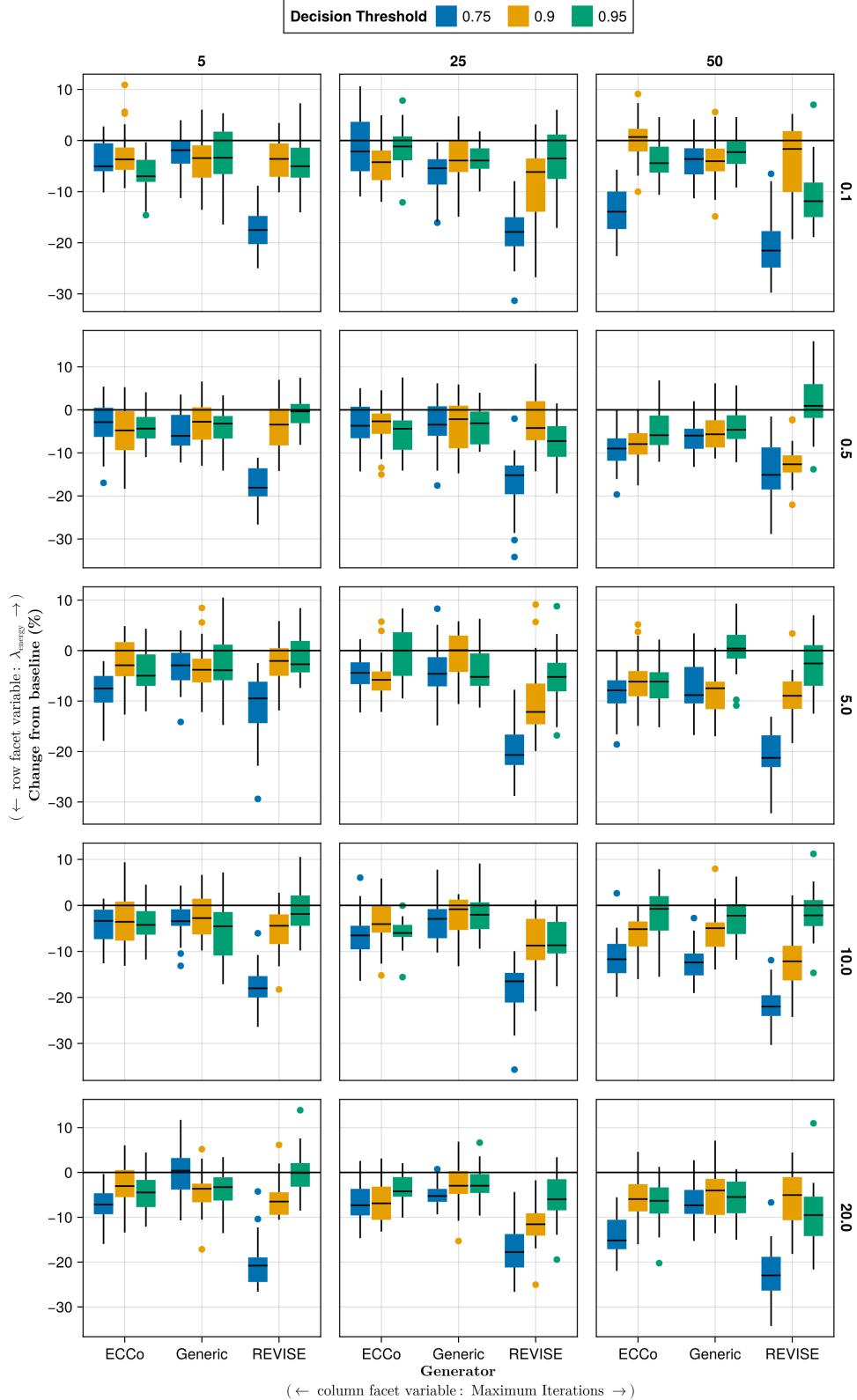


Figure 16 Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Overlapping.

D.3.1 Predictive Performance

Predictive performance measures for this grid search are shown in Table 8.

Table 8 Predictive performance measures by dataset and objective averaged across training-phase parameters (Note 5) and evaluation-phase parameters (Note 6).

Dataset	Variable	Objective	Mean	Se
Circ	Accuracy	Full	0.99	0.01
Circ	Accuracy	Vanilla	1.0	0.0
Circ	F1-score	Full	0.99	0.01
Circ	F1-score	Vanilla	1.0	0.0
LS	Accuracy	Full	1.0	0.01
LS	Accuracy	Vanilla	1.0	0.0
LS	F1-score	Full	1.0	0.01
LS	F1-score	Vanilla	1.0	0.0
Moon	Accuracy	Full	0.99	0.04
Moon	Accuracy	Vanilla	1.0	0.01
Moon	F1-score	Full	0.99	0.04
Moon	F1-score	Vanilla	1.0	0.01
OL	Accuracy	Full	0.91	0.02
OL	Accuracy	Vanilla	0.92	0.0
OL	F1-score	Full	0.91	0.02
OL	F1-score	Vanilla	0.92	0.0

D.3.2 Plausibility

The results with respect to the plausibility measure are shown in Figure 17 to Figure 20.

D.3.3 Cost

The results with respect to the cost measure are shown in Figure 21 to Figure 24.

D.4 Other Parameters

The hyperparameter grid with other varying training parameters is shown in Note 7. The corresponding evaluation grid used for these experiments is shown in Note 8.

Note 7space Training Phase

- Generator: `ecco`, `generic`, `revise`
- Model: `mlp`
- Training Parameters:
 - Burnin: 0.0, 0.5
 - No. Counterfactuals: 100, 1000
 - No. Epochs: 50, 100
 - Objective: `full`, `vanilla`

Note 8space Evaluation Phase

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

D.4.1 Predictive Performance

Predictive performance measures for this grid search are shown in Table 9.

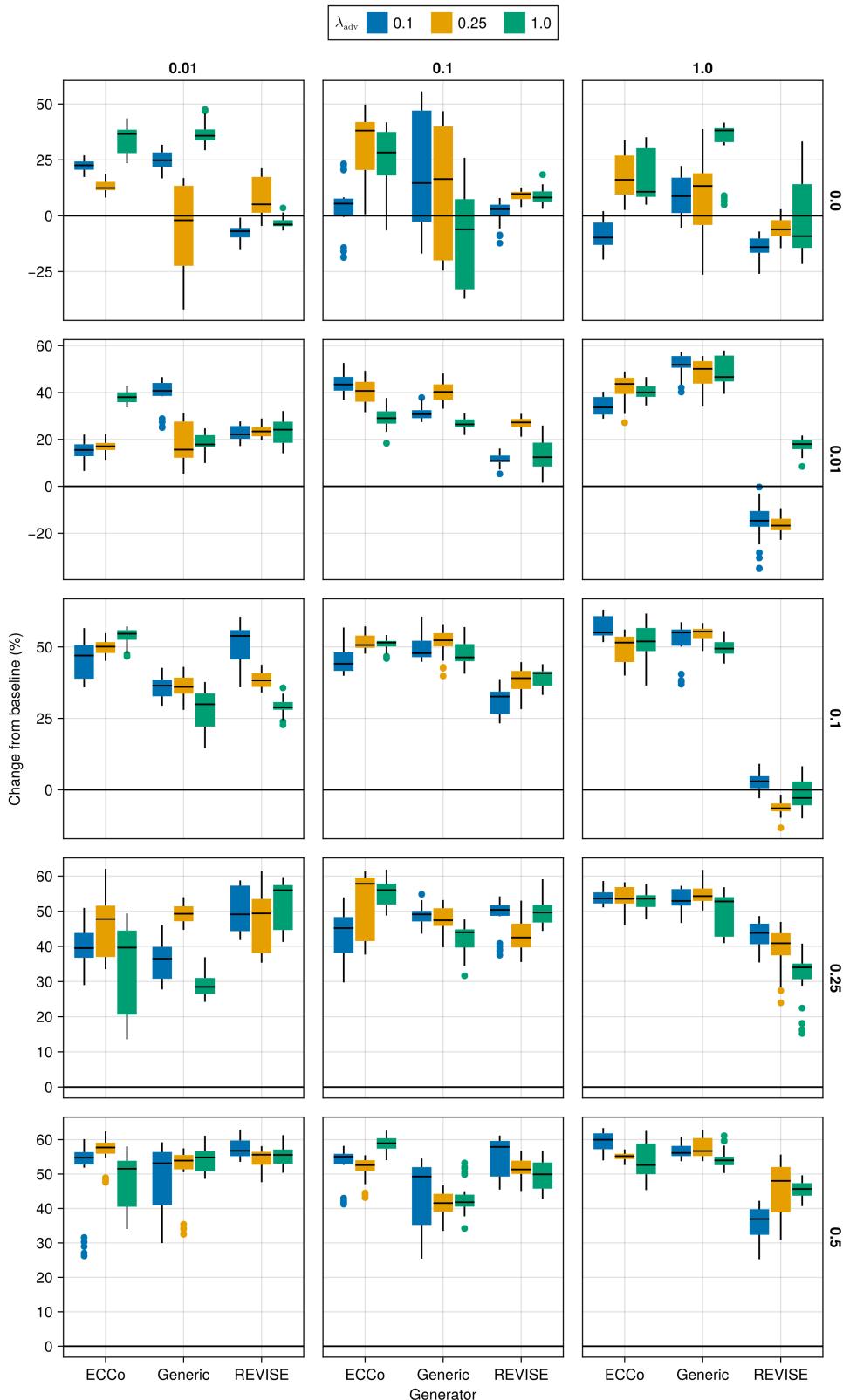


Figure 17 Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Circles.

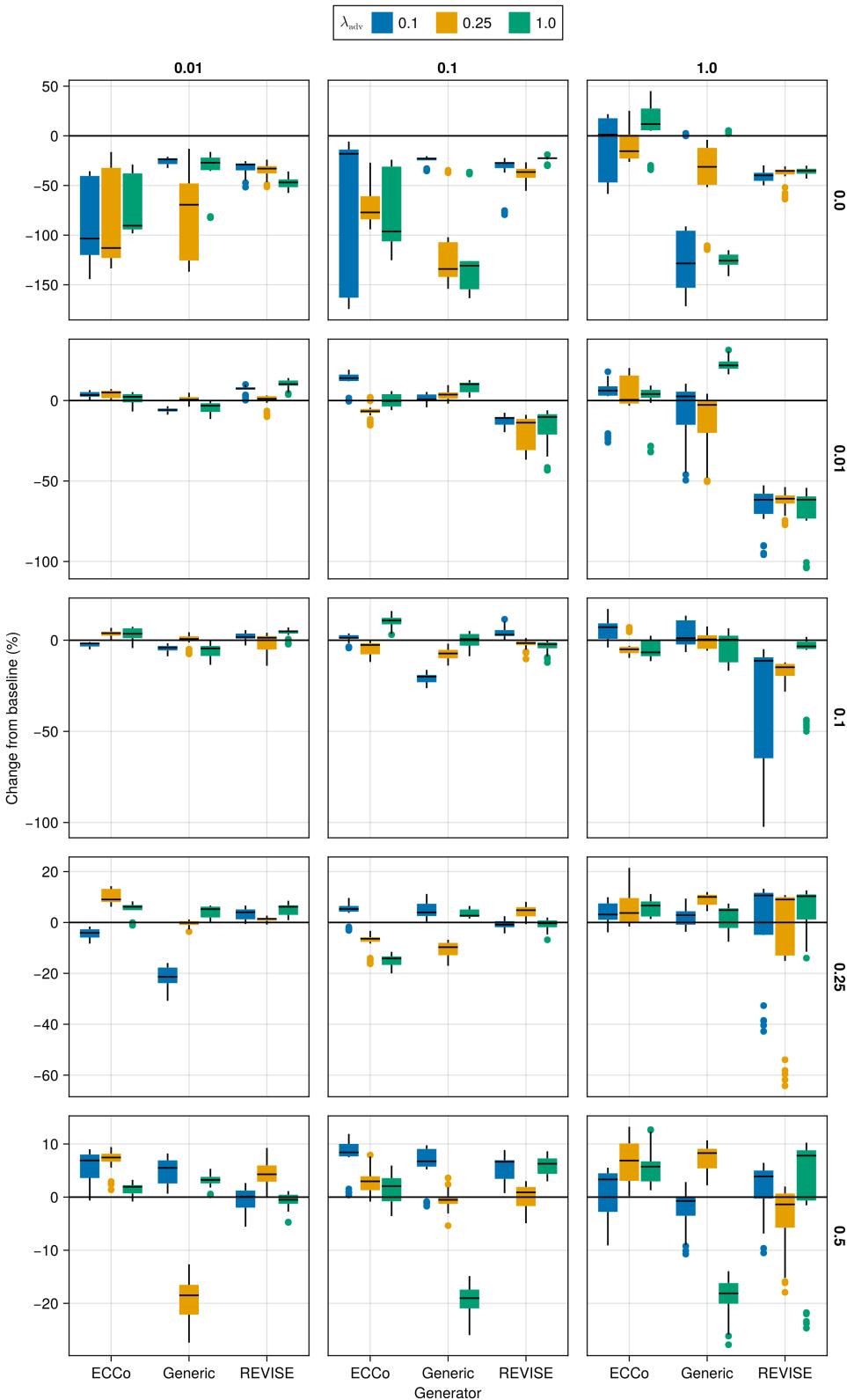


Figure 18 Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Linearly Separable.

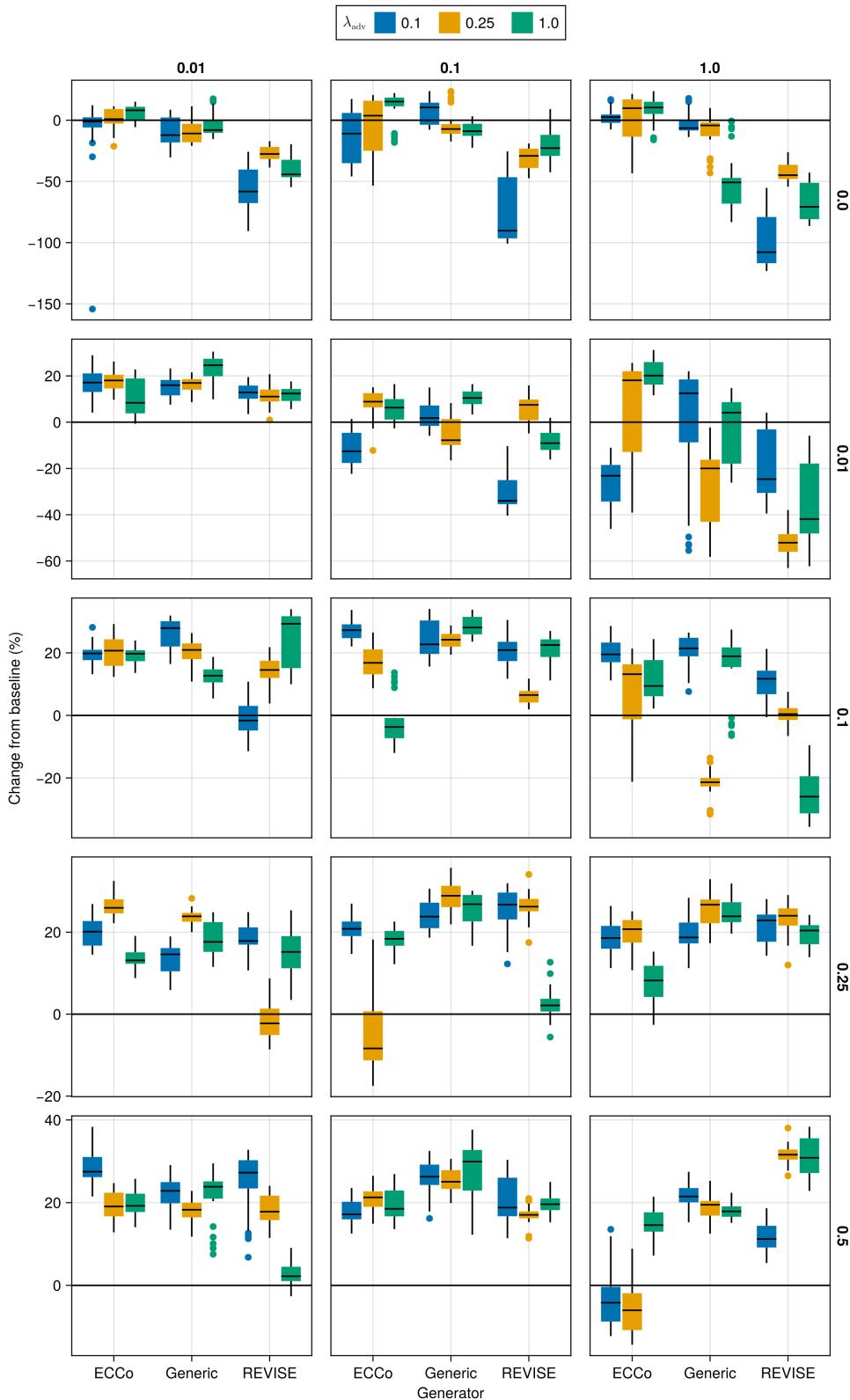


Figure 19 Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Moons.

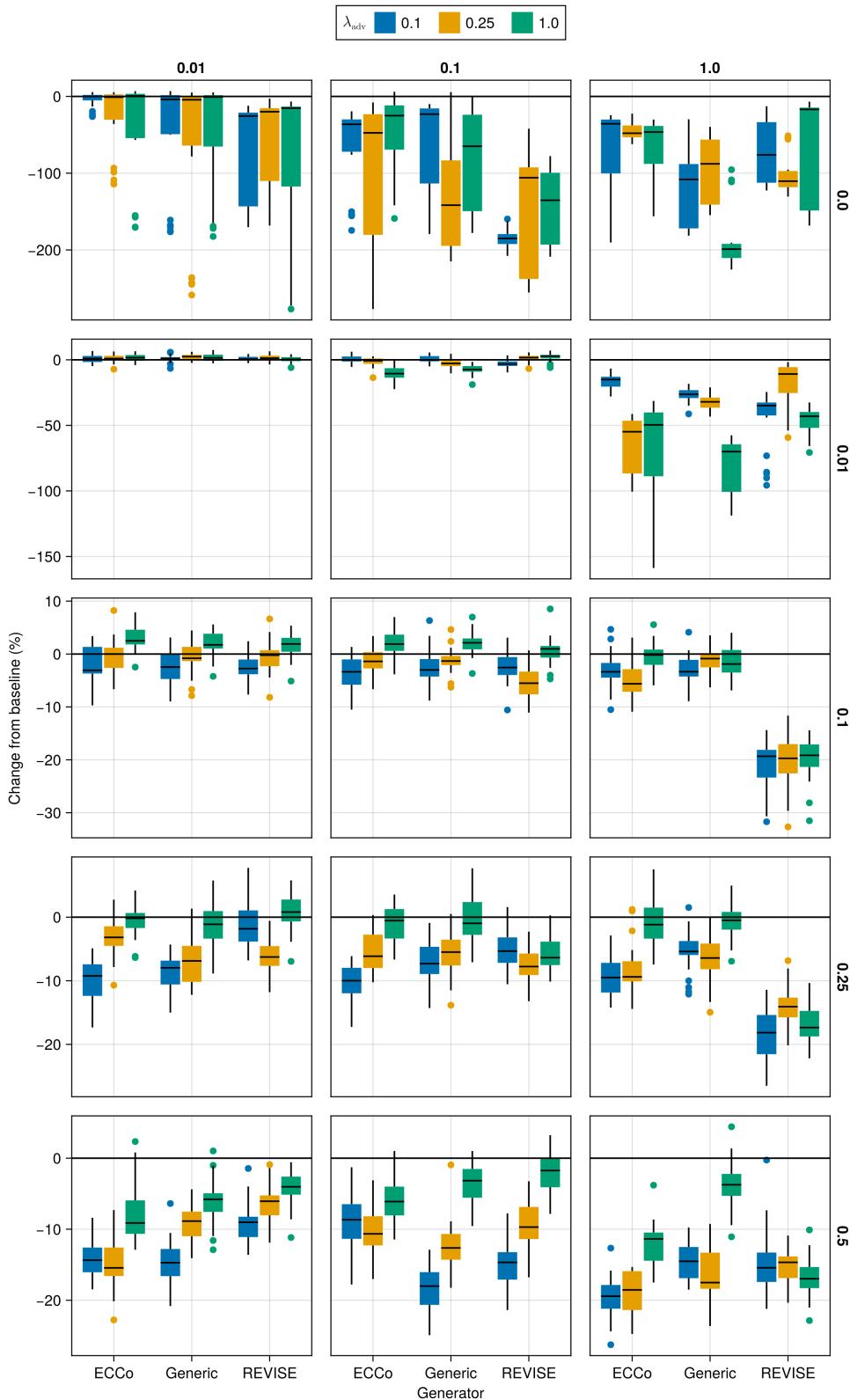


Figure 20 Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

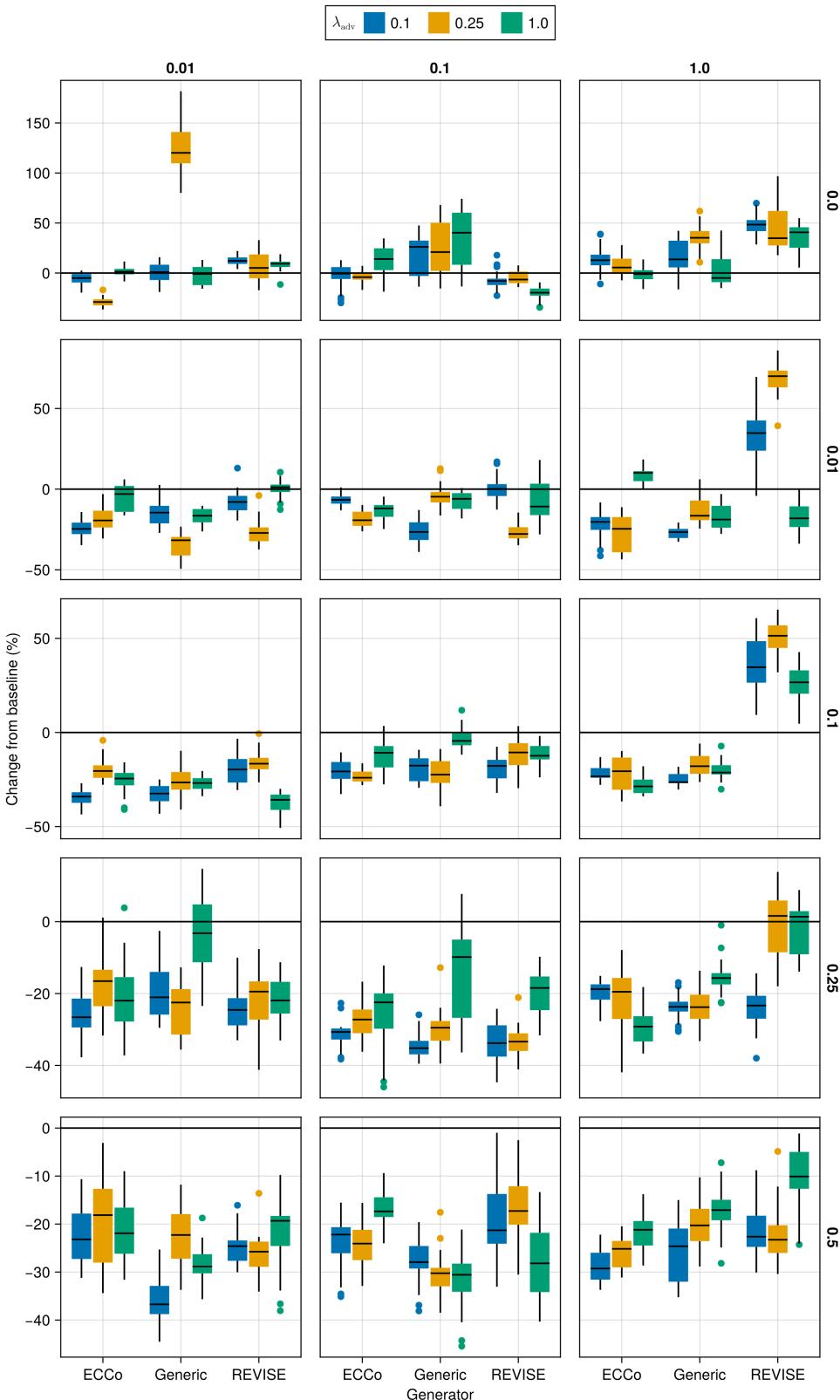


Figure 21 Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Circles.

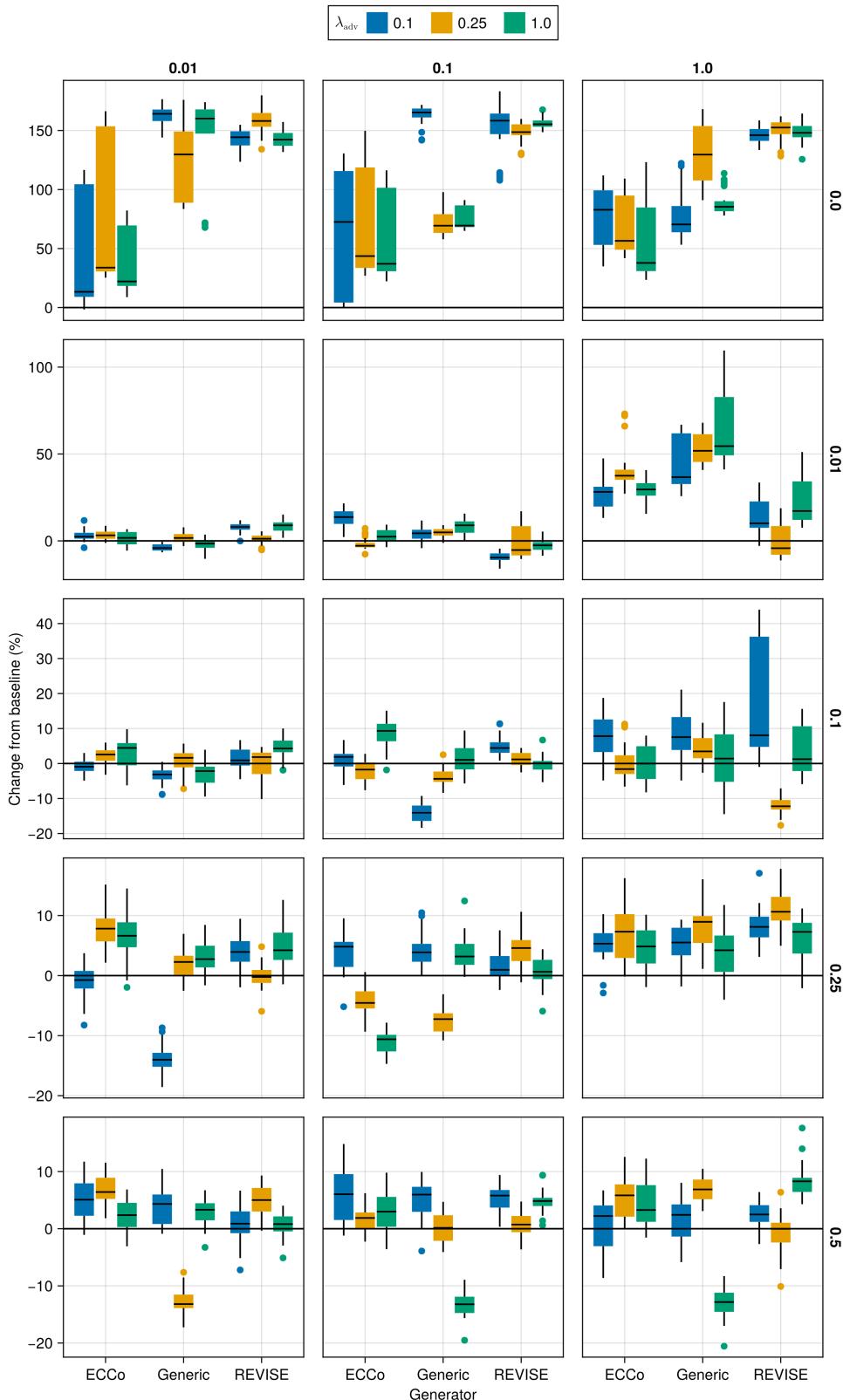


Figure 22 Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Linearly Separable.

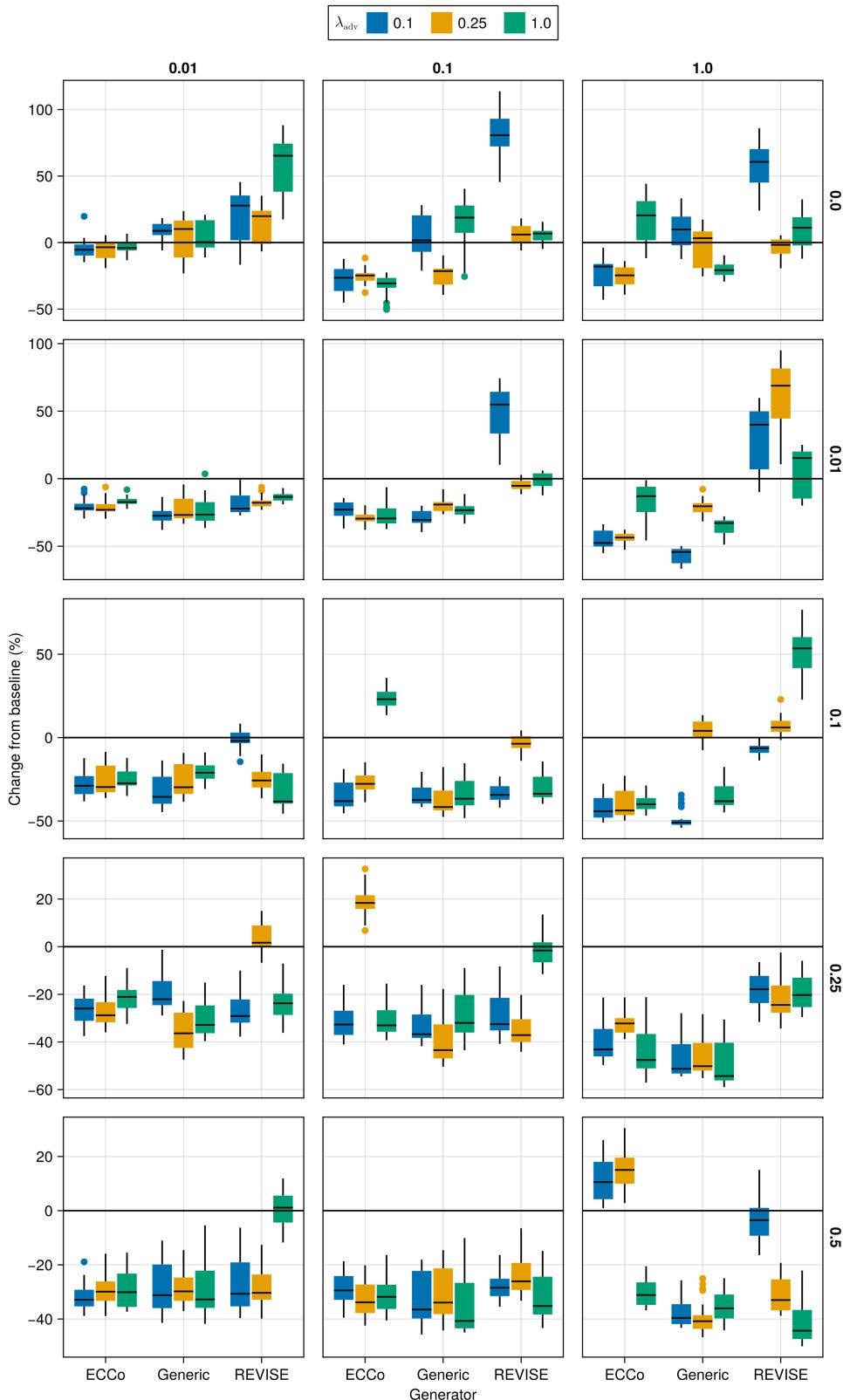


Figure 23 Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Moons.

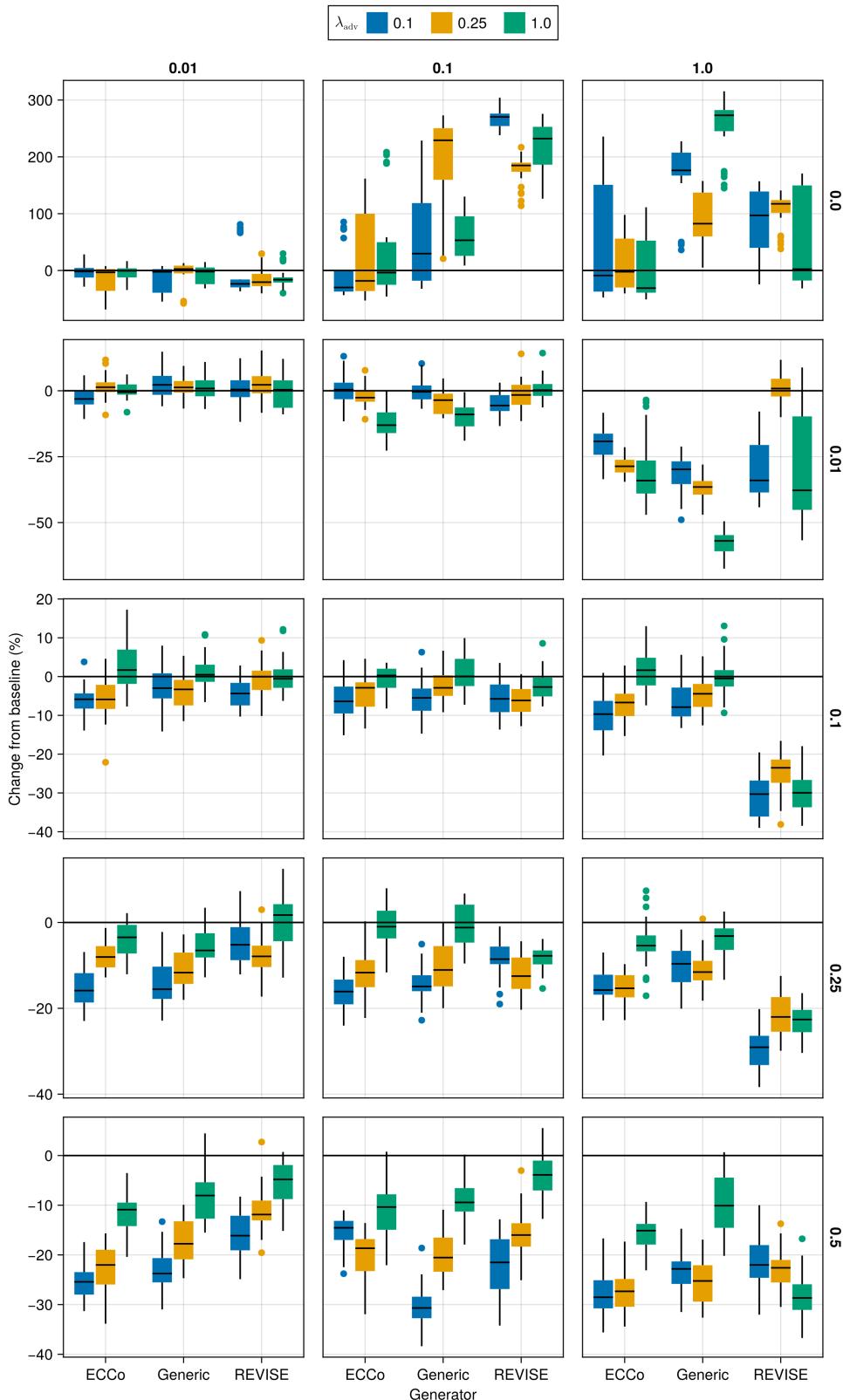


Figure 24 Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

Table 9 Predictive performance measures by dataset and objective averaged across training-phase parameters (Note 7) and evaluation-phase parameters (Note 8).

Dataset	Variable	Objective	Mean	Se
Circ	Accuracy	Full	0.99	0.0
Circ	Accuracy	Vanilla	1.0	0.0
Circ	F1-score	Full	0.99	0.0
Circ	F1-score	Vanilla	1.0	0.0
LS	Accuracy	Full	1.0	0.0
LS	Accuracy	Vanilla	1.0	0.0
LS	F1-score	Full	1.0	0.0
LS	F1-score	Vanilla	1.0	0.0
Moon	Accuracy	Full	1.0	0.01
Moon	Accuracy	Vanilla	0.99	0.02
Moon	F1-score	Full	1.0	0.01
Moon	F1-score	Vanilla	0.99	0.02
OL	Accuracy	Full	0.91	0.01
OL	Accuracy	Vanilla	0.92	0.0
OL	F1-score	Full	0.91	0.01
OL	F1-score	Vanilla	0.92	0.0

D.4.2 Plausibility

The results with respect to the plausibility measure are shown in Figure 25 to Figure 28.

D.4.3 Cost

The results with respect to the cost measure are shown in Figure 29 to Figure 32.

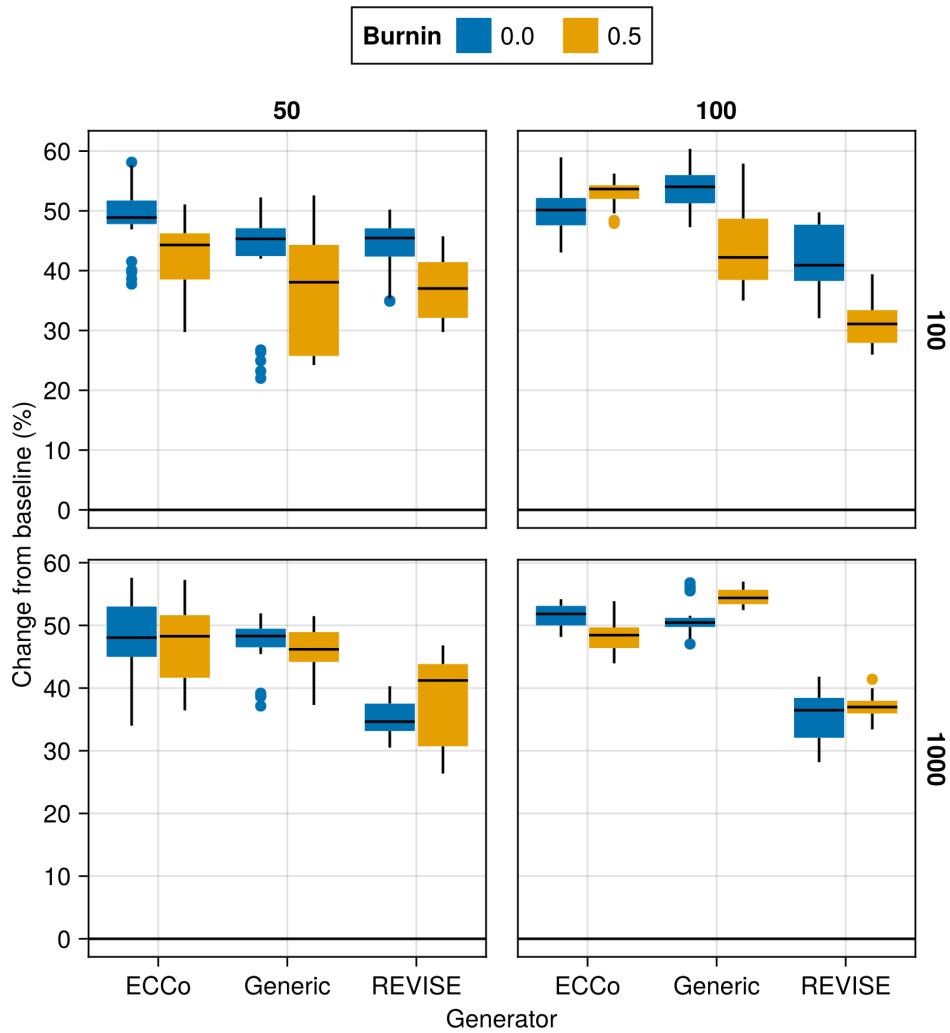


Figure 25 Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Circles.

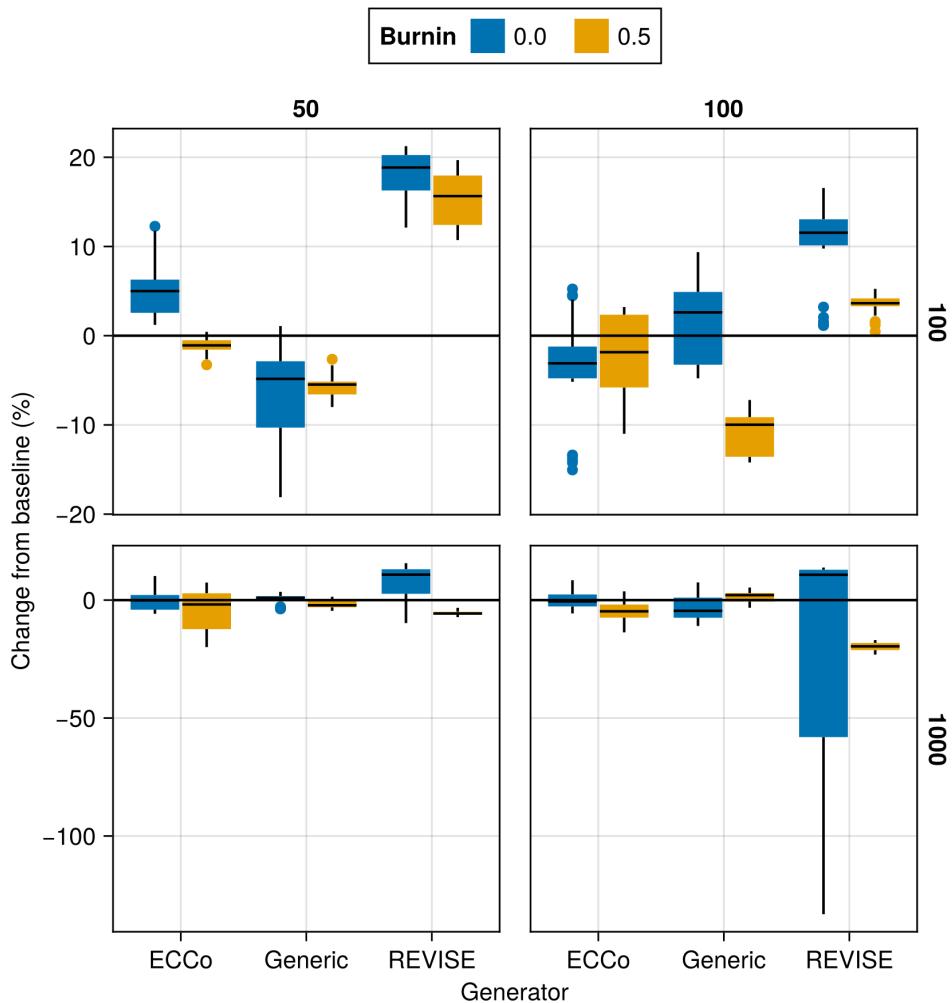


Figure 26 Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Linearly Separable.

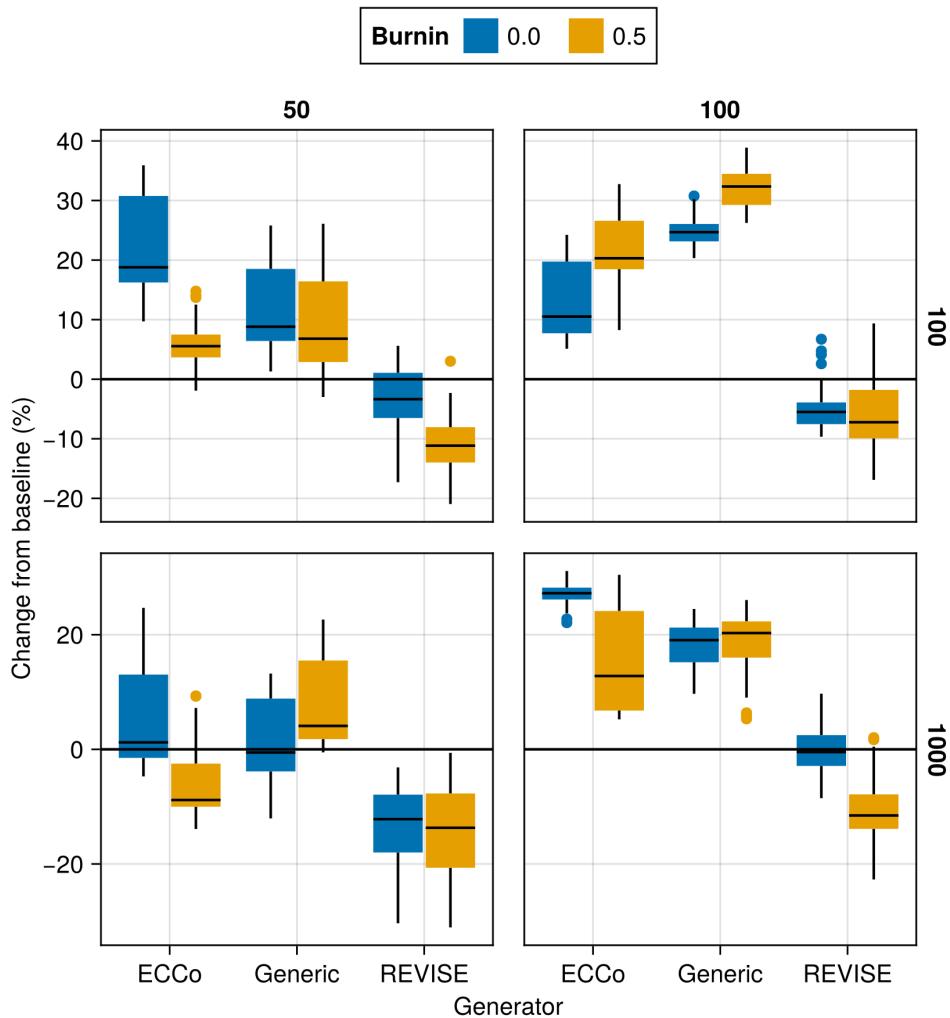


Figure 27 Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Moons.

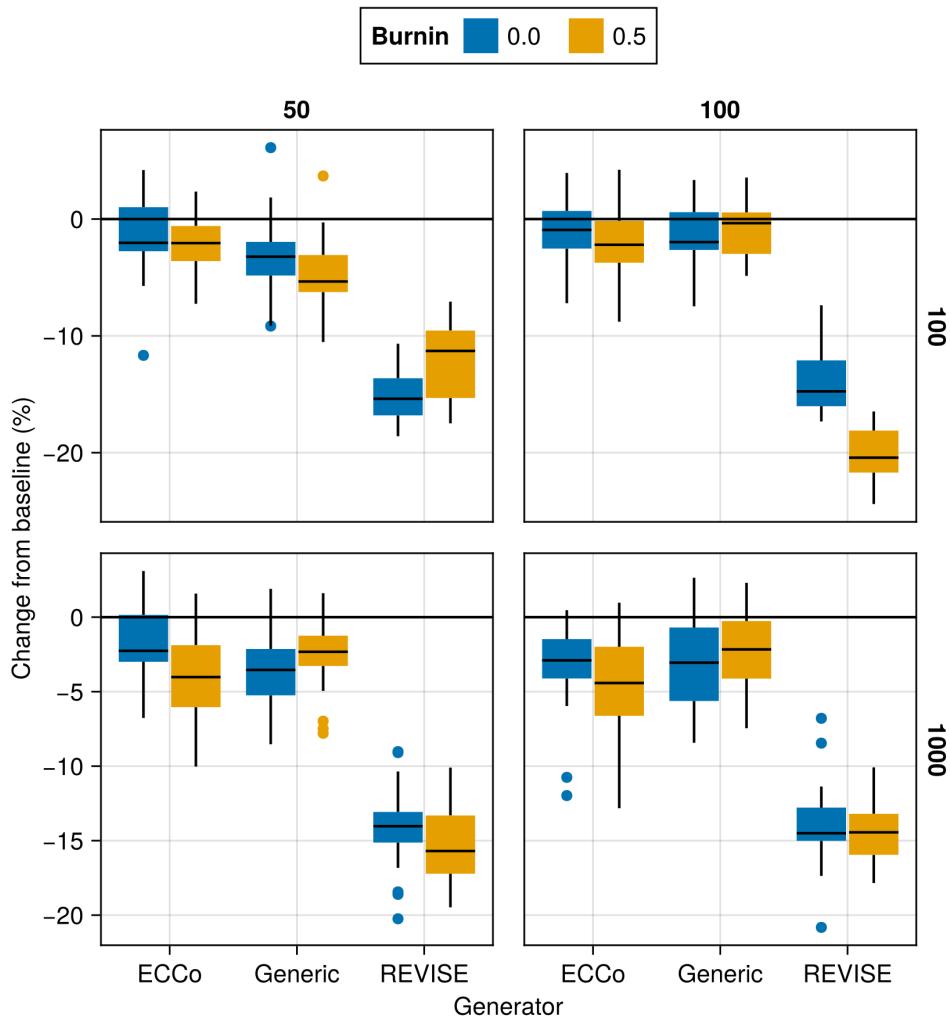


Figure 28 Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

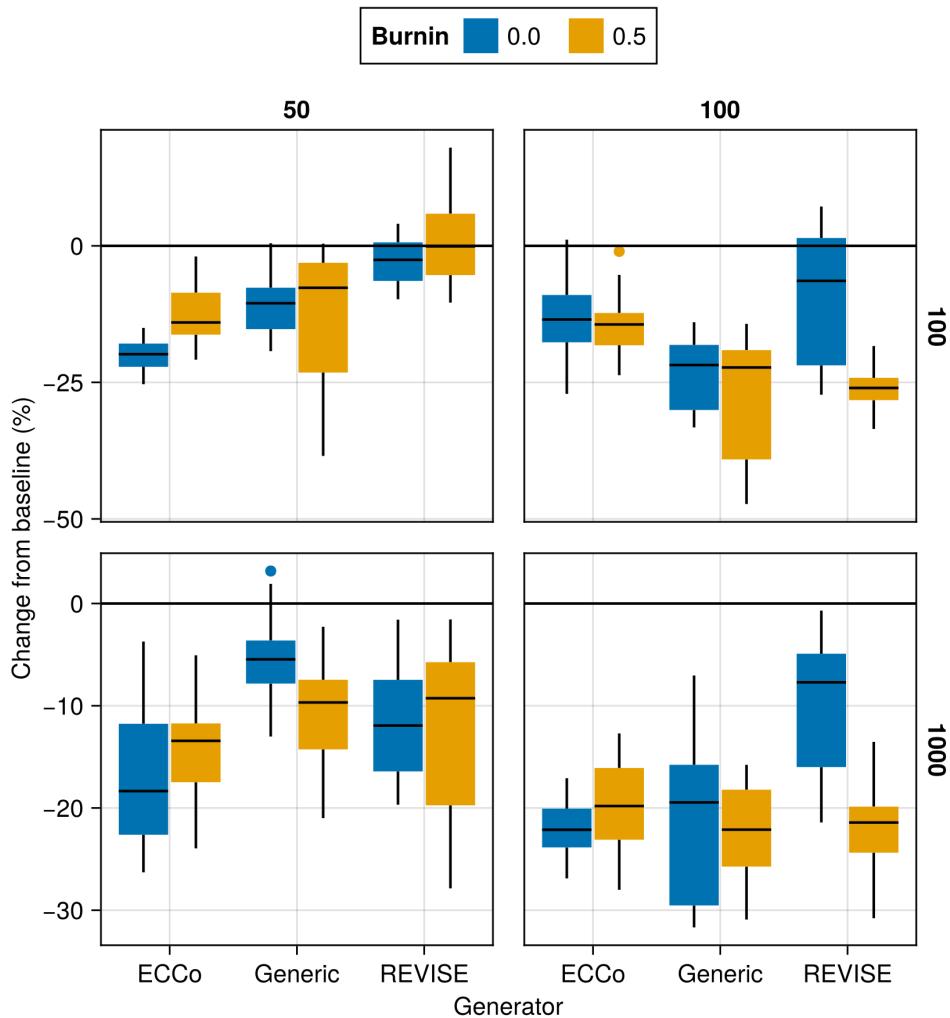


Figure 29 Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Circles.

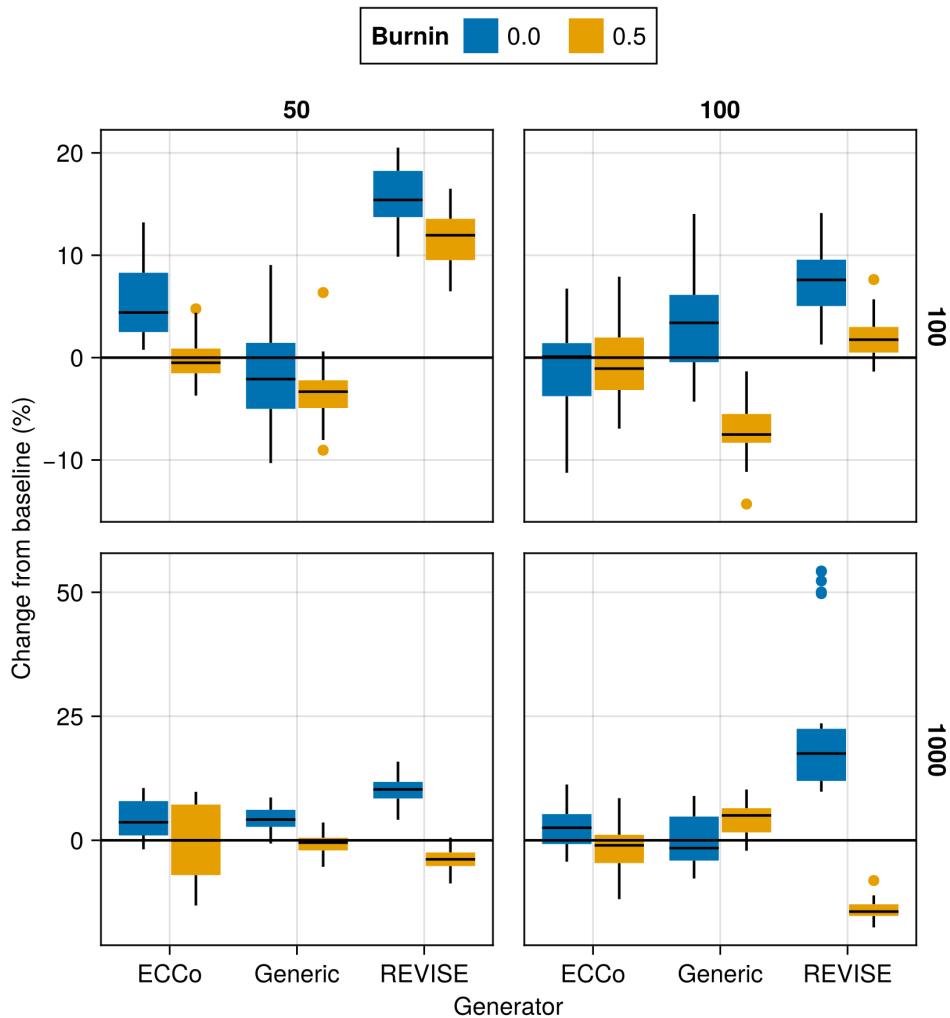


Figure 30 Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Linearly Separable.

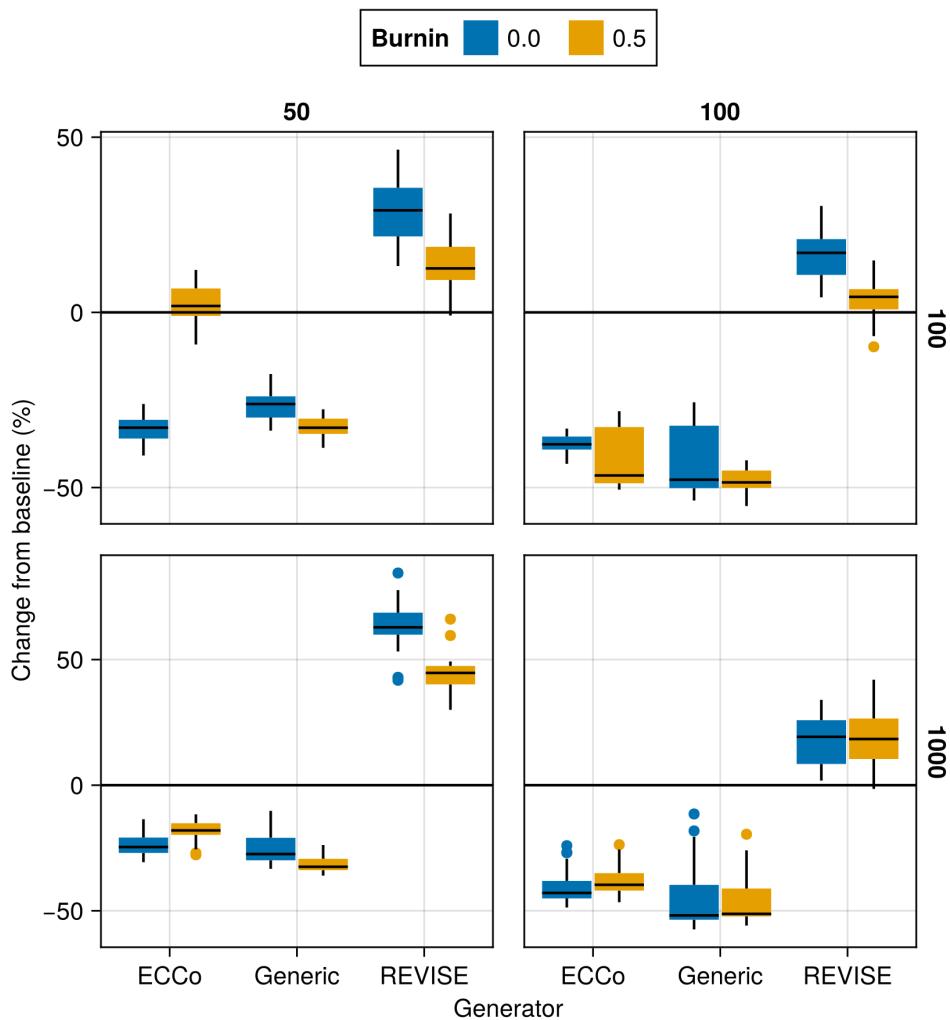


Figure 31 Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Moons.

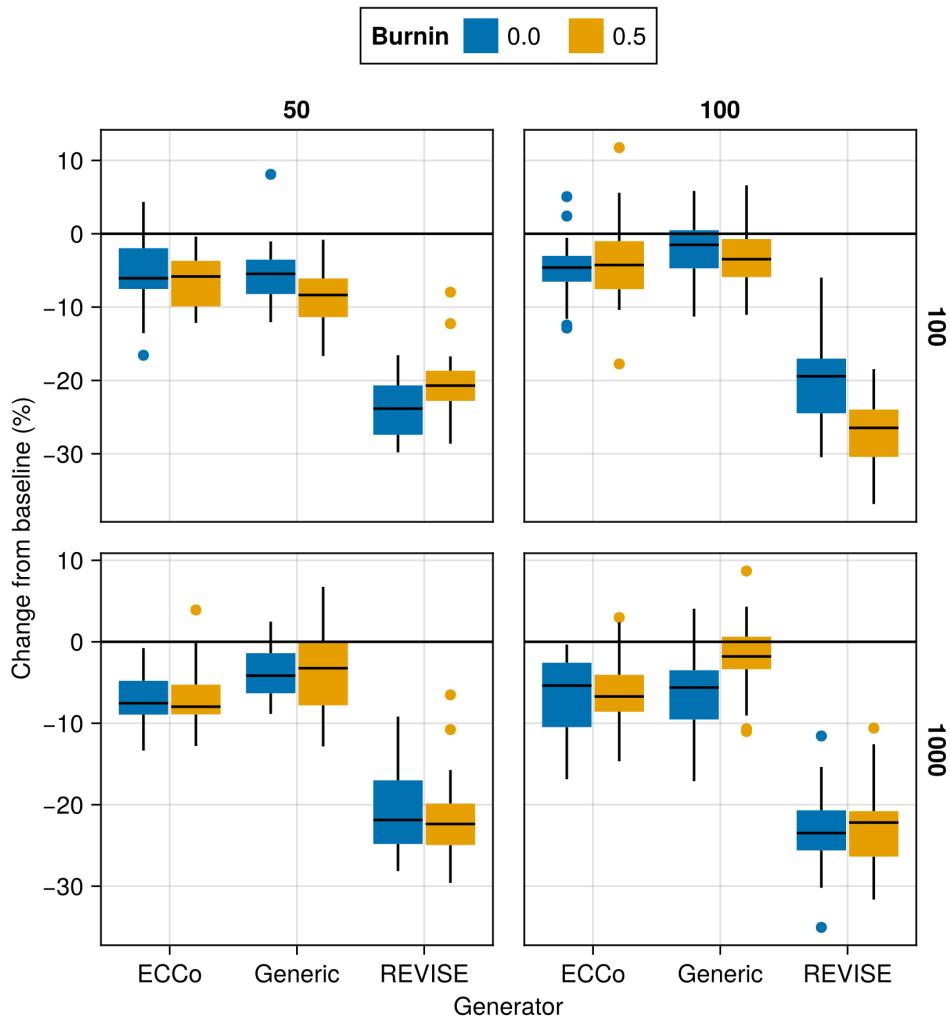


Figure 32 Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

Appendix E Tuning Key Parameters

Based on the findings from our initial large grid searches (Section D), we tune selected hyperparameters for all datasets: namely, the decision threshold τ and the strength of the energy regularization λ_{reg} . The final hyperparameter choices for each dataset are presented in Table 4 in Section C. Detailed results for each data set are shown in Figure 33 to Figure 50. From Table 4, we notice that the same decision threshold of $\tau = 0.5$ is optimal for all but one dataset. We attribute this to the fact that a low decision threshold results in a higher share of mature counterfactuals and hence more opportunities for the model to learn from examples (Figure 42 to Figure 50). This has played a role in particular for our real-world tabular datasets and MNIST, which suffered from low levels of maturity for higher decision thresholds. In cases where maturity is not an issue, as for *Moons*, higher decision thresholds lead to better outcomes, which may have to do with the fact that the resulting counterfactuals are more faithful to the model. Concerning the regularization strength, we find somewhat high variation across datasets. Most notably, we find that relatively low levels of regularization are optimal for MNIST. We hypothesize that this finding may be attributed to the uniform scaling of all input features (digits).

Finally, to increase the proportion of mature counterfactuals for some datasets, we have also investigated the effect on the learning rate η for the counterfactual search and even smaller regularization strengths for a fixed decision threshold of 0.5 (Figure 51 to Figure 59). For the given low decision threshold, we find that the learning rate has no discernible impact on the proportion of mature counterfactuals (Figure 60 to Figure 68). We do notice, however, that the results for MNIST are much improved when using a low value λ_{reg} , the strength for the energy regularization: plausibility is increased by up to $\sim 10\%$ (Figure 57) and the proportion of mature counterfactuals reaches 100%.

One consideration worth exploring is to combine high decision thresholds with high learning rates, which we have not investigated here.

E.1 Key Parameters

The hyperparameter grid for tuning key parameters is shown in Note 9. The corresponding evaluation grid used for these experiments is shown in Note 10.

Note 9space Training Phase

- Generator Parameters:
 - Decision Threshold: 0.5, 0.75, 0.9
- Model: mlp
- Training Parameters:
 - λ_{reg} : 0.1, 0.25, 0.5
 - Objective: full, vanilla

Note 10space Evaluation Phase

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

E.1.1 Plausibility

The results with respect to the plausibility measure are shown in Figure 33 to Figure 41.

E.1.2 Proportion of Mature CE

The results with respect to the proportion of mature counterfactuals in each epoch are shown in Figure 42 to Figure 50.

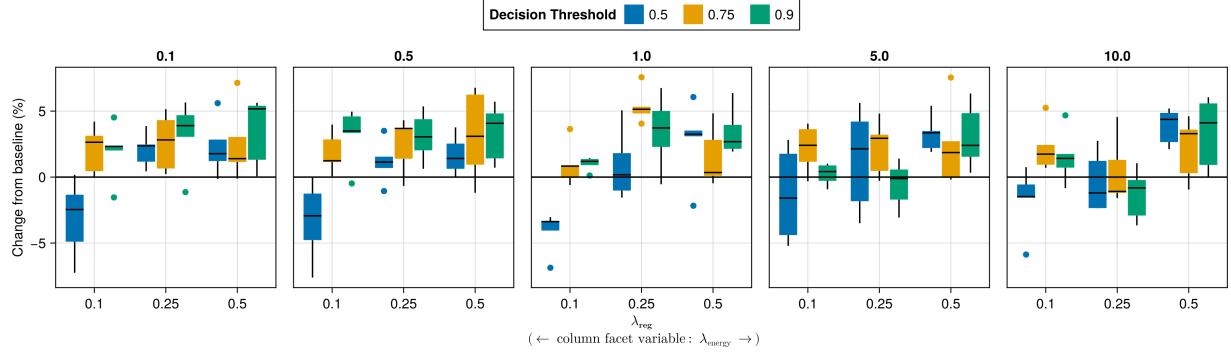


Figure 33 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Adult.

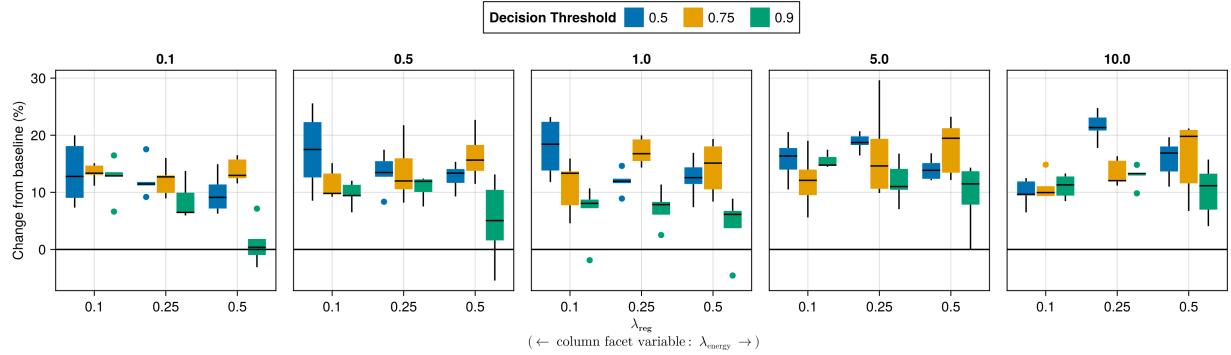


Figure 34 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: California Housing.

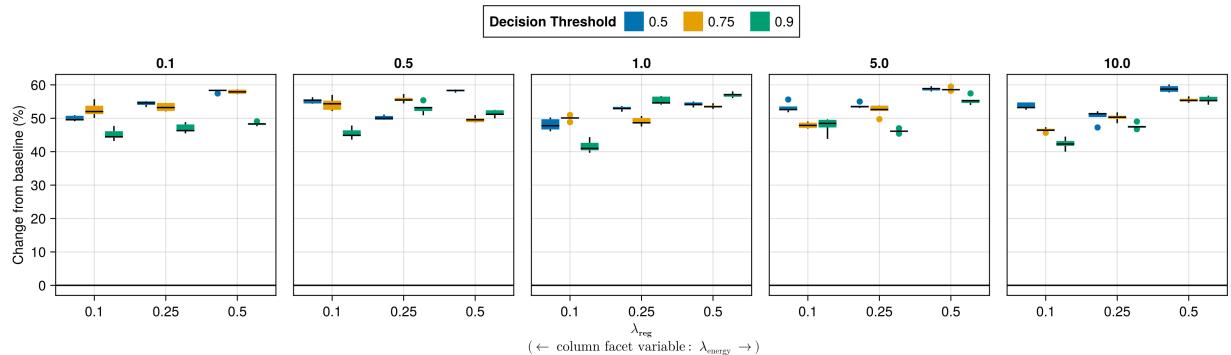


Figure 35 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Circles.

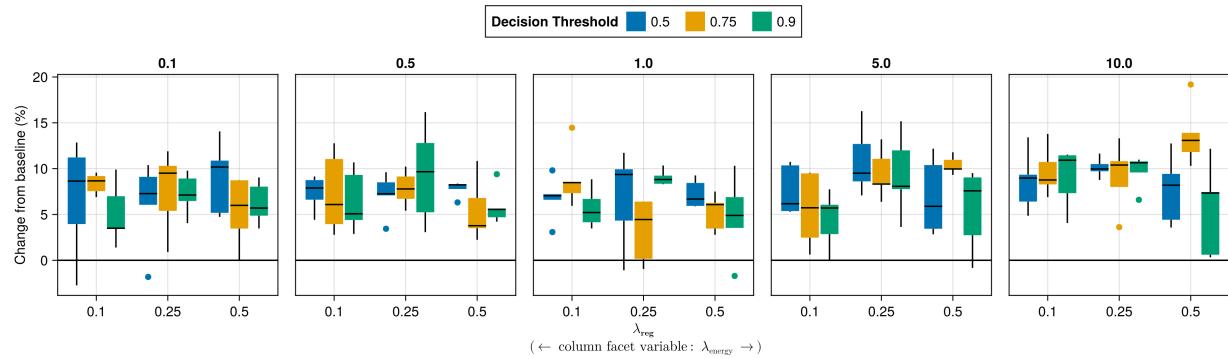


Figure 36 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for ECCCo). Data: Credit.

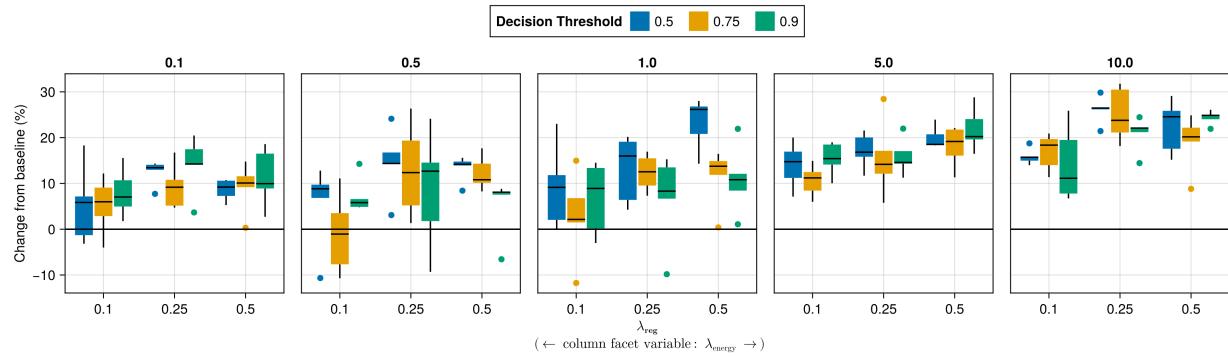


Figure 37 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for ECCCo). Data: GMSC.

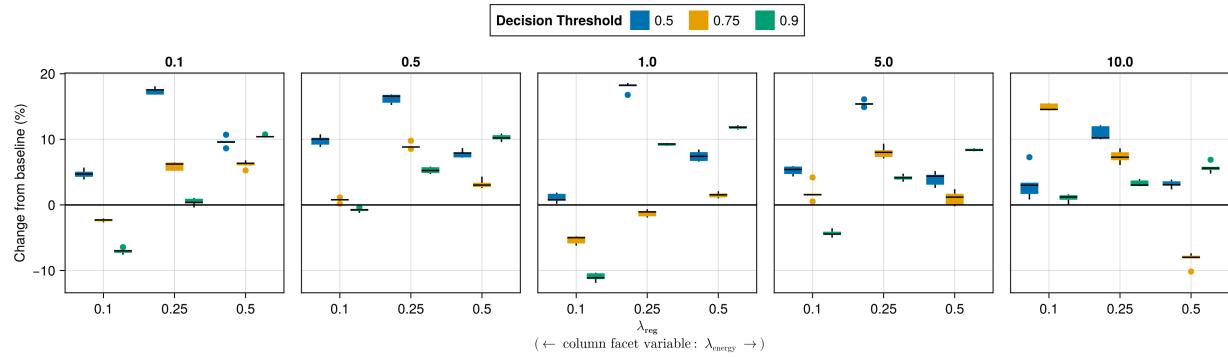


Figure 38 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for ECCCo). Data: Linearly Separable.

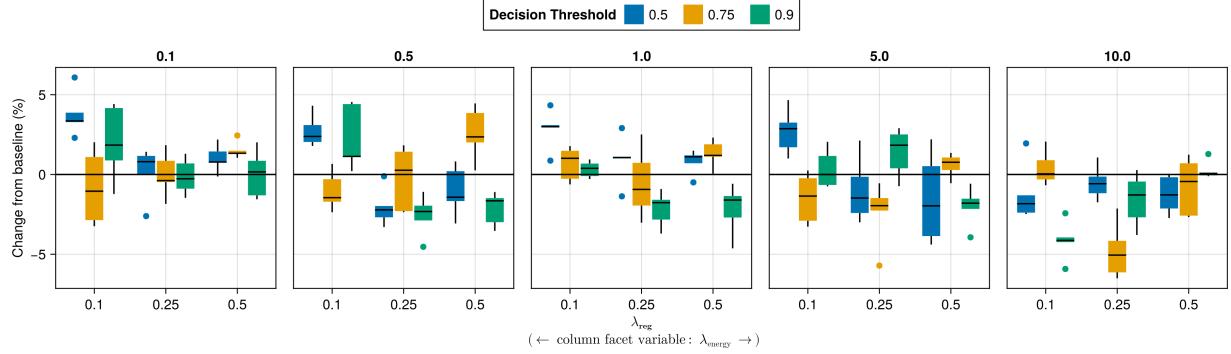


Figure 39 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: MNIST.

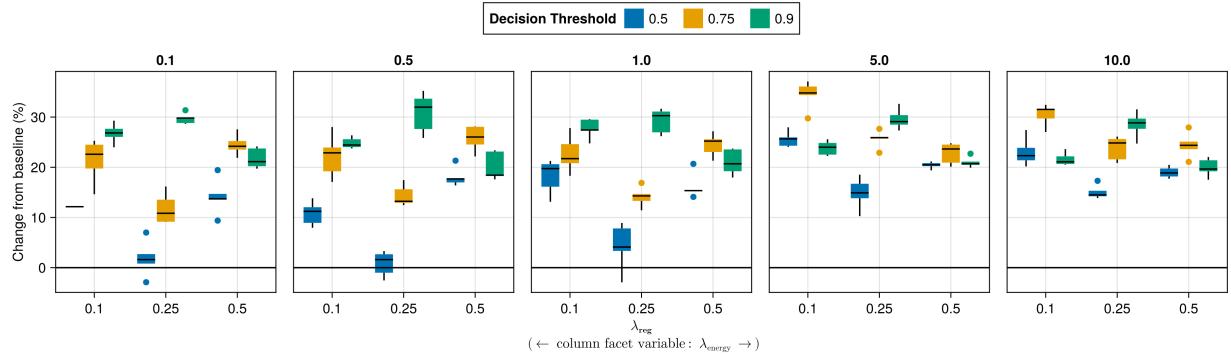


Figure 40 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Moons.

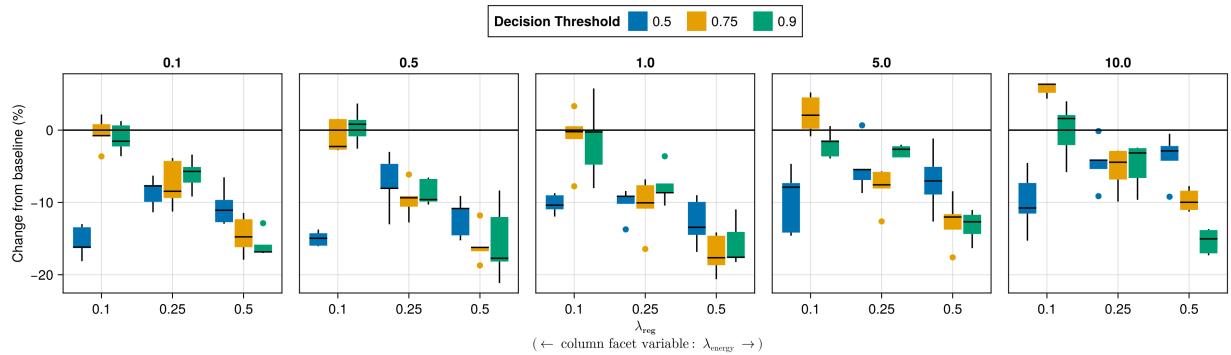


Figure 41 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Overlapping.

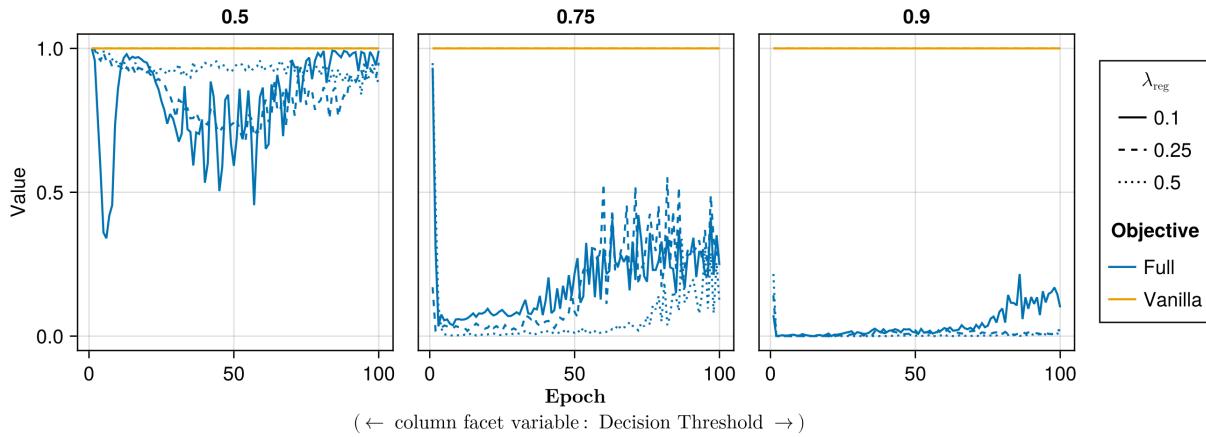


Figure 42 Proportion of mature counterfactuals in each epoch. Data: Adult.

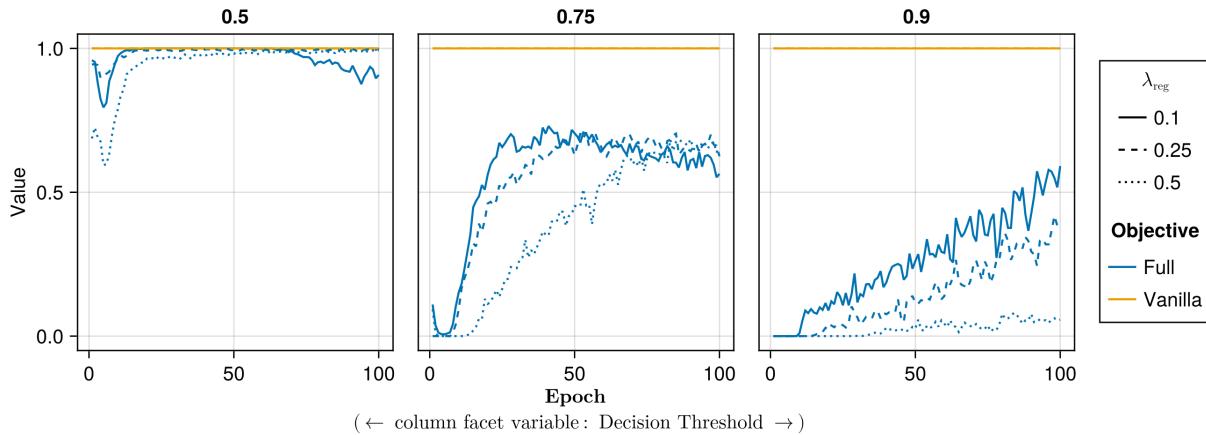


Figure 43 Proportion of mature counterfactuals in each epoch. Data: California Housing.

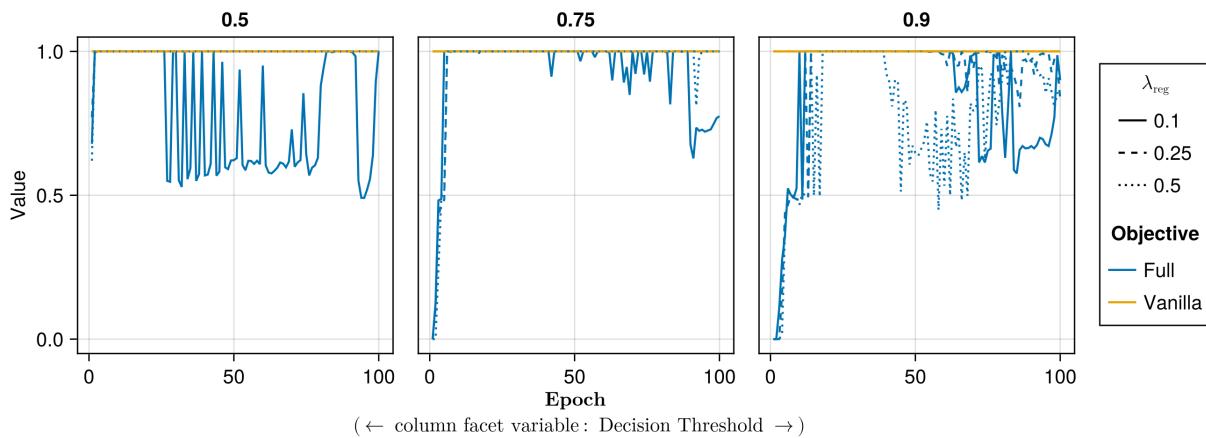


Figure 44 Proportion of mature counterfactuals in each epoch. Data: Circles.

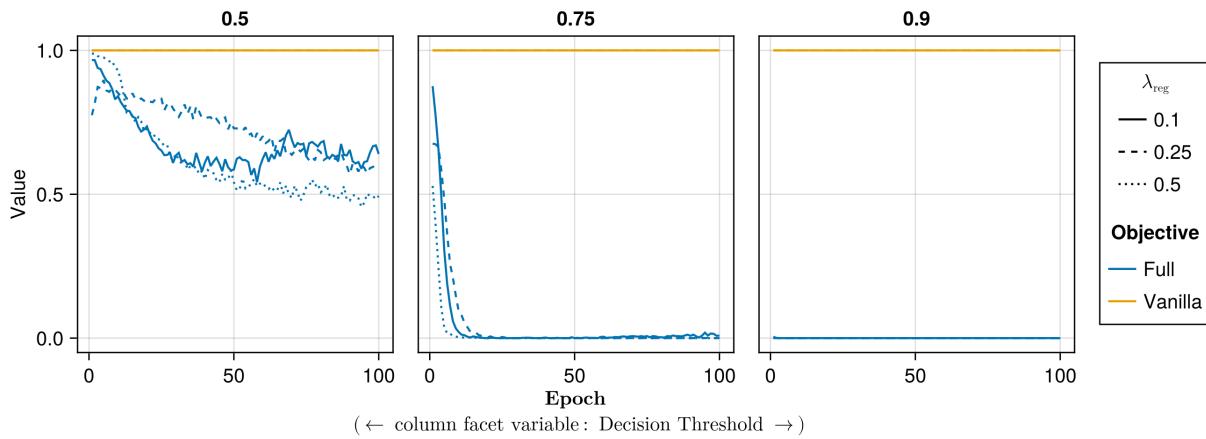


Figure 45 Proportion of mature counterfactuals in each epoch. Data: Credit.

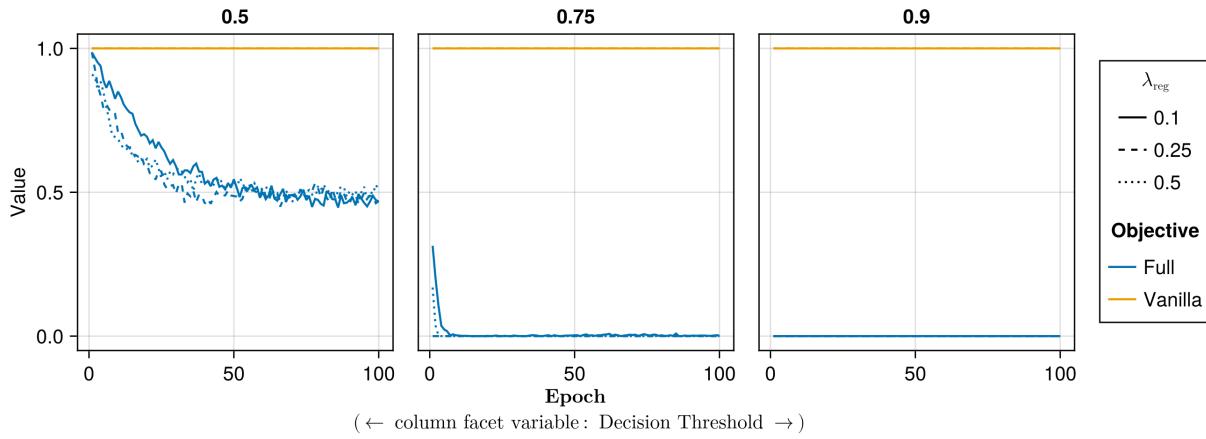


Figure 46 Proportion of mature counterfactuals in each epoch. Data: GMSC.

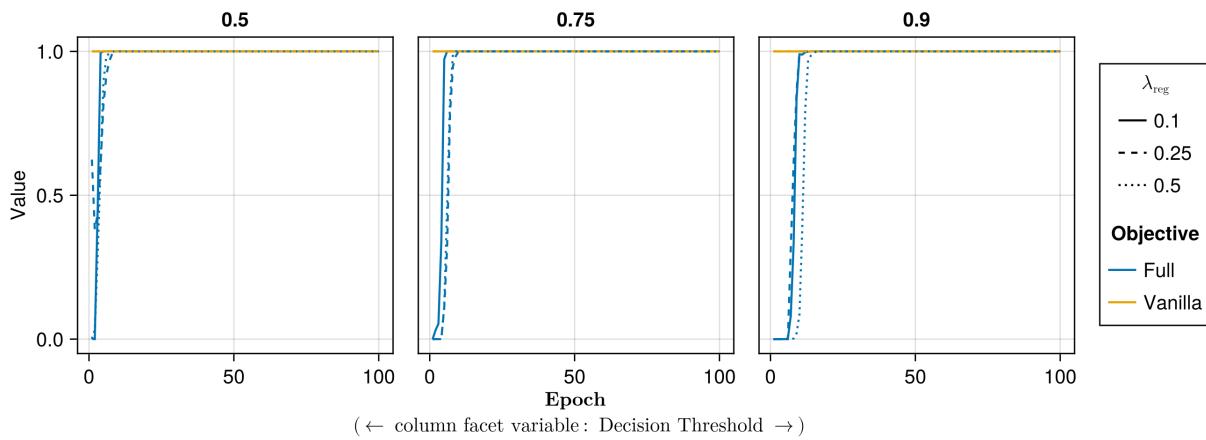


Figure 47 Proportion of mature counterfactuals in each epoch. Data: Linearly Separable.

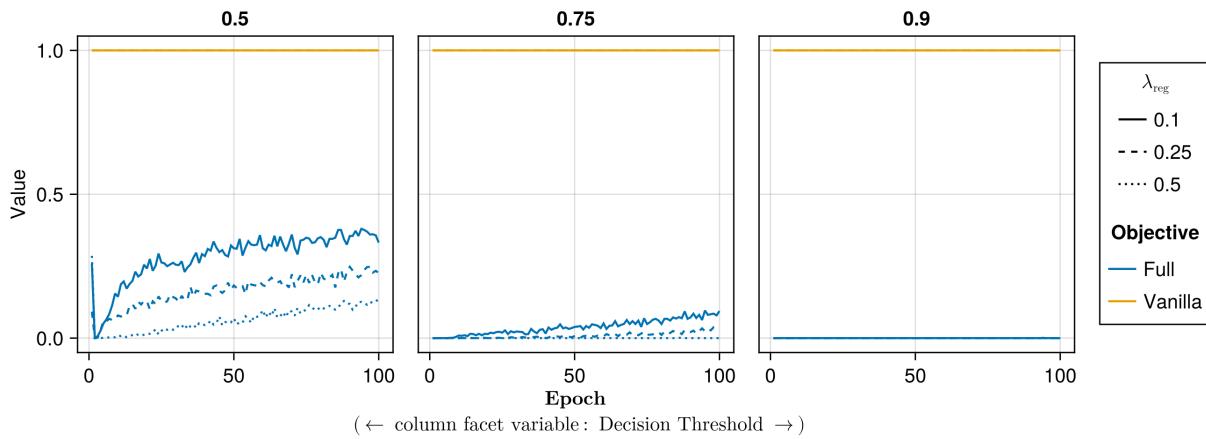


Figure 48 Proportion of mature counterfactuals in each epoch. Data: MNIST.

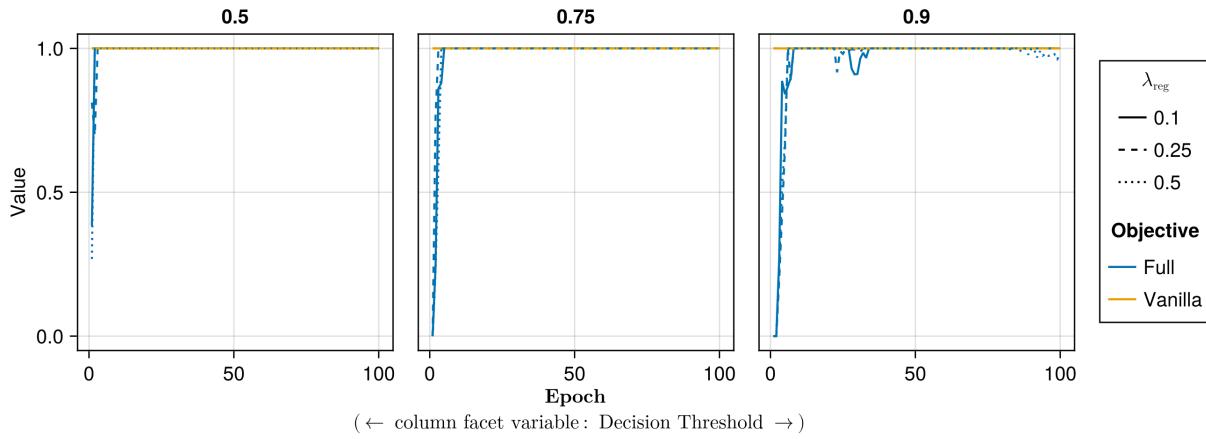


Figure 49 Proportion of mature counterfactuals in each epoch. Data: Moons.

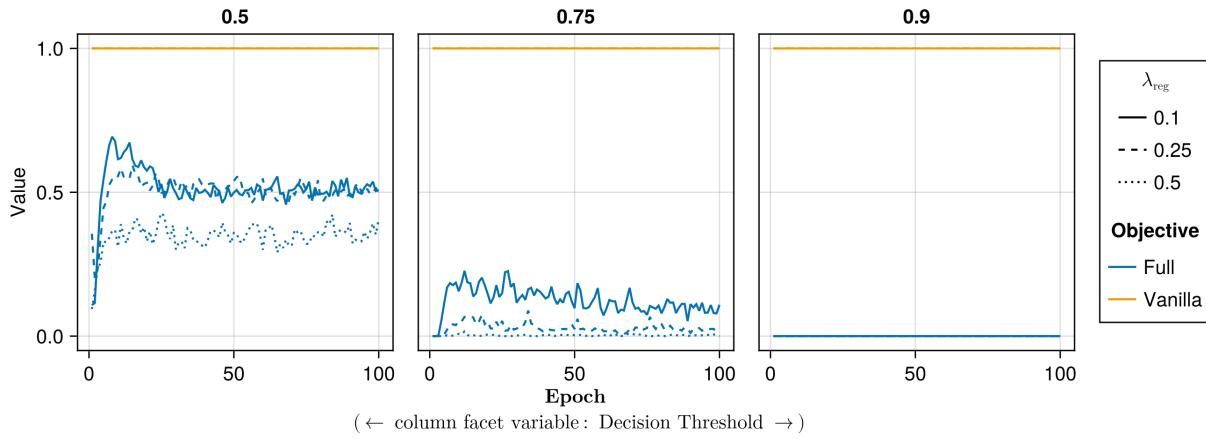


Figure 50 Proportion of mature counterfactuals in each epoch. Data: Overlapping.

E.2 Learning Rate

The hyperparameter grid for tuning the learning rate is shown in Note 11. The corresponding evaluation grid used for these experiments is shown in Note 12.

Note 11space Training Phase

- Generator Parameters:
 - Learning Rate: 0.1, 0.5, 1.0
- Model: mlp
- Training Parameters:
 - λ_{reg} : 0.01, 0.1, 0.5
 - Objective: full, vanilla

Note 12space Evaluation Phase

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

E.2.1 Plausibility

The results with respect to the plausibility measure are shown in Figure 51 to Figure 59.

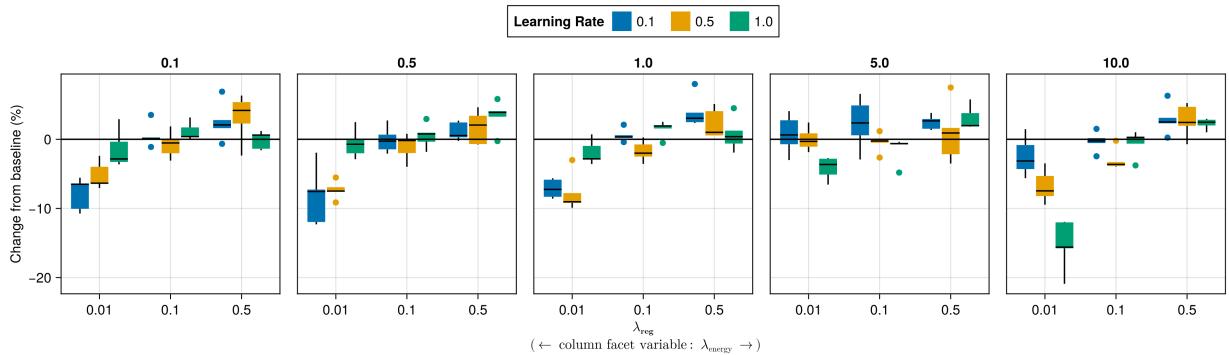


Figure 51 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for ECCCo). Data: Adult.

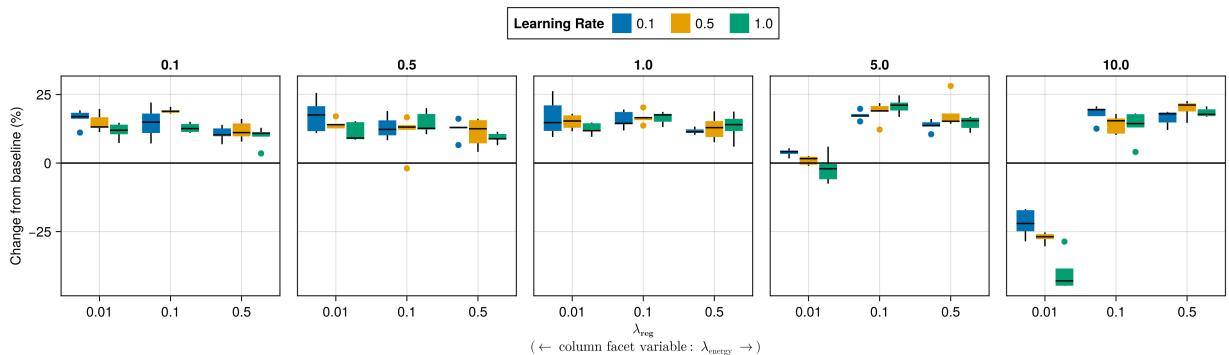


Figure 52 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for ECCCo). Data: California Housing.

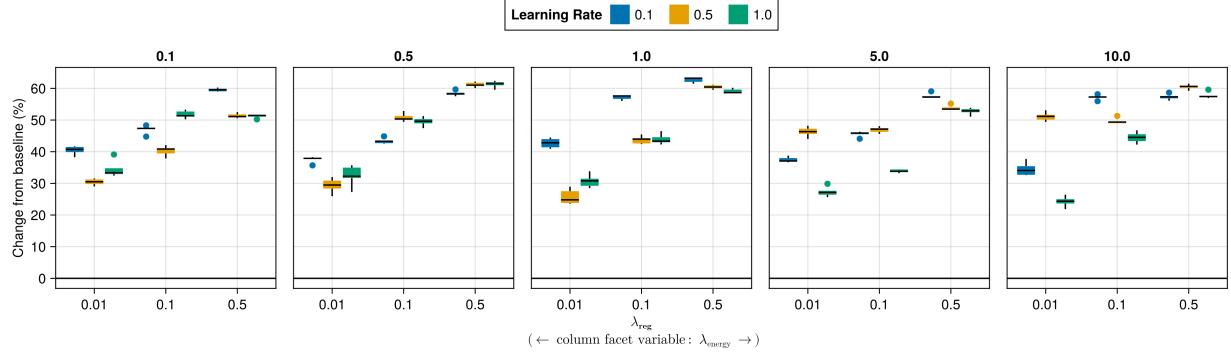


Figure 53 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Circles.

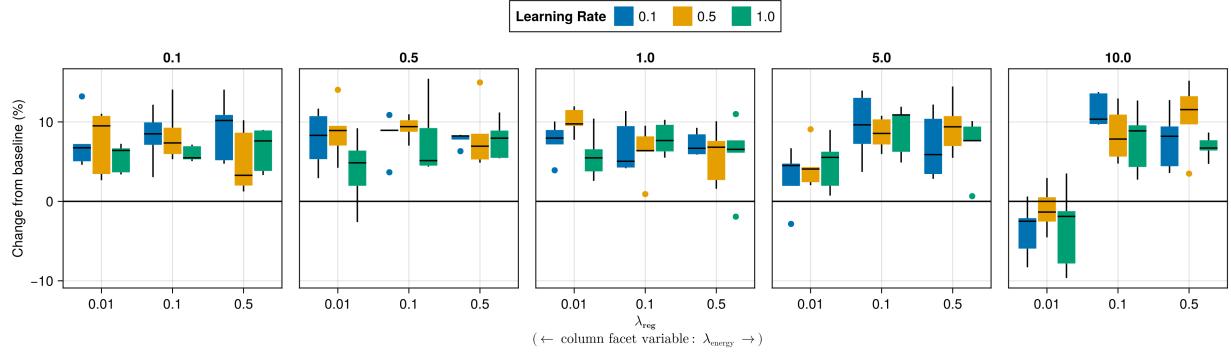


Figure 54 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Credit.

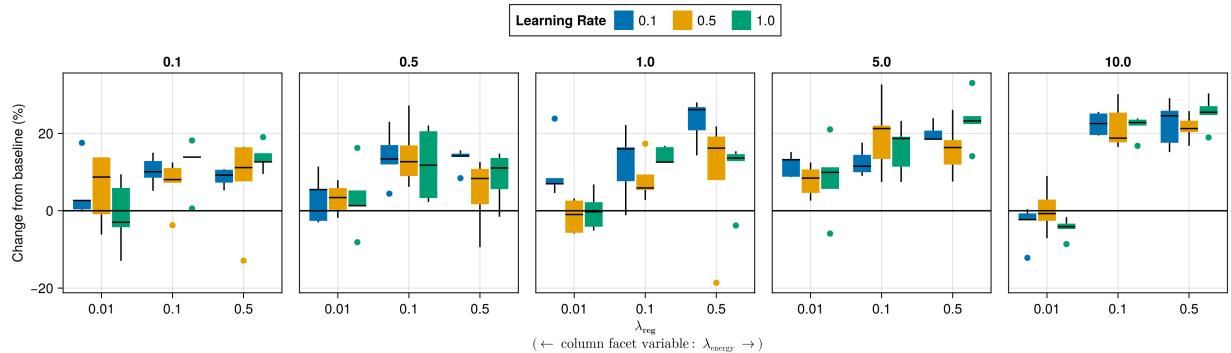


Figure 55 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: GMSC.

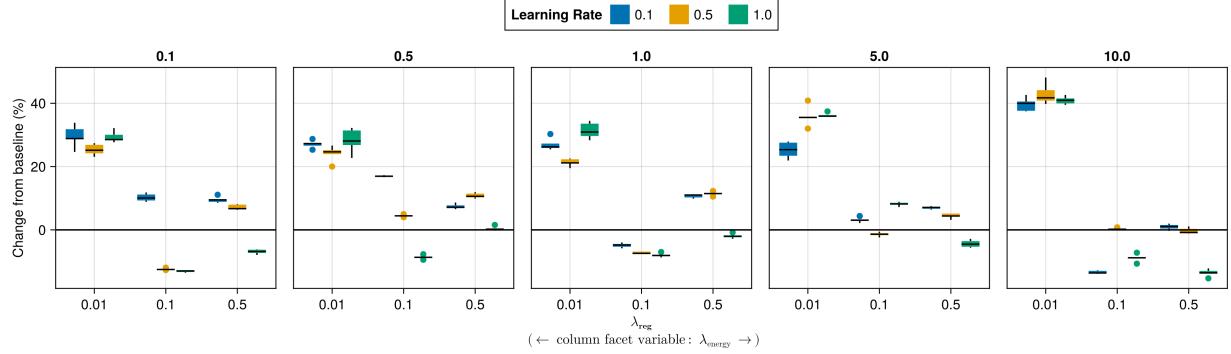


Figure 56 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for $ECCCo$). Data: Linearly Separable.

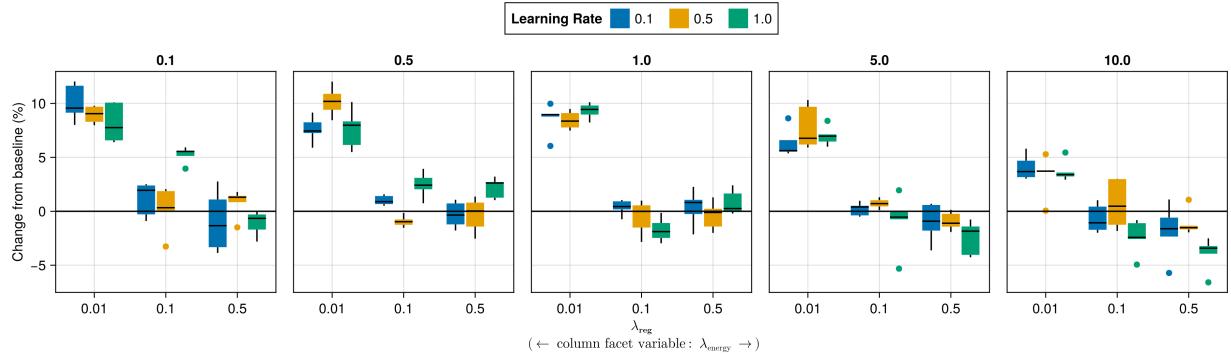


Figure 57 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for $ECCCo$). Data: MNIST.

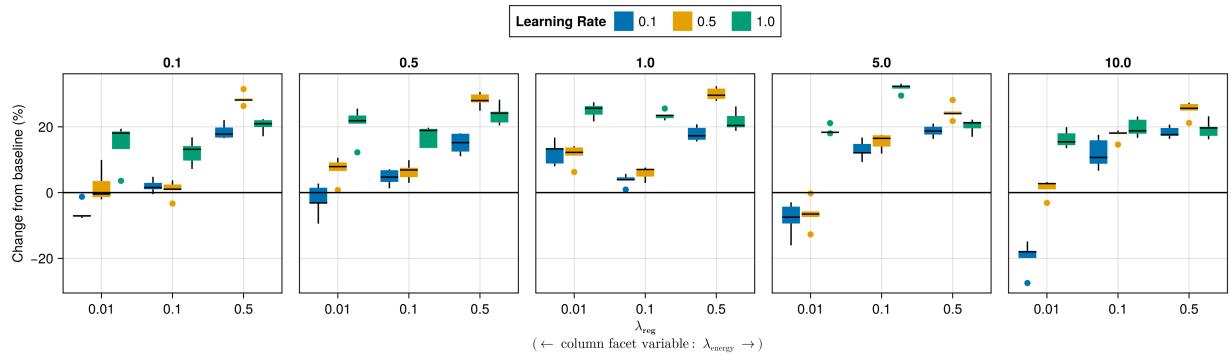


Figure 58 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for $ECCCo$). Data: Moons.

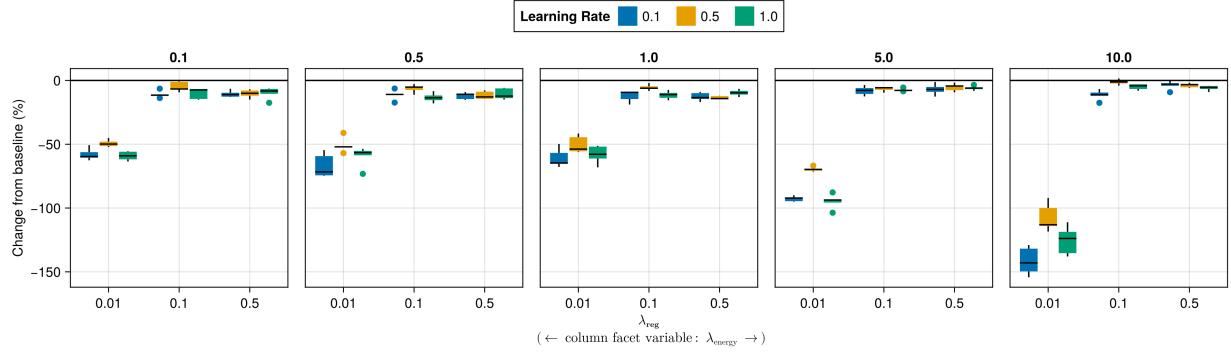


Figure 59 Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (IP). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Overlapping.

E.2.2 Proportion of Mature CE

The results with respect to the proportion of mature counterfactuals in each epoch are shown in Figure 60 to Figure 68.

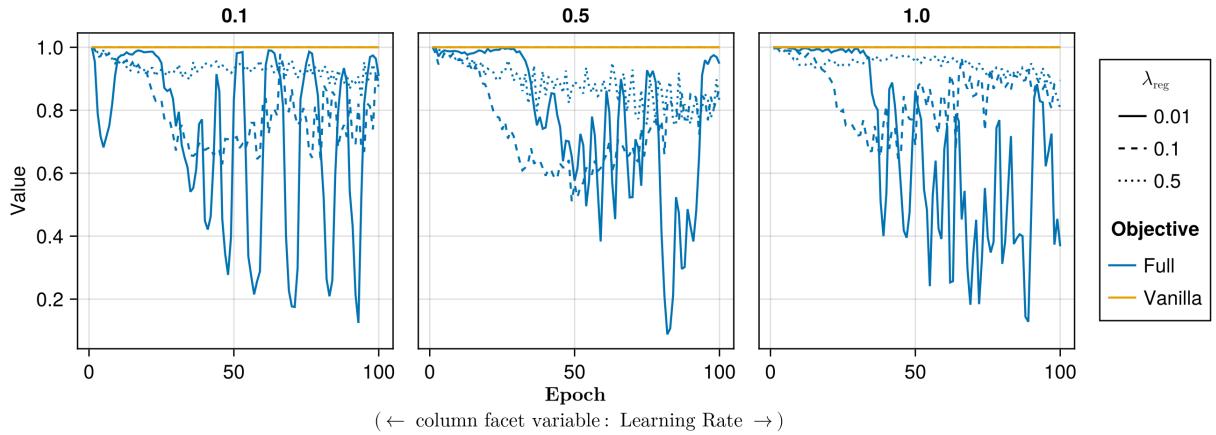


Figure 60 Proportion of mature counterfactuals in each epoch. Data: Adult.

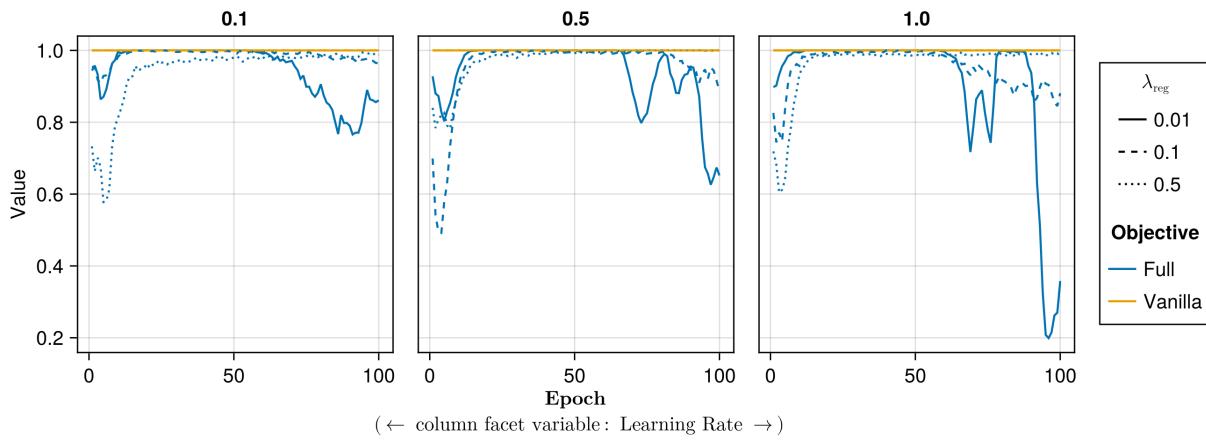


Figure 61 Proportion of mature counterfactuals in each epoch. Data: California Housing.

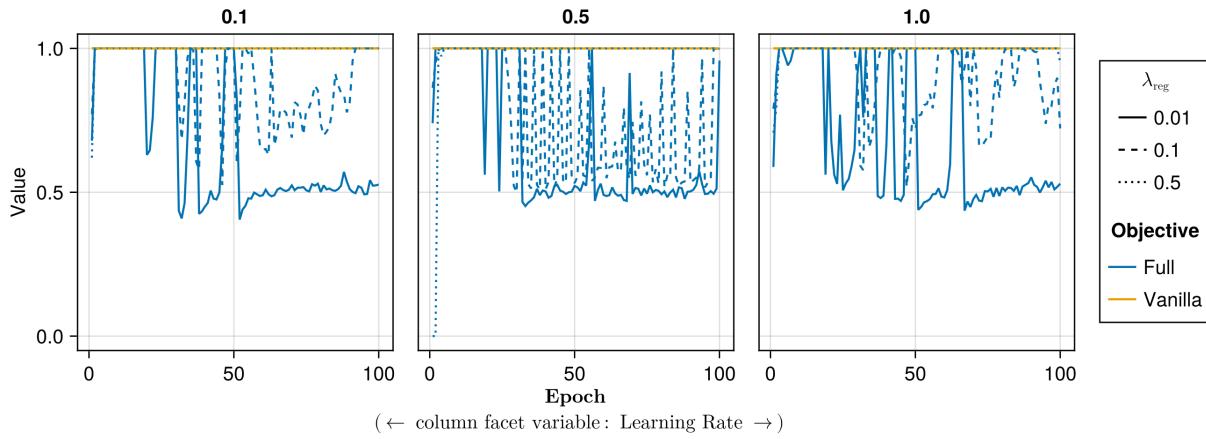


Figure 62 Proportion of mature counterfactuals in each epoch. Data: Circles.

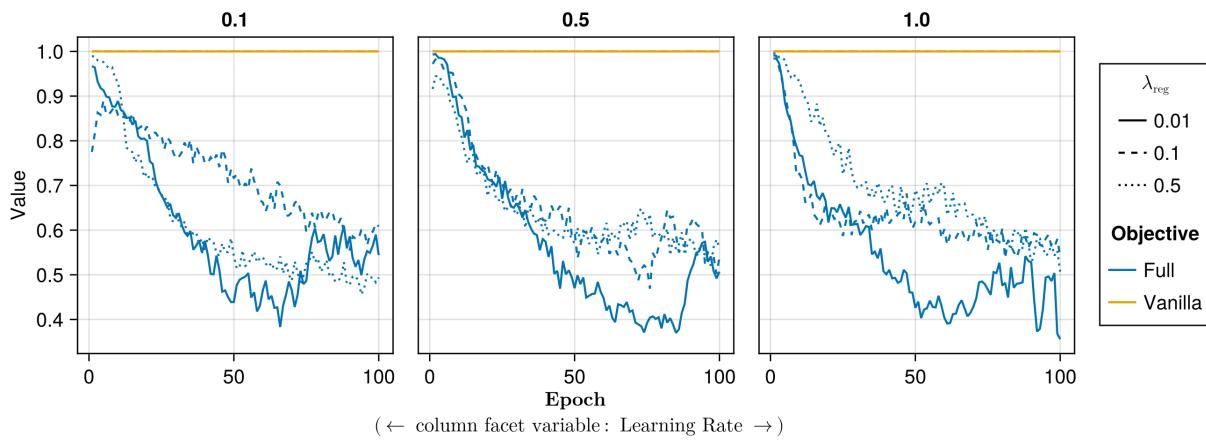


Figure 63 Proportion of mature counterfactuals in each epoch. Data: Credit.

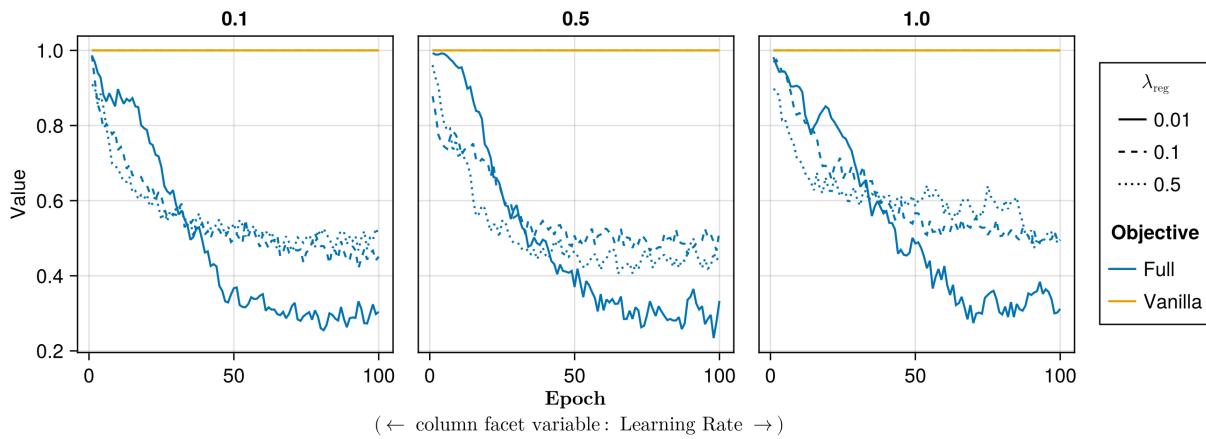


Figure 64 Proportion of mature counterfactuals in each epoch. Data: GMSC.

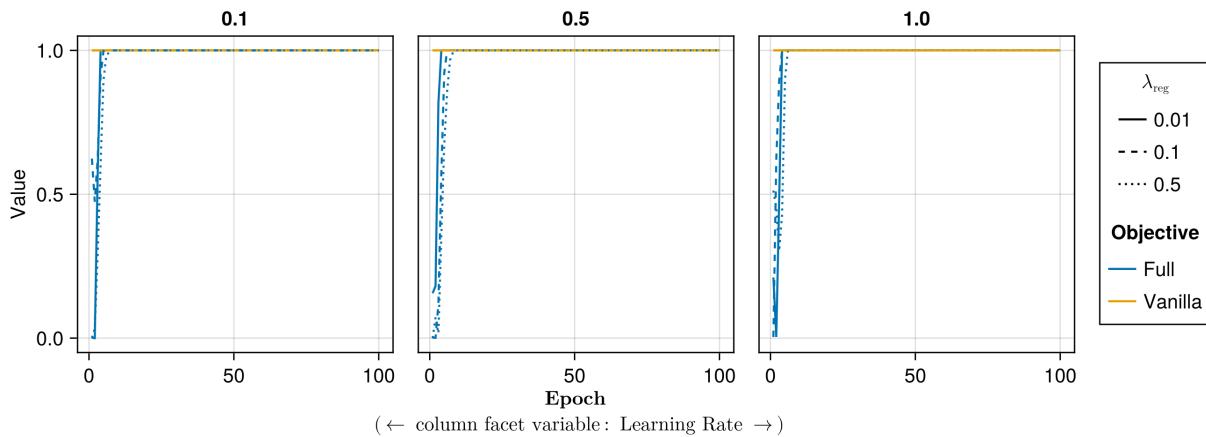


Figure 65 Proportion of mature counterfactuals in each epoch. Data: Linearly Separable.

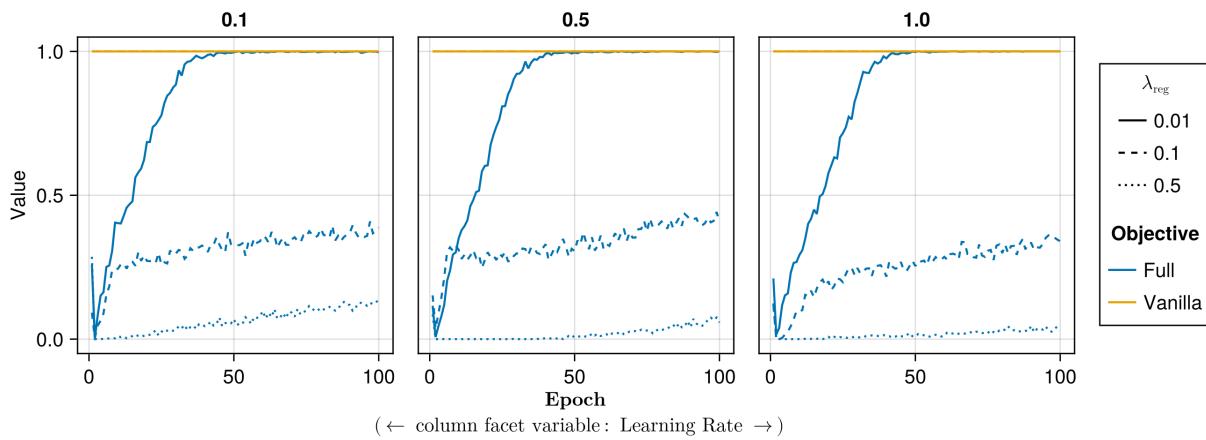


Figure 66 Proportion of mature counterfactuals in each epoch. Data: MNIST.

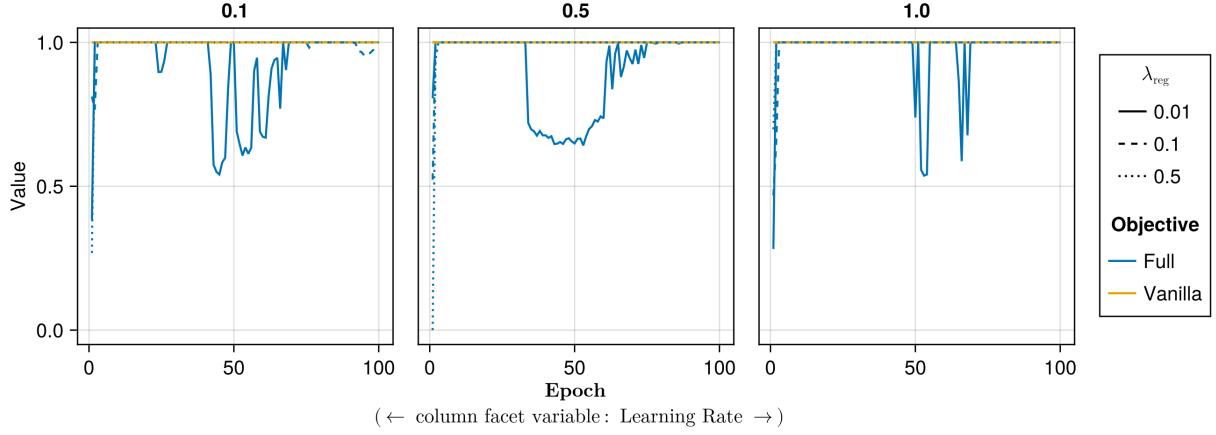


Figure 67 Proportion of mature counterfactuals in each epoch. Data: Moons.

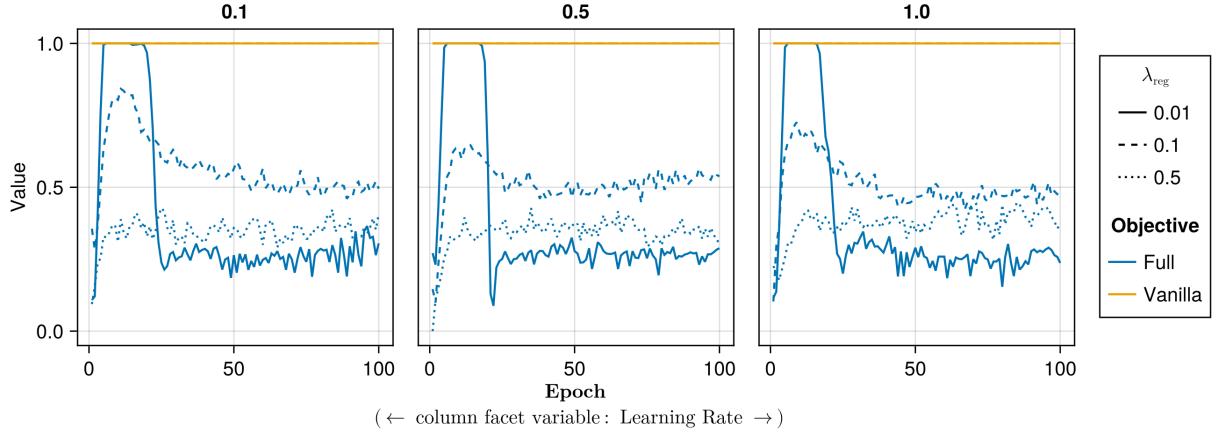


Figure 68 Proportion of mature counterfactuals in each epoch. Data: Overlapping.

Appendix F Computation Details

F.1 Hardware

We performed our experiments on a high-performance cluster ([\(DHPC\) 2022](#)). 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* used.

F.1.1 Grid Searches

Model training for the largest grid searches with 270 unique parameter combinations was parallelized across 34 CPUs with 2GB memory each. The time to completion varied by dataset: 0h49m (*Moons*), 1h4m (*Linearly Separable*), 1h49m (*Circles*), 3h52m (*Overlapping*). Model evaluations for large grid searches were parallelized across 20 CPUs with 3GB memory each. Evaluations for all data sets took less than one hour (<1h) to complete but were generally more memory-intensive (see Section F.2 for additional details)

F.1.2 Tuning

For tuning of selected hyperparameters, we distributed the task of generating counterfactuals during training across 40 CPUs with 2GB memory each for all tabular datasets. Except for the *Adult* dataset, all training runs were completed in less than half an hour (<0h30m). The *Adult* dataset took around 0h35m to complete. Evaluations across 20 CPUs with 3GB memory each generally took less than 0h30m to complete. For *MNIST*, we relied on 100 CPUs with 2GB memory each. For the *MLP*, training of all models could be completed in 1h30m, while the evaluation across 20 CPUs

(6GB memory) took 4h12m. For the *CNN*, training of all models took ~8h, with conventionally trained models taking ~0h15m each and model with CT taking ~0h30m-0h45m each.

F.2 Software

Our code has been open-sourced on GitHub as Julia package: [CounterfactualTraining.jl](#). All computations were performed in the Julia Programming Language (Bezanson et al. 2017). We have developed a package for counterfactual training that leverages and extends the functionality provided by several existing packages, most notably [CounterfactualExplanations.jl](#) (Altmeyer, Deursen, and Liem 2023) and the [Flux.jl](#) library for deep learning (Michael Innes et al. 2018; Mike Innes 2018). We chose to work with [CounterfactualExplanations.jl](#) because it currently appears to be the most comprehensive and extensible package for counterfactual explanations. Despite its good interplay with Flux.jl, the package is not, however, optimized to be used in training. This has caused some issues with memory management and bottlenecked performance. The code is commented with clearly marked references to the paper (look for `# ----- PAPER REF -----`).

For data-wrangling and presentation-ready tables we relied on [DataFrames.jl](#) (Bouchet-Valat and Kamiski 2023) and [PrettyTables.jl](#) (Chagas et al. 2024), respectively. For plots and visualizations we used both [Plots.jl](#) (Christ et al. 2023) and [Makie.jl](#) (Danisch and Krumbiegel 2021), in particular [AlgebraOfGraphics.jl](#). To distribute computational tasks across multiple processors, we have relied on [MPI.jl](#) (Byrne, Wilcox, and Churavy 2021).

F.3 Reproducibility

We have taken care to set random seeds for reproducibility using Julia’s Random.jl package from the standard library. A global seed and (if applicable or wanted) dataset-specific seeds can be specified in TOML configuration files, environment variables or in interactive Julia sessions. Additional details can be found in the code base.