
COUNTERFACTUAL TRAINING: TEACHING MODELS PLAUSIBLE AND ACTIONABLE EXPLANATIONS

A PREPRINT

Patrick Altmeyer 

Faculty of Electrical Engineering, Mathematics and Computer Science
Delft University of Technology

p.altmeyer@tudelft.nl

Aleksander Buszydlik

Faculty of Electrical Engineering, Mathematics and Computer Science
Delft University of Technology

Arie van Deursen

Faculty of Electrical Engineering, Mathematics and Computer Science
Delft University of Technology

Cynthia C. S. Liem

Faculty of Electrical Engineering, Mathematics and Computer Science
Delft University of Technology

March 15, 2025

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-word decision-making systems, counterfactuals should be (1) plausible with respect to the underlying data and (2) actionable with respect to the user-defined 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 ad-hoc 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 explanations while promoting robustness and preserving high predictive performance.

¹⁴ **Keywords** Counterfactual Training • Counterfactual Explanations • Algorithmic Recourse • Explainable AI • Representation Learning

16 1 Introduction

17 Today’s prominence of artificial intelligence (AI) has largely been driven by **representation learning**: instead of relying
 18 on features and rules that are carefully hand-crafted by humans, modern machine learning (ML) models are tasked
 19 with learning representations directly from data, guided by narrow objectives such as predictive accuracy (Goodfellow,
 20 Bengio, and Courville 2016). Modern advances in computing have made it possible to provide such models with
 21 ever-growing degrees of freedom to achieve that task, which frequently allows them to outperform traditionally more
 22 parsimonious models. Unfortunately, in doing so, models learn increasingly complex and highly sensitive representa-
 23 tions that humans can no longer easily interpret.

24 The trend towards complexity for the sake of performance has come under serious scrutiny in recent years. At the
 25 very cusp of the deep learning (DL) revolution, Szegedy et al. (2014) showed that artificial neural networks (ANN)
 26 are sensitive to adversarial examples (AEs): perturbed versions of data instances that yield vastly different model
 27 predictions despite being “imperceptible” in that they are semantically indifferent from their factual counterparts.
 28 Even though some partially effective mitigation strategies have been proposed—most notably **adversarial training**
 29 (Goodfellow, Shlens, and Szegedy 2015)—truly robust deep learning remains unattainable even for models that are
 30 considered “shallow” by today’s standards (Kolter 2023).

31 Part of the problem is that the high degrees of freedom provide room for many solutions that are locally optimal with
 32 respect to narrow objectives (Wilson 2020).¹ Indeed, recent work on the so-called “lottery ticket hypothesis” suggests
 33 that modern neural networks can be pruned by up to 90% while preserving their predictive performance (Frankle
 34 and Carbin 2019). Similarly, Zhang et al. (2021) showed that state-of-the-art neural networks are expressive enough
 35 to fit randomly labeled data. Thus, looking at the predictive performance alone, the solutions may seem to provide
 36 compelling explanations for the data, when in fact they are based on purely associative, semantically meaningless
 37 patterns. This poses two challenges. Firstly, there is no dependable way to verify if representations correspond to
 38 meaningful, plausible explanations. Secondly, even if we could resolve the first challenge, it remains undecided how
 39 to ensure that models can *only* learn valuable explanations.

40 The first challenge has attracted an abundance of research on **explainable AI** (XAI), a paradigm that focuses on
 41 the development of tools to derive (post-hoc) explanations from complex model representations. Such explanations
 42 should mitigate a scenario in which practitioners deploy opaque models and blindly rely on their predictions. On
 43 countless occasions, this has happened in practice and caused real harms to people who were adversely and unfairly
 44 affected by automated decision-making (ADM) systems involving opaque models (see, e.g., (O’Neil 2016)). Effective
 45 XAI tools can aid us in monitoring models and providing recourse to individuals to turn negative outcomes (e.g.,
 46 “loan application rejected”) into positive ones (e.g., “application accepted”). In line with this, our work builds upon
 47 **counterfactual explanations** (CE) proposed by Wachter, Mittelstadt, and Russell (2017) as an effective approach to
 48 achieve this goal. CEs prescribe minimal changes for factual inputs that, if implemented, would prompt some fitted
 49 model to produce a desired output.

50 To our surprise, the second challenge has not yet attracted major research interest. Specifically, there has been no
 51 concerted effort towards improving the “explanatory capacity” of models, i.e., the degree to which learned represen-
 52 tations correspond to explanations that are **interpretable** and deemed **plausible** by humans (see Def. 3.1). Instead,
 53 the choice has generally been to improve the ability of XAI tools to identify the subset of explanations that are both
 54 plausible and valid for any given model, independent of whether the learned representations are in fact compatible
 55 with plausible explanations (Altmeyer et al. 2024). Fortunately, recent findings indicate that improved explanatory
 56 capacity can arise as a consequence of regularization techniques aimed at other training objectives such as generative
 57 capacity, generalization, or robustness (Altmeyer et al. 2024; Augustin, Meinke, and Hein 2020; Schut et al. 2021).
 58 As further discussed in Section 2, our work consolidates these findings within a single objective.

59 **Specifically, we introduce Counterfactual Training (CT):** a novel training regime explicitly geared towards improv-
 60 ing the explainability of models. In high-level terms, we define this concept as as the extent to which valid explanations
 61 derived for an opaque model are also deemed plausible with respect to the underlying data and the global actionability
 62 constraints. To the best of our knowledge, our framework represents the first attempt to address this challenge by
 63 employing counterfactual explanations already in the training phase.

64 The remainder of this manuscript is structured as follows. Section 2 presents related work, focusing on the link between
 65 AEs and CEs. Then follow our two principal contributions. In Section 3, we introduce our methodological framework
 66 and show theoretically that it can be employed to enforce global actionability constraints. In Section 4, through
 67 extensive experiments, we demonstrate that CT substantially improves explainability and positively contributes to
 68 the robustness of trained models without sacrificing predictive performance. Finally, in Section 5, we discuss the
 69 challenges and conclude that CT is a promising approach towards making opaque models more trustworthy.

¹We follow the standard ML convention, where “degrees of freedom” refer to the number of parameters estimated from data.

70 2 Related Literature

71 To make the desiderata for our framework more concrete, we follow Augustin, Meinke, and Hein (2020) in tying the
 72 concept of explainability to the quality of CEs that can be generated for a given model. The authors show that CEs—
 73 understood as minimal input perturbations that yield some desired model prediction—tend to be more meaningful if the
 74 underlying model is more robust to adversarial examples. We can make intuitive sense of this finding when looking
 75 at adversarial training (AT) through the lens of representation learning with high degrees of freedom. As argued
 76 before, learned representations may be sensitive to producing implausible explanations and mispredicting for worst-
 77 case counterfactuals (i.e., AEs). Thus, by inducing models to “unlearn” susceptibility to such examples, adversarial
 78 training can effectively remove implausible explanations from the solution space.

79 2.1 Adversarial Examples are Counterfactual Explanations

80 This interpretation of the link between explainability through counterfactuals on one side and robustness to adversarial
 81 examples on the other is backed by empirical evidence. Sauer and Geiger (2021) demonstrates that using counter-
 82 factual images during classifier training improves model robustness. Similarly, Abbasnejad et al. (2020) argues that
 83 counterfactuals represent potentially useful training data in machine learning, especially in supervised settings where
 84 inputs may be reasonably mapped to multiple outputs. They, too, demonstrate that augmenting the training data of
 85 image classifiers can improve generalization. Finally, Teney, Abbasnejad, and Hengel (2020) proposes an approach
 86 using counterfactuals in training that does not rely on data augmentation: they argue that counterfactual pairs typically
 87 already exist in training datasets. Specifically, their approach relies on identifying similar input samples with different
 88 annotations and ensuring that the gradient of the classifier aligns with the vector between such pairs of counterfactual
 89 inputs using the cosine distance as the loss function.

90 In the natural language processing (NLP) domain, CEs have also been used to improve models through data augmen-
 91 tation. Wu et al. (2021) proposes *Polyjuice*, a general-purpose counterfactual generator for language models. The
 92 authors empirically demonstrate that the augmentation of training data with their method improves robustness in a
 93 number of NLP tasks. Balashankar et al. (2023) similarly uses *Polyjuice* to augment NLP datasets through diverse
 94 counterfactuals and show that classifier robustness improves by up to 20%. Finally, Luu and Inoue (2023) introduces
 95 Counterfactual Adversarial Training (CAT) that also aims to improve generalization and robustness of language mod-
 96 els through a three-step procedure: the authors identify training samples that are subject to high predictive uncertainty,
 97 generate CEs for them, and fine-tune the language model on a dataset augmented with the CEs.

98 There have also been several attempts at formalizing the relationship between counterfactual explanations and adver-
 99 sarial examples. Pointing to clear similarities in how CEs and AEs are generated, Freiesleben (2022) makes the case
 100 for jointly studying the opaqueness and robustness problems in representation learning. Formally, AEs can be seen as
 101 the subset of CEs for which misclassification is achieved (Freiesleben 2022). Similarly, Pawelczyk et al. (2022) shows
 102 that CEs and AEs are equivalent under certain conditions.

103 Two recent works are closely related to ours in that they use counterfactuals during training with the explicit goal of
 104 affecting certain properties of the post-hoc counterfactual explanations. Firstly, Ross, Lakkaraju, and Bastani (2024)
 105 proposes a way to train models that guarantee individual recourse to some positive target class with high probability.
 106 Their approach builds on adversarial training by explicitly inducing susceptibility to targeted adversarial examples for
 107 the positive class. Additionally, the proposed method allows for imposing a set of actionability constraints ex-ante.
 108 For example, users can specify that certain features are immutable. Secondly, Guo, Nguyen, and Yadav (2023) is the
 109 first to propose an end-to-end training pipeline that includes CEs as part of the training procedure. In particular, they
 110 propose a specific network architecture that includes a predictor and CE generator network, where the parameters of
 111 the CE generator network are learnable. Counterfactuals are generated during each training iteration and fed back
 112 to the predictor network. In contrast to Guo, Nguyen, and Yadav (2023), we impose no restrictions on the ANN
 113 architecture at all.

114 2.2 Aligning Representations with Plausible Explanations

115 Improving the adversarial robustness of models is not the only path towards aligning representations with plausible
 116 explanations. In a work closely related to this one, Altmeyer et al. (2024) shows that explainability can be improved
 117 through model averaging and refined model objectives. The authors propose a way to generate counterfactuals that
 118 are maximally faithful to the model in that they are consistent with what the model has learned about the underlying
 119 data. Formally, they rely on tools from energy-based modelling (Teh et al. 2003) to minimize the divergence between
 120 the distribution of counterfactuals and the conditional posterior over inputs learned by the model. Their proposed
 121 counterfactual explainer, *ECCCo*, yields plausible explanations if and only if the underlying model has learned repre-
 122 sentations that align with them. The authors find that both deep ensembles (Lakshminarayanan, Pritzel, and Blundell
 123 2017) and joint energy-based models (JEMs) (Grathwohl et al. 2020) tend to do well in this regard.

Once again it helps to look at these findings through the lens of representation learning with high degrees of freedom. Deep ensembles are approximate Bayesian model averages, which are most called for when models are underspecified by the available data (Wilson 2020). Averaging across solutions mitigates the aforementioned risk of relying on a single locally optimal representations that corresponds to semantically meaningless explanations for the data. Previous work of Schut et al. (2021) similarly found that generating plausible (“interpretable”) CEs is almost trivial for deep ensembles that have also undergone adversarial training. The case for JEMs is even clearer: they involve a hybrid objective that induces both high predictive performance and generative capacity (Grathwohl et al. 2020). This is closely related to the idea of aligning models with plausible explanations and has inspired our CT objective.

3 Counterfactual Training

In this section we propose our novel counterfactual training objective. In CT, we combine ideas from adversarial training, counterfactual explanations, and energy-based modelling with the explicit goal of aligning representations with plausible explanations that comply with user-defined actionability constraints.

In the context of counterfactual explanations, plausibility has broadly been defined as the degree to which counterfactuals comply with the underlying data-generating process (Altmeyer et al. 2024; Guidotti 2022; Poyiadzi et al. 2020). Plausibility is a necessary but insufficient condition for using CEs to provide algorithmic recourse (AR) to individuals (negatively) affected by opaque models. To be actionable, AR recommendations must also be attainable. A plausible CE for a rejected 20-year-old loan applicant, for example, might reveal that their application would have been accepted, if only they were 20 years older. Ignoring all other features, this would comply with the definition of plausibility if 40-year-old individuals were in fact more credit-worthy on average than young adults. But of course this CE does not qualify for providing actionable recourse to the applicant since *age* is not a (directly) mutable feature. Counterfactual training aims to improve model explainability by aligning models with counterfactuals that meet both desiderata: plausibility and actionability. Formally, we define explainability as follows:

Definition 3.1 (Model Explainability). Let $M_\theta : \mathcal{X} \mapsto \mathcal{Y}$ denote a supervised classification model that maps from the D -dimensional input space \mathcal{X} to representations $\phi(\mathbf{x}; \theta)$ and finally to the K -dimensional output space \mathcal{Y} . Assume that for any given input-output pair $\{\mathbf{x}, \mathbf{y}\}_i$ there exists a counterfactual $\mathbf{x}' = \mathbf{x} + \Delta : M_\theta(\mathbf{x}') = \mathbf{y}^+ \neq \mathbf{y} = M_\theta(\mathbf{x})$ where $\arg \max_y \mathbf{y}^+ = y^+$ and y^+ denotes the index of the target class.

We say that M_θ is **explainable** to the extent that faithfully generated counterfactuals are plausible and actionable. We define these properties as follows:

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

where $p_\theta(\mathbf{x} | \mathbf{y}^+)$ denotes the conditional posterior over inputs.

The characterization of faithfulness and plausibility in Def. 3.1 is the same as in Altmeyer et al. (2024), with adapted notation. Intuitively, plausible counterfactuals are consistent with the data and faithful counterfactuals are consistent with what the model has learned about input data. Actionability constraints in Def. 3.1 vary and depend on the context in which M_θ is deployed. In this work, we focus on domain and mutability constraints for individual features x_d for $d = 1, \dots, D$. We limit ourselves to classification tasks for reasons discussed in Section 5.

3.1 Our Proposed Objective

Let \mathbf{x}'_t for $t = 0, \dots, T$ denote a counterfactual explanation generated through gradient descent over T iterations as initially proposed by Wachter, Mittelstadt, and Russell (2017). For our purposes, we let T vary and consider the counterfactual search as converged as soon as the predicted probability for the target class has reached a pre-determined threshold, $\tau : \mathcal{S}(M_\theta(\mathbf{x}'))[y^+] \geq \tau$, where \mathcal{S} is the softmax function.²

To train models with high explainability as defined in Def. 3.1, we propose to leverage counterfactuals in the following objective:

$$\begin{aligned} \min_{\theta} & \text{yloss}(M_\theta(\mathbf{x}), \mathbf{y}) + \lambda_{\text{div}} \text{div}(\mathbf{x}^+, \mathbf{x}'_T, y; \theta) + \lambda_{\text{adv}} \text{advloss}(M_\theta(\mathbf{x}'_{t \leq T}), \mathbf{y}) \\ & + \lambda_{\text{reg}} \text{ridge}(\mathbf{x}^+, \mathbf{x}'_T, y; \theta) \end{aligned} \quad (1)$$

²For detailed background information on gradient-based counterfactual search and convergence see supplementary appendix.

168 where $y\text{loss}(\cdot)$ is a classification loss that induces discriminative performance (e.g., cross-entropy). The second and
 169 third terms are explained in detail below. For now, they can be summarized as inducing explainability directly and
 170 indirectly by penalizing: (1) the contrastive divergence, $\text{div}(\cdot)$, between mature counterfactuals \mathbf{x}'_T and observed
 171 samples $\mathbf{x}^+ \in \mathcal{X}^+ = \{\mathbf{x} : y = y^+\}$ in the target class y^+ , and, (2) the adversarial loss, $\text{advloss}(\cdot)$, with respect to
 172 nascent counterfactuals $\mathbf{x}'_{t \leq T}$. Finally, $\text{ridge}(\cdot)$ denotes a Ridge penalty (ℓ_2 -norm) that regularizes the magnitude of
 173 the energy terms involved in $\text{div}(\cdot)$ (Du and Mordatch 2020). The trade-off between the components can be governed
 174 through penalties λ_{div} , λ_{adv} and λ_{reg} .

175 3.2 Directly Inducing Explainability with Contrastive Divergence

176 As observed by Grathwohl et al. (2020), any classifier can be re-interpreted as a joint energy-based model (JEM)
 177 that learns to discriminate output classes conditional on the observed (training) samples from $p(\mathbf{x})$ and the generated
 178 samples from $p_\theta(\mathbf{x})$. The authors show that JEMs can be trained to perform well at both tasks by directly maximizing
 179 the joint log-likelihood factorized as $\log p_\theta(\mathbf{x}, \mathbf{y}) = \log p_\theta(\mathbf{y}|\mathbf{x}) + \log p_\theta(\mathbf{x})$. The first term can be optimized using
 180 conventional cross-entropy as in Equation 1. To optimize $\log p_\theta(\mathbf{x})$, Grathwohl et al. (2020) minimizes the contrastive
 181 divergence between the observed samples from $p(\mathbf{x})$ and generated samples from $p_\theta(\mathbf{x})$.

182 A key empirical finding in Altmeyer et al. (2024) was that JEMs tend to do well with respect to the plausibility
 183 objective in Def. 3.1. This follows directly if we consider samples drawn from $p_\theta(\mathbf{x})$ as counterfactuals because
 184 the JEM objective effectively minimizes the divergence between the conditional posterior and $p(\mathbf{x}|y^+)$. To generate
 185 samples, Grathwohl et al. (2020) relies on Stochastic Gradient Langevin Dynamics (SGLD) using an uninformative
 186 prior for initialization but we depart from their methodology. Instead of SGLD, we propose to use counterfactual
 187 explainers to generate counterfactuals of observed training samples. Specifically, we have:

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

188 where $\mathcal{E}_\theta(\cdot)$ denotes the energy function defined as $\mathcal{E}_\theta(\mathbf{x}, y) = -\mathbf{M}_\theta(\mathbf{x})[y^+]$, with y^+ denoting the index of the
 189 randomly drawn target class, $y^+ \sim p(y)$. Conditional on the target class y^+ , \mathbf{x}'_T denotes a mature counterfactual for a
 190 randomly sampled factual from a non-target class generated with a gradient-based CE generator for up to T iterations.
 191 Mature counterfactuals are ones that have either reached convergence wrt. the decision threshold τ or exhausted T .

192 Intuitively, the gradient of Equation 2 decreases the energy of observed training samples (positive samples) while
 193 increasing the energy of counterfactuals (negative samples) (Du and Mordatch 2020). As the counterfactuals get more
 194 plausible (Def. 3.1) during training, these opposing effects gradually balance each other out (Lippe 2024).

195 The departure from SGLD of (Grathwohl et al. 2020) allows us to tap into the vast repertoire of explainers that have
 196 been proposed in the literature to meet different desiderata. For example, many methods facilitate the imposition of
 197 domain and mutability constraints. In principle, any existing approach for generating counterfactual explanations is
 198 viable, so long as it does not violate the faithfulness condition. Like JEMs (Murphy 2022), CT can be considered a
 199 form of contrastive representation learning.

200 3.3 Indirectly Inducing Explainability with Adversarial Robustness

201 Based on our analysis in Section 2, counterfactuals \mathbf{x}' can be repurposed as additional training samples (Balashankar
 202 et al. 2023; Luu and Inoue 2023) or adversarial examples (Freiesleben 2022; Pawelczyk et al. 2022). This leaves some
 203 flexibility wrt. the choice for $\text{advloss}(\cdot)$ in Equation 1. An intuitive functional form, but likely not the only sensible
 204 choice, is inspired by adversarial training:

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

205 Under this choice, we consider nascent counterfactuals $\mathbf{x}'_{t \leq T}$ as AEs as long as the magnitude of the perturbation to
 206 any single feature is at most ε . This is closely aligned with Szegedy et al. (2014) that defines an adversarial attack as
 207 an “imperceptible non-random perturbation”. Thus, we choose to work with a different distinction between CE and
 208 AE than Freiesleben (2022) that considers misclassification as the key distinguishing feature of AE. One of the key
 209 observations of this work is that we can leverage CEs during training and get adversarial examples essentially for free,
 210 which can be used to reap the aforementioned benefits of adversarial training.

211 3.4 Encoding Actionability Constraints

212 Many existing counterfactual explainers support domain and mutability constraints out-of-the-box. In fact, both types
 213 of constraints can be implemented for any counterfactual explainer that relies on gradient descent in the feature space
 214 for optimization (Altmeyer, Deursen, and Liem 2023). In this context, domain constraints can be imposed by simply

215 projecting counterfactuals back to the specified domain, if the previous gradient step resulted in updated feature values
 216 that were out-of-domain. Mutability constraints can similarly be enforced by setting partial derivatives to zero to
 217 ensure that features are only perturbed in the allowed direction, if at all.

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

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

226 From a Bayesian perspective, setting $\mathbf{x}^+[d] := \mathbf{x}'[d]$ can be understood as assuming a point mass prior for $p(\mathbf{x}^+)$
 227 with respect to feature d . Intuitively, we think of this simply in terms ignoring implausibility costs with respect
 228 to immutable features, which effectively forces the model to instead seek plausibility with respect to the remaining
 229 features. This in turn results in lower overall sensitivity to immutable features, which we demonstrate empirically for
 230 different classifiers in Section 4. Under certain conditions, this result holds theoretically:³

231 **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
 232 softmax activation \mathcal{S} where $y \in \{1, \dots, K\} = \mathcal{K}$ and $\mathbf{x} \in \mathbb{R}^D$. If we assume multivariate Gaussian class densities with
 233 common diagonal covariance matrix $\Sigma_k = \Sigma$ for all $k \in \mathcal{K}$, then protecting an immutable feature from the contrastive
 234 divergence penalty will result in lower classifier sensitivity to that feature relative to the remaining features, provided
 235 that at least one of those is discriminative and mutable.*

236 It is worth highlighting that Proposition 3.1 assumes independence of features. This raises a valid concern about
 237 the effect of protecting immutable features in the presence of proxies that remain unprotected. We address this in
 238 Section 5.

239 3.5 Example (Prediction of Consumer Credit Default)

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

244 We have simulated this scenario using synthetic data with independent *age* and *debt* features, and Gaussian class-
 245 conditional densities in Figure 1. The four panels show the outcomes for different training procedures using the same
 246 model architectures (a linear classifier). In panels (a) and (c) we have trained the models conventionally, while in
 247 panels (b) and (d) we used CT.

248 In all cases, all counterfactuals (stars) are valid—they have crossed the decision boundary (green)—but their quality
 249 differs. In panel (a), they are not plausible: they do not comply with the distribution of the factuals in y^+ to the point
 250 where they form a clearly distinguishable cluster. In panel (b), they are highly plausible, meeting the first objective
 251 of Def. 3.1. In panel (c), the CEs involve substantial reductions in *debt* for younger applicants. By comparison,
 252 counterfactual paths are shorter on average in panel (d) where we have protected the immutable *age* as described in
 253 Section 3.4. Due to the classifier's lower sensitivity to *age*, recommendations with respect to *debt* are much more
 254 homogenous and do not unfairly punish younger individuals. These counterfactuals are also plausible with respect to
 255 the mutable feature. Thus, we consider the model in panel (d) as the most explainable according to Def. 3.1.

256 4 Experiments

257 In our experiments we seek to answer the following three research questions:

- 258 1. To what extent does the counterfactual training objective as it is defined in Equation 1 induce models to learn
 259 plausible explanations?
 260 2. To what extent does the CT objective produce more favorable algorithmic recourse outcomes in the presence
 261 of actionability constraints?
 262 3. What are the effects of hyperparameter selection wrt. the CT objective?

³For the proof, see the supplementary appendix.

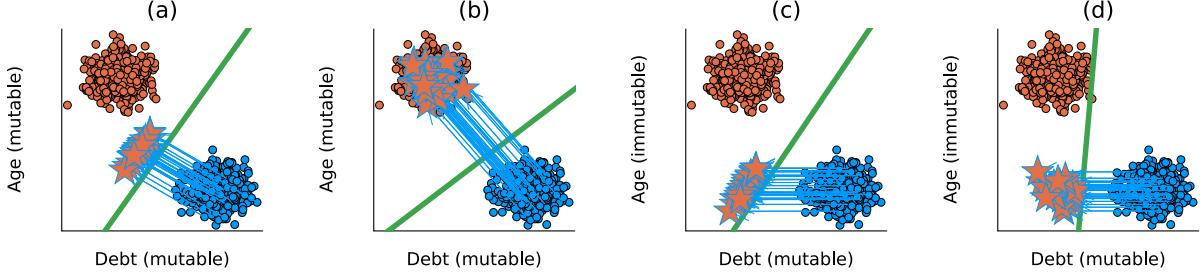


Figure 1: Illustration of how CT improves model explainability: (a) conventional training, all mutable; (b) CT, all mutable; (c) conventional, *age* immutable; (d) CT, *age* immutable. The linear decision boundary is shown in green along with training data colored according to their ground-truth label: $y^- = 1$ in blue and $y^+ = 2$ in orange. Stars indicate counterfactuals in the target class.

263 4.1 Experimental Setup

264 Our key outcome of interest is improvement in explainability (Def. 3.1). To this end, we focus primarily on the
 265 plausibility and cost of faithfully generated counterfactuals at test time. To measure the cost, we follow the standard
 266 convention of using distances (ℓ_1 -norm) between factuals and counterfactuals as a proxy. For plausibility, we assess
 267 how similar CEs are to the observed samples in the target domain, $\mathbf{X}' \subset \mathcal{X}^+$. We rely on the distance-based metric
 268 used in Altmeyer et al. (2024),

$$\text{IP}(\mathbf{x}', \mathbf{X}') = \frac{1}{|\mathbf{X}'|} \sum_{\mathbf{x} \in \mathbf{X}'} \text{dist}(\mathbf{x}', \mathbf{x}) \quad (4)$$

269 and introduce a novel divergence metric,

$$\text{IP}^*(\mathbf{X}', \mathbf{X}') = \text{MMD}(\mathbf{X}', \mathbf{X}') \quad (5)$$

270 where \mathbf{X}' denotes a collection of counterfactuals and $\text{MMD}(\cdot)$ is an unbiased estimate of the squared population
 271 maximum mean discrepancy (Gretton et al. 2012). The metric in Equation 5 is equal to zero iff the two distributions
 272 are the same, $\mathbf{X}' = \mathbf{X}^+$.

273 In addition to cost and plausibility, we compute other standard metrics to evaluate counterfactuals including validity
 274 and redundancy. Finally, we also assess the predictive performance of models using standard metrics, including robust
 275 accuracy estimated on adversarially perturbed data using FGSM (Goodfellow, Shlens, and Szegedy 2015).

276 We run the experiments with three gradient-based generators: *Generic* of Wachter, Mittelstadt, and Russell (2017) as
 277 a simple baseline approach, *REVISE* (Joshi et al. 2019) that aims to generate plausible counterfactuals using a surro-
 278 gate Variational Autoencoder (VAE), and *ECCo*—the generator of Altmeyer et al. (2024) but without the conformal
 279 prediction component—as a method that directly targets both faithfulness and plausibility of the counterfactuals.

280 We make use of nine classification datasets common in the CE/AR literature. Four of them are synthetic with two
 281 classes and different characteristics: linearly separable clusters (*LS*), overlapping clusters (*OL*), concentric circles
 282 (*Circ*), and interlocking moons (*Moon*). These datasets are generated using the library of (Altmeyer, Deursen, and
 283 Liem 2023) and we present them in the supplementary appendix. Next, we have four real-world binary tabular datasets
 284 from the domain of economics: *Adult* (a.k.a. Census data) of (Becker and Kohavi 1996), California housing (*CH*) of
 285 (Pace and Barry 1997), Default of Credit Card Clients (*Cred*) of (Yeh 2016), and Give Me Some Credit (*GMSC*) of
 286 (Kaggle 2011). Finally, for the convenience of illustration, we use of the 10-class *MNIST* vision dataset (LeCun 1998).

287 To assess CT, we investigate the improvements in performance metrics when using it on top of a weak baseline (BL):
 288 a multilayer perceptron (*MLP*). This is the best way to get a clear picture of the effectiveness of CT, and it is consistent
 289 with how assessment is done in the related literature (Goodfellow, Shlens, and Szegedy 2015; Ross, Lakkaraju, and
 290 Bastani 2024; Teney, Abbasnejad, and Hengel 2020).

291 4.2 Experimental Results

292 4.2.1 Plausibility

293 Table 1 presents our main empirical findings. For all datasets except *OL* and across all test settings, the average
 294 distance of CEs from observed samples in the target class is reduced, indicating improved plausibility. The magnitude
 295 of improvements varies. For the simple synthetic datasets, distance reductions range from around 20-40% (*LS*, *Moon*)

Table 1: Key performance metrics across all datasets (column 1). **Plausibility**: Columns 2-6 show the percentage reduction in implausibility (IP) for varying degrees of the energy penalty λ_{egy} used for *ECCo* at test time; column 7 shows the reduction in IP* (MMD), aggregated across all test specifications. **Accuracy** (columns 8-11): test accuracies and robust accuracies (Acc.*[†]) for CT and the baseline (BL). **Actionability** (column 12): average reduction in costs when imposing mutability constraints reported for the four datasets for which we could identify meaningful features to protect.

Data	IP (-%)	IP (-%)	IP (-%)	IP (-%)	IP (-%)	IP* (-%) (agg.)	Acc. (CT)	Acc. (BL)	Acc.* (CT)	Acc.* (BL)	Cost (-%)
λ_{egy}	0.1	0.5	1.0	5.0	10.0						
Adult	2.9	3.4	3.5	2.9	3.2	34.8	0.85	0.85	0.83	0.41	
CH	9.6	9.3	10.4	11.9	14.6	66.6	0.79	0.85	0.76	0.75	
Circ	56.5	57.1	56.5	58.5	49.3	93.4	1.0	1.0	0.99	1.0	35.0
Cred	6.7	6.2	7.2	7.0	7.8	51.6	0.71	0.71	0.7	0.52	
GMSC	11.0	13.4	13.4	21.4	27.9	77.9	0.61	0.75	0.58	0.42	
LS	27.1	26.7	26.6	27.1	38.6	54.5	1.0	1.0	1.0	1.0	26.3
MNIST	9.1	8.3	8.1	6.1	3.5	-2.3	0.9	0.92	0.84	0.78	
Moon	20.4	21.4	21.6	19.0	19.8	27.6	1.0	1.0	1.0	1.0	23.4
OL	-6.7	-6.2	-6.1	-2.8	-1.4	-25.5	0.92	0.91	0.91	0.91	15.5

296 to almost 60% (*Circ*). For the real-world tabular datasets, improvements tend to be smaller but still substantial, with
 297 around 10-15% for *CH*, 11-28% for *GMSC*, 7-8% for *Cred*, and around 3% for *Adult*. For the vision dataset (*MNIST*),
 298 distances are reduced by up to 9%. The results for our proposed divergence metric are qualitatively similar, but
 299 generally even more pronounced: for the *Circ* dataset, implausibility is reduced by almost 94% to virtually zero as
 300 we verified by the absolute outcome. Improvements for other datasets range from 28% (*Moon*) up to 78% (*GMSC*).
 301 For *OL* the reduction is negative, consistent with the distance-based metric. *MNIST* is the only dataset for which the
 302 distance and divergence metrics disagree. Upon visual inspection of the image counterfactuals we find that CT clearly
 303 improves plausibility (see supplementary appendix for images).

304 4.2.2 Predictive Performance

305 Test accuracy for CT is virtually identical to the baseline for *Adult*, *Circ*, *LS*, *Moon*, and *OL*, and even slightly improved
 306 for *Cred*. Exceptions to this general pattern are *MNIST*, *CH*, and *GMSC*, for which we observe a reduction in test
 307 accuracy of 2, 5, and 15 percentage points respectively. When looking at robust test accuracies (Acc.*[†]) for these
 308 datasets in particular, we find that CT strongly outperforms the baseline. In fact, we find that CT improves adversarial
 309 robustness on all datasets.

310 4.2.3 Actionability

311 In Section 3, we show that our proposed way for encoding mutability constraints leads to lower classifier sensitivity
 312 wrt. immutable features for linear models, tilting the decision boundary in favour of mutable features instead. For
 313 binding constraints at test time, this leads to shorter counterfactual paths and hence smaller average costs (ℓ_1 -norm) to
 314 individuals. To extend this to the non-linear case, we test the effect of imposing mutability constraints empirically for
 315 our synthetic data using the same evaluation scheme as above. The final row in Table 1 reports the average reduction in
 316 costs for CT compared to the “vanilla” baseline, when imposing that either the first or the second feature is immutable.
 317 In all cases, costs are reduced substantially, indicating that classifiers trained with CT are indeed more sensitive to
 318 mutable features.

319 4.2.4 Impact of hyperparameter settings.

320 We test the impact of three types of hyperparameters; our complete results are in the supplementary appendix.
 321 We note that CT is highly sensitive to the choice of a CE generator and its hyperparameters but (a) there are manageable
 322 patterns and (b) we can typically identify settings that improve either plausibility or cost, and commonly both of them
 323 at the same time. For example, *REVISE* tends to perform the worst, most likely because it uses a surrogate VAE to
 324 generate counterfactuals which impedes faithfulness (Altmeyer et al. 2024). Increasing T , the maximum number of
 325 steps, generally yields better outcomes because more CEs can mature in each training epoch. The impact of τ , the
 326 required decision threshold is more difficult to predict. On “harder” datasets it may be difficult to satisfy high τ for any
 327 given sample (i.e., also factuals) and so increasing this threshold does not seem to correlate with better outcomes. In
 328 fact, the choice of $\tau = 0.5$ generally leads to optimal results because it is associated with high proportions of mature
 329 counterfactuals.

330 The strength of the energy regularization, λ_{reg} is highly impactful and leads to poor performance in terms of decreased

331 plausibility and increased costs if insufficiently high. The sensitivity with respect to λ_{div} and λ_{adv} is much less evident.
 332 While high values of λ_{reg} may increase the variability in outcomes when combined with high values of λ_{div} or λ_{adv} ,
 333 this effect is not very pronounced.

334 The effectiveness and stability of CT is positively associated with the number of counterfactuals generated during
 335 each training epoch. We also confirm that a higher number of training epochs is beneficial. Interestingly, we observed
 336 desired improvements when CT was combined with conventional training and applied only for the final 50% of epochs
 337 of the complete training process. Put differently, CT can improve the explainability of models in a fine-tuning manner.

338 5 Conclusions

339 As our results indicate, counterfactual training produces models that are more explainable. Nonetheless, it brings
 340 about three important limitations.

341 *CT increases the training time of models.* CT can be more time-consuming than conventional training regimes. While
 342 higher numbers of CEs per iteration positively impact the quality of solutions, they also increase the amount of
 343 computations. Relatively small grids with 270 settings can take almost four hours for more demanding datasets on
 344 a high-performance computing cluster with 34 2GB CPUs.⁴ Three factors attenuate this effect: (1) CT amortizes
 345 the cost of CEs for the training samples; (2) it can retain its value when used as a “fine-tuning” technique for
 346 conventionally-trained models; and (3) it yields itself to parallel execution, which we have leveraged for our own
 347 experiments.

348 *Immutable features may have proxies.* We propose an approach to protect immutable features and thus increase the
 349 actionability of the generated CEs. However, it requires that model owners define the mutability constraints for
 350 (all) features considered by the model. Even if all immutable features are protected, there may exist proxies that
 351 are mutable (and hence should not be protected) but preserve enough information about the principals to hinder the
 352 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.

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

361
 362 We also highlight several important directions for future research. Firstly, it is an interesting challenge to ex-
 363 tend CT beyond classification settings. Our formulation relies on the distinction between non-target class(es) y^- and
 364 target class(es) y^+ to generate counterfactuals through Equation 1. While y^- and y^+ can be arbitrarily defined, CT
 365 requires the output space \mathcal{Y} to be discrete. Thus, it does not apply to ML tasks where the change in outcome cannot be
 366 readily quantified. Focus on classification models is a common restriction in research on CEs and AR. Other settings
 367 have attracted some interest (e.g., regression in ([Spooner et al. 2021](#))), but there is little consensus how to robustly
 368 extend the notion of CEs.

369 Secondly, our approach is susceptible to training instabilities. This problem has been recognized for JEMs ([Grathwohl et al. 2020](#)) and even though we depart from the SGLD-based sampling, we still encounter considerable variability
 370 in the outcomes. CT is exposed to two potential sources of instabilities: (1) the energy-based contrastive divergence
 371 term in Equation 2, and (2) the underlying counterfactual explainers. We find several promising ways to mitigate this
 372 problem: regularizing energy (λ_{reg}), generating sufficiently many counterfactuals during each epoch, and including
 373 only mature counterfactuals for contrastive divergence.

375 Finally, we believe that it is possible to substantially improve hyperparameter selection procedures. Our method
 376 benefits from the tuning of certain key hyperparameters (see Section 4.2.4). In this work, we have relied exclusively
 377 on grid search for this task. Future work on CT could benefit from investigating more sophisticated approaches.
 378 Notably, CT is iterative which makes methods such as Bayesian or gradient-based optimization applicable (see, e.g.,
 379 ([Bischl et al. 2023](#))).

380

381 To conclude, state-of-the-art machine learning models are prone to learning complex representations that cannot be
 382 interpreted by humans and existing post-hoc explainability approaches cannot guarantee that the explanations agree
 383 with the model’s learned representation of data. As a step towards addressing this challenge, we introduced coun-
 384 terfactual training, a novel training regime that incentivizes highly-explainable models. Our approach leads to expla-
 385 nations that are both plausible—compliant with the underlying data-generating process—and actionable—compliant

⁴See supplementary appendix for computational details.

386 with user-specified mutability constraints—and thus meaningful to their recipients. Through extensive experiments
 387 we demonstrate that CT satisfies its objective while promoting robustness and preserving the predictive performance
 388 of the models. It can also be used to fine-tune conventionally-trained models and achieve similar gains. Finally, this
 389 work showcases that it is practical to improve models *and* their explanations at the same time.

390 References

- 391 Abbasnejad, Ehsan, Damien Teney, Amin Parvaneh, Javen Shi, and Anton van den Hengel. 2020. “Counterfactual
 392 Vision and Language Learning.” In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition
 393 (CVPR)*, 10041–51. <https://doi.org/10.1109/CVPR42600.2020.01006>.
- 394 Altmeyer, Patrick, Arie van Deursen, and Cynthia C. S. Liem. 2023. “Explaining Black-Box Models through Coun-
 395 terfactuals.” In *Proceedings of the JuliaCon Conferences*, 1:130.
- 396 Altmeyer, Patrick, Mojtaba Farmanbar, Arie van Deursen, and Cynthia C. S. Liem. 2024. “Faithful Model Ex-
 397 planations through Energy-Constrained Conformal Counterfactuals.” In *Proceedings of the Thirty-Eighth AAAI
 398 Conference on Artificial Intelligence*, 38:10829–37. 10. <https://doi.org/10.1609/aaai.v38i10.28956>.
- 399 Augustin, Maximilian, Alexander Meinke, and Matthias Hein. 2020. “Adversarial Robustness on In- and Out-
 400 Distribution Improves Explainability.” In *Computer Vision – ECCV 2020*, edited by Andrea Vedaldi, Horst Bischof,
 401 Thomas Brox, and Jan-Michael Frahm, 228–45. Cham: Springer.
- 402 Balashankar, Ananth, Xuezhi Wang, Yao Qin, Ben Packer, Nithum Thain, Ed Chi, Jilin Chen, and Alex Beutel. 2023.
 403 “Improving Classifier Robustness through Active Generative Counterfactual Data Augmentation.” In *Findings of
 404 the Association for Computational Linguistics: EMNLP 2023*, 127–39. ACL. <https://doi.org/10.18653/v1/2023.f>
 405 indings-emnlp.10.
- 406 Becker, Barry, and Ronny Kohavi. 1996. “Adult.” UCI Machine Learning Repository.
- 407 Bell, Andrew, Joao Fonseca, Carlo Abate, Francesco Bonchi, and Julia Stoyanovich. 2024. “Fairness in Algorithmic
 408 Recourse Through the Lens of Substantive Equality of Opportunity.” <https://arxiv.org/abs/2401.16088>.
- 409 Bezanson, Jeff, Alan Edelman, Stefan Karpinski, and Viral B Shah. 2017. “Julia: A Fresh Approach to Numerical
 410 Computing.” *SIAM Review* 59 (1): 65–98. <https://doi.org/10.1137/141000671>.
- 411 Bischl, Bernd, Martin Binder, Michel Lang, Tobias Pielok, Jakob Richter, Stefan Coors, Janek Thomas, et al. 2023.
 412 “Hyperparameter optimization: Foundations, algorithms, best practices, and open challenges.” *WIREs Data Min-
 413 ing and Knowledge Discovery* 13 (2): e1484. <https://doi.org/10.1002/widm.1484>.
- 414 Bouchet-Valat, Milan, and Bogumi Kamiski. 2023. “DataFrames.jl: Flexible and Fast Tabular Data in Julia.” *Journal
 415 of Statistical Software* 107 (4): 1–32. <https://doi.org/10.18637/jss.v107.i04>.
- 416 Byrne, Simon, Lucas C. Wilcox, and Valentin Churavy. 2021. “MPI.jl: Julia Bindings for the Message Passing
 417 Interface.” *Proceedings of the JuliaCon Conferences* 1 (1): 68. <https://doi.org/10.21105/jcon.00068>.
- 418 Chagas, Ronan Arraes Jardim, Ben Baumgold, Glen Hertz, Hendrik Ranocha, Mark Wells, Nathan Boyer, Nicholas
 419 Ritchie, et al. 2024. “Ronisbr/PrettyTables.jl: V2.4.0.” Zenodo. <https://doi.org/10.5281/zenodo.1383553>.
- 420 Christ, Simon, Daniel Schwabeneder, Christopher Rackauckas, Michael Krabbe Borregaard, and Thomas Breloff.
 421 2023. “Plots.jl – a User Extendable Plotting API for the Julia Programming Language.” <https://doi.org/https://doi.org/10.5334/jors.431>.
- 422 Danisch, Simon, and Julius Krumbiegel. 2021. “Makie.jl: Flexible High-Performance Data Visualization for Julia.”
 423 *Journal of Open Source Software* 6 (65): 3349. <https://doi.org/10.21105/joss.03349>.
- 424 Du, Yilun, and Igor Mordatch. 2020. “Implicit Generation and Generalization in Energy-Based Models.” <https://arxiv.org/abs/1903.08689>.
- 425 Frankle, Jonathan, and Michael Carbin. 2019. “The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural
 426 Networks.” In *International Conference on Learning Representations*.
- 427 Freiesleben, Timo. 2022. “The Intriguing Relation Between Counterfactual Explanations and Adversarial Examples.”
 428 *Minds and Machines* 32 (1): 77–109.
- 429 Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press.
- 430 Goodfellow, Ian, Jonathon Shlens, and Christian Szegedy. 2015. “Explaining and Harnessing Adversarial Examples.”
 431 <https://arxiv.org/abs/1412.6572>.
- 432 Grathwohl, Will, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swer-
 433 sky. 2020. “Your Classifier Is Secretly an Energy Based Model and You Should Treat It Like One.” In *International
 434 Conference on Learning Representations*.
- 435 Gretton, Arthur, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. 2012. “A Kernel
 436 Two-Sample Test.” *The Journal of Machine Learning Research* 13 (1): 723–73.
- 437 Guidotti, Riccardo. 2022. “Counterfactual Explanations and How to Find Them: Literature Review and Benchmark-
 438 ing.” *Data Mining and Knowledge Discovery* 38 (5): 2770–2824. <https://doi.org/10.1007/s10618-022-00831-6>.
- 439 Guo, Hangzhi, Thanh H. Nguyen, and Amulya Yadav. 2023. “CounterNet: End-to-End Training of Prediction Aware
 440 Counterfactual Explanations.” In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery*

- 443 and Data Mining, 577--589. KDD '23. New York, NY, USA: Association for Computing Machinery. <https://doi.org/10.1145/3580305.3599290>.
- 444 Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2009. *The Elements of Statistical Learning*. Springer New
445 York. <https://doi.org/10.1007/978-0-387-84858-7>.
- 446 Innes, Michael, Elliot Saba, Keno Fischer, Dhairyा Gandhi, Marco Concetto Rudilosso, Neethu Mariya Joy, Tejan
447 Karmali, Avik Pal, and Viral Shah. 2018. "Fashionable Modelling with Flux." <https://arxiv.org/abs/1811.01457>.
- 448 Innes, Mike. 2018. "Flux: Elegant Machine Learning with Julia." *Journal of Open Source Software* 3 (25): 602.
449 <https://doi.org/10.21105/joss.00602>.
- 450 Joshi, Shalmali, Oluwasanmi Koyejo, Warut Vigitbenjaronk, Been Kim, and Joydeep Ghosh. 2019. "Towards Realistic
451 Individual Recourse and Actionable Explanations in Black-Box Decision Making Systems." <https://arxiv.org/abs/1907.09615>.
- 452 Kaggle. 2011. "Give Me Some Credit, Improve on the State of the Art in Credit Scoring by Pre-
453 dicting the Probability That Somebody Will Experience Financial Distress in the Next Two Years." <https://www.kaggle.com/c/GiveMeSomeCredit>; Kaggle.
- 454 Kolter, Zico. 2023. "Keynote Addresses: SaTML 2023 ." In *2023 IEEE Conference on Secure and Trustworthy
455 Machine Learning (SaTML)*. Los Alamitos, CA, USA: IEEE Computer Society. [https://doi.org/10.1109/SaTML5
4575.2023.00009](https://doi.org/10.1109/SaTML5
456 4575.2023.00009).
- 457 Lakshminarayanan, Balaji, Alexander Pritzel, and Charles Blundell. 2017. "Simple and Scalable Predictive Un-
458 certainty Estimation Using Deep Ensembles." In *Proceedings of the 31st International Conference on Neural
459 Information Processing Systems*, 6405–16. NIPS'17. Red Hook, NY, USA: Curran Associates Inc.
- 460 LeCun, Yann. 1998. "The MNIST database of handwritten digits." <http://yann.lecun.com/exdb/mnist/>.
- 461 Lippe, Phillip. 2024. "UvA Deep Learning Tutorials." <https://uvadlc-notebooks.readthedocs.io/en/latest/>.
- 462 Luu, Hoai Linh, and Naoya Inoue. 2023. "Counterfactual Adversarial Training for Improving Robustness of Pre-
463 trained Language Models." In *Proceedings of the 37th Pacific Asia Conference on Language, Information and
464 Computation*, 881–88. ACL. <https://aclanthology.org/2023.paclic-1.88>.
- 465 Murphy, Kevin P. 2022. *Probabilistic Machine Learning: An Introduction*. MIT Press.
- 466 O'Neil, Cathy. 2016. *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*.
467 Crown.
- 468 Pace, R Kelley, and Ronald Barry. 1997. "Sparse Spatial Autoregressions." *Statistics & Probability Letters* 33 (3):
469 291–97. [https://doi.org/10.1016/s0167-7152\(96\)00140-x](https://doi.org/10.1016/s0167-7152(96)00140-x).
- 470 Pawelczyk, Martin, Chirag Agarwal, Shalmali Joshi, Sohini Upadhyay, and Himabindu Lakkaraju. 2022. "Exploring
471 Counterfactual Explanations Through the Lens of Adversarial Examples: A Theoretical and Empirical Analysis." In
472 *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, edited by Gustau
473 Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, 151:4574–94. Proceedings of Machine Learning Research.
474 PMLR. <https://proceedings.mlr.press/v151/pawelczyk22a.html>.
- 475 Poyiadzi, Rafael, Kacper Sokol, Raul Santos-Rodriguez, Tijl De Bie, and Peter Flach. 2020. "FACE: Feasible and
476 Actionable Counterfactual Explanations." In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*,
477 344–50.
- 478 Ross, Alexis, Himabindu Lakkaraju, and Osbert Bastani. 2024. "Learning Models for Actionable Recourse." In
479 *Proceedings of the 35th International Conference on Neural Information Processing Systems*. NIPS '21. Red
480 Hook, NY, USA: Curran Associates Inc.
- 481 Sauer, Axel, and Andreas Geiger. 2021. "Counterfactual Generative Networks." <https://arxiv.org/abs/2101.06046>.
- 482 Schut, Lisa, Oscar Key, Rory McGrath, Luca Costabello, Bogdan Sacaleanu, Yarin Gal, et al. 2021. "Generating
483 Interpretable Counterfactual Explanations By Implicit Minimisation of Epistemic and Aleatoric Uncertainties." In
484 *International Conference on Artificial Intelligence and Statistics*, 1756–64. PMLR.
- 485 Sharma, Shubham, Jette Henderson, and Joydeep Ghosh. 2020. "CERTIFAI: A Common Framework to Provide
486 Explanations and Analyse the Fairness and Robustness of Black-box Models." In *Proceedings of the AAAI/ACM
487 Conference on AI, Ethics, and Society*, 166–72. AIES '20. New York, NY, USA: Association for Computing
488 Machinery. <https://doi.org/10.1145/3375627.3375812>.
- 489 Spooner, Thomas, Danial Dervovic, Jason Long, Jon Shepard, Jiahao Chen, and Daniele Magazzeni. 2021. "Counter-
490 factual Explanations for Arbitrary Regression Models." <https://arxiv.org/abs/2106.15212>.
- 491 Szegedy, Christian, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus.
492 2014. "Intriguing Properties of Neural Networks." <https://arxiv.org/abs/1312.6199>.
- 493 Teh, Yee Whye, Max Welling, Simon Osindero, and Geoffrey E. Hinton. 2003. "Energy-Based Models for Sparse
494 Overcomplete Representations." *J. Mach. Learn. Res.* 4 (null): 1235–60.
- 495 Teney, Damien, Ehsan Abbasnejad, and Anton van den Hengel. 2020. "Learning What Makes a Difference from
496 Counterfactual Examples and Gradient Supervision." In *Computer Vision - ECCV 2020*, 580–99. Berlin, Heidelberg:
497 Springer-Verlag. https://doi.org/10.1007/978-3-030-58607-2_34.

- 501 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,
 502 NY, USA: Association for Computing Machinery. <https://doi.org/10.1145/3351095.3372876>.
- 503 Wachter, Sandra, Brent Mittelstadt, and Chris Russell. 2017. “Counterfactual Explanations Without Opening the Black
 504 Box: Automated Decisions and the GDPR.” *Harv. JL & Tech.* 31: 841. <https://doi.org/10.2139/ssrn.3063289>.
- 505 Wilson, Andrew Gordon. 2020. “The Case for Bayesian Deep Learning.” <https://arxiv.org/abs/2001.10995>.
- 506 Wu, Tongshuang, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2021. “Polyjuice: Generating Counterfactuals
 507 for Explaining, Evaluating, and Improving Models.” In *Proceedings of the 59th Annual Meeting of the Association
 508 for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing
 509 (Volume 1: Long Papers)*, 6707–23. ACL. <https://doi.org/10.18653/v1/2021.acl-long.523>.
- 510 Yeh, I-Cheng. 2016. “Default of Credit Card Clients.” UCI Machine Learning Repository.
- 511 Zhang, Chiyuan, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. 2021. “Understanding Deep
 512 Learning (Still) Requires Rethinking Generalization.” *Commun. ACM* 64 (3): 107–15. <https://doi.org/10.1145/3446776>.
- 513
- 514

515 **Appendix A Notation**

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

- 517 • y^+ : The target class and also the index of the target class.
- 518 • y^- : The non-target class and also the index of non-the target class.
- 519 • \mathbf{x} : a single training sample.
- 520 • \mathbf{x}' : a counterfactual.
- 521 • \mathbf{x}^+ : a training sample in the target class (ground-truth).
- 522 • \mathbf{y}^+ : The one-hot encoded output vector for the target class.
- 523 • θ : Model parameters (unspecified).
- 524 • Θ : Matrix of parameters.
- 525 • $\mathbf{M}(\cdot)$: linear predictions (logits) of the classifier.

526 **A.1 Other Technical Details**

527 Maximum mean discrepancy is defined as follows,

$$\begin{aligned} \text{MMD}(X', \tilde{X}') &= \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(x_i, x_j) \\ &\quad + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n k(\tilde{x}_i, \tilde{x}_j) \\ &\quad - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, \tilde{x}_j) \end{aligned} \tag{6}$$

528 where $k(\cdot, \cdot)$ is a kernel function (Gretton et al. 2012). We make use of a Gaussian kernel with a constant length-scale
 529 parameter of 0.5. In our implementation, Equation 6 is by default applied to the entire subset of the training data for
 530 which $y = y^+$.

531 **Appendix B Technical Details of Our Approach**

532 **B.1 Generating Counterfactuals through Gradient Descent**

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

535 **B.1.1 Background**

536 Gradient-based counterfactual search was originally proposed by Wachter, Mittelstadt, and Russell (2017). It generally
 537 solves the following unconstrained objective,

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

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

542 For all generators used in this work we use standard logit crossentropy loss for $\text{ylloss}(\cdot)$. All generators also penalize
 543 the distance (ℓ_1 -norm) of counterfactuals from their original factual state. For *Generic* and *ECCo*, we have $\mathcal{Z} := \mathcal{X}$
 544 and $g(\mathbf{z}) = g(\mathbf{z})^{-1} = \mathbf{z}$, that is counterfactual are searched directly in the feature space. Conversely, *REVISE* traverses
 545 the latent space of a variational autoencoder (VAE) fitted to the training data, where $g(\cdot)$ corresponds to the decoder
 546 (Joshi et al. 2019). In addition to the distance penalty, *ECCo* uses an additional penalty component that regularizes
 547 the energy associated with the counterfactual, \mathbf{x}' (Altmeyer et al. 2024).

548 **B.1.2 Convergence**

549 An important consideration when generating counterfactual explanations using gradient-based methods is how to
 550 define convergence. Two common choices are to 1) perform gradient descent over a fixed number of iterations T , or
 551 2) conclude the search as soon as the predicted probability for the target class has reached a pre-determined threshold,

552 $\tau: \mathcal{S}(\mathbf{M}_\theta(\mathbf{x}'))[y^+] \geq \tau$. We prefer the latter for our purposes, because it explicitly defines convergence in terms of the
 553 black-box model, $\mathbf{M}(\mathbf{x})$.

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

558 B.2 Protecting Mutability Constraints with Linear Classifiers

559 In Section 3.4 we explain that to avoid penalizing implausibility that arises due to mutability constraints, we impose a
 560 point mass prior on $p(\mathbf{x})$ for the corresponding feature. We argue in Section 3.4 that this approach induces models to
 561 be less sensitive to immutable features and demonstrate this empirically in Section 4. Below we derive the analytical
 562 results in Prp.~3.1.

563 *Proof.* Let d_{mtbl} and d_{immmtbl} denote some mutable and immutable feature, respectively. Suppose that $\mu_{y^-, d_{\text{immmtbl}}} <$
 564 $\mu_{y^+, d_{\text{immmtbl}}}$ and $\mu_{y^-, d_{\text{mtbl}}} > \mu_{y^+, d_{\text{mtbl}}}$, where $\mu_{k,d}$ denotes the conditional sample mean of feature d in class k . In words,
 565 we assume that the immutable feature tends to take lower values for samples in the non-target class y^- than in the
 566 target class y^+ . We assume the opposite to hold for the mutable feature.

567 Assuming multivariate Gaussian class densities with common diagonal covariance matrix $\Sigma_k = \Sigma$ for all $k \in \mathcal{K}$, we
 568 have for the log likelihood ratio between any two classes $k, m \in \mathcal{K}$ (Hastie, Tibshirani, and Friedman 2009):

$$\log \frac{p(k|\mathbf{x})}{p(m|\mathbf{x})} = \mathbf{x}^\top \Sigma^{-1} (\mu_k - \mu_m) + \text{const} \quad (7)$$

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

$$\log \frac{p(k|\mathbf{x})}{p(m|\mathbf{x})} = \sum_{d=1}^D \frac{\mu_{k,d} - \mu_{m,d}}{\sigma_d^2} x_d + \text{const} \quad (8)$$

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

$$\log \frac{p(k|\mathbf{x})}{p(m|\mathbf{x})} = \sum_{d=1}^D (\theta_{k,d} - \theta_{m,d}) x_d + \text{const} \quad (9)$$

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

572 Based on Equation 8 and Equation 9 we can identify that $(\mu_{k,d} - \mu_{m,d}) \propto (\theta_{k,d} - \theta_{m,d})$ under the assumptions we
 573 made above. Hence, we have that $(\theta_{y^-, d_{\text{immmtbl}}} - \theta_{y^+, d_{\text{immmtbl}}}) < 0$ and $(\theta_{y^-, d_{\text{mtbl}}} - \theta_{y^+, d_{\text{mtbl}}}) > 0$

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

$$\frac{\partial}{\partial \theta_{y^+, d}} (\text{div}(\mathbf{x}^+, \mathbf{x}', \mathbf{y}; \theta)) = \frac{\partial}{\partial \theta_{y^+, d}} ((-\mathbf{M}_\theta(\mathbf{x}^+)[y^+]) - (-\mathbf{M}_\theta(\mathbf{x}') [y^+])) = x'_d - x_d^+ \quad (10)$$

577 and equal to zero everywhere else.

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

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

584 By the same rationale, the contrastive divergence penalty can generally be expected to increase $\theta_{y^-, d_{\text{mtbl}}}$ exacerbating
 585 $(\theta_{y^-, d_{\text{mtbl}}} - \theta_{y^+, d_{\text{mtbl}}}) > 0$. In words, this has the effect of making the classifier more sensitive to the mutable feature, in
 586 that it would be more likely to predict class y^- as opposed to y^+ for higher values of d_{mtbl} .

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

590 B.3 Domain Constraints

591 We apply domain constraints on counterfactuals during training and evaluation. There are at least two good reasons for
 592 doing so. Firstly, within the context of explainability and algorithmic recourse, real-world attributes are often domain
 593 constrained: the *age* feature, for example, is lower bounded by zero and upper bounded by the maximum human
 594 lifespan. Secondly, domain constraints help mitigate training instabilities commonly associated with energy-based
 595 modelling (Grathwohl et al. 2020; Altmeyer et al. 2024).

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

$$\begin{aligned} x_d^{\text{LB}} &= \arg \min_{x_d} \{\mu_d - n_{\sigma_d} \sigma_d, \arg \min_{x_d} x_d\} \\ x_d^{\text{UB}} &= \arg \max_{x_d} \{\mu_d + n_{\sigma_d} \sigma_d, \arg \max_{x_d} x_d\} \end{aligned} \quad (11)$$

599 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
 600 higher values and hence wider bounds may be appropriate depending on the application.

601 B.4 Training Hyperparameters

602 Note 1 presents the default hyperparameters used during training.

Note 1: 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

603

- Objective: full
- Optimizer: adam

604

605 **B.5 Evaluation Details**

606 For all of our evaluations, we proceed as follows: for each experiment setting we generate multiple counterfactuals
 607 (“No. Counterfactuals”), randomly choosing the factual and target class each time (Note 2). We do this across multiple
 608 rounds (“No. Runs”) with different random seeds to account for stochasticity (Note 2). This is in line with standard
 609 practice in the related literature on CE. Note 2 presents the default hyperparameters used during evaluation. For
 610 our final results presented in the main paper, we rely on held out test sets to sample factuals (and outputs for our
 611 performance metrics). For tuning purposes we rely on training or validation sets.

612 **B.5.1 Robust Accuracy**

613 To evaluate robust accuracy (Acc.*), we use the Fast Gradient Sign Method (FGSM) to perturb test samples (Goodfellow,
 614 Shlens, and Szegedy 2015). For the main results, we have set the perturbation size to $\epsilon = 0.03$. We have also
 615 tested other perturbation sizes, as well as randomly perturbed data. Although not reported here, we have consistently
 616 found strong outperformance of CT compared to the weak baseline.

Note 2: Evaluation Phase

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

617

618 **Appendix C Details on Main Experiments**619 **C.1 Final Hyperparameters**

620 As discussed Section 4, CT is sensitive to certain hyperparameter choices. We study the effect of many hyperparameters
 621 extensively in Section D. For the main results, we tune a small set of key hyperparameters (Section E). The final
 622 choices for the main results are presented for each data set in Table 2 along with training, test and batch sizes.

Table 2: Final hyperparameters used for the main results presented in Section 4. 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}
Adult	26049	5010	1000	none	0.75	5000	0.25
CH	16504	3101	1000	none	0.5	5000	0.25
Circ	3600	600	30	none	0.5	1000	0.5
Cred	10617	1923	1000	none	0.5	5000	0.25
GMSC	13371	2474	1000	none	0.5	5000	0.5
LS	3600	600	30	none	0.5	1000	0.01
MNIST	11000	2000	1000	(-1.0, 1.0)	0.5	5000	0.01
Moon	3600	600	30	none	0.9	1000	0.25
OL	3600	600	30	none	0.5	1000	0.25

623 **C.2 Qualitative Findings for Image Data**

624 Figure 2 shows much more plausible (faithful) counterfactuals for a model with CT than the model with conventional
 625 training (Figure 3).

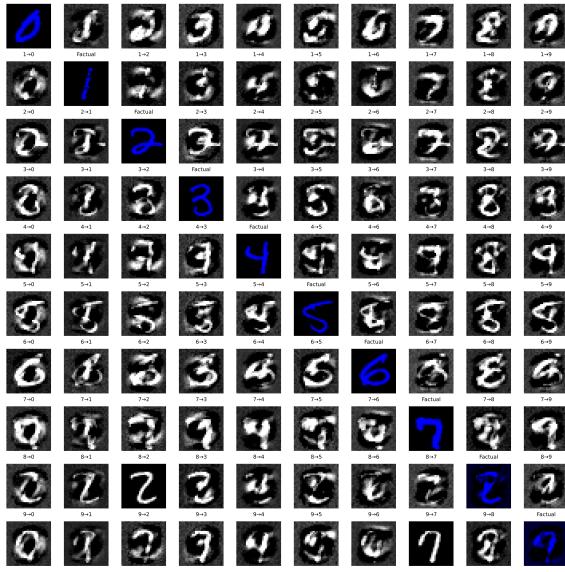


Figure 2: Counterfactual images for *MLP* with counterfactual training. Factual images are shown on the diagonal, with the corresponding counterfactual for each target class (columns) in that same row. The underlying generator, *ECCo*, aims to generate counterfactuals that are faithful to the model (Altmeier et al. 2024).

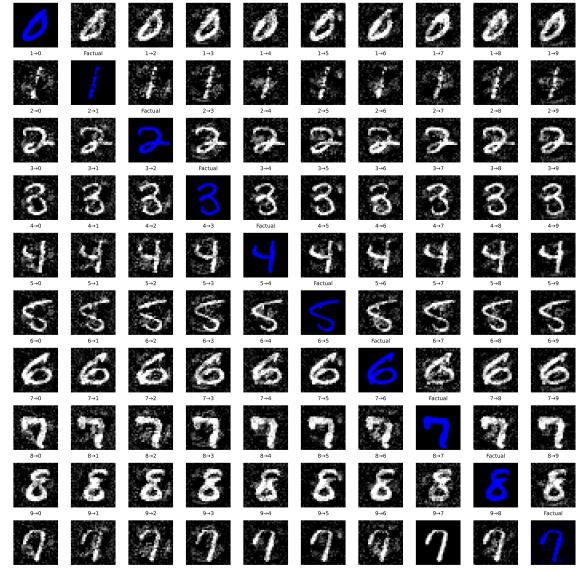


Figure 3: The same setup, factuals, model architecture and generator as in Figure 2, but the model was trained with CT.

626 **Appendix D Grid Searches**

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

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

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

638 **D.1 Evaluation Details**

639 To measure predictive performance, we compute the accuracy and F1-score for all models on test data (Table 3,
 640 Table 4, Table 5). With respect to explanatory performance, we report here our findings for the (im)plausibility and
 641 cost of counterfactuals at test time. Since the computation of our proposed divergence metric (Equation 5) is memory-
 642 intensive, we rely on the distance-based metric for the grid searches. For the counterfactual evaluation, we draw factual
 643 samples from the training data for the grid searches to avoid data leakage with respect to our final results reported in
 644 the body of the paper. Specifically, we want to avoid choosing our default hyperparameters based on results on the
 645 test data. Since we are optimizing for explainability, not predictive performance, we still present test accuracy and
 646 F1-scores.

647 **D.1.1 Predictive Performance**

648 We find that CT is associated with little to no decrease in average predictive performance for our synthetic datasets: test
 649 accuracy and F1-scores decrease by at most ~1 percentage point, but generally much less (Table 3, Table 4, Table 5).
 650 Variation across hyperparameters is negligible as indicated by small standard deviations for these metrics across the
 651 board.

652 **D.1.2 Counterfactual Outcomes**

653 Overall, we find that counterfactual training (CT) achieves its key objectives consistently across all hyperparameter
 654 settings and also broadly across datasets: plausibility is improved by up to ~60 percent (%) for the *Circles* data (e.g.
 655 Figure 4), ~25-30% for the *Moons* data (e.g. Figure 6) and ~10-20% for the *Linearly Separable* data (e.g. Figure 5). At
 656 the same time, the average costs of faithful counterfactuals are reduced in many cases by around ~20-25% for *Circles*
 657 (e.g. Figure 8) and up to ~50% for *Moons* (e.g. Figure 10). For the *Linearly Separable* data, costs are generally
 658 increased although typically by less than 10% (e.g. Figure 9), which reflects a common tradeoff between costs and
 659 plausibility (Altmeyer et al. 2024).

660 We do observe strong sensitivity to certain hyperparameters, with clear manageable patterns. Concerning generator
 661 parameters, we firstly find that using *REVISE* to generate counterfactuals during training typically yields the worst
 662 outcomes out of all generators, often leading to a substantial decrease in plausibility. This finding can be attributed to
 663 the fact that *REVISE* effectively assigns the task of learning plausible explanations from the model itself to a surrogate
 664 VAE. In other words, counterfactuals generated by *REVISE* are less faithful to the model than *ECCo* and *Generic*, and
 665 hence we would expect them to be a less effective and, in fact, potentially detrimental role in our training regime.
 666 Secondly, we observe that allowing for a higher number of maximum steps T for the counterfactual search generally
 667 yields better outcomes. This is intuitive, because it allows more counterfactuals to reach maturity in any given iteration.
 668 Looking in particular at the results for *Linearly Separable*, it seems that higher values for T in combination with higher
 669 decision thresholds (τ) yields the best results when using *ECCo*. But depending on the degree of class separability
 670 of the underlying data, a high decision-threshold can also affect results adversely, as evident from the results for the
 671 *Overlapping* data (Figure 7): here we find that CT generally fails to achieve its objective because only a tiny proportion
 672 of counterfactuals ever reaches maturity.

673 Regarding penalty strengths, we find that the strength of the energy regularization, λ_{reg} is a key hyperparameter, while
 674 sensitivity with respect to λ_{div} and λ_{adv} is much less evident. In particular, we observe that not regularizing energy
 675 enough or at all typically leads to poor performance in terms of decreased plausibility and increased costs, in particular
 676 for *Circles* (Figure 12), *Linearly Separable* (Figure 13) and *Overlapping* (Figure 15). High values of λ_{reg} can increase
 677 the variability in outcomes, in particular when combined with high values for λ_{div} and λ_{adv} , but this effect is less
 678 pronounced.

679 Finally, concerning other hyperparameters we observe that the effectiveness and stability of CT is positively associated
 680 with the number of counterfactuals generated during each training epoch, in particular for *Circles* (Figure 20) and
 681 *Moons* (Figure 22). We further find that a higher number of training epochs is beneficial as expected, where we tested
 682 training models for 50 and 100 epochs. Interestingly, we find that it is not necessary to employ CT during the entire
 683 training phase to achieve the desired improvements in explainability: specifically, we have tested training models
 684 conventionally during the first half of training before switching to CT after this initial burn-in period.

685 **D.2 Generator Parameters**

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

Note 3: 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*

688

Note 4: Evaluation Phase

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

689

690 **D.2.1 Predictive Performance**

691 Predictive performance measures for this grid search are shown in Table 3.

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

Dataset	Variable	Objective	Mean	Std
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

692 **D.2.2 Plausibility**

693 The results with respect to the plausibility measure are shown in Figure 4 to Figure 7.

694 **D.2.3 Cost**

695 The results with respect to the cost measure are shown in Figure 8 to Figure 11.

696 **D.3 Penalty Strengths**

697 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 5: 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`

699

Note 6: Evaluation Phase

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

700

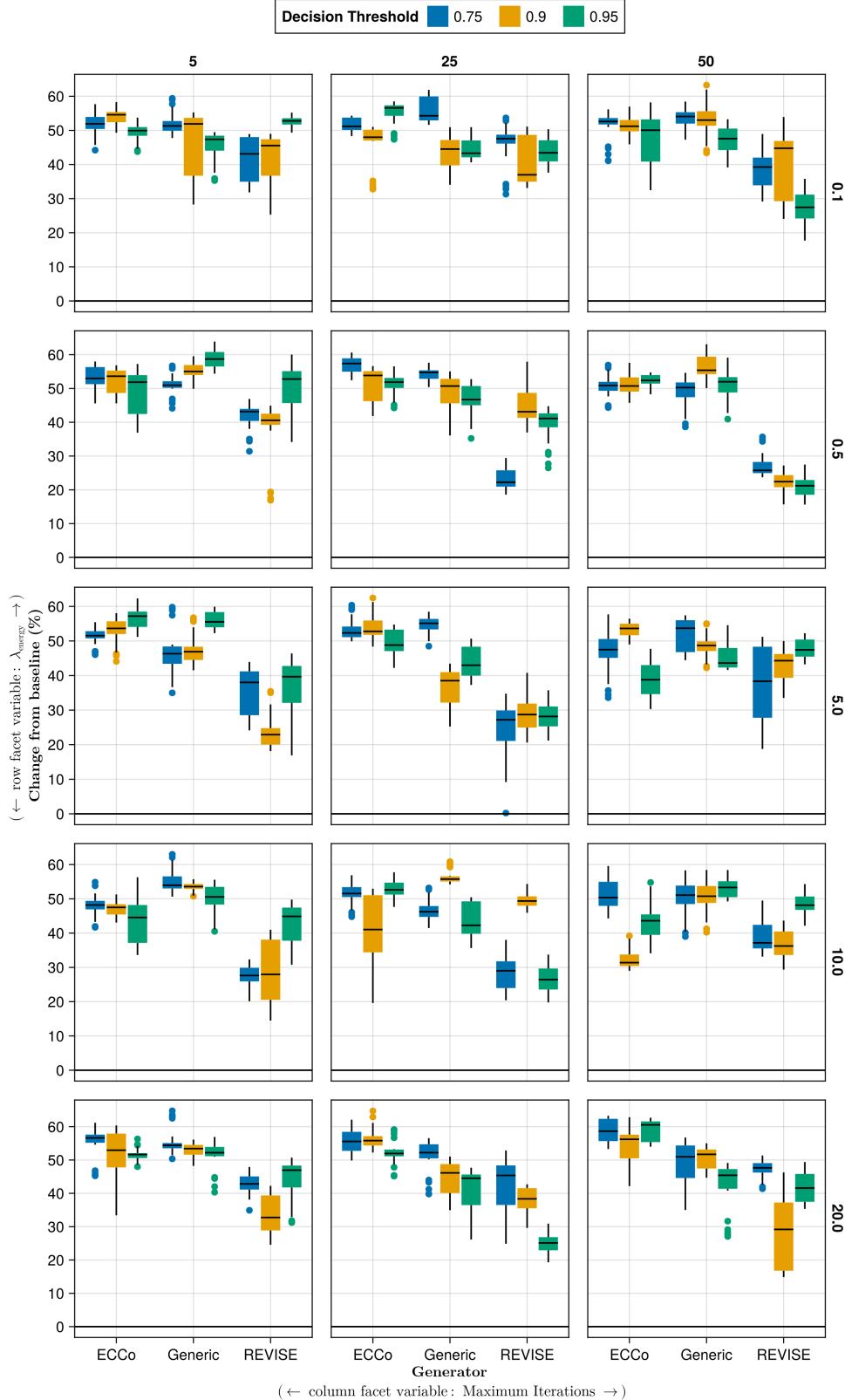


Figure 4: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Circles.

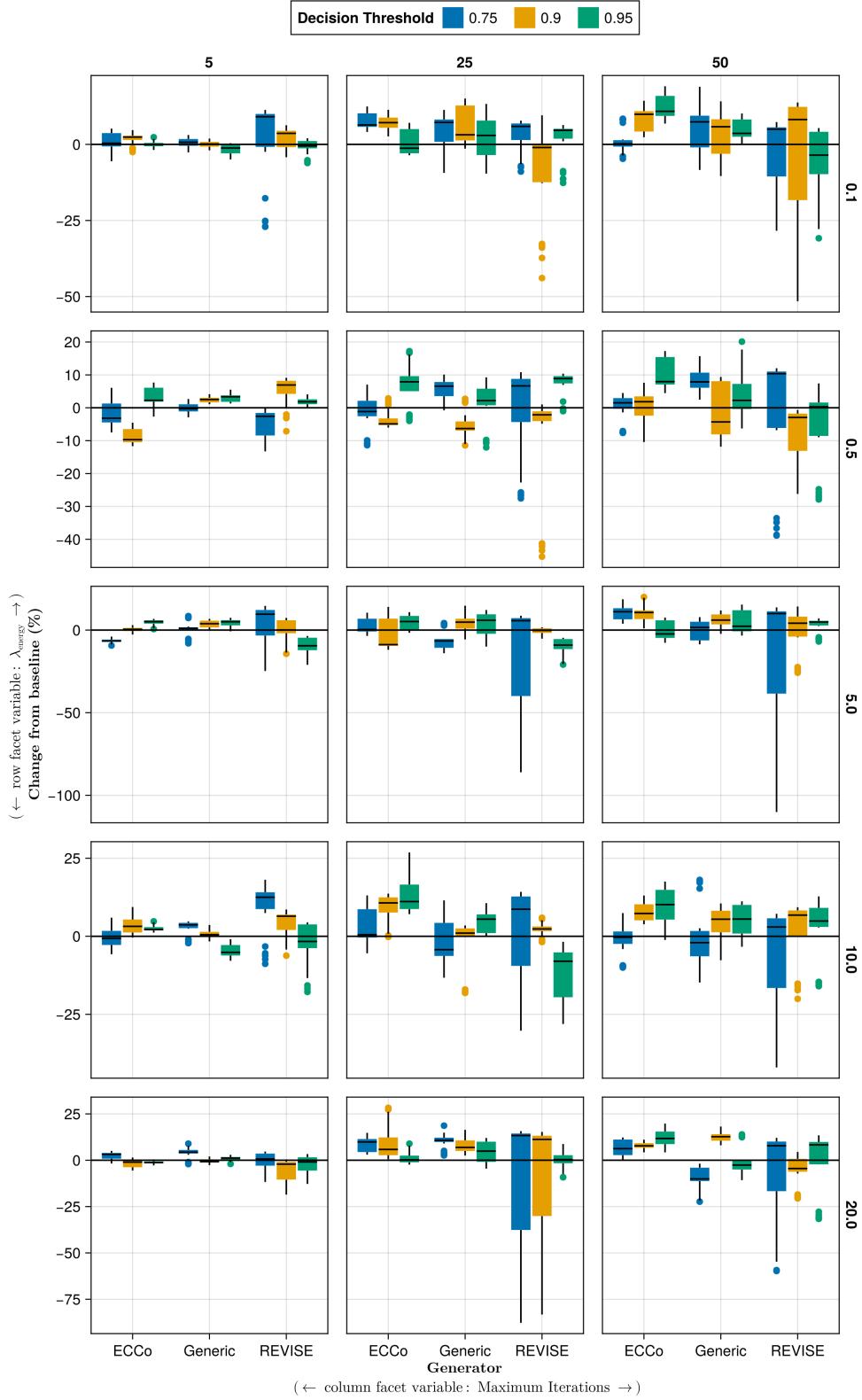


Figure 5: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Linearly Separable.

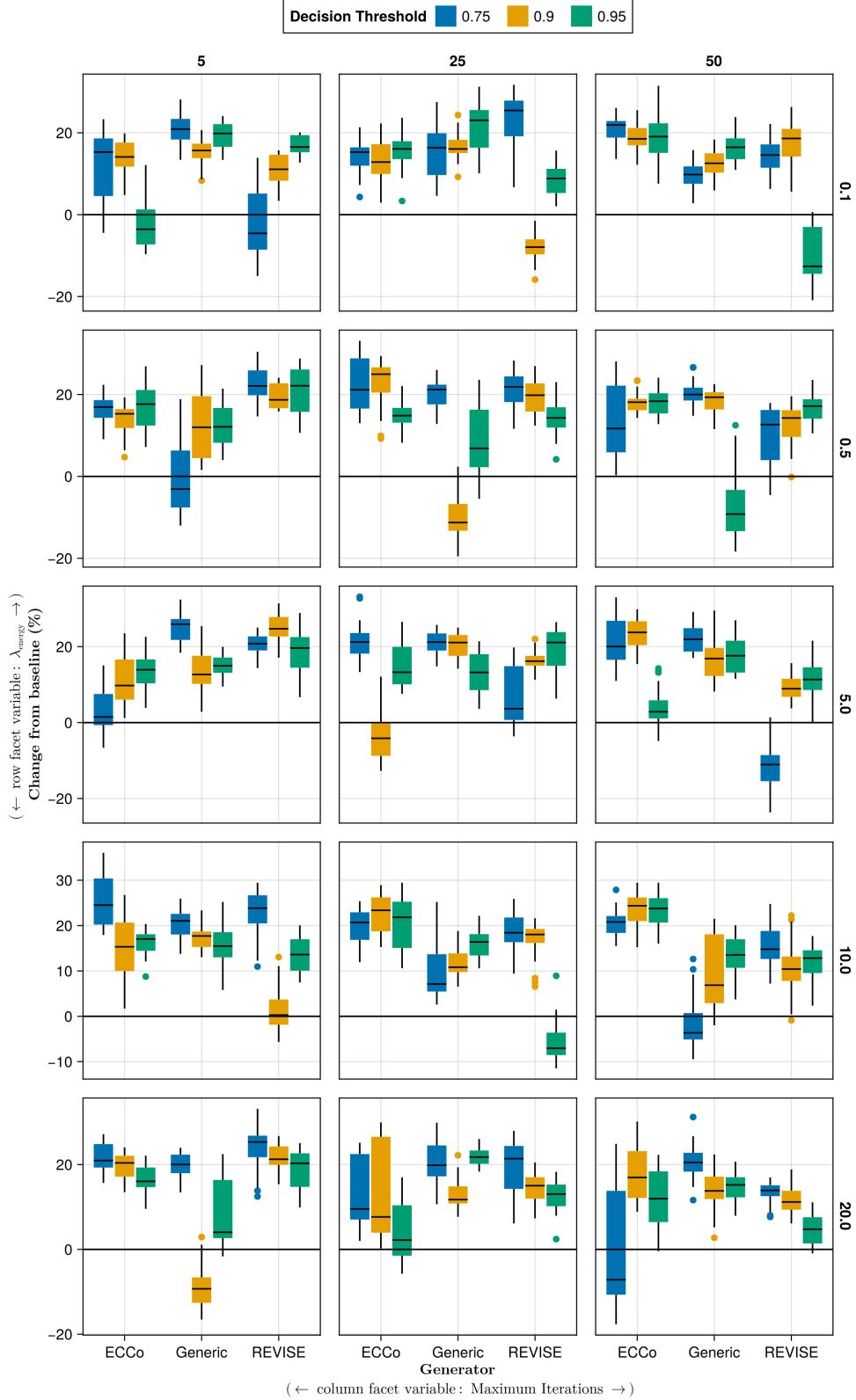


Figure 6: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Moons.

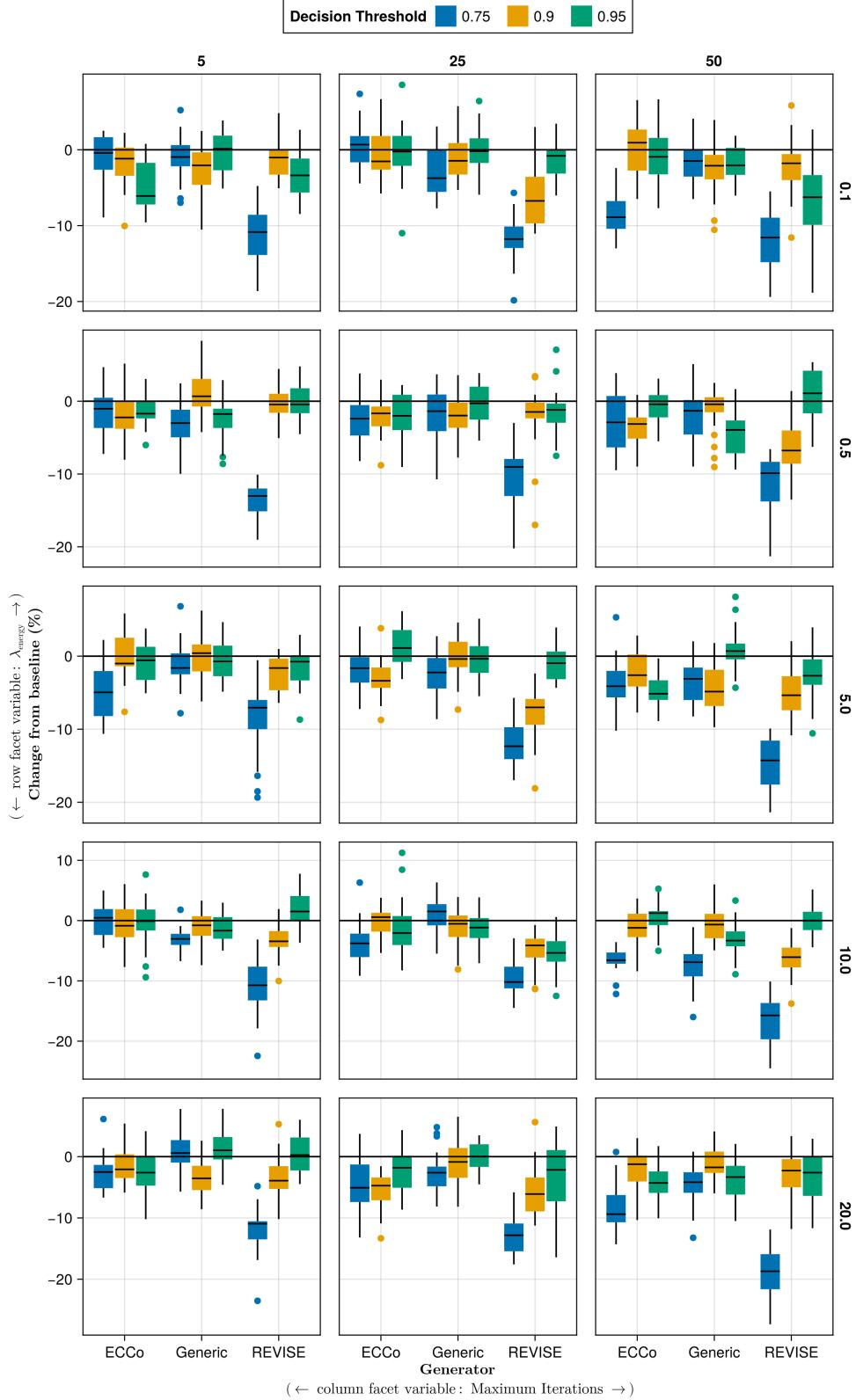


Figure 7: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

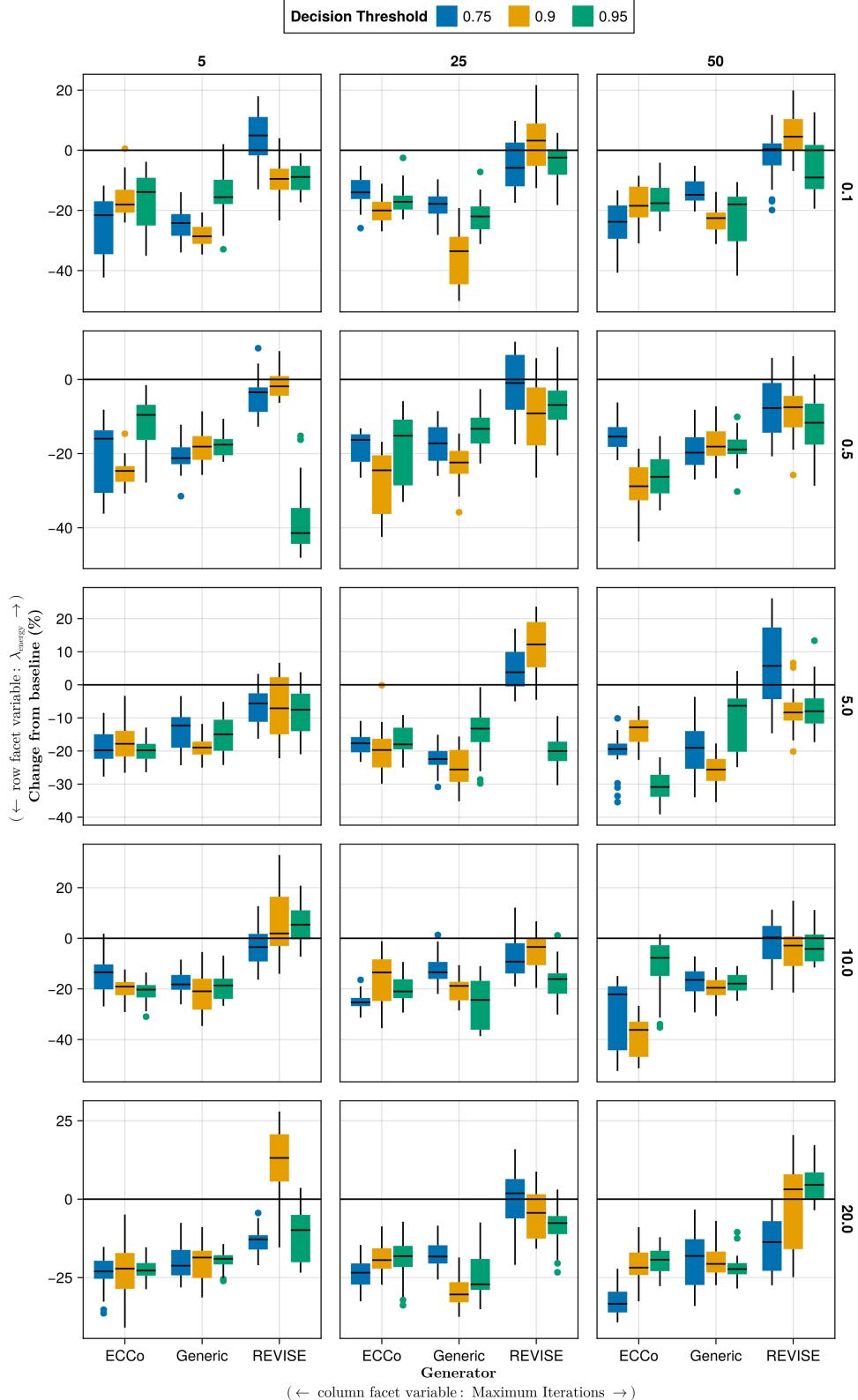


Figure 8: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Circles.

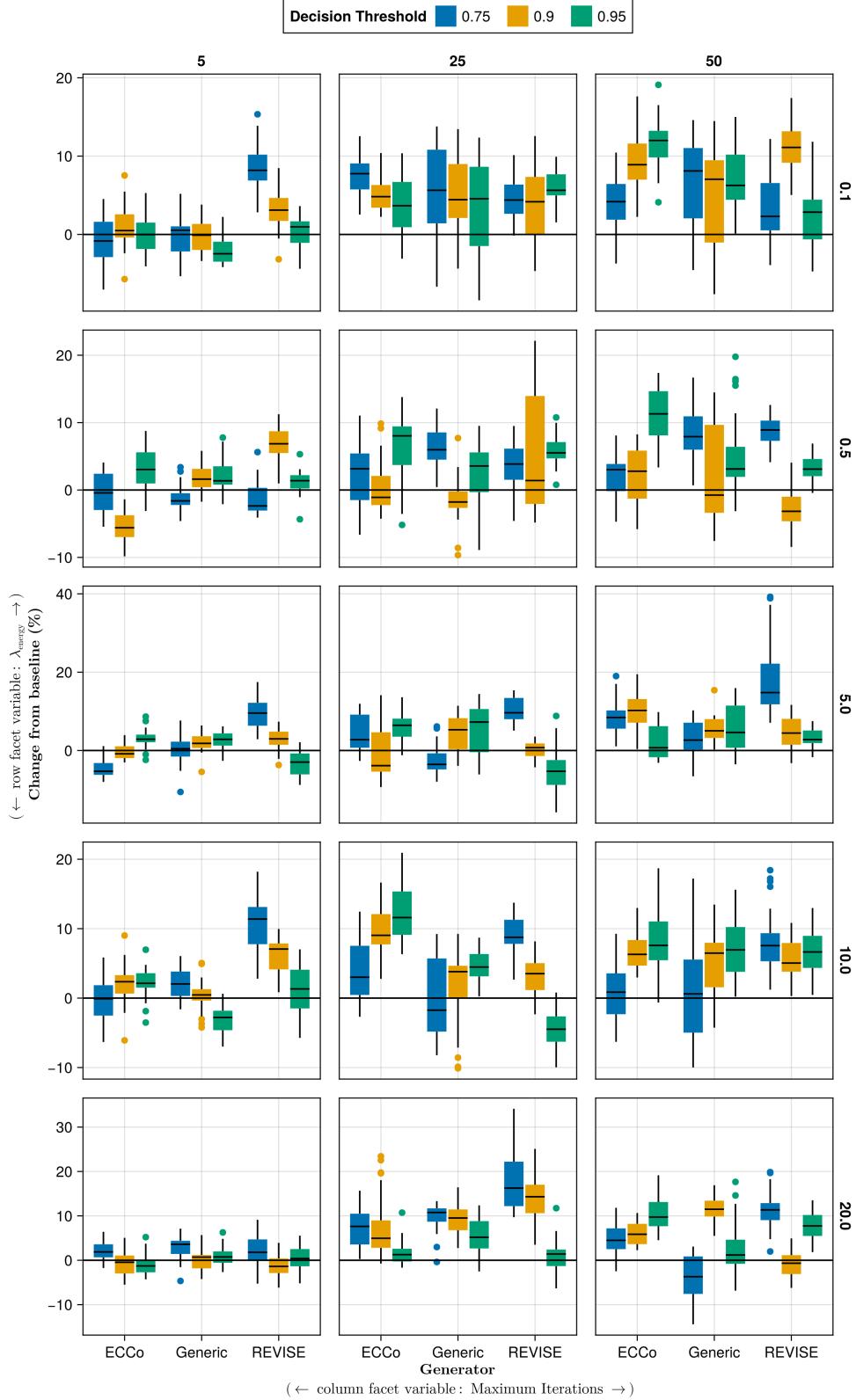


Figure 9: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Linearly Separable.

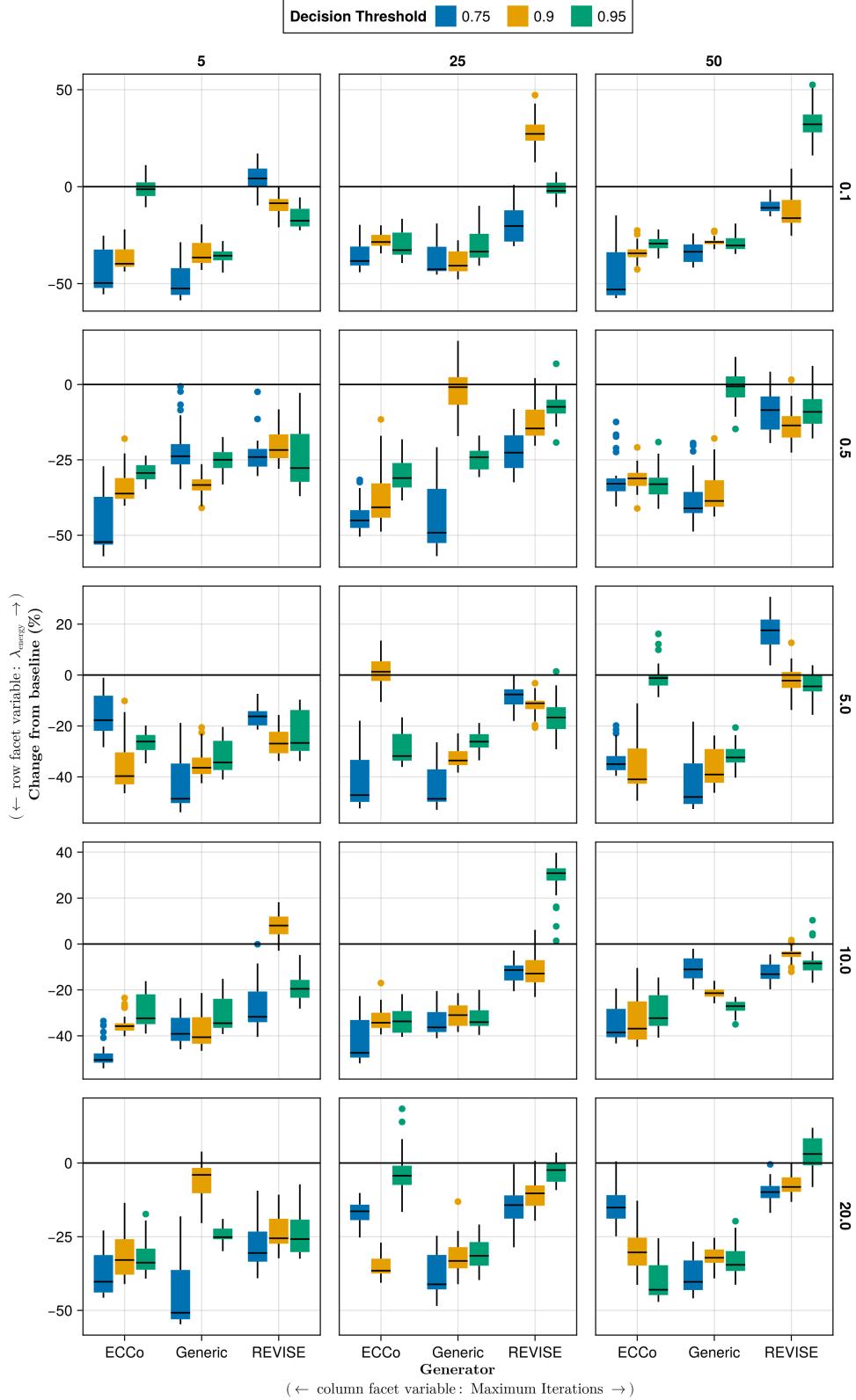


Figure 10: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric (Wachter, Mittelstadt, and Russell 2017). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Moons.

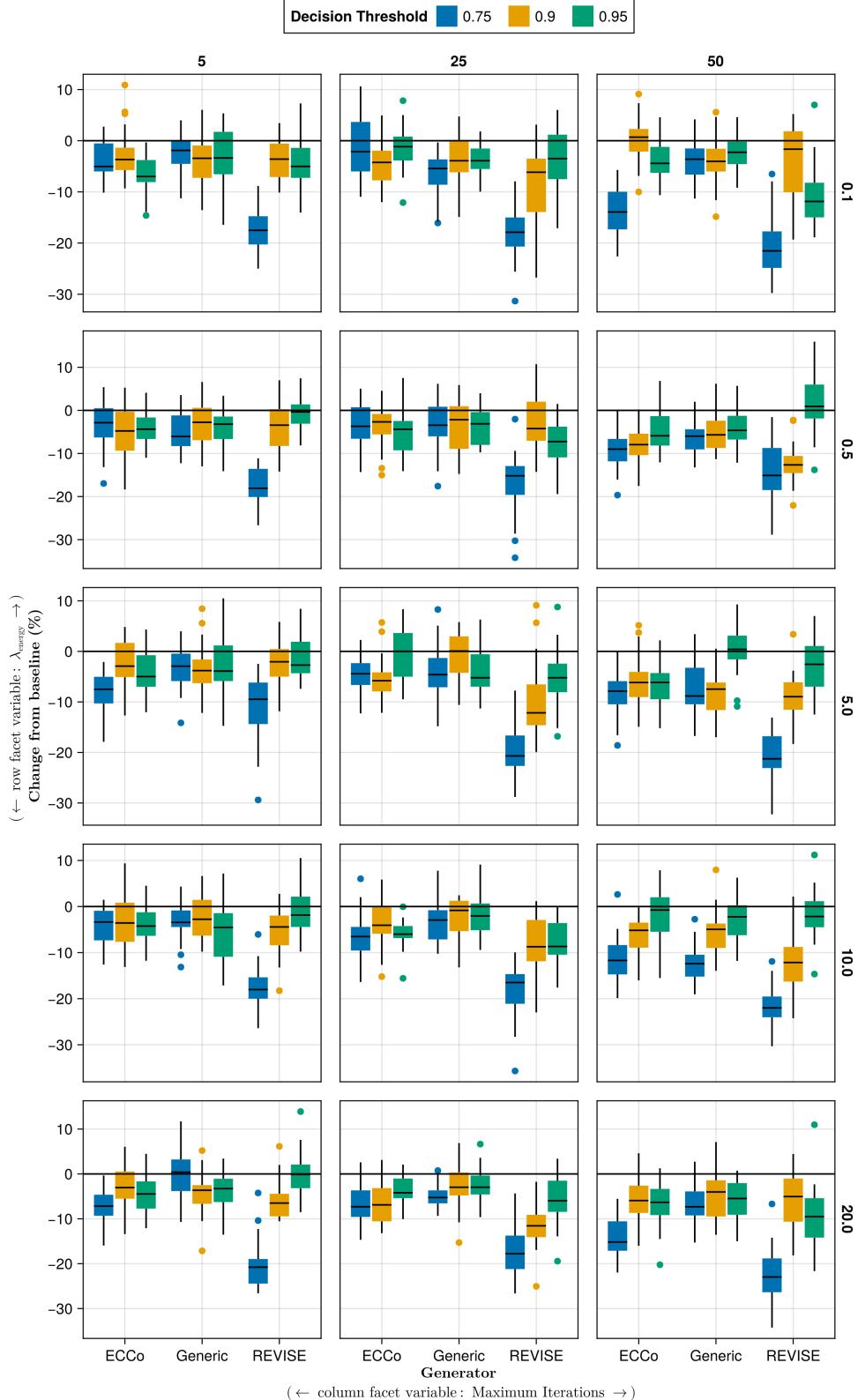


Figure 11: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric (Wachter, Mittelstadt, and Russell 2017). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

701 **D.3.1 Predictive Performance**

702 Predictive performance measures for this grid search are shown in Table 4.

Table 4: Predictive performance measures by dataset and objective averaged across training-phase parameters (Note 5) and evaluation-phase parameters (Note 6).

Dataset	Variable	Objective	Mean	Std
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

703 **D.3.2 Plausibility**

704 The results with respect to the plausibility measure are shown in Figure 12 to Figure 15.

705 **D.3.3 Cost**

706 The results with respect to the cost measure are shown in Figure 16 to Figure 19.

707 **D.4 Other Parameters**708 The hyperparameter grid with other varying training parameters is shown in Note 7. The corresponding evaluation
709 grid used for these experiments is shown in Note 8.

Note 7: 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`

710

Note 8: Evaluation Phase

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

711

712 **D.4.1 Predictive Performance**

713 Predictive performance measures for this grid search are shown in Table 5.

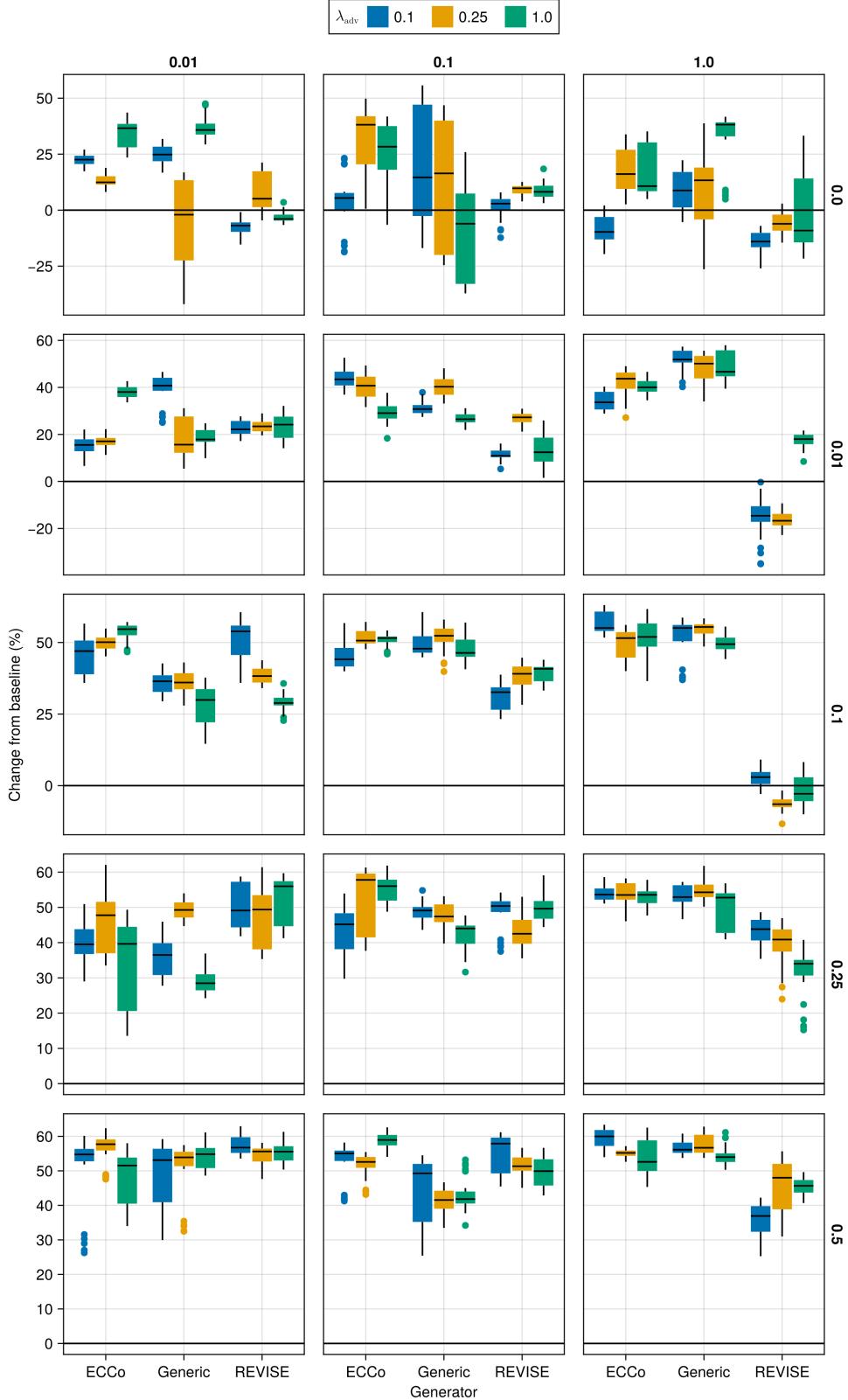


Figure 12: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Circles.

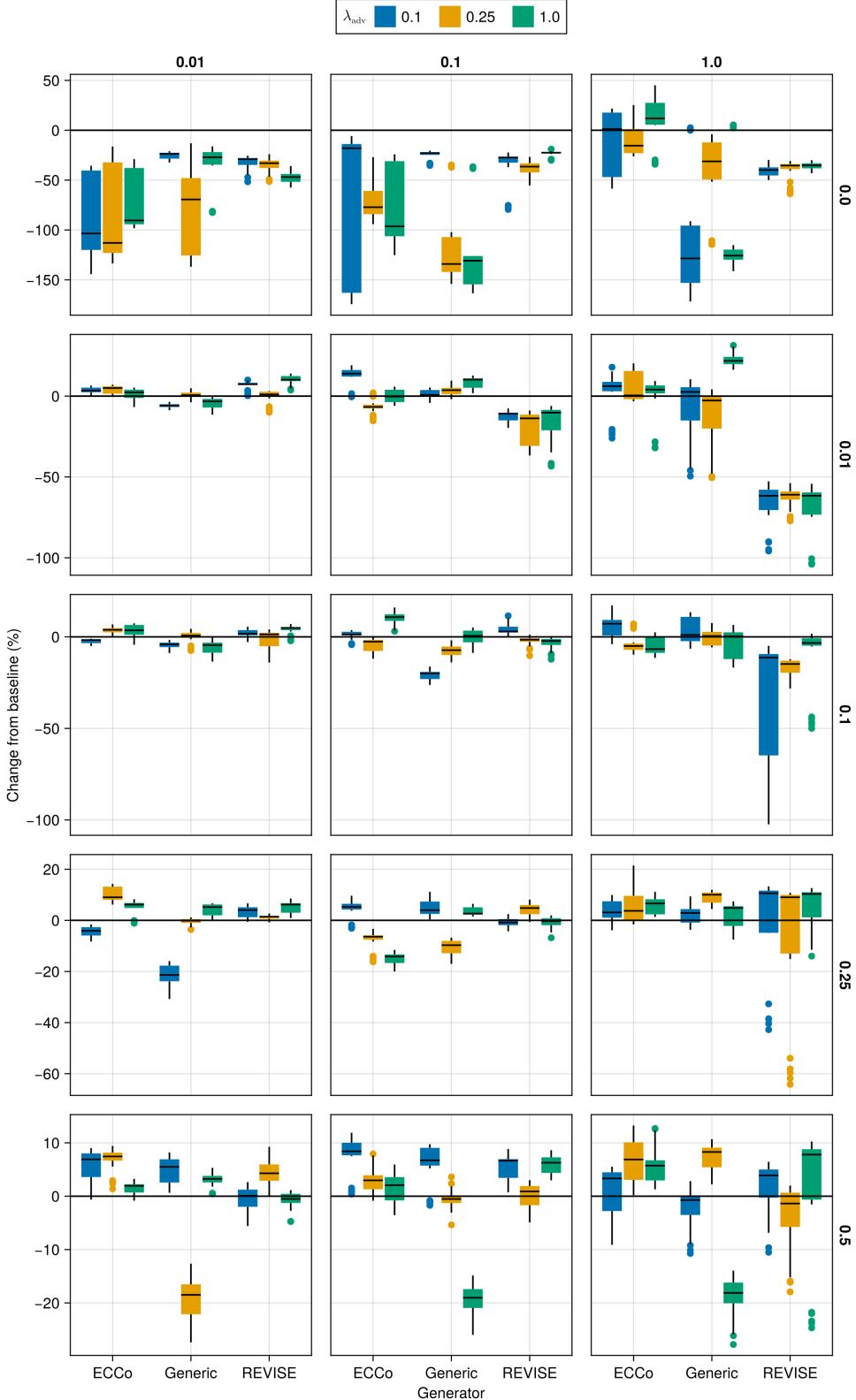


Figure 13: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Linearly Separable.

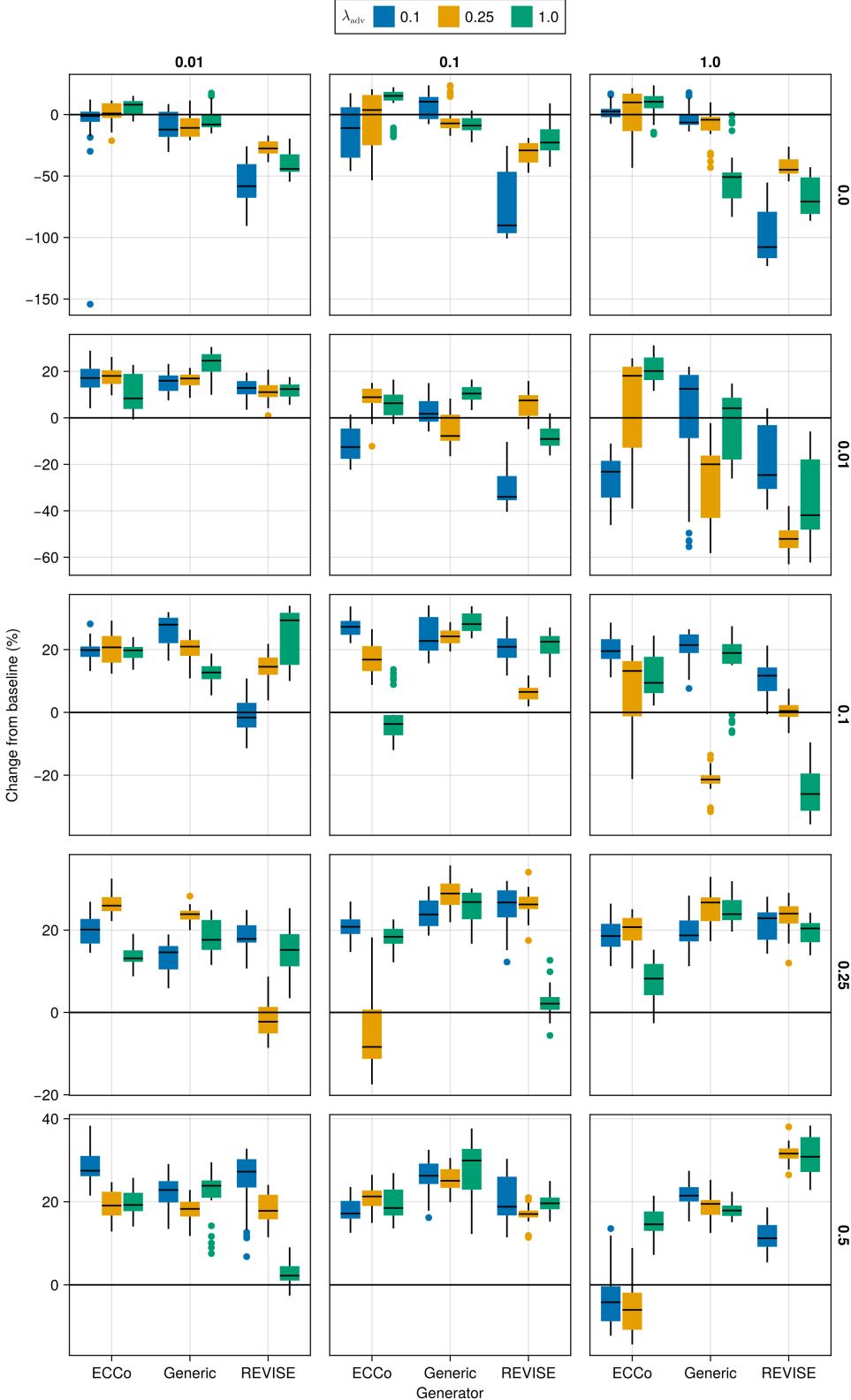


Figure 14: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Moons.

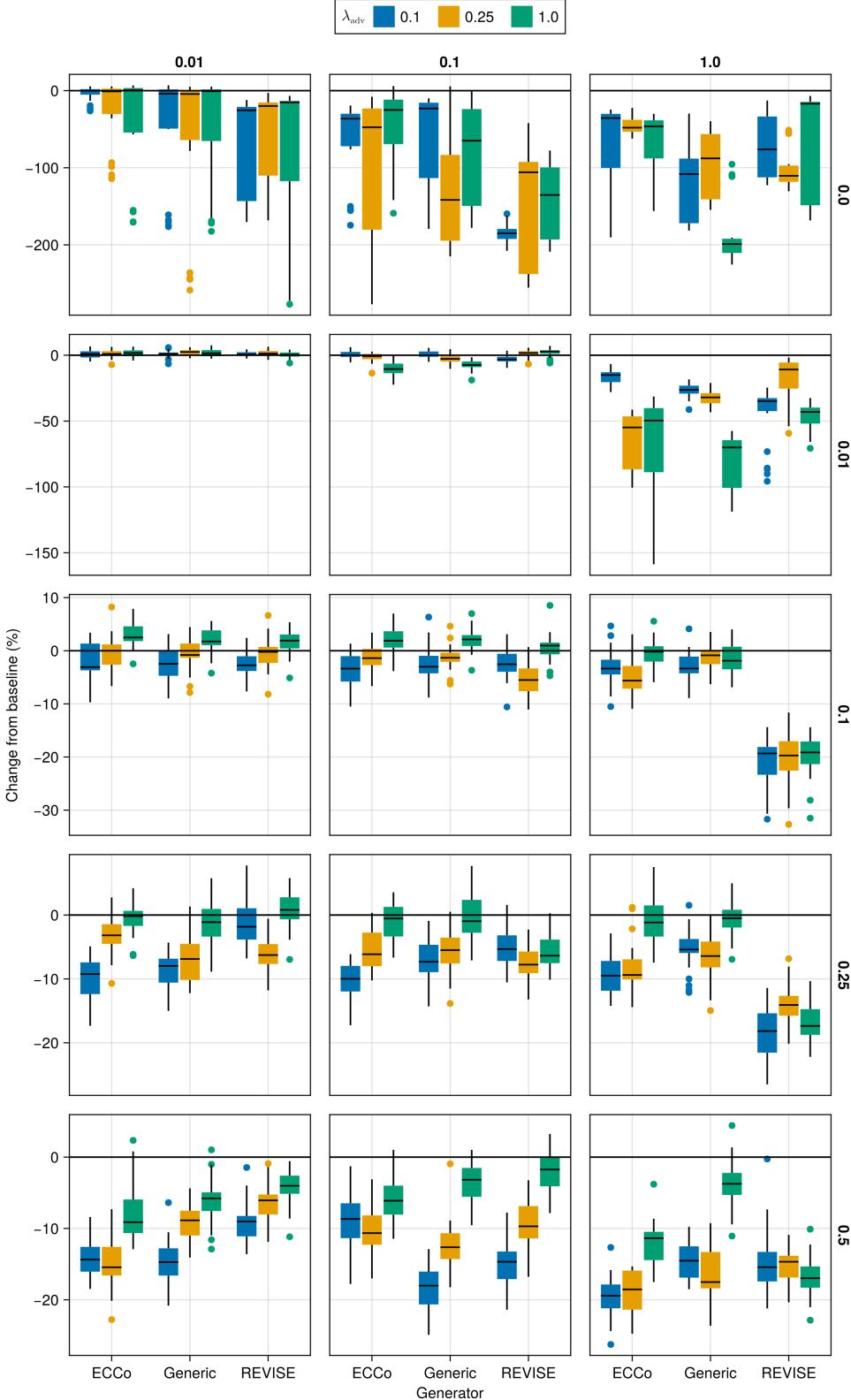


Figure 15: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

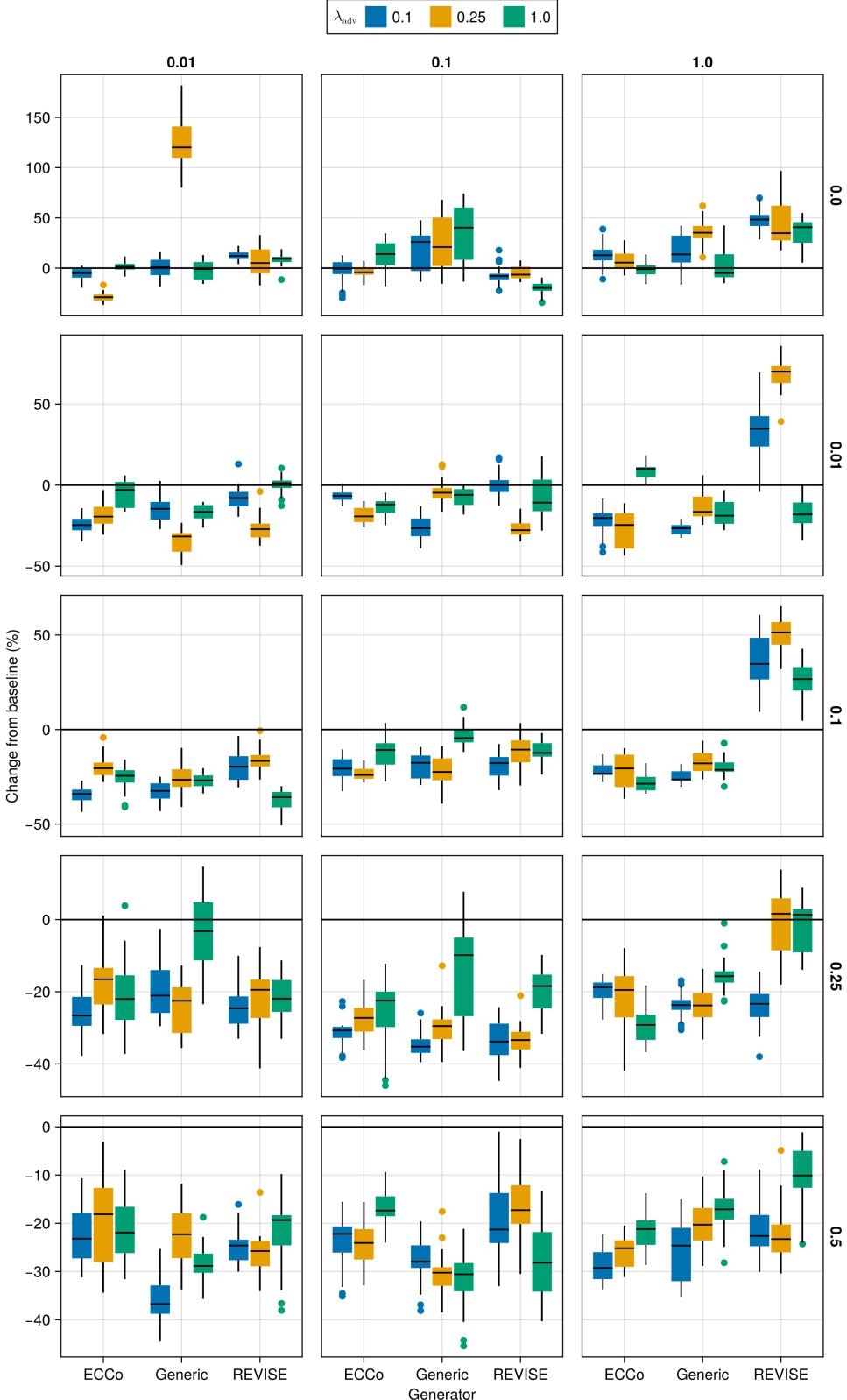


Figure 16: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

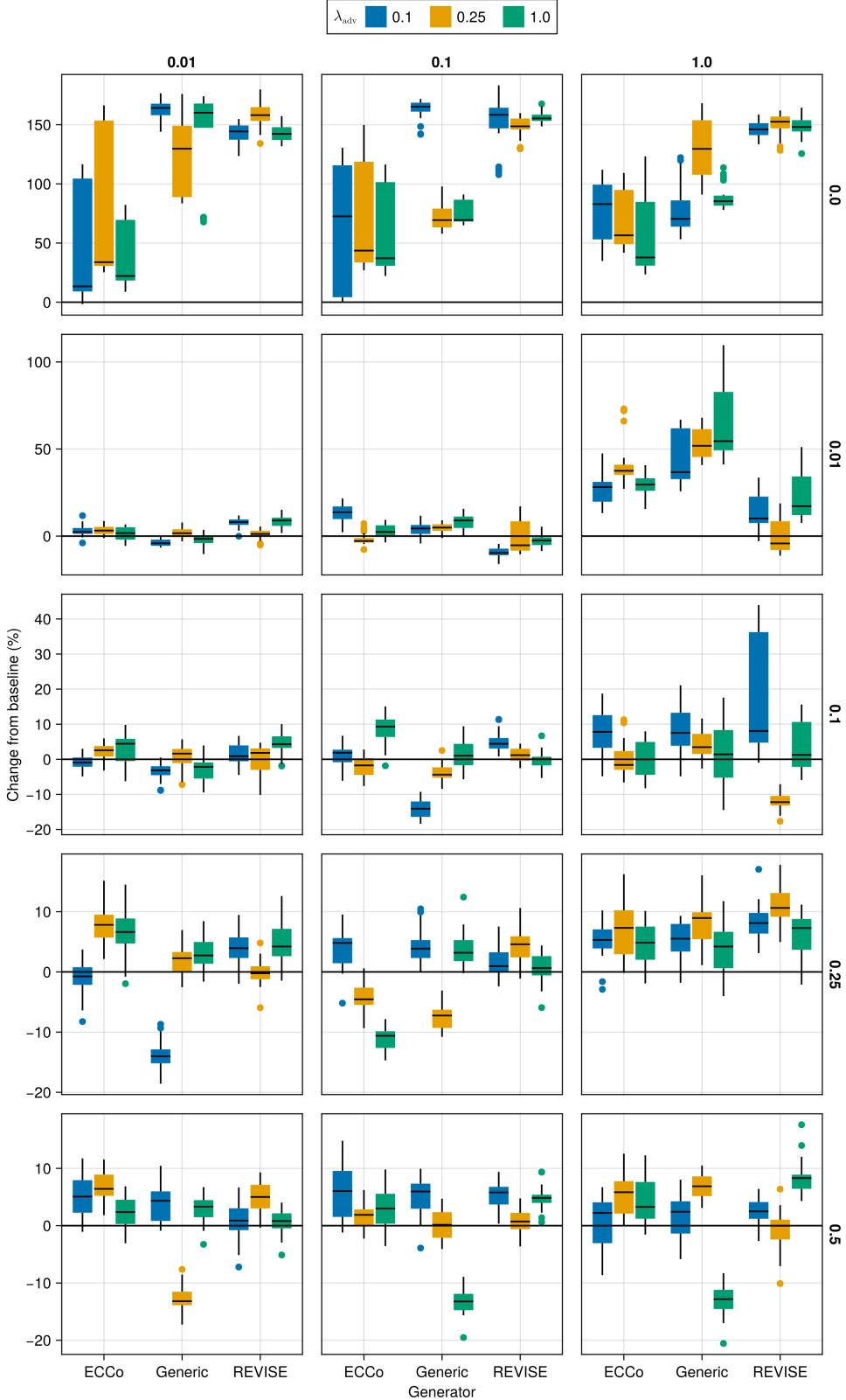


Figure 17: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Linearly Separable.

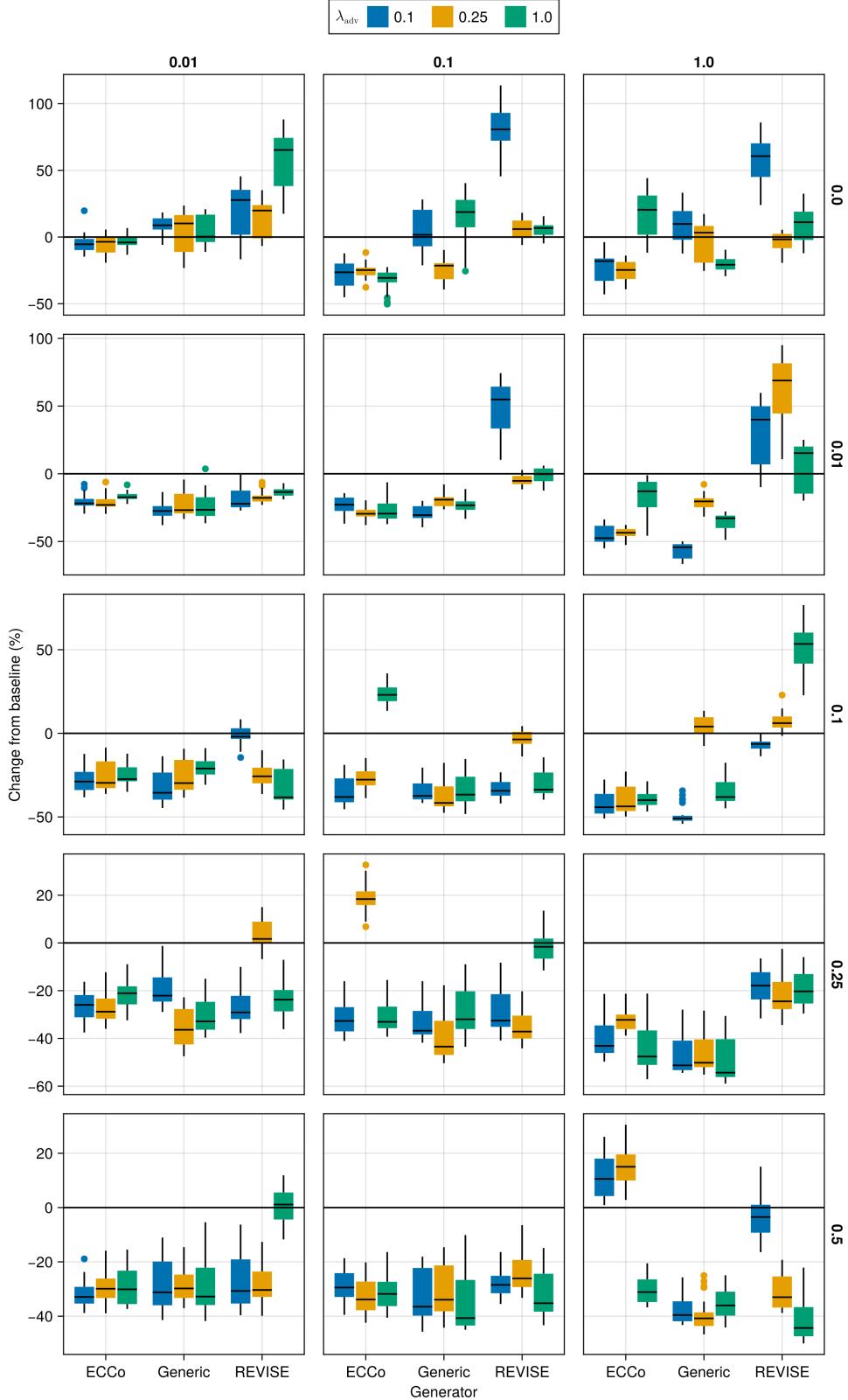


Figure 18: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Moons.

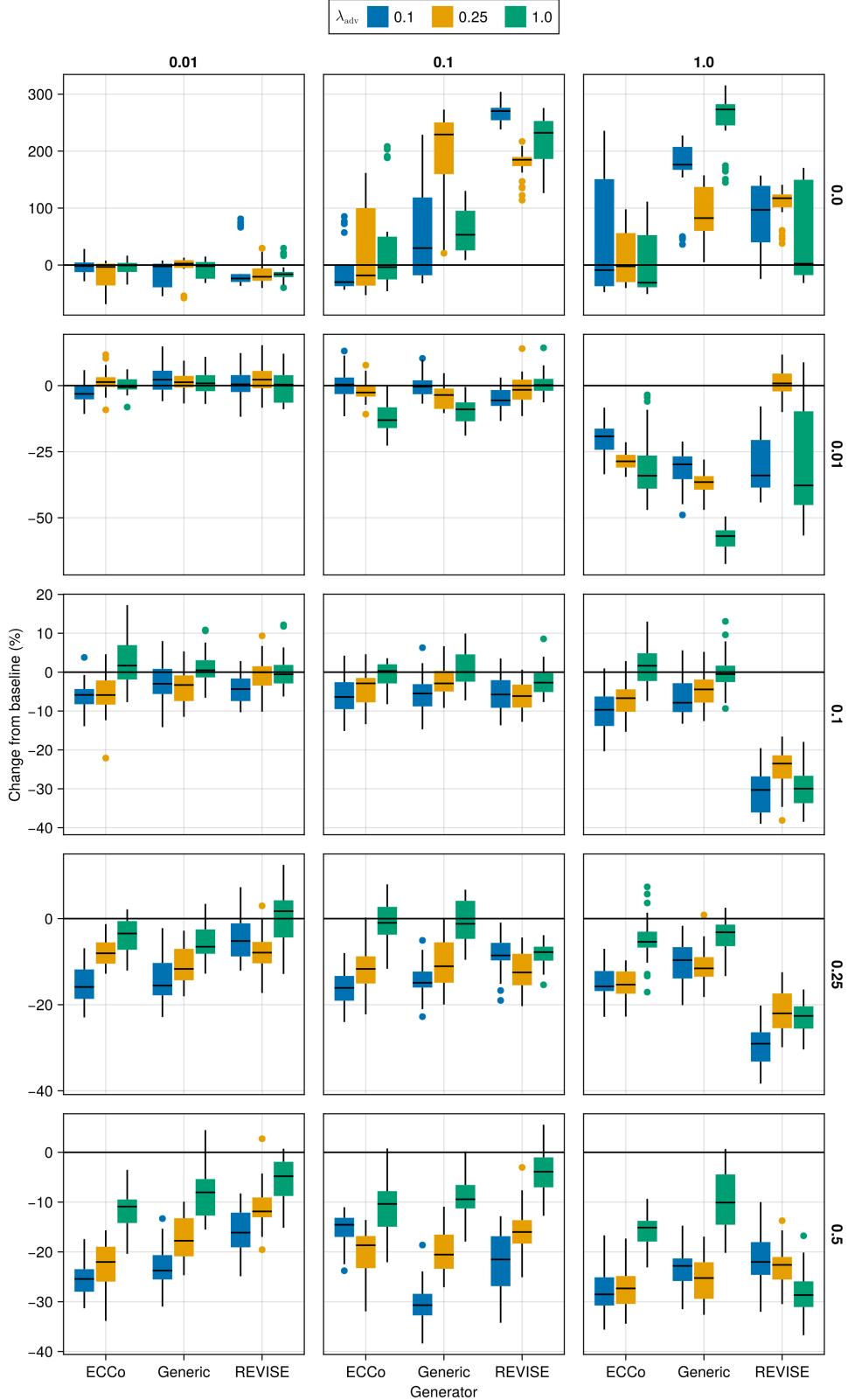


Figure 19: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

Table 5: Predictive performance measures by dataset and objective averaged across training-phase parameters (Note 7) and evaluation-phase parameters (Note 8).

Dataset	Variable	Objective	Mean	Std
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

⁷¹⁴ **D.4.2 Plausibility**

⁷¹⁵ The results with respect to the plausibility measure are shown in Figure 20 to Figure 23.

⁷¹⁶ **D.4.3 Cost**

⁷¹⁷ The results with respect to the cost measure are shown in Figure 24 to Figure 27.

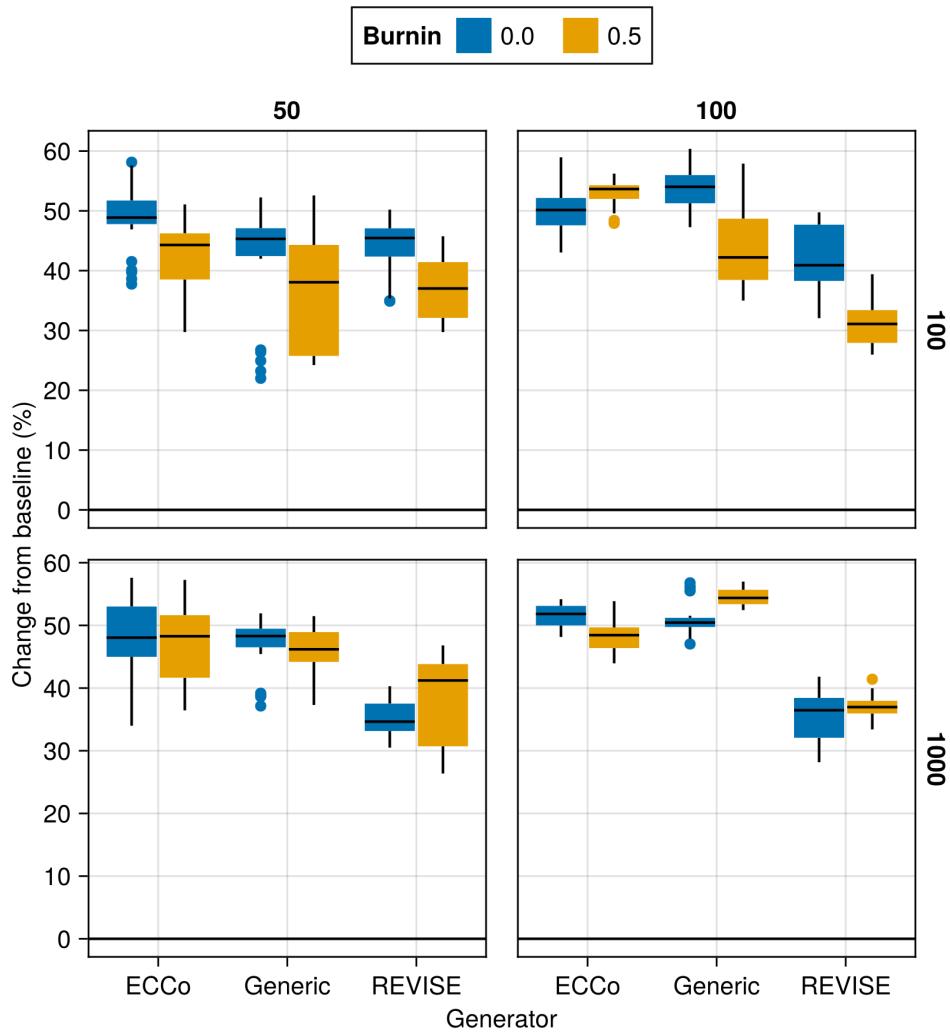


Figure 20: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Circles.

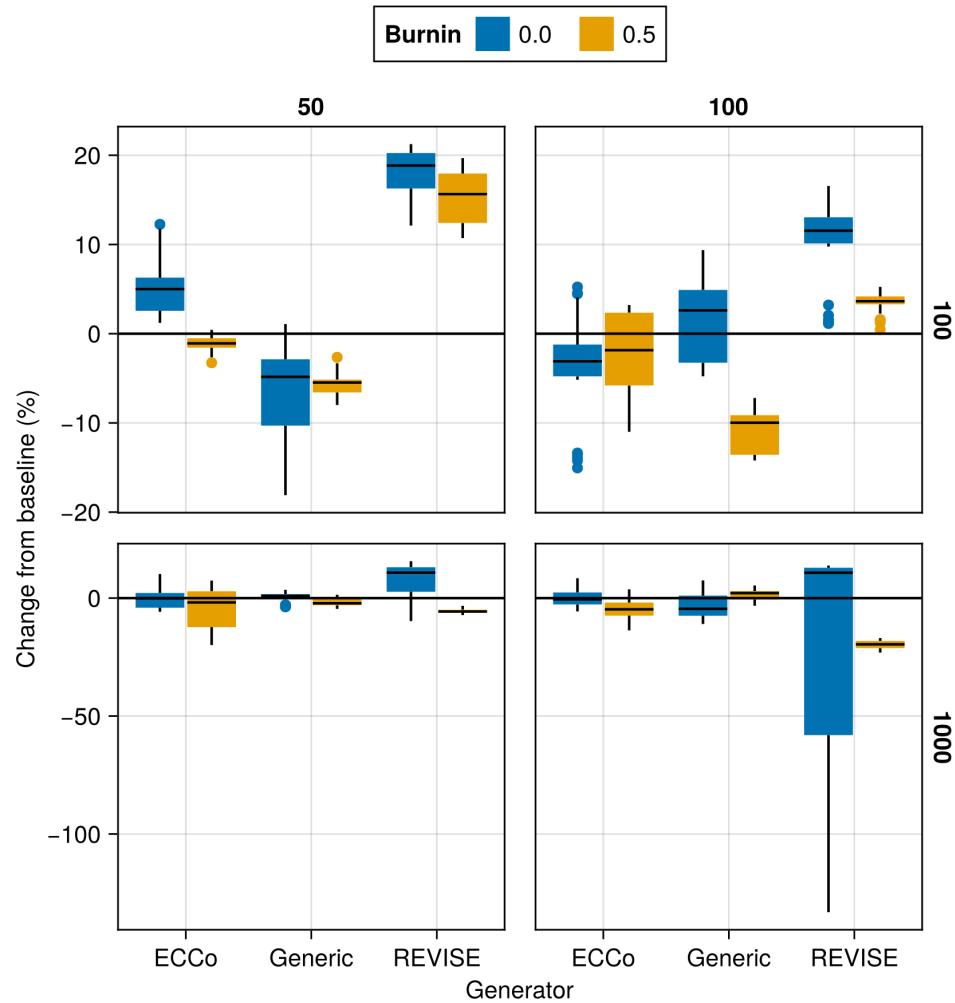


Figure 21: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Linearly Separable.

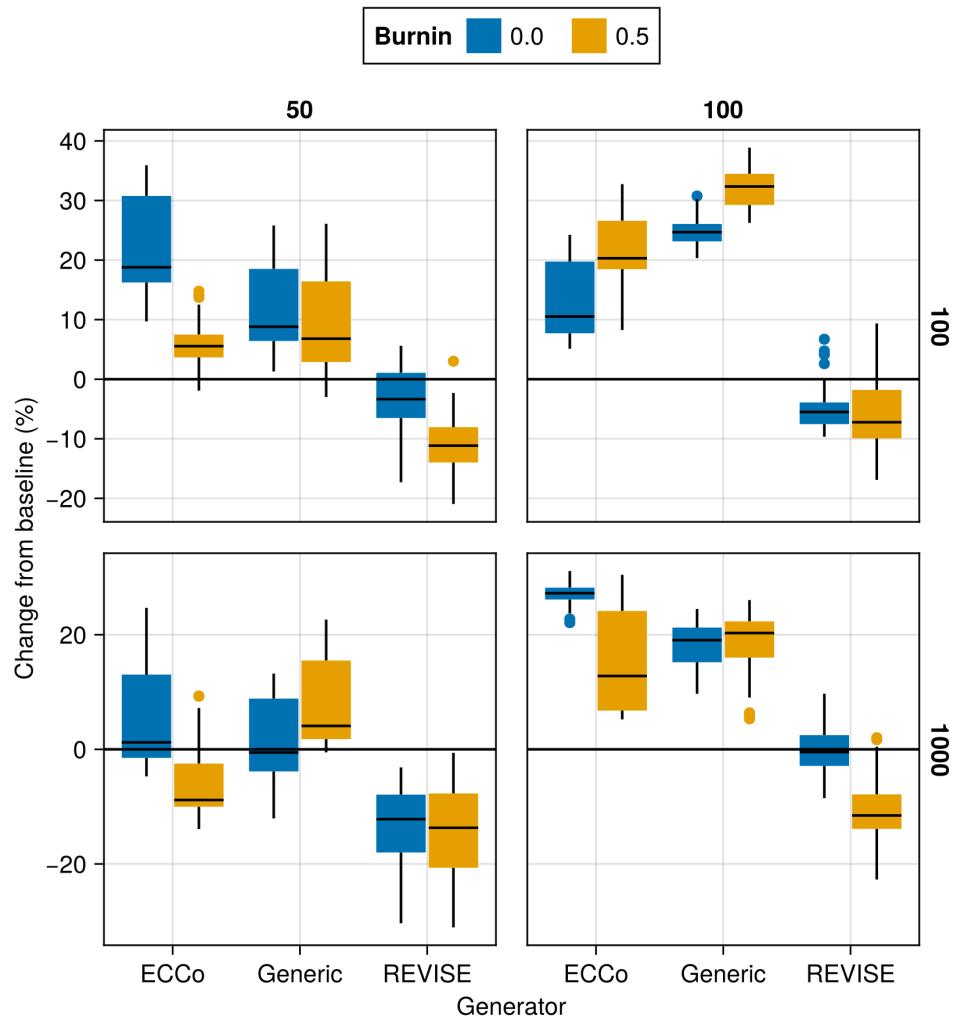


Figure 22: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Moons.

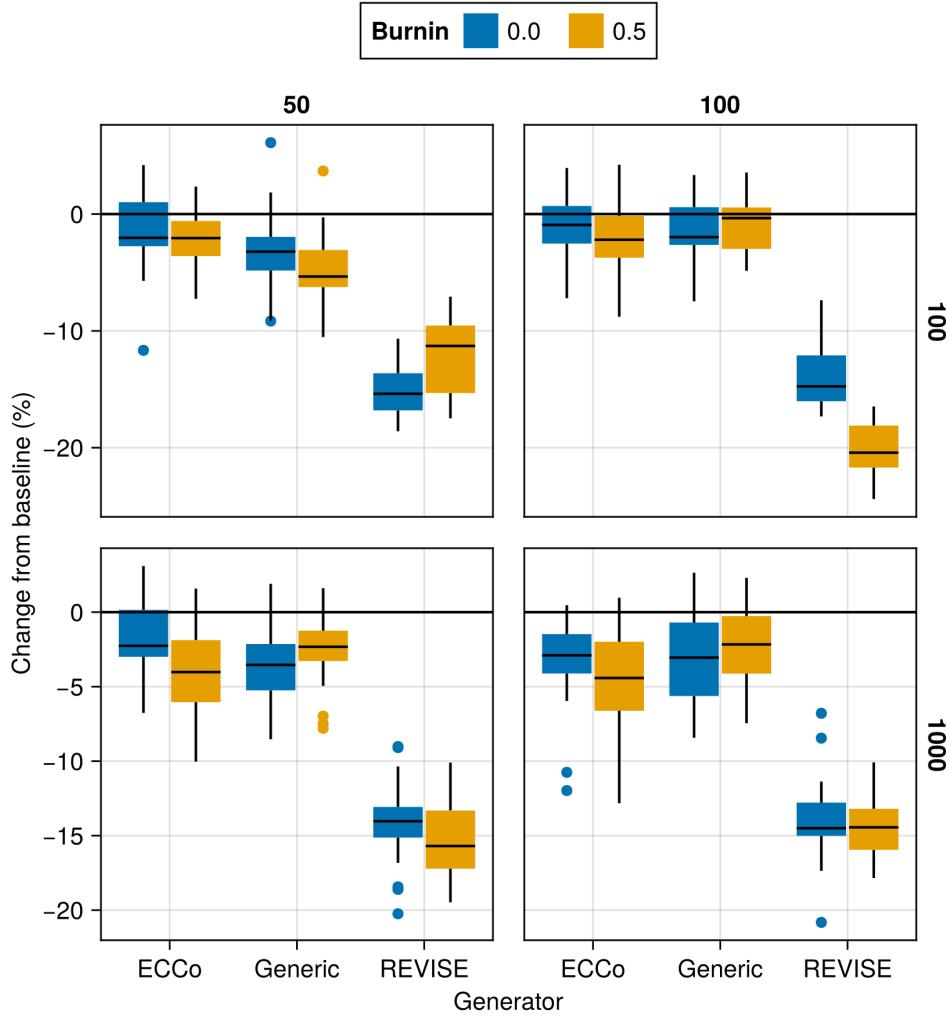


Figure 23: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

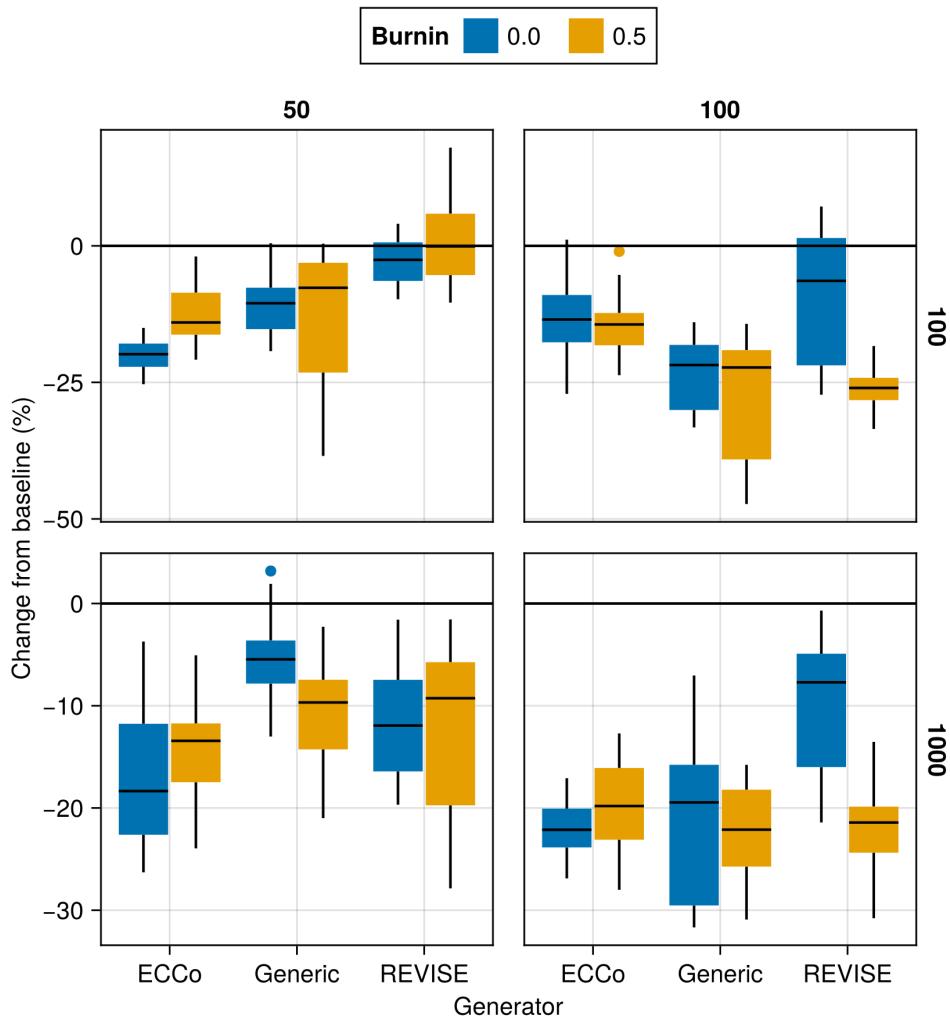


Figure 24: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

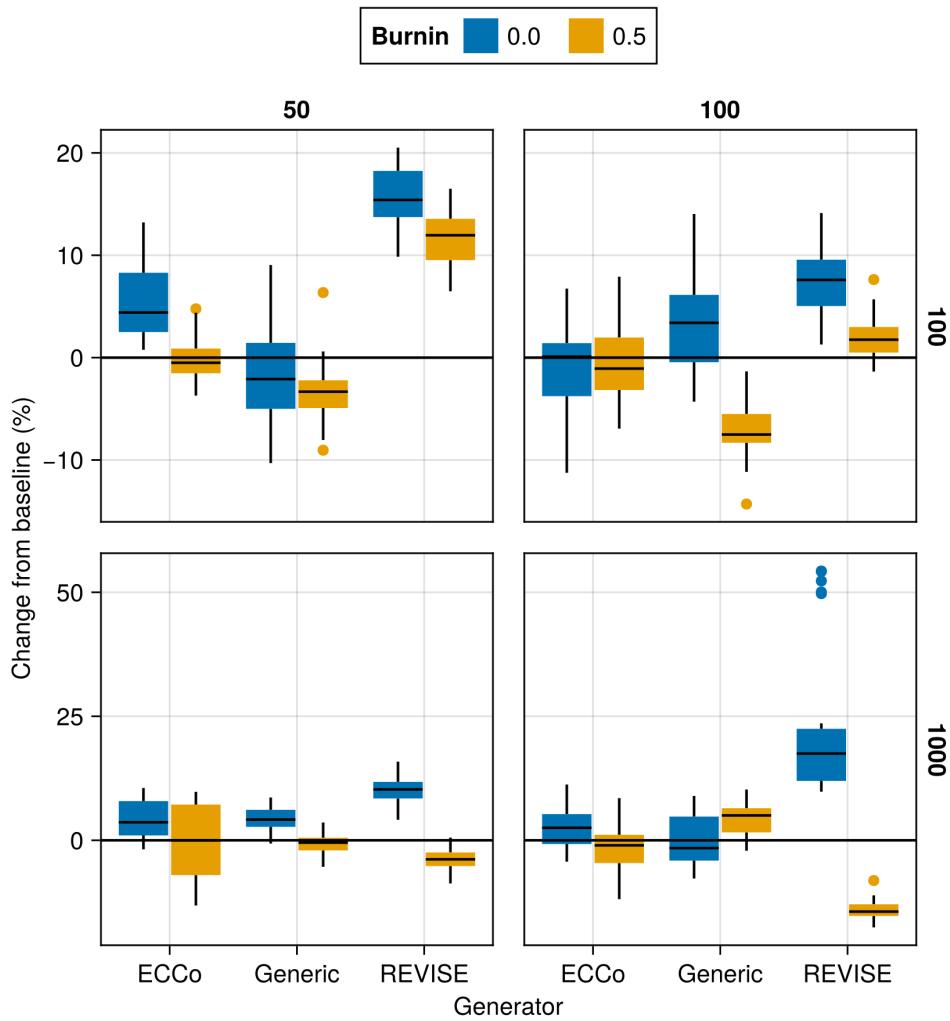


Figure 25: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Linearly Separable.

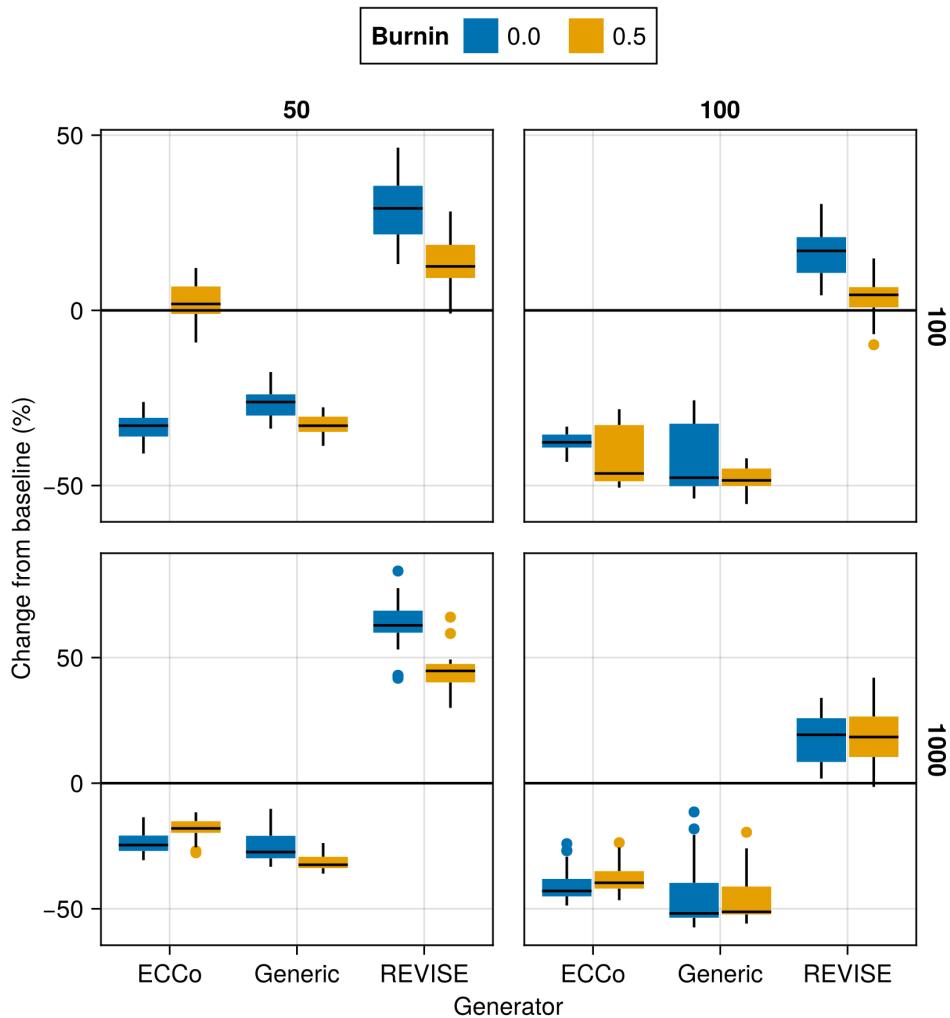


Figure 26: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Moons.

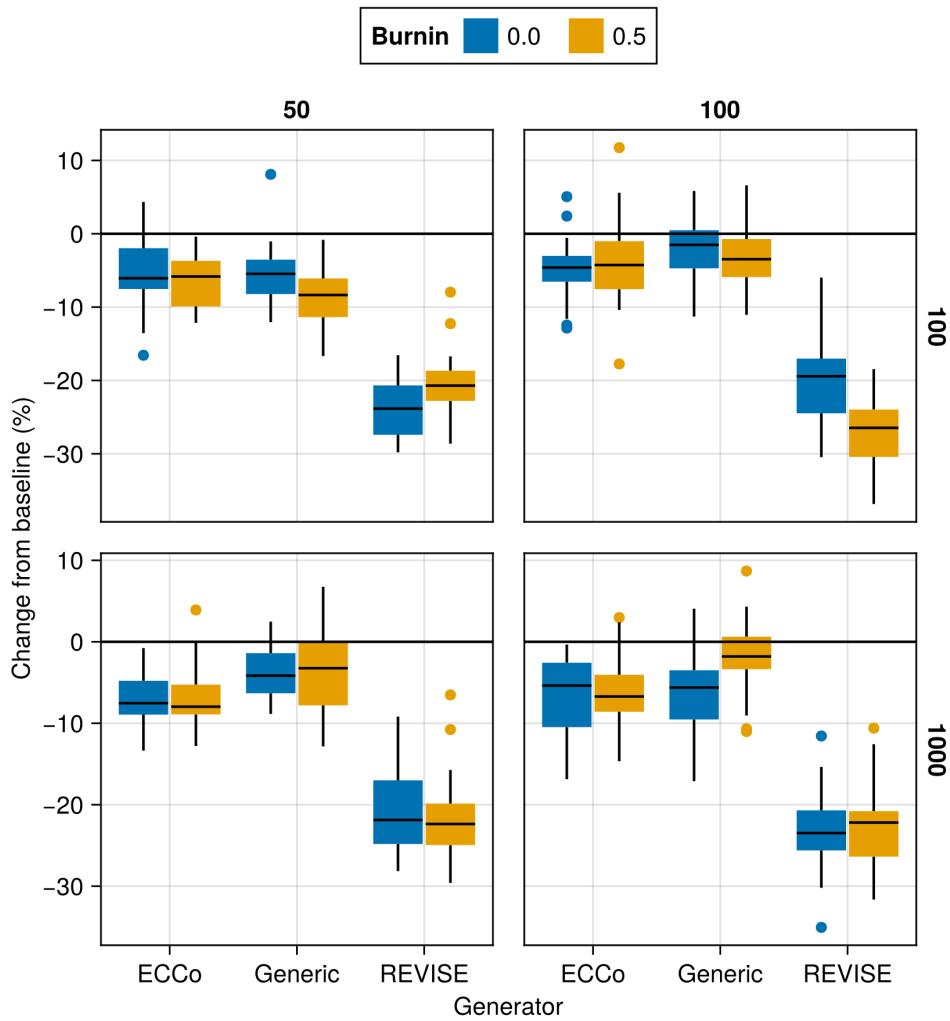


Figure 27: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

718 **Appendix E Tuning Key Parameters**

719 Based on the findings from our initial large grid searches (Section D), we tune selected hyperparameters for all datasets:
 720 namely, the decision threshold τ and the strength of the energy regularization λ_{reg} . The final hyperparameter choices
 721 for each dataset are presented in Table 2 in Section C. Detailed results for each data set are shown in Figure 28 to
 722 Figure 45. From Table 2, we notice that the same decision threshold of $\tau = 0.5$ is optimal for all but one dataset. We
 723 attribute this to the fact that a low decision threshold results in a higher share of mature counterfactuals and hence more
 724 opportunities for the model to learn from examples (Figure 37 to Figure 45). This has played a role in particular for
 725 our real-world tabular datasets and MNIST, which suffered from low levels of maturity for higher decision thresholds.
 726 In cases where maturity is not an issue, as for *Moons*, higher decision thresholds lead to better outcomes, which may
 727 have to do with the fact that the resulting counterfactuals are more faithful to the model. Concerning the regularization
 728 strength, we find somewhat high variation across datasets. Most notably, we find that relatively low levels of regulariza-
 729 tion are optimal for MNIST. We hypothesize that this finding may be attributed to the uniform scaling of all input
 730 features (digits).

731 Finally, to increase the proportion of mature counterfactuals for some datasets, we have also investigated the effect on
 732 the learning rate η for the counterfactual search and even smaller regularization strengths for a fixed decision threshold
 733 of 0.5 (Figure 46 to Figure 51). For the given low decision threshold, we find that the learning rate has no discernable
 734 impact on the proportion of mature counterfactuals (Figure 52 to Figure 57). We do notice, however, that the results
 735 for MNIST are much improved when using a low value λ_{reg} , the strength for the energy regularization: plausibility is
 736 increased by up to ~10% (Figure 50) and the proportion of mature counterfactuals reaches 100%.

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

Package Version (Reproducibility)

Tuning was run using v1.1.3 of TaijaData. The follow-up version v1.1.4 introduced an option to split real-world tabular datasets into train and test set, ensuring that pre-processing steps like standardization is fit on the training set only. If you are rerunning the tuning experiments with a version of TaijaData that is higher than v1.1.3, than for the default parameters specified in the configuration files, you may end up with slightly different results, although we would not expect any changes in terms of qualitative findings. For exact reproducibility, please use v1.1.3.

739

740 **E.1 Key Parameters**

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

Note 9: 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

743

Note 10: Evaluation Phase

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

744

745 **E.1.1 Plausibility**

746 The results with respect to the plausibility measure are shown in Figure 28 to Figure 36.

747 **E.1.2 Proportion of Mature CE**

748 The results with respect to the proportion of mature counterfactuals in each epoch are shown in Figure 37 to Figure 45.

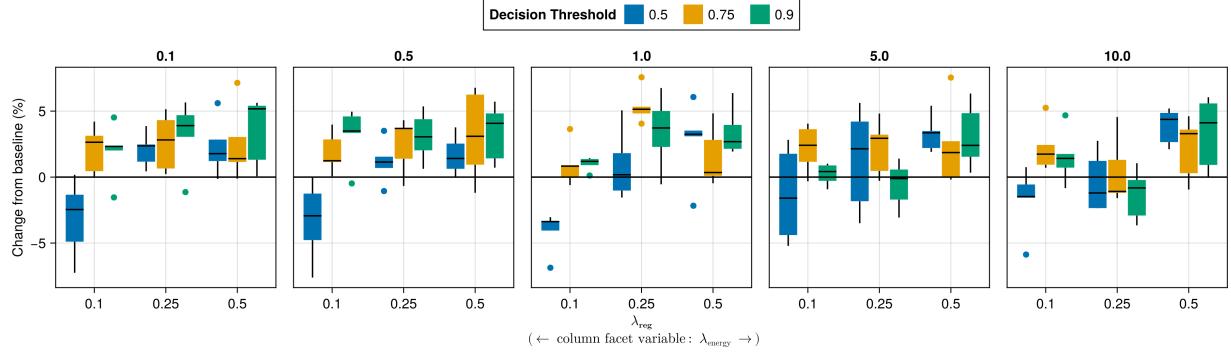


Figure 28: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Adult.

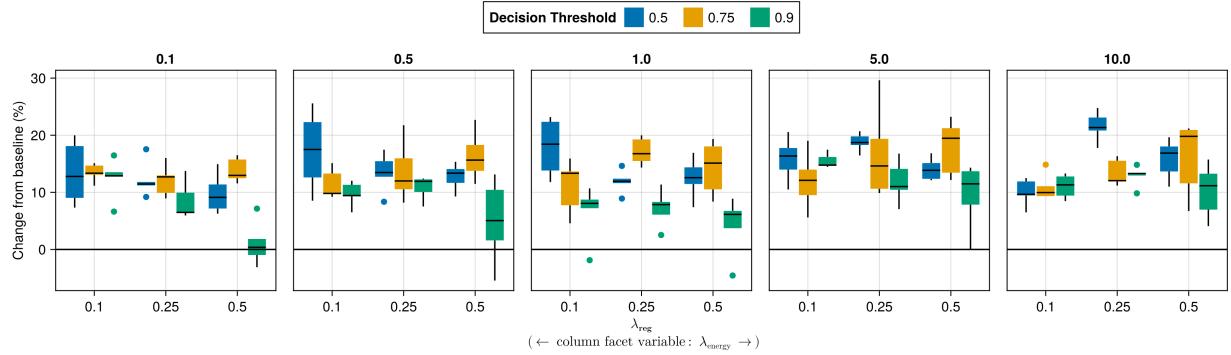


Figure 29: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: California Housing.

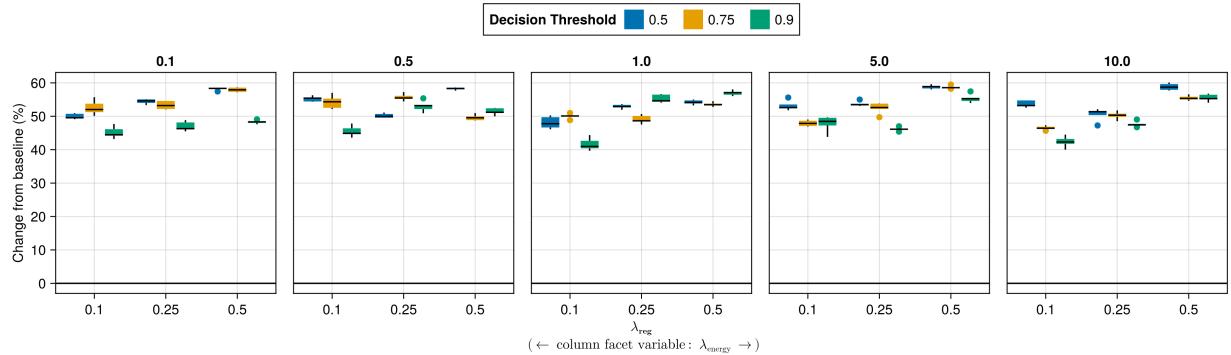


Figure 30: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Circles.

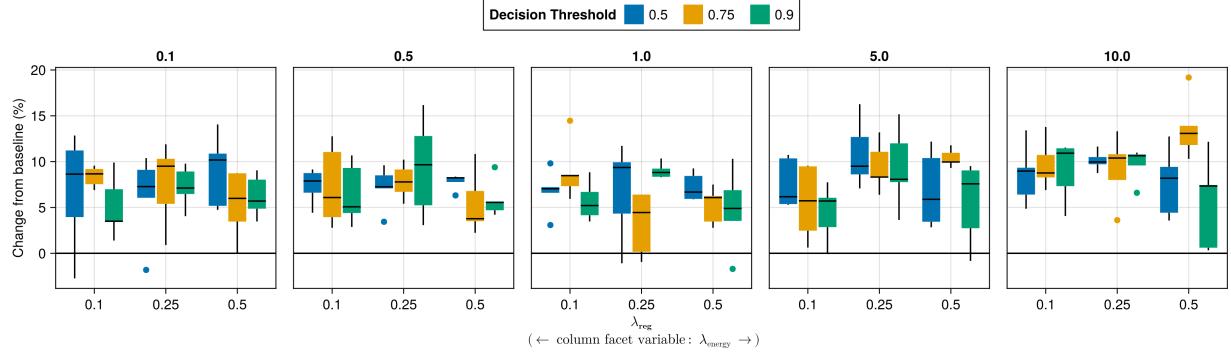


Figure 31: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Credit.

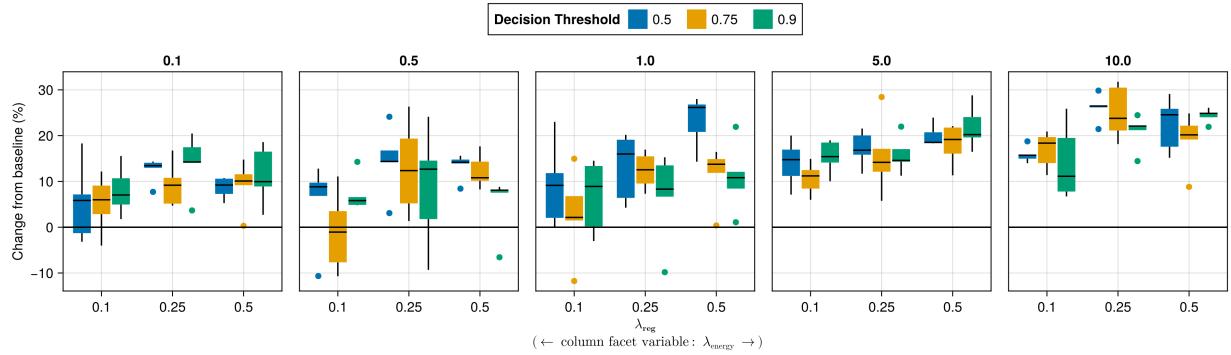


Figure 32: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: GMSC.

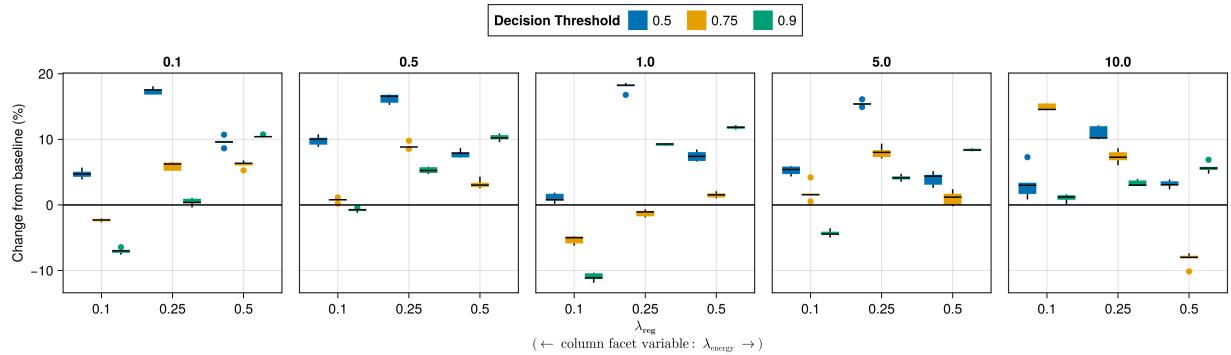


Figure 33: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Linearly Separable.

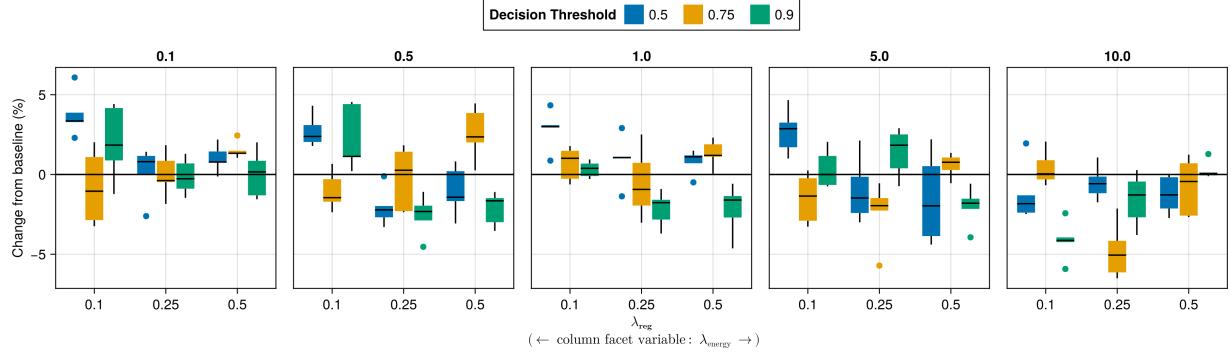


Figure 34: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: MNIST.

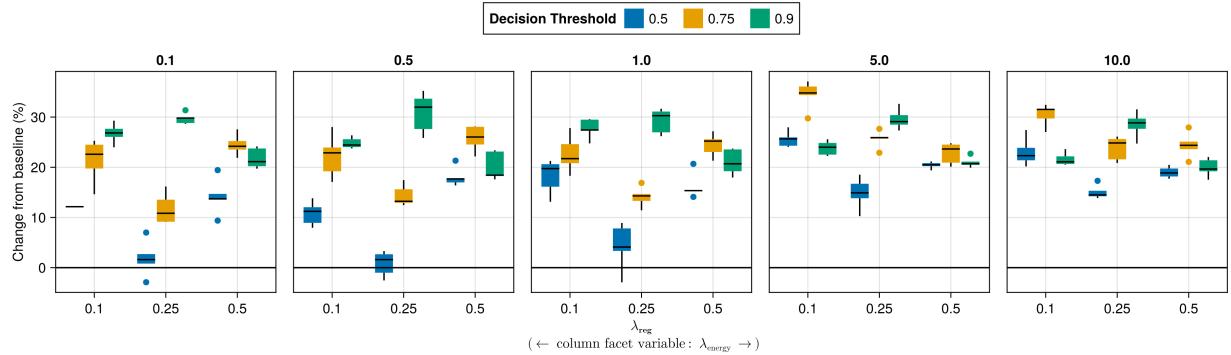


Figure 35: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Moons.

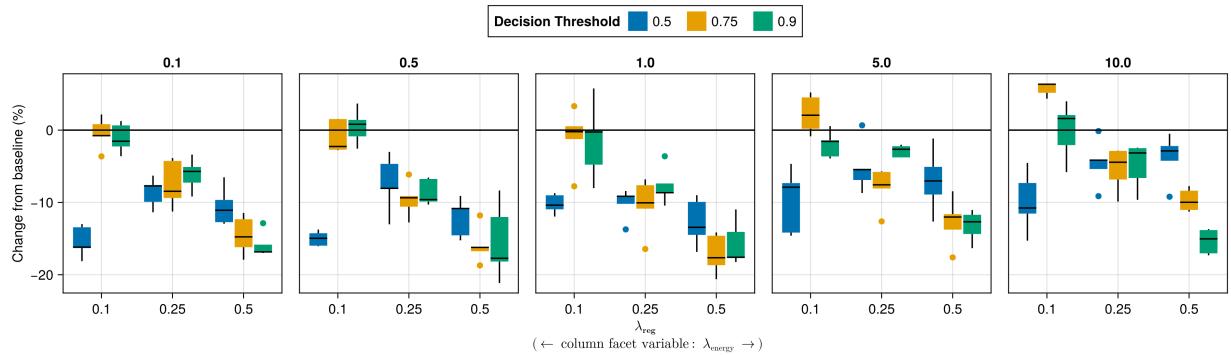


Figure 36: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

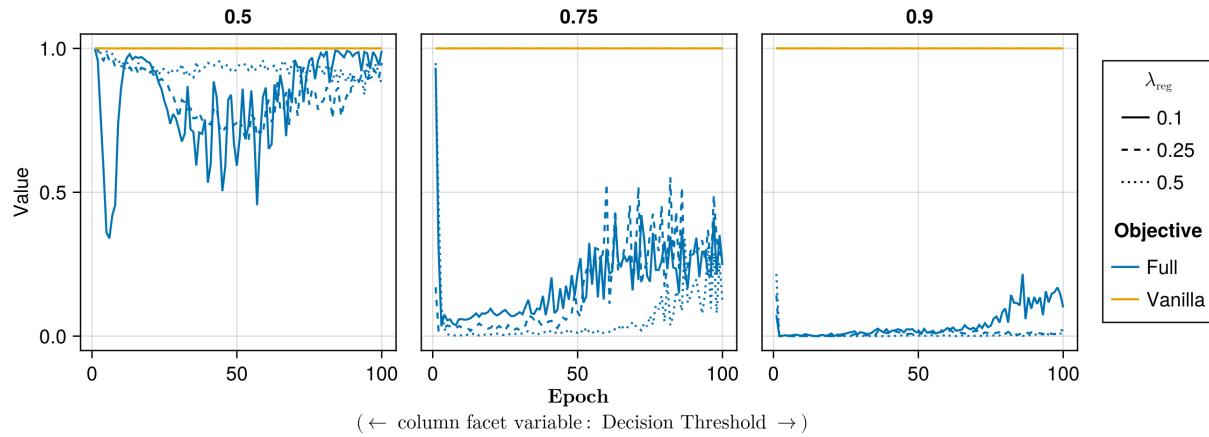


Figure 37: Proportion of mature counterfactuals in each epoch. Data: Adult.

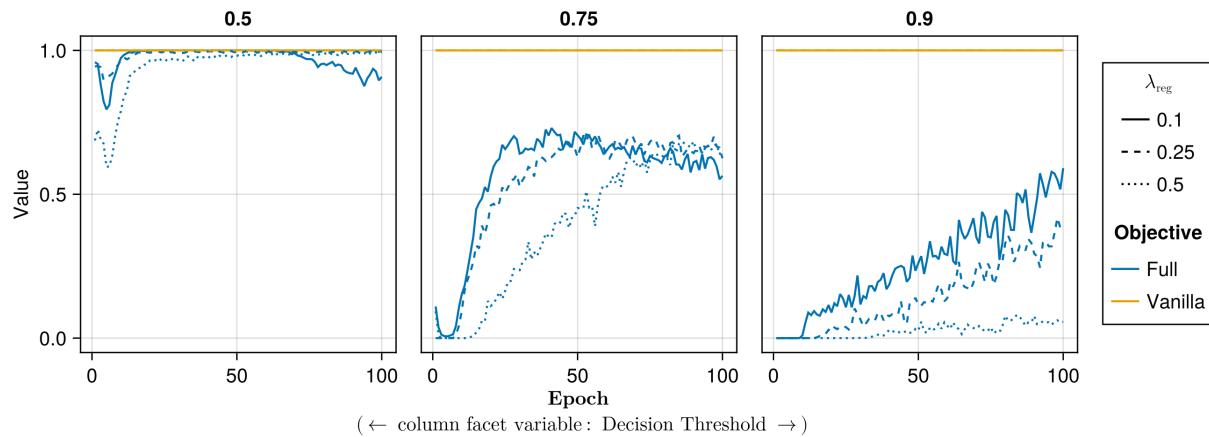


Figure 38: Proportion of mature counterfactuals in each epoch. Data: California Housing.

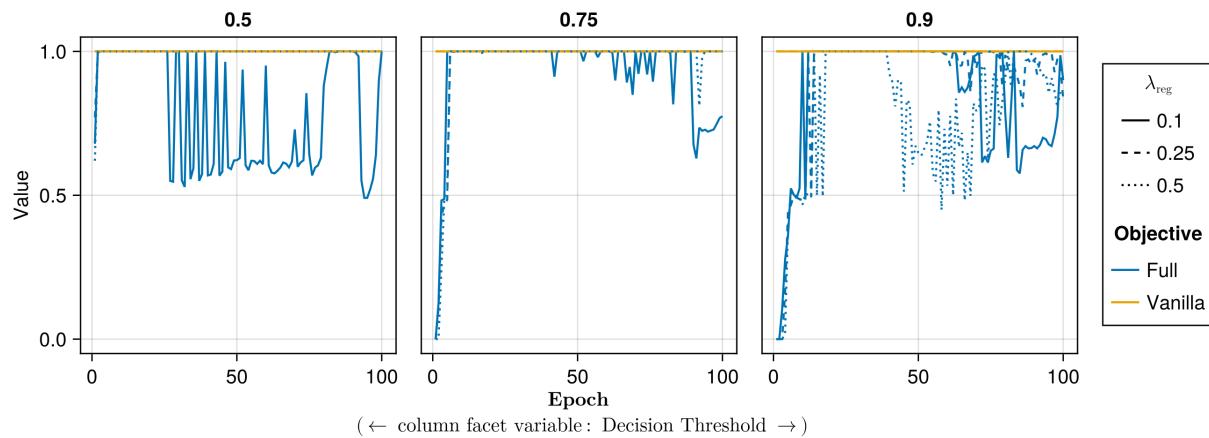


Figure 39: Proportion of mature counterfactuals in each epoch. Data: Circles.

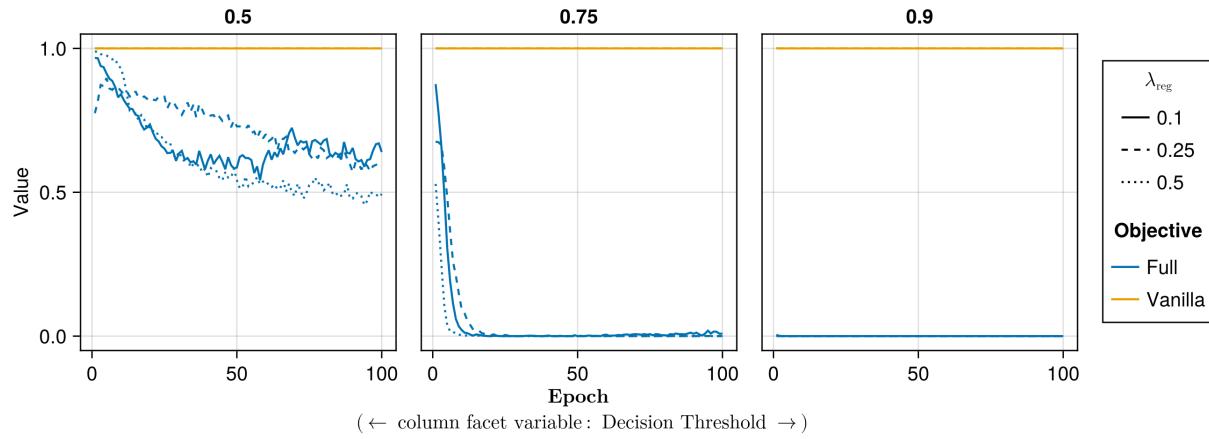


Figure 40: Proportion of mature counterfactuals in each epoch. Data: Credit.

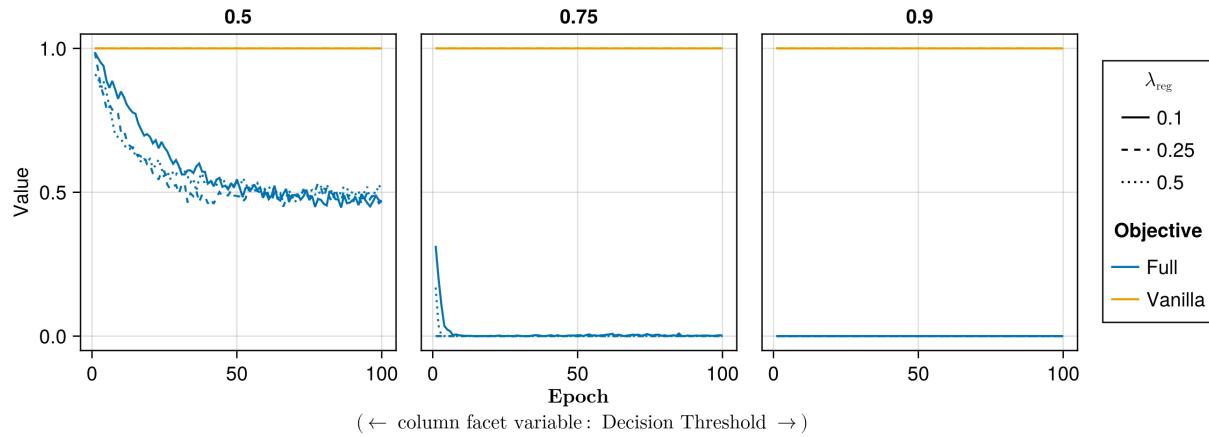


Figure 41: Proportion of mature counterfactuals in each epoch. Data: GMSC.

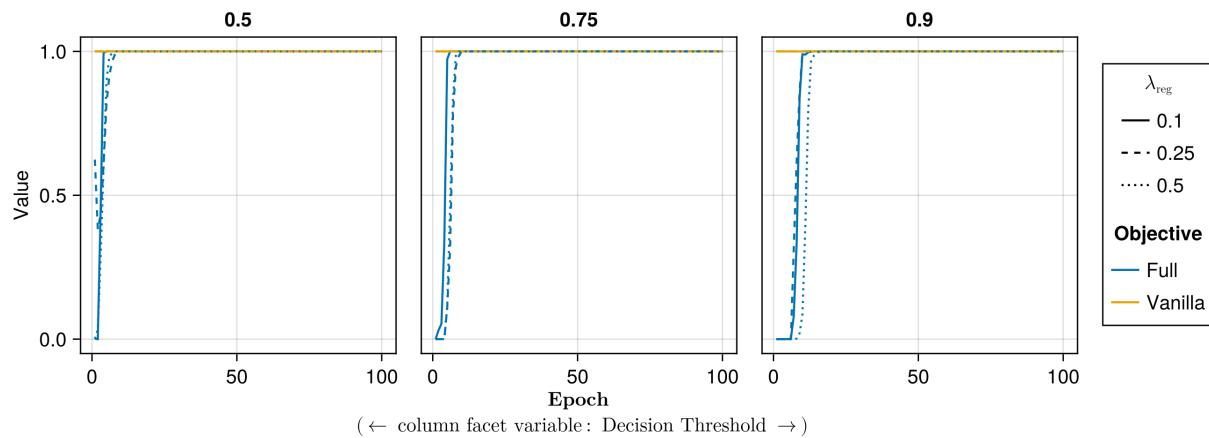


Figure 42: Proportion of mature counterfactuals in each epoch. Data: Linearly Separable.

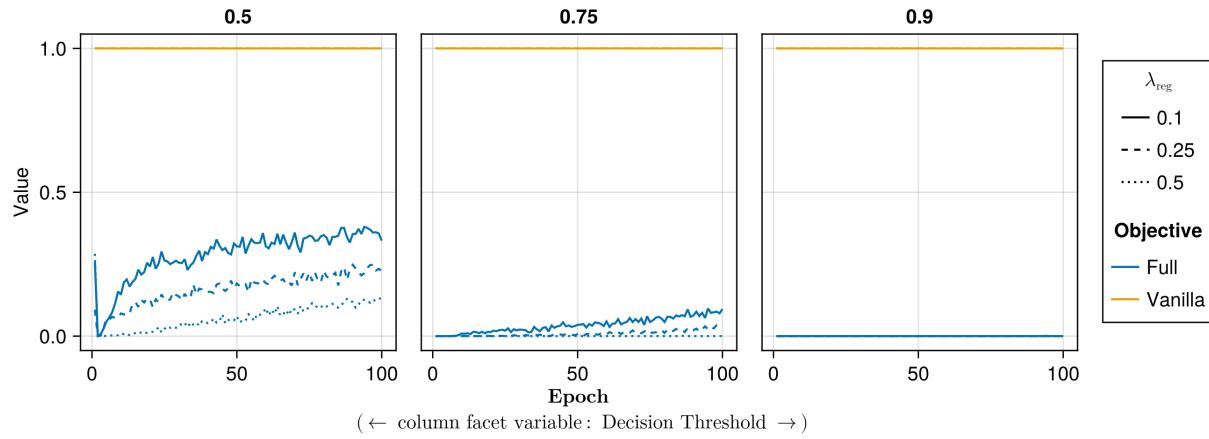


Figure 43: Proportion of mature counterfactuals in each epoch. Data: MNIST.

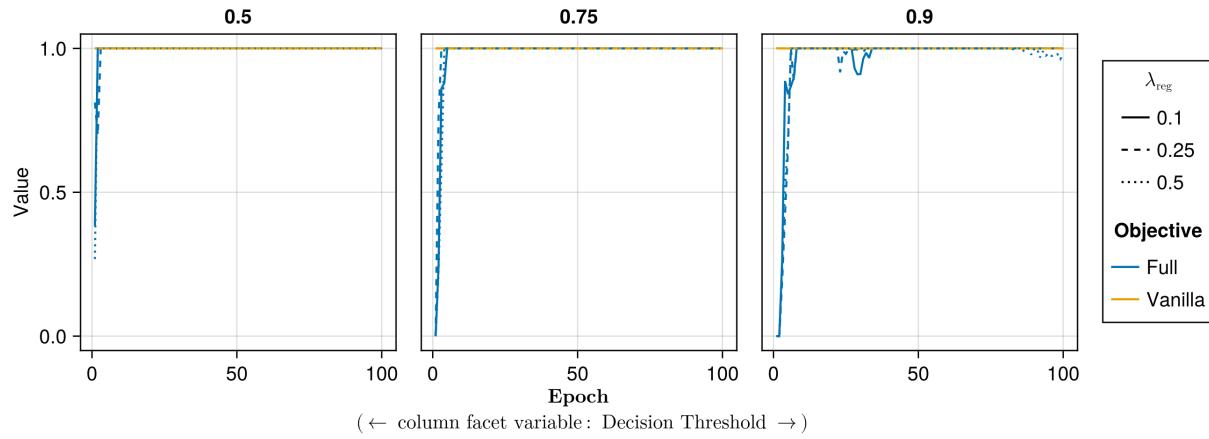


Figure 44: Proportion of mature counterfactuals in each epoch. Data: Moons.

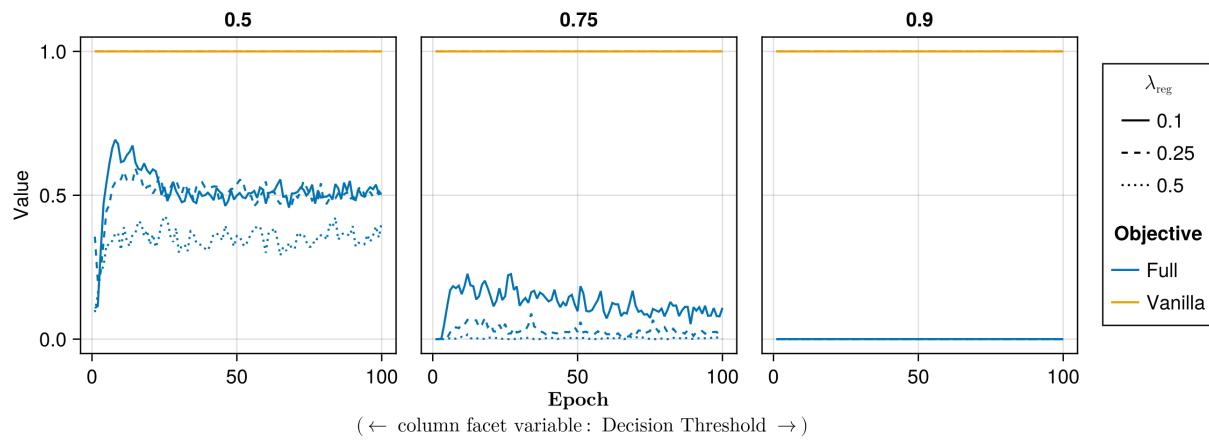


Figure 45: Proportion of mature counterfactuals in each epoch. Data: Overlapping.

749 **E.2 Learning Rate**

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

Note 11: 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

752

Note 12: Evaluation Phase

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

753

754 **E.2.1 Plausibility**

755 The results with respect to the plausibility measure are shown in Figure 46 to Figure 51.

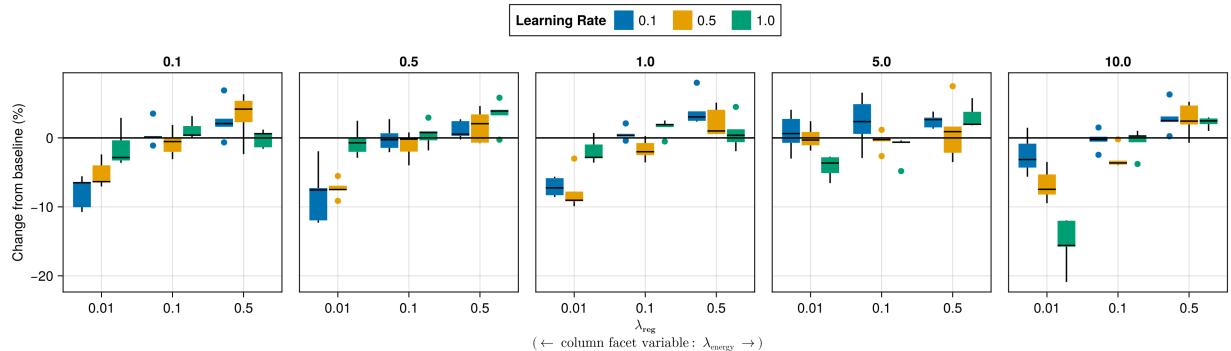


Figure 46: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Adult.

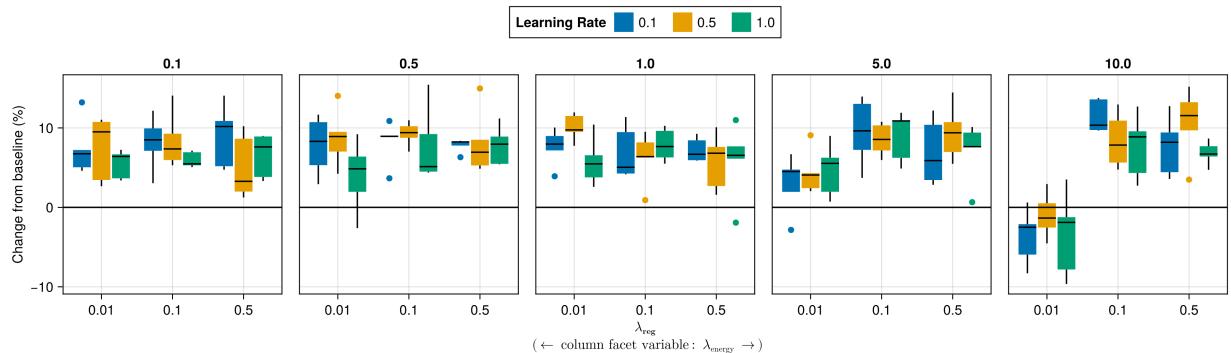


Figure 47: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Credit.

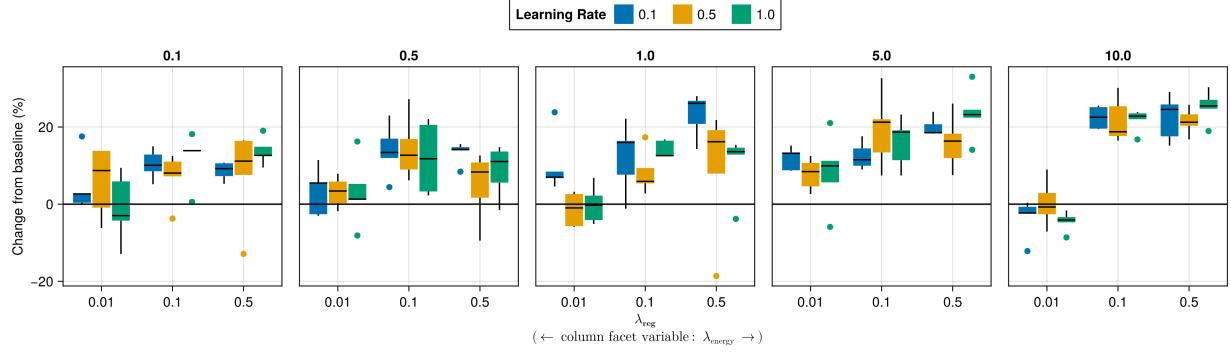


Figure 48: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for $ECCo$). Data: GMSC.

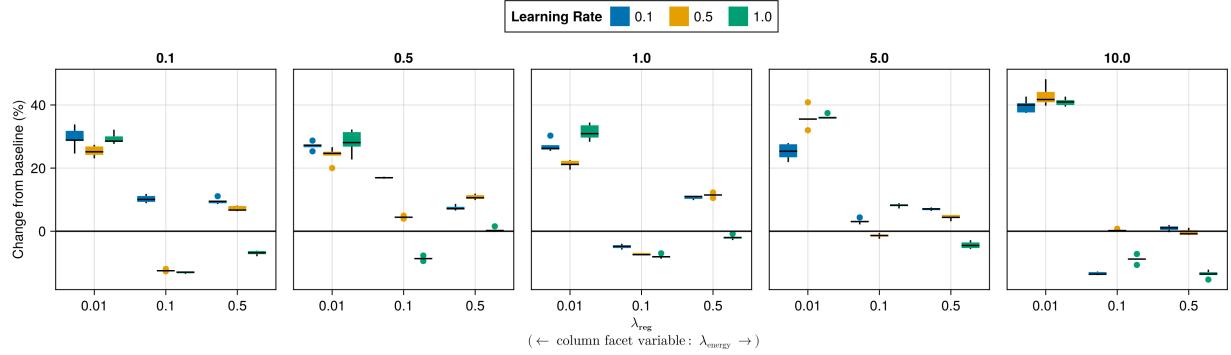


Figure 49: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for $ECCo$). Data: Linearly Separable.

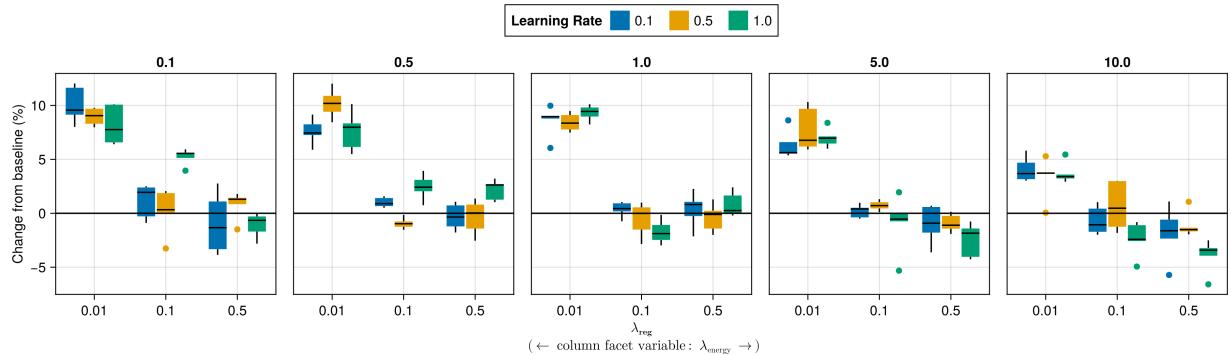


Figure 50: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for $ECCo$). Data: MNIST.

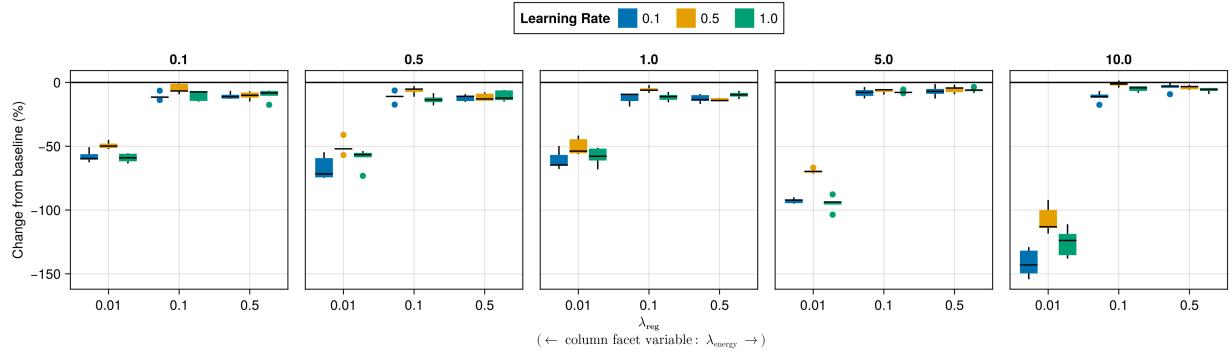


Figure 51: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (Equation 4). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

756 E.2.2 Proportion of Mature CE

757 The results with respect to the proportion of mature counterfactuals in each epoch are shown in Figure 52 to Figure 57.

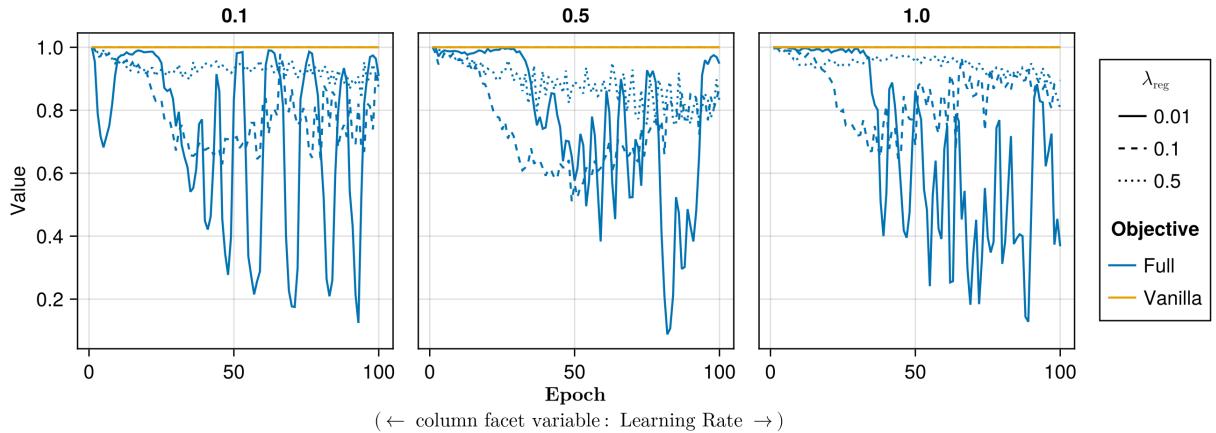


Figure 52: Proportion of mature counterfactuals in each epoch. Data: Adult.

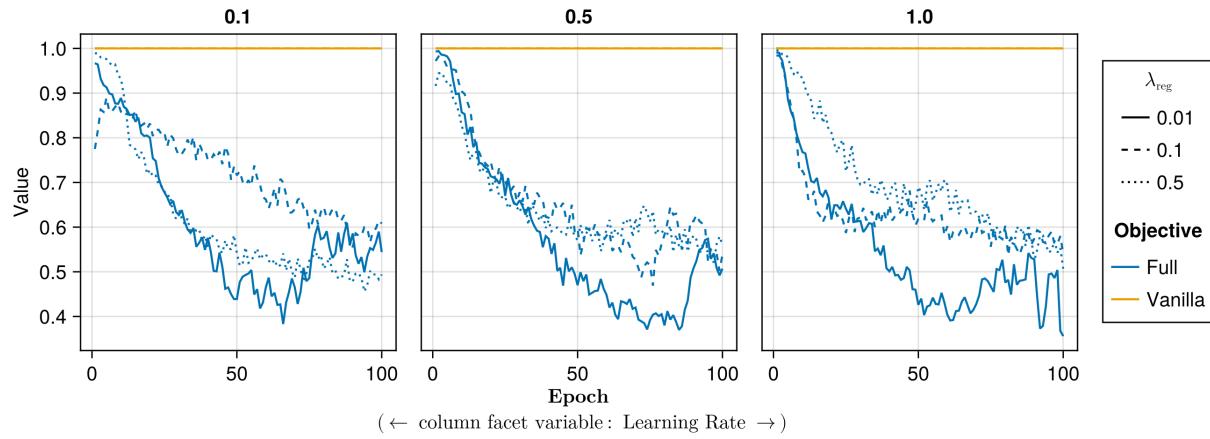


Figure 53: Proportion of mature counterfactuals in each epoch. Data: Credit.

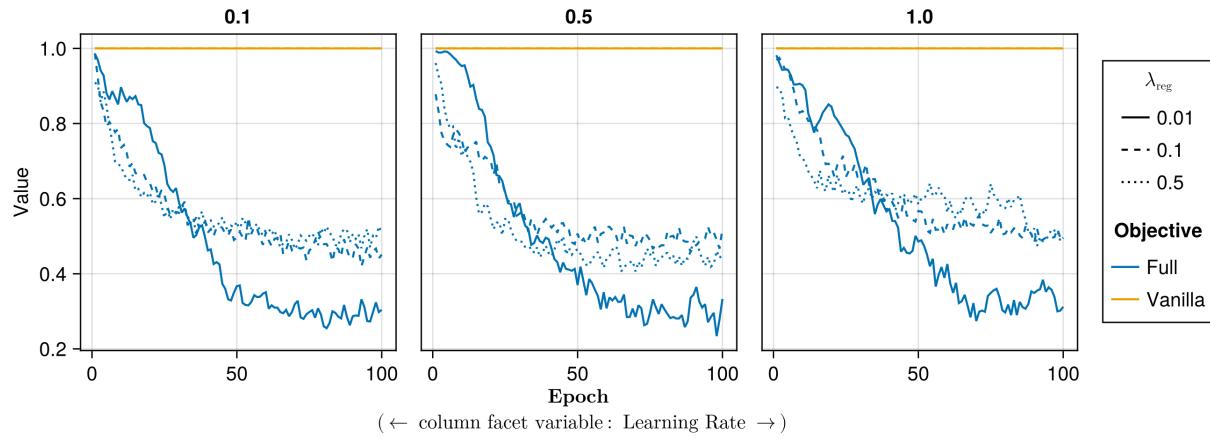


Figure 54: Proportion of mature counterfactuals in each epoch. Data: GMSC.

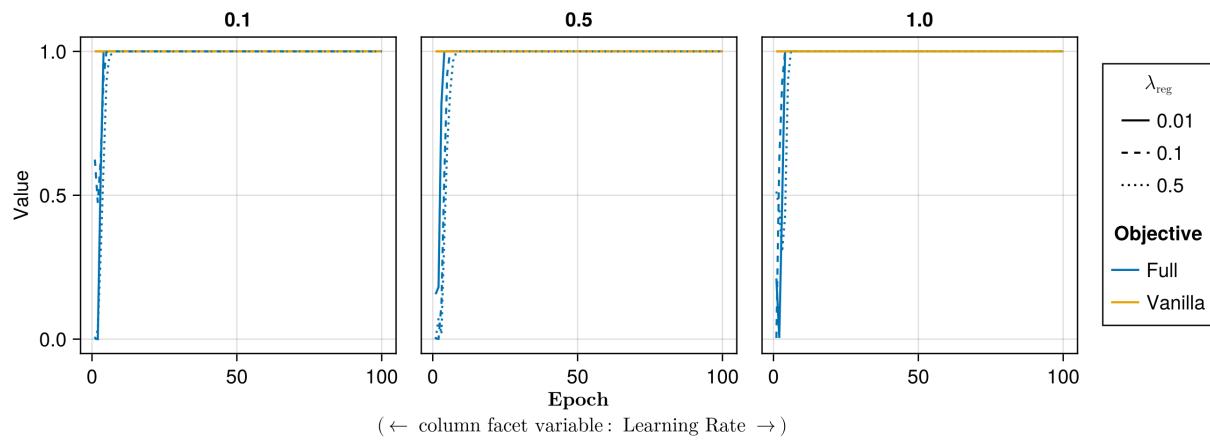


Figure 55: Proportion of mature counterfactuals in each epoch. Data: Linearly Separable.

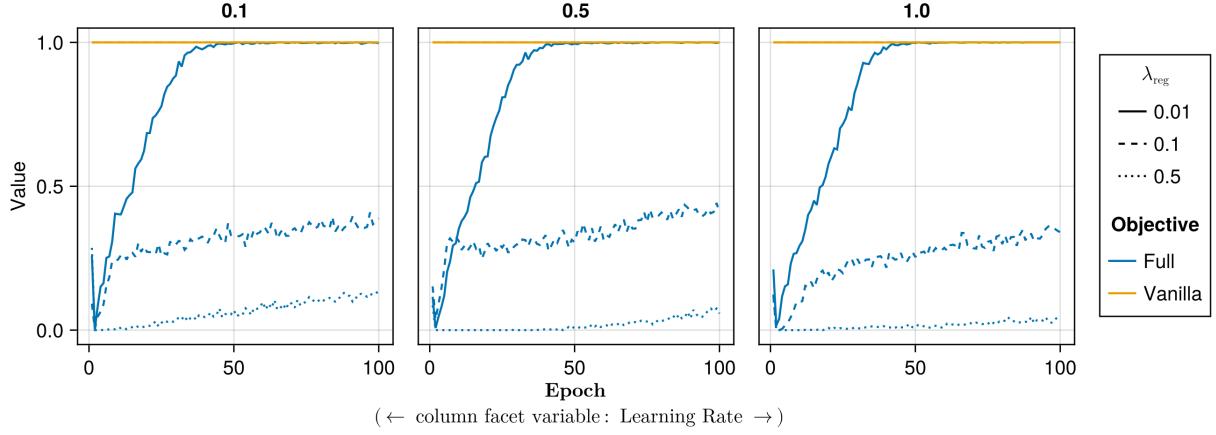


Figure 56: Proportion of mature counterfactuals in each epoch. Data: MNIST.

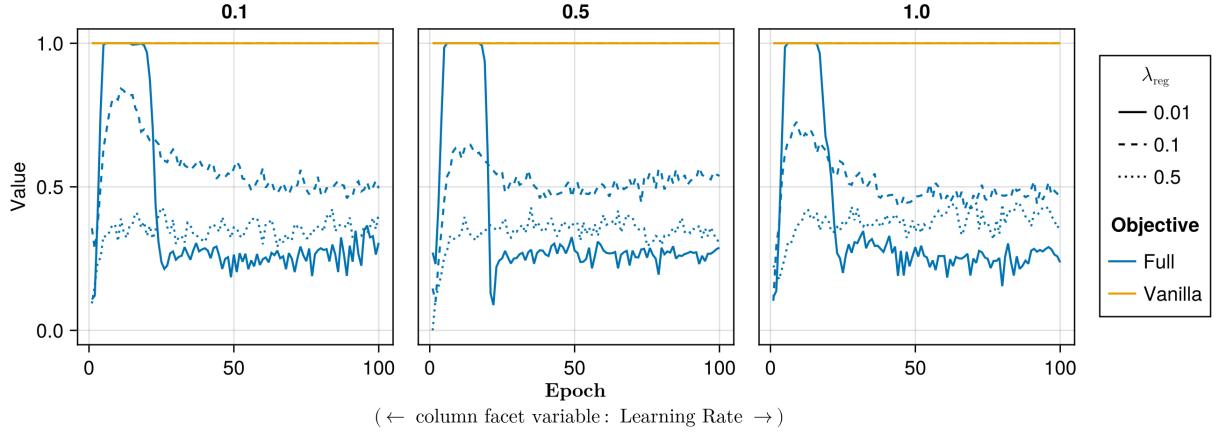


Figure 57: Proportion of mature counterfactuals in each epoch. Data: Overlapping.

758 Appendix F Computation Details

759 F.1 Hardware

760 We performed our experiments on a high-performance cluster. Details about the cluster will be disclosed upon publication to avoid revealing information that might interfere with the double-blind review process. Since our experiments 761 involve highly parallel tasks and rather small models by today’s standard, we have relied on distributed computing 762 across multiple central processing units (CPU). Graphical processing units (GPU) were not required. 763

764 F.1.1 Grid Searches

765 Model training for the largest grid searches with 270 unique parameter combinations was parallelized across 34 CPUs 766 with 2GB memory each. The time to completion varied by dataset for reasons discussed in Section 5: 0h49m (*Moons*), 767 1h4m (*Linearly Separable*), 1h49m (*Circles*), 3h52m (*Overlapping*). Model evaluations for large grid searches were 768 parallelized across 20 CPUs with 3GB memory each. Evaluations for all data sets took less than one hour (<1h) to 769 complete. 770

F.1.2 Tuning

771 For tuning of selected hyperparameters, we distributed the task of generating counterfactuals during training across 40 772 CPUs with 2GB memory each for all tabular datasets. Except for the *Adult* dataset, all training runs were completed 773 in less than half an hour (<0h30m). The *Adult* dataset took around 0h35m to complete. Evaluations across 20 CPUs 774 with 3GB memory each generally took less than 0h30m to complete. For *MNIST*, we relied on 100 CPUs with 2GB 775 memory each. For the *MLP*, training of all models could be completed in 1h30m, while the evaluation across 20 CPUs

776 (6GB memory) took 4h12m. For the *CNN*, training of all models took ~8h, with conventionally trained models taking
777 ~0h15m each and model with CT taking ~0h30m-0h45m each.

778 **F.2 Software**

779 All computations were performed in the Julia Programming Language ([Bezanson et al. 2017](#)). We have developed a
780 package for counterfactual training that leverages and extends the functionality provided by several existing packages,
781 most notably [CounterfactualExplanations.jl](#) ([Altmeyer, Deursen, and Liem 2023](#)) and the [Flux.jl](#) library for deep
782 learning ([Michael Innes et al. 2018; Mike Innes 2018](#)). For data-wrangling and presentation-ready tables we relied on
783 [DataFrames.jl](#) ([Bouchet-Valat and Kamiski 2023](#)) and [PrettyTables.jl](#) ([Chagas et al. 2024](#)), respectively. For plots and
784 visualizations we used both [Plots.jl](#) ([Christ et al. 2023](#)) and [Makie.jl](#) ([Danisch and Krumbiegel 2021](#)), in particular
785 [AlgebraOfGraphics.jl](#). To distribute computational tasks across multiple processors, we have relied on [MPI.jl](#) ([Byrne,
786 Wilcox, and Churavy 2021](#)).