
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 14, 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 plausible with respect to the underlying data and actionable with respect to the stakeholder requirements. 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 maintaining high predictive performance.

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

1 Introduction

Today's prominence of artificial intelligence (AI) has largely been driven by **representation learning**: instead of relying on features and rules that are carefully hand-crafted by humans, modern machine learning (ML) models are tasked

18 with learning representations directly from data, guided by narrow objectives such as predictive accuracy ([I. Goodfellow, Bengio, and Courville 2016](#)). Modern advances in computing have made it possible to provide such models
 19 with ever-growing degrees of freedom to achieve that task, which frequently allows them to outperform traditionally
 20 more parsimonious models. Unfortunately, in doing so, models learn increasingly complex and highly sensitive
 21 representations that humans can no longer easily interpret.

22

23 The trend towards complexity for the sake of performance has come under serious scrutiny in recent years. At the
 24 very cusp of the deep learning (DL) revolution, Szegedy et al. ([2013](#)) showed that artificial neural networks (ANN)
 25 are sensitive to adversarial examples (AEs): perturbed versions of data instances that yield vastly different model
 26 predictions despite being “imperceptible” in that they are semantically indifferent from their factual counterparts.
 27 Even though some partially effective mitigation strategies have been proposed—most notably **adversarial training** ([I.
 28 J. Goodfellow, Shlens, and Szegedy 2014](#))—truly robust deep learning remains unattainable even for models that are
 29 considered “shallow” by today’s standards ([Kolter 2023](#)).

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

39 The first challenge has attracted an abundance of research on **explainable AI** (XAI), a paradigm that focuses on the
 40 development of tools to derive (post-hoc) explanations from complex model representations. Such explanations should
 41 mitigate a scenario in which practitioners deploy opaque models and blindly rely on their predictions. On countless
 42 occasions, this has happened in practice and caused real harms to people who were adversely and unfairly affected
 43 by automated decision-making (ADM) systems involving opaque models ([O’Neil 2016; McGregor 2021](#)). Effective
 44 XAI tools can aid us in monitoring models and providing recourse to individuals to turn negative outcomes (e.g.,
 45 “loan application rejected”) into positive ones (e.g., “application accepted”). Our work builds upon **counterfactual
 46 explanations** (CE) proposed by Wachter, Mittelstadt, and Russell ([2017](#)) as an effective approach to achieve this goal.
 47 CEs prescribe minimal changes for factual inputs that, if implemented, would prompt some fitted model to produce a
 48 desired output.

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

58 **Specifically, we propose counterfactual training (CT):** a novel training regime that aligns learned representations
 59 with plausible explanations compliant with user requirements. The remainder of this paper is structured as follows.
 60 Section [2](#) presents related work, focusing on the link between adversarial examples and counterfactual explanations.
 61 Then follow our main contributions:

- 62 1. In Section [3](#), we introduce our methodological framework and show theoretically that it can be used to
 63 enforce global actionability constraints.
- 64
- 65 2. In Section [4](#), through extensive experiments we demonstrate that CT substantially improves explainability
 66 without sacrificing predictive performance.

67 We discuss the challenges in Section [5](#) and conclude in Section [6](#) that CT is a promising approach towards making
 68 opaque models more trustworthy.

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

69 2 Related Literature

70 To the best of our knowledge, the proposed framework for counterfactual training represents the first attempt to use
 71 counterfactual explanations during training to improve model explainability. In high-level terms, we define model
 72 explainability as the extent to which valid explanations derived for an opaque model are also deemed plausible with
 73 respect to the underlying data and (global) actionability constraints. To make the desiderata for our framework more
 74 concrete, we follow Augustin, Meinke, and Hein (2020) in tying the concept of explainability to the quality of CEs that
 75 can be generated for a given model. The authors show that CEs—understood as minimal input perturbations that yield
 76 some desired model prediction—tend to be more meaningful if the underlying model is more robust to adversarial
 77 examples. We can make intuitive sense of this finding when looking at adversarial training (AT) through the lens of
 78 representation learning with high degrees of freedom. As argued before, learned representations may be sensitive to
 79 producing implausible explanations and mispredicting for worst-case counterfactuals (i.e., AEs). Thus, by inducing
 80 models to “unlearn” susceptibility to such examples, AT can effectively remove implausible explanations from the
 81 solution space.

82 2.1 Adversarial Examples are Counterfactual Explanations

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

93 In the natural language processing (NLP) domain, counterfactuals have similarly been used to improve models through
 94 data augmentation. Wu et al. (2021) propose *Polyjuice*, a general-purpose counterfactual generator for language mod-
 95 els. They demonstrate empirically that the augmentation of training data through *Polyjuice* counterfactuals improves
 96 robustness in a number of NLP tasks. Balashankar et al. (2023) similarly use *Polyjuice* to augment NLP datasets
 97 through diverse counterfactuals and show that classifier robustness improves by up to 20%. Finally, Luu and Inoue
 98 (2023) introduce Counterfactual Adversarial Training (CAT), which also aims at improving generalization and robust-
 99 ness of language models through a three-step procedure. First, the authors identify training samples that are subject
 100 to high predictive uncertainty. Second, they generate counterfactual explanations for those samples. Finally, they
 101 fine-tune the given language model on the augmented dataset that includes the generated counterfactuals.

102 There have also been several attempts at formalizing the relationship between counterfactual explanations and adver-
 103 sarial examples. Pointing to clear similarities in how CEs and AEs are generated, Freiesleben (2022) makes the case
 104 for jointly studying the opaqueness and robustness problems in representation learning. Formally, AEs can be seen as
 105 the subset of CEs for which misclassification is achieved (Freiesleben 2022). Similarly, Pawelczyk et al. (2022) show
 106 that CEs and AEs are equivalent under certain conditions and derive theoretical upper bounds on distances between
 107 them.

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

119 2.2 Beyond Robustness

120 Improving the adversarial robustness of models is not the only path towards aligning representations with plausible
 121 explanations. In a work closely related to this one, Altmeyer et al. (2024) show that explainability can be improved
 122 through model averaging and refined model objectives. The authors propose a way to generate counterfactuals that
 123 are maximally faithful to the model in that they are consistent with what the model has learned about the underlying

124 data. Formally, they rely on tools from energy-based modelling to minimize the divergence between the distribution
 125 of counterfactuals and the conditional posterior over inputs learned by the model. Their proposed counterfactual
 126 explainer, *ECCCo*, yields plausible explanations if and only if the underlying model has learned representations that
 127 align with them. The authors find that both deep ensembles ([Lakshminarayanan, Pritzel, and Blundell 2017](#)) and joint
 128 energy-based models (JEMs) ([Grathwohl et al. 2020](#)) tend to do well in this regard.

129 Once again it helps to look at these findings through the lens of representation learning with high degrees of freedom.
 130 Deep ensembles are approximate Bayesian model averages, which are most called for when models are underspecified
 131 by the available data ([Wilson 2020](#)). Averaging across solutions mitigates the aforementioned risk of relying on a
 132 single locally optimal representations that corresponds to semantically meaningless explanations for the data. Previous
 133 work by Schut et al. ([2021](#)) similarly found that generating plausible (“interpretable”) counterfactual explanations is
 134 almost trivial for deep ensembles that have also undergone adversarial training. The case for JEMs is even clearer:
 135 they involve a hybrid objective that induces both high predictive performance and generative capacity ([Grathwohl et al.
 136 2020](#)). This is closely related to the idea of aligning models with plausible explanations and has inspired our proposed
 137 CT objective, as we explain in Section 3.

138 3 Counterfactual Training

139 Counterfactual training combines ideas from adversarial training, energy-based modelling and counterfactuals explana-
 140 tions with the explicit goal of aligning representations with plausible explanations that comply with user requirements.
 141 In the context of CEs, plausibility has broadly been defined as the degree to which counterfactuals comply with the
 142 underlying data-generating process ([Poyiadzi et al. 2020; Guidotti 2022; Altmeyer et al. 2024](#)). Plausibility is a neces-
 143 sary but insufficient condition for using CEs to provide algorithmic recourse (AR) to individuals (negatively) affected
 144 by opaque models. For AR recommendations to be actionable, they need to not only result in plausible counterfactuals
 145 but also be attainable. A plausible CE for a rejected 20-year-old loan applicant, for example, might reveal that their
 146 application would have been accepted, if only they were 20 years older. Ignoring all other features, this would comply
 147 with the definition of plausibility if 40-year-old individuals were in fact more credit-worthy on average than young
 148 adults. But of course this CE does not qualify for providing actionable recourse to the applicant since *age* is not a
 149 (directly) mutable feature. CT aims to improve model explainability by aligning models with counterfactuals that meet
 150 both desiderata: plausibility and actionability. Formally, we define explainability as follows:

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

155 We say that \mathbf{M}_θ is **explainable** to the extent that faithfully generated counterfactuals are plausible and actionable. We
 156 define these properties as follows:

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

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

161 The characterization of faithfulness and plausibility in Def. 3.1 is the same as in Altmeyer et al. ([2024](#)), with adapted
 162 notation. Intuitively, plausible counterfactuals are consistent with the data and faithful counterfactuals are consistent
 163 with what the model has learned about input data. Actionability constraints in Def. 3.1 vary and depend on the context
 164 in which \mathbf{M}_θ is deployed. In this work, we focus on domain and mutability constraints for individual features x_d for
 165 $d = 1, \dots, D$. We limit ourselves to classification tasks for reasons discussed in Section 5.

166 3.1 Our Proposed Objective

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

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

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

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

173 where $\text{yloss}(\cdot)$ is a classification loss that induces discriminative performance (e.g., cross-entropy). The second and
 174 third terms are explained in detail below. For now, they can be summarized as inducing explainability directly and
 175 indirectly by penalizing: (1) the contrastive divergence, $\text{div}(\cdot)$, between mature counterfactuals \mathbf{x}'_T and observed
 176 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
 177 nascent counterfactuals $\mathbf{x}'_{t \leq T}$. Finally, $\text{ridge}(\cdot)$ denotes a Ridge penalty (ℓ_2 -norm) that regularizes the magnitude of
 178 the energy terms involved in $\text{div}(\cdot)$ (Du and Mordatch 2020). The trade-off between the components can be governed
 179 through penalties λ_{div} , λ_{adv} and λ_{reg} .

180 3.2 Directly Inducing Explainability with Contrastive Divergence

181 Grathwohl et al. (2020) observe that any classifier can be re-interpreted as a joint energy-based model (JEM) that
 182 learns to discriminate output classes conditional on the observed (training) samples from $p(\mathbf{x})$ and the generated
 183 samples from $p_{\theta}(\mathbf{x})$. The authors show that JEMs can be trained to perform well at both tasks by directly maximizing
 184 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
 185 conventional cross-entropy as in Equation 1. Then, to optimize $\log p_{\theta}(\mathbf{x})$ Grathwohl et al. (2020) minimize the
 186 contrastive divergence between these observed samples from $p(\mathbf{x})$ and generated samples from $p_{\theta}(\mathbf{x})$.

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

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

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

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

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

205 3.3 Indirectly Inducing Explainability with Adversarial Robustness

206 Based on our analysis in Section 2, counterfactuals \mathbf{x}' can be repurposed as additional training samples (Luu and Inoue
 207 2023; Balashankar et al. 2023) or AEs (Freiesleben 2022; Pawelczyk et al. 2022). This leaves some flexibility with
 208 respect to the choice for $\text{advloss}(\cdot)$ in Equation 1. An intuitive functional form, but likely not the only sensible choice,
 209 is inspired by adversarial training:

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

210 Under this choice, we consider nascent counterfactuals $\mathbf{x}'_{t \leq T}$ as AEs as long as the magnitude of the perturbation to
 211 any single feature is at most ε . This is closely aligned with Szegedy et al. (2013) who define an adversarial attack as
 212 an “imperceptible non-random perturbation”. Thus, we choose to work with a different distinction between CE and
 213 AE than Freiesleben (2022) who consider misclassification as the key distinguishing feature of AE. One of the key
 214 observations in this work is that we can leverage CEs during training and get adversarial examples essentially for free.

215 **3.4 Encoding Actionability Constraints**

216 Many existing counterfactual explainers support domain and mutability constraints out-of-the-box. In fact, both types
 217 of constraints can be implemented for any counterfactual explainer that relies on gradient descent in the feature space
 218 for optimization (Altmeyer, Deursen, et al. 2023). In this context, domain constraints can be imposed by simply
 219 projecting counterfactuals back to the specified domain, if the previous gradient step resulted in updated feature values
 220 that were out-of-domain. Mutability constraints can similarly be enforced by setting partial derivatives to zero to
 221 ensure that features are only perturbed in the allowed direction, if at all.

222 Since such actionability constraints are binding at test time, we should also impose them when generating \mathbf{x}' during
 223 each training iteration to inform model representations. Through their effect on \mathbf{x}' , both types of constraints influence
 224 model outcomes via Equation 2. Here it is crucial that we avoid penalizing implausibility that arises due to mutability
 225 constraints. For any mutability-constrained feature d this can be achieved by enforcing $\mathbf{x}^+[d] - \mathbf{x}'[d] := 0$ whenever
 226 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:

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

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

235 **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
 236 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
 237 common diagonal covariance matrix $\Sigma_k = \Sigma$ for all $k \in \mathcal{K}$, then protecting an immutable feature from the contrastive
 238 divergence penalty will result in lower classifier sensitivity to that feature relative to the remaining features, provided
 239 that at least one of those is discriminative and mutable.*

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

242 **3.5 Example (Prediction of Consumer Credit Default)**

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

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

253 In panel (a), we have trained our model conventionally, and we do not impose mutability constraints at test time.
 254 The generated counterfactuals are all valid, but not plausible: they do not comply with the distribution of the factual
 255 samples in the target class to the point where they are clearly distinguishable from the ground-truth data. In panel (b),
 256 we have trained our model with CT, once again without any mutability constraints. We observe that the counterfactuals
 257 are highly plausible, meeting the first objective of Def. 3.1.

258 In panel (c), we have used conventional training again, this time imposing the mutability constraint on *age* at test time.
 259 Counterfactuals are valid but involve some substantial reductions in *debt* for some individuals (very young applicants).
 260 By comparison, counterfactual paths are shorter on average in panel (d), where we have used CT and protected the
 261 immutable feature as described in Section 3.4. We observe that due to the classifier's lower sensitivity to *age*, recourse
 262 recommendations with respect to *debt* are much more homogenous and do not disproportionately punish younger
 263 individuals. The counterfactuals are also plausible with respect to the mutable feature. Thus, we consider the model
 264 in panel (d) as the most explainable according to Def. 3.1.

³For the proof, see the supplementary appendix.

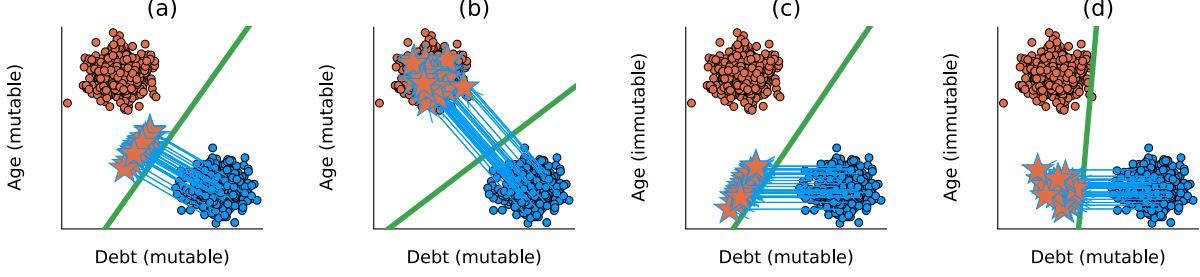


Figure 1: Illustration of how CT improves model explainability.

265 4 Experiments

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

- 267 1. To what extent does our proposed counterfactual training objective in Equation 1 induce models to learn
268 plausible explanations?
269 2. To what extent does our proposed counterfactual training objective in Equation 1 yield more favorable algo-
270 rithmic recourse outcomes in the presence of actionability constraints?
271 3. What are the effects of hyperparameter selection with respect to Equation 1?

272 **4.1 Experimental Setup**

273 **4.1.1 Evaluation**

274 Our key outcome of interest is how well do models perform with respect to explainability (Def. 3.1). To this end, we
275 focus primarily on the plausibility and cost of faithfully generated counterfactuals at test time. To measure the cost of
276 counterfactuals, we follow the standard convention of using distances (ℓ_1 -norm) between factuals and counterfactuals
277 as a proxy. For plausibility, we assess how similar counterfactuals are to observed samples in the target domain. We
278 rely on the distance-based metric used by Altmeyer et al. (2024),

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

279 and introduce a novel divergence metric,

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

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

283 In addition to cost and plausibility, we also compute other standard metrics to evaluate counterfactuals at test time in-
284 cluding validity and redundancy. Finally, we also assess the predictive performance of models using standard metrics.

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

289 **4.2 Experimental Results**

290 **4.2.1 Plausibility**

291 Table 1 presents our main empirical findings. The top five rows show the percentage reduction in implausibility
292 according to Equation 4 for varying degrees of the energy penalty used for *ECCo* at test time. The following row shows
293 the reduction in implausibility as measured by Equation 5 and aggregated across all test specifications of *ECCo*. The
294 final two rows show the test accuracies for the model trained with CT and conventionally trained models (“vanilla”).

295 We observe that for all datasets except *OL* and across all test settings, the average distance of counterfactuals from
296 observed samples in the target class is reduced, indicating improved plausibility. The magnitude of improvements
297 varies by dataset: for the simple synthetic datasets, distance reductions range from around 20-40% (*LS*, *Moon*) to
298 almost 60% (*Circ*). For the real-world tabular datasets, improvements are generally smaller but still substantial in

Table 1: Key plausibility and predictive performance metrics for all datasets. The top five rows show the percentage reduction in implausibility according to Equation 4 for varying degrees of the energy penalty used for *ECCo* at test time. The following row shows the reduction in implausibility as measured by Equation 5 and aggregated across all test specifications of *ECCo*. The final two rows show the test accuracies for the model trained with CT and conventionally trained models (“vanilla”).

Measure	λ_{egy}	Adult	CH	Circ	Cred	GMSC	LS	MNIST	Moon	OL
IP ($-\Delta\%$)	0.1	2.93	9.59	56.5	6.7	11	27.1	9.11	20.4	-6.72
IP ($-\Delta\%$)	0.5	3.4	9.26	57.1	6.18	13.4	26.7	8.26	21.4	-6.19
IP ($-\Delta\%$)	1	3.53	10.4	56.5	7.19	13.4	26.6	8.07	21.6	-6.1
IP ($-\Delta\%$)	5	2.88	11.9	58.5	7.01	21.4	27.1	6.1	19	-2.77
IP ($-\Delta\%$)	10	3.15	14.6	49.3	7.78	27.9	38.6	3.53	19.8	-1.44
IP* ($-\Delta\%$) (agg.)		34.8	66.6	93.4	51.6	77.9	54.5	-2.28	27.6	-25.5
Acc. (CT)		0.848	0.794	0.997	0.712	0.608	1	0.902	0.999	0.918
Acc. (vanilla)		0.854	0.85	0.999	0.706	0.751	1	0.922	1	0.914

many cases with around 10-15% for *CH*, 11-28% for *GMSC*, 7-8% for *Cred* and around 3% for *Adult*. For our only vision dataset (*MNIST*), distances are reduced by up to 9%. The results for our proposed divergence metric are qualitatively similar, but generally even more pronounced: for the *Circ* dataset, implausibility is reduced by almost 94% to virtually zero as we verified by looking at the absolute outcome. Improvements for other datasets range from 28% (*Moon*) to 78% (*GMSC*). For *OL* the reduction is negative, consistent with the distance-based metric. *MNIST* is the only dataset for which the two metrics disagree.

These broad and substantial improvements in plausibility generally do not come at the cost of decreased predictive performance: test accuracy for CT is virtually identical to the baseline for *Adult*, *Circ*, *LS*, *Moon* and *OL*, and even slightly improved for *Cred*. Exceptions to this general pattern are *MNIST*, *CH* and *GMSC*, for which we observe reduction in test accuracy of 2, 5 and 15 percentage points, respectively. We note in this context, that we have not optimized our models for predictive performance at all and worked with very small networks. In summary, we find that CT can substantially improve the quality of explanations learned by models without sacrificing predictive accuracy.

4.2.2 Actionability

4.2.3 Impact of hyperparameter settings

We test in-depth the impact of three types of hyperparameters; our complete results are in the appendix.

Hyperparameters of the CE generators. First, we observe that CT is highly sensitive to hyperparameter settings but (a) there are manageable patterns and (b) we can typically identify settings that improve either plausibility or cost, and commonly both of them at the same time. Second, we note that the choice of a CE generator has a major impact on the results. For example, *REVISE* tends to perform the worst, most likely because it uses a surrogate VAE to generate counterfactuals which impedes faithfulness (Altmeyer et al. 2024). Third, increasing T , the maximum number of steps, generally yields better outcomes because more CEs can mature in each training epoch. Fourth, the impact of τ , the required decision threshold is more difficult to predict. On “harder” datasets it may be difficult to satisfy high τ for any given sample (i.e., also factuals) and so increasing this threshold does not seem to correlate with better outcomes. In fact, we have generally found that a choice of $\tau = 0.5$ leads to optimal results because it is associated with high proportions of mature counterfactuals.

Hyperparameters for penalties. We find that the strength of the energy regularization, λ_{reg} is highly impactful; energy must be sufficiently regularized to avoid poor performance in terms of decreased plausibility and increased costs. The sensitivity with respect to λ_{div} and λ_{adv} is much less evident. While high values of λ_{reg} may increase the variability in outcomes when combined with high values of λ_{div} or λ_{adv} , this effect is not very pronounced.

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

5 Discussion

We first address the direct extensions of CT in Section 5.1. Then, we look at its limitations and challenges in Section 5.2.

336 **5.1 Future Research**

337 ***CT is defined only for classification settings.*** Our formulation relies on the distinction between non-target class(es)
 338 y^- and target class(es) y^+ to generate counterfactuals through Equation 1. While y^- and y^+ can be arbitrarily defined,
 339 CT requires the output space \mathcal{Y} to be discrete. Thus, it does not apply to ML tasks where the change in outcome
 340 cannot be readily quantified. Focus on classification models is a common restriction in research on CEs and AR. Other
 341 settings have attracted some interest (e.g., regression in (Spooner et al. 2021; Zhao, Broelemann, and Kasneci 2023)),
 342 but there is little consensus how to robustly extend the notion of counterfactuals.

343 ***CT is subject to training instabilities.*** JEMs are susceptible to instabilities during training (Grathwohl et al. 2020)
 344 and even though we depart from the SGLD-based sampling, we still encounter major variability in the outcomes. CT
 345 is exposed to two potential sources of instabilities: (1) the energy-based contrastive divergence term in Equation 2,
 346 and (2) the underlying counterfactual explainers. Still, we find that training instabilities can be successfully mitigated
 347 by regularizing energy (λ_{reg}), generating sufficiently many counterfactuals during each training epoch, and including
 348 only mature counterfactuals for contrastive divergence.

349 ***CT is sensitive to hyperparameter selection.*** Our method benefits from the tuning of certain key hyperparameters
 350 (see Section 4.2.3). In this work, we have relied exclusively on grid search for this task. Future work on CT could
 351 benefit from investigating more sophisticated approaches towards hyperparameter tuning. Notably, CT is iterative
 352 which makes a variety of methods applicable, including Bayesian (e.g., Snoek, Larochelle, and Adams 2012) or
 353 gradient-based (e.g., Franceschi et al. 2017) optimization.

354 **5.2 Current Limitations**

355 ***CT increases the training time of models.*** CT can be more time-consuming than conventional training regimes.
 356 While higher numbers of CEs per iteration positively impact the quality of solutions, they also increase the amount of
 357 computations. Relatively small grids with 270 settings can take almost four hours for more demanding datasets on a
 358 high-performance computing cluster with 34 2GB CPUs.⁴ Three factors attenuate this effect. First, CT amortizes the
 359 cost of CEs for the training samples. Second, we find that it can retain its value when used as a “fine-tuning” technique
 360 for conventionally-trained models. Third, it yields itself to parallel execution, which we have leveraged for our own
 361 experiments.

362 ***Immutable features may have proxies.*** We propose an approach to protect immutable features and thus increase the
 363 actionability of the generated CEs. However, it requires that model owners define the mutability constraints for (all)
 364 features considered by the model. Even if all immutable features are protected, there may exist proxies that are mutable
 365 (and hence should not be protected) but preserve enough information about the principals to hinder the protections.
 366 Delineating actionability is a major undecided challenge in the AR literature (see, e.g., Venkatasubramanian and
 367 Alfano 2020) impacting the capacity of CT to increase the explainability of the model.

368 ***Interventions on features may impact fairness.*** We provide a tool that allows practitioners to modify the sensitivity
 369 of a model with respect to certain features, which may have implication for the fair and equitable treatment of decision
 370 subjects. As protecting a set of features leads the model to assign higher relative importance to unprotected features,
 371 model owners could misuse our solution by enforcing explanations based on features that are more difficult to modify
 372 by some (group of) individuals. For example, consider the Adult dataset used in our experiments, where *workclass* or
 373 *education* may be more difficult to change for underprivileged groups. When applied irresponsibly, CT could result
 374 in an unfairly assigned burden of recourse (e.g., Sharma, Henderson, and Ghosh 2020), threatening the equality of
 375 opportunity in the system (Bell et al. 2024). Still, these phenomena are not specific to CT.

376 **6 Conclusion**

377 State-of-the-art machine learning models are prone to learning complex representations that cannot be interpreted by
 378 humans and existing post-hoc explainability approaches cannot guarantee that the explanations agree with the model’s
 379 learned representation of data. As a step towards addressing this challenge, we introduced counterfactual training, a
 380 novel training regime that incentivizes highly-explainable models. Our approach leads to explanations that are both
 381 plausible—compliant with the underlying data-generating process—and actionable—compliant with user-specified
 382 mutability constraints—and thus meaningful to their recipients. Through extensive experiments we demonstrate that
 383 CT satisfies its objectives while preserving the predictive performance of the models. Our approach can also be used
 384 to fine-tune conventionally-trained models and achieve similar gains in explainability. Finally, this work showcases
 385 that it is practical to improve models *and* their explanations at the same time.

⁴See supplementary appendix for computational details.

386 **References**

- 387 Abbasnejad, Ehsan, Damien Teney, Amin Parvaneh, Javen Shi, and Anton van den Hengel. 2020. “Counterfactual
 388 Vision and Language Learning.” In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition
 389 (CVPR)*, 10041–51. <https://doi.org/10.1109/CVPR42600.2020.01006>.
- 390 Altmeyer, Patrick, Arie van Deursen, et al. 2023. “Explaining Black-Box Models Through Counterfactuals.” In
 391 *Proceedings of the JuliaCon Conferences*, 1:130. 1.
- 392 Altmeyer, Patrick, Mojtaba Farmanbar, Arie van Deursen, and Cynthia CS Liem. 2024. “Faithful Model Explanations
 393 Through Energy-Constrained Conformal Counterfactuals.” In *Proceedings of the AAAI Conference on Artificial
 394 Intelligence*, 38:10829–37. 10.
- 395 Augustin, Maximilian, Alexander Meinke, and Matthias Hein. 2020. “Adversarial Robustness on in-and Out-
 396 Distribution Improves Explainability.” In *European Conference on Computer Vision*, 228–45. Springer.
- 397 Balashankar, Ananth, Xuezhi Wang, Yao Qin, Ben Packer, Nithum Thain, Ed Chi, Jilin Chen, and Alex Beutel. 2023.
 398 “Improving Classifier Robustness Through Active Generative Counterfactual Data Augmentation.” In *Findings of
 399 the Association for Computational Linguistics: EMNLP 2023*, 127–39.
- 400 Bell, Andrew, Joao FONSECA, Carlo Abrate, Francesco Bonchi, and Julia Stoyanovich. 2024. “Fairness in Algorithmic
 401 Recourse Through the Lens of Substantive Equality of Opportunity.” <https://arxiv.org/abs/2401.16088>.
- 402 Bezanson, Jeff, Alan Edelman, Stefan Karpinski, and Viral B Shah. 2017. “Julia: A Fresh Approach to Numerical
 403 Computing.” *SIAM Review* 59 (1): 65–98. <https://doi.org/10.1137/141000671>.
- 404 Bouchet-Valat, Milan, and Bogumi Kamiski. 2023. “DataFrames.jl: Flexible and Fast Tabular Data in Julia.” *Journal
 405 of Statistical Software* 107 (4): 1–32. <https://doi.org/10.18637/jss.v107.i04>.
- 406 Byrne, Simon, Lucas C. Wilcox, and Valentin Churavy. 2021. “MPI.jl: Julia Bindings for the Message Passing
 407 Interface.” *Proceedings of the JuliaCon Conferences* 1 (1): 68. <https://doi.org/10.21105/jcon.00068>.
- 408 Chagas, Ronan Arraes Jardim, Ben Baumgold, Glen Hertz, Hendrik Ranocha, Mark Wells, Nathan Boyer, Nicholas
 409 Ritchie, et al. 2024. “Ronisbr/PrettyTables.jl: V2.4.0.” Zenodo. <https://doi.org/10.5281/zenodo.1383553>.
- 410 Christ, Simon, Daniel Schwabeneder, Christopher Rackauckas, Michael Krabbe Borregaard, and Thomas Breloff.
 411 2023. “Plots.jl – a User Extendable Plotting API for the Julia Programming Language.” <https://doi.org/https://doi.org/10.5334/jors.431>.
- 412 Danisch, Simon, and Julius Krumbiegel. 2021. “Makie.jl: Flexible High-Performance Data Visualization for Julia.”
 413 *Journal of Open Source Software* 6 (65): 3349. <https://doi.org/10.21105/joss.03349>.
- 414 Du, Yilun, and Igor Mordatch. 2020. “Implicit Generation and Generalization in Energy-Based Models.” <https://arxiv.org/abs/1903.08689>.
- 415 Franceschi, Luca, Michele Donini, Paolo Frasconi, and Massimiliano Pontil. 2017. “Forward and Reverse Gradient-
 416 Based Hyperparameter Optimization.” In *Proceedings of the 34th International Conference on Machine Learning*,
 417 edited by Doina Precup and Yee Whye Teh, 70:1165–73. Proceedings of Machine Learning Research. PMLR.
 418 <https://proceedings.mlr.press/v70/franceschi17a.html>.
- 419 Frankle, Jonathan, and Michael Carbin. 2019. “The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural
 420 Networks.” In *International Conference on Learning Representations*. <https://openreview.net/forum?id=rJl-b3RcF7>.
- 421 Freiesleben, Timo. 2022. “The Intriguing Relation Between Counterfactual Explanations and Adversarial Examples.”
 422 *Minds and Machines* 32 (1): 77–109.
- 423 Goodfellow, Ian J, Jonathon Shlens, and Christian Szegedy. 2014. “Explaining and Harnessing Adversarial Examples.”
 424 <https://arxiv.org/abs/1412.6572>.
- 425 Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press.
- 426 Grathwohl, Will, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swer-
 427 sky. 2020. “Your Classifier Is Secretly an Energy Based Model and You Should Treat It Like One.” In *International
 428 Conference on Learning Representations*.
- 429 Gretton, Arthur, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. 2012. “A Kernel
 430 Two-Sample Test.” *The Journal of Machine Learning Research* 13 (1): 723–73.
- 431 Guidotti, Riccardo. 2022. “Counterfactual Explanations and How to Find Them: Literature Review and Benchmark-
 432 ing.” *Data Mining and Knowledge Discovery*, 1–55.
- 433 Guo, Hangzhi, Thanh H. Nguyen, and Amulya Yadav. 2023. “CounterNet: End-to-End Training of Prediction Aware
 434 Counterfactual Explanations.” In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery
 435 and Data Mining*, 577–89. KDD ’23. New York, NY, USA: Association for Computing Machinery. <https://doi.org/10.1145/3580305.3599290>.
- 436 Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2009. *The Elements of Statistical Learning*. Springer New
 437 York. <https://doi.org/10.1007/978-0-387-84858-7>.
- 438 Innes, Michael, Elliot Saba, Keno Fischer, Dhairya Gandhi, Marco Conchetto Rudilosso, Neethu Mariya Joy, Tejan
 439 Karmali, Avik Pal, and Viral Shah. 2018. “Fashionable Modelling with Flux.” <https://arxiv.org/abs/1811.01457>.

- 444 Innes, Mike. 2018. "Flux: Elegant Machine Learning with Julia." *Journal of Open Source Software* 3 (25): 602.
 445 <https://doi.org/10.21105/joss.00602>.
- 446 Joshi, Shalmali, Oluwasanmi Koyejo, Warut Vigitbenjaronk, Been Kim, and Joydeep Ghosh. 2019. "Towards Realistic
 447 Individual Recourse and Actionable Explanations in Black-Box Decision Making Systems." <https://arxiv.org/abs/1907.09615>.
- 448 Kolter, Zico. 2023. "Keynote Addresses: SaTML 2023 ." In *2023 IEEE Conference on Secure and Trustworthy
 449 Machine Learning (SaTML)*, xvi–. Los Alamitos, CA, USA: IEEE Computer Society. <https://doi.org/10.1109/SaTML54575.2023.00009>.
- 450 Lakshminarayanan, Balaji, Alexander Pritzel, and Charles Blundell. 2017. "Simple and Scalable Predictive Uncer-
 451 tainty Estimation Using Deep Ensembles." *Advances in Neural Information Processing Systems* 30.
- 452 Lippe, Phillip. 2024. "UvA Deep Learning Tutorials." <https://uvadlc-notebooks.readthedocs.io/en/latest/>.
- 453 Luu, Hoai Linh, and Naoya Inoue. 2023. "Counterfactual Adversarial Training for Improving Robustness of Pre-
 454 Trained Language Models." In *Proceedings of the 37th Pacific Asia Conference on Language, Information and
 455 Computation*, 881–88.
- 456 McGregor, Sean. 2021. "Preventing repeated real world AI failures by cataloging incidents: The AI incident database." In
 457 *Proceedings of the AAAI Conference on Artificial Intelligence*, 35:15458–63. 17.
- 458 Murphy, Kevin P. 2022. *Probabilistic Machine Learning: An Introduction*. MIT Press.
- 459 O'Neil, Cathy. 2016. *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*.
 460 Crown.
- 461 Pawelczyk, Martin, Chirag Agarwal, Shalmali Joshi, Sohini Upadhyay, and Himabindu Lakkaraju. 2022. "Exploring
 462 Counterfactual Explanations Through the Lens of Adversarial Examples: A Theoretical and Empirical Analysis." In
 463 *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, edited by Gustau
 464 Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, 151:4574–94. Proceedings of Machine Learning Research.
 465 PMLR. <https://proceedings.mlr.press/v151/pawelczyk22a.html>.
- 466 Poyiadzi, Rafael, Kacper Sokol, Raul Santos-Rodriguez, Tijl De Bie, and Peter Flach. 2020. "FACE: Feasible and
 467 Actionable Counterfactual Explanations." In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*,
 468 344–50.
- 469 Ross, Alexis, Himabindu Lakkaraju, and Osbert Bastani. 2024. "Learning Models for Actionable Recourse." In
 470 *Proceedings of the 35th International Conference on Neural Information Processing Systems*. NIPS '21. Red
 471 Hook, NY, USA: Curran Associates Inc.
- 472 Sauer, Axel, and Andreas Geiger. 2021. "Counterfactual Generative Networks." <https://arxiv.org/abs/2101.06046>.
- 473 Schut, Lisa, Oscar Key, Rory Mc Grath, Luca Costabello, Bogdan Sacaleanu, Yarin Gal, et al. 2021. "Generating
 474 Interpretable Counterfactual Explanations By Implicit Minimisation of Epistemic and Aleatoric Uncertainties." In
 475 *International Conference on Artificial Intelligence and Statistics*, 1756–64. PMLR.
- 476 Sharma, Shubham, Jette Henderson, and Joydeep Ghosh. 2020. "CERTIFAI: A Common Framework to Provide
 477 Explanations and Analyse the Fairness and Robustness of Black-Box Models." In *Proceedings of the AAAI/ACM
 478 Conference on AI, Ethics, and Society*, 166–72. AIES '20. New York, NY, USA: Association for Computing
 479 Machinery. <https://doi.org/10.1145/3375627.3375812>.
- 480 Snoek, Jasper, Hugo Larochelle, and Ryan P. Adams. 2012. "Practical Bayesian Optimization of Machine Learning
 481 Algorithms." In *Advances in Neural Information Processing Systems*, edited by F. Pereira, C. J. Burges, L. Bottou,
 482 and K. Q. Weinberger. Vol. 25. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2012/file05311655a15b75fab86956663e1819cd-Paper.pdf.
- 483 Spooner, Thomas, Danial Dervovic, Jason Long, Jon Shepard, Jiahao Chen, and Daniele Magazzeni. 2021. "Counter-
 484 factual Explanations for Arbitrary Regression Models." *CoRR* abs/2106.15212. <https://arxiv.org/abs/2106.15212>.
- 485 Szegedy, Christian, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus.
 486 2013. "Intriguing Properties of Neural Networks." <https://arxiv.org/abs/1312.6199>.
- 487 Teney, Damien, Ehsan Abbasnedjad, and Anton van den Hengel. 2020. "Learning What Makes a Difference from
 488 Counterfactual Examples and Gradient Supervision." In *Computer Vision–ECCV 2020: 16th European Conference,
 489 Glasgow, UK, August 23–28, 2020, Proceedings, Part x 16*, 580–99. Springer.
- 490 Venkatasubramanian, Suresh, and Mark Alfano. 2020. "The Philosophical Basis of Algorithmic Recourse." In
 491 *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 284–93. FAT* '20. New York,
 492 NY, USA: Association for Computing Machinery. <https://doi.org/10.1145/3351095.3372876>.
- 493 Wachter, Sandra, Brent Mittelstadt, and Chris Russell. 2017. "Counterfactual Explanations Without Opening the Black
 494 Box: Automated Decisions and the GDPR." *Harv. JL & Tech.* 31: 841. <https://doi.org/10.2139/ssrn.3063289>.
- 495 Wilson, Andrew Gordon. 2020. "The Case for Bayesian Deep Learning." <https://arxiv.org/abs/2001.10995>.
- 496 Wu, Tongshuang, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2021. "Polyjuice: Generating Counterfactuals
 497 for Explaining, Evaluating, and Improving Models." In *Proceedings of the 59th Annual Meeting of the Association
 498 for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*
 499 500

- 502 (Volume 1: Long Papers), edited by Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, 6707–23. Online:
503 Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.acl-long.523>.
- 504 Zhang, Chiyuan, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. 2021. “Understanding Deep
505 Learning (Still) Requires Rethinking Generalization.” *Commun. ACM* 64 (3): 107–15. <https://doi.org/10.1145/3446776>.
- 506 Zhao, Xuan, Klaus Broelemann, and Gjergji Kasneci. 2023. “Counterfactual Explanation for Regression via Disentan-
507 glement in Latent Space.” In *2023 IEEE International Conference on Data Mining Workshops (ICDMW)*, 976–84.
508 Los Alamitos, CA, USA: IEEE Computer Society. <https://doi.org/10.1109/ICDMW60847.2023.00130>.

510 **G Notation**

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

516 **G.1 Other Technical Details**

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

517 **H Technical Details of Our Approach**

518 **H.1 Generating Counterfactuals through Gradient Descent**

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

521 **H.1.1 Background**

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

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

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

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

534 **H.1.2 Convergence**

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

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

544 **H.2 Protecting Mutability Constraints with Linear Classifiers**

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

549 *Proof.* Let d_{mtbl} and d_{immtbl} denote some mutable and immutable feature, respectively. Suppose that $\mu_{y^-, d_{\text{immtbl}}} <$
 550 $\mu_{y^+, d_{\text{immtbl}}} \text{ and } \mu_{y^-, d_{\text{mtbl}}} > \mu_{y^+, d_{\text{mtbl}}}$, where $\mu_{k,d}$ denotes the conditional sample mean of feature d in class k . In words,
 551 we assume that the immutable feature tends to take lower values for samples in the non-target class y^- than in the
 552 target class y^+ . We assume the opposite to hold for the mutable feature.

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

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

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

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

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

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

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

558 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
 559 made above. Hence, we have that $(\theta_{y^-, d_{\text{immtbl}}} - \theta_{y^+, d_{\text{immtbl}}}) < 0$ and $(\theta_{y^-, d_{\text{mtbl}}} - \theta_{y^+, d_{\text{mtbl}}}) > 0$

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

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

563 and equal to zero everywhere else.

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

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

570 By the same rationale, the contrastive divergence penalty can generally be expected to increase $\theta_{y^-, d_{\text{mtbl}}}$ exacerbating
 571 $(\theta_{y^-, d_{\text{mtbl}}} - \theta_{y^+, d_{\text{mtbl}}}) > 0$. In words, this has the effect of making the classifier more sensitive to the mutable feature, in
 572 that it would be more likely to predict class y^- as opposed to y^+ for higher values of d_{mtbl} .

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

576 H.3 Domain Constraints

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

Table A2: Final hyperparameters used for the main results for the different datasets.

Data	No. Train	No. Test	Batchsize	Domain	Decision Threshold	No. Counterfactuals	λ_{reg}
Adult	$2.6 \cdot 10^4$	$5.01 \cdot 10^3$	$1 \cdot 10^3$	none	0.75	$5 \cdot 10^3$	0.25
CH	$1.65 \cdot 10^4$	$3.1 \cdot 10^3$	$1 \cdot 10^3$	none	0.5	$5 \cdot 10^3$	0.25
Circ	$3.6 \cdot 10^3$	600	30	none	0.5	$1 \cdot 10^3$	0.5
Cred	$1.06 \cdot 10^4$	$1.92 \cdot 10^3$	$1 \cdot 10^3$	none	0.5	$5 \cdot 10^3$	0.25
GMSC	$1.34 \cdot 10^4$	$2.47 \cdot 10^3$	$1 \cdot 10^3$	none	0.5	$5 \cdot 10^3$	0.5
LS	$3.6 \cdot 10^3$	600	30	none	0.5	$1 \cdot 10^3$	0.01
MNIST	$1.1 \cdot 10^4$	$2 \cdot 10^3$	$1 \cdot 10^3$	(-1.0, 1.0)	0.5	$5 \cdot 10^3$	0.01
Moon	$3.6 \cdot 10^3$	600	30	none	0.9	$1 \cdot 10^3$	0.25
OL	$3.6 \cdot 10^3$	600	30	none	0.5	$1 \cdot 10^3$	0.25

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

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

585 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
 586 higher values and hence wider bounds may be appropriate depending on the application.

587 H.4 Training Details

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

590 I Details on Main Experiments

591 I.1 Final Hyperparameters

592 As discussed Section 4, CT is sensitive to certain hyperparameter choices. We study the effect of many hyperparame-
 593 ters extensively in Section J. For the main results, we tune a small set of key hyperparameters (Section K). The final
 594 choices for the main results are presented for each data set in Table A2 along with training, test and batch sizes.

595 I.2 Qualitative Findings for Image Data

Note

Figure A2 shows much more plausible (faithful) counterfactuals for a model with CT than the model with conventional training (Figure A3). In fact, this is not even using ECCo+ and still showing better results than the best results we achieved in our AAAI paper for JEM ensembles.

596

597 J Grid Searches

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

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

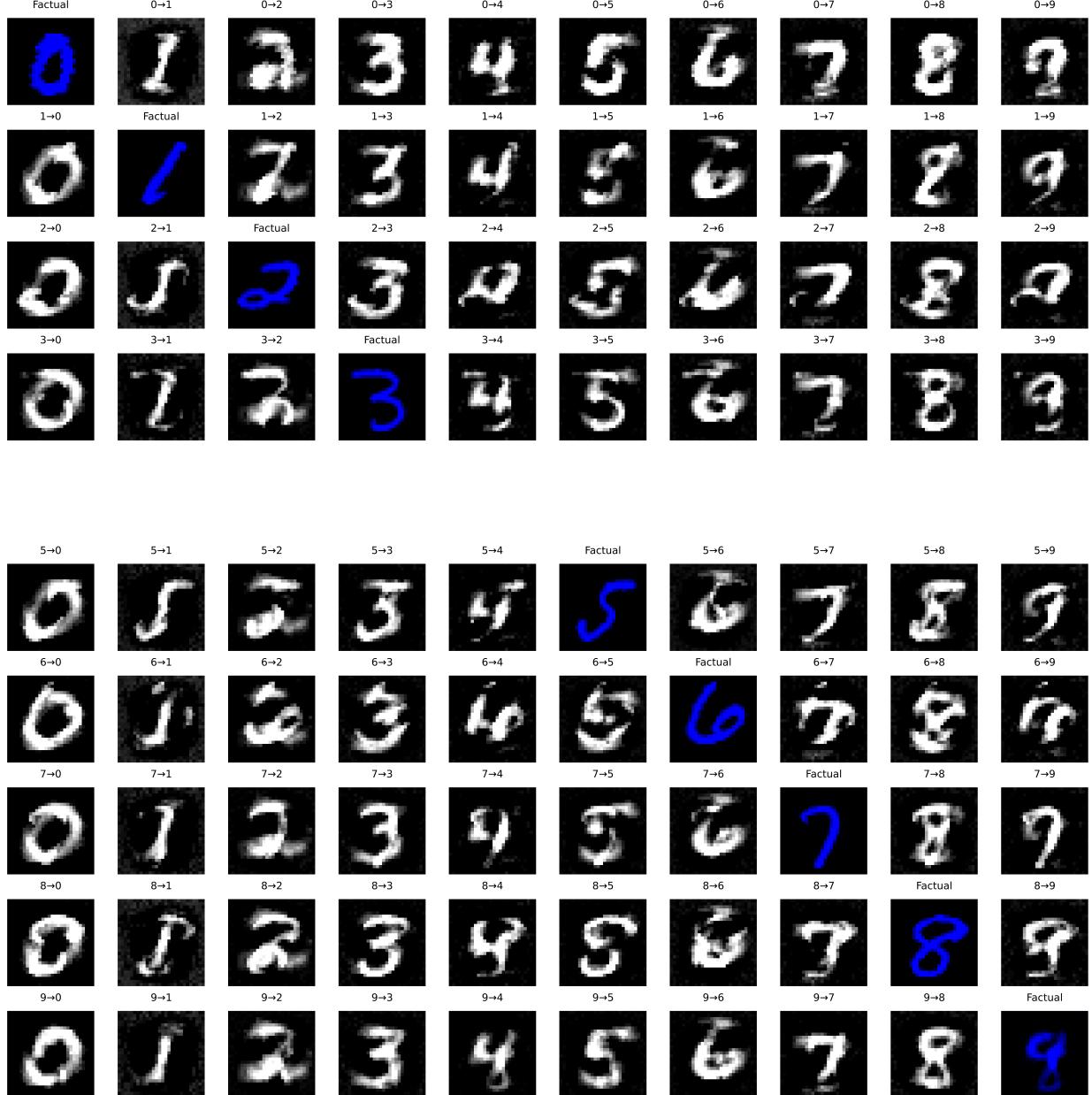


Figure A2: Counterfactual images for *MLP* with counterfactual training. The underlying generator, *ECCo*, aims to generate counterfactuals that are faithful to the model (Altmeyer et al. 2024).

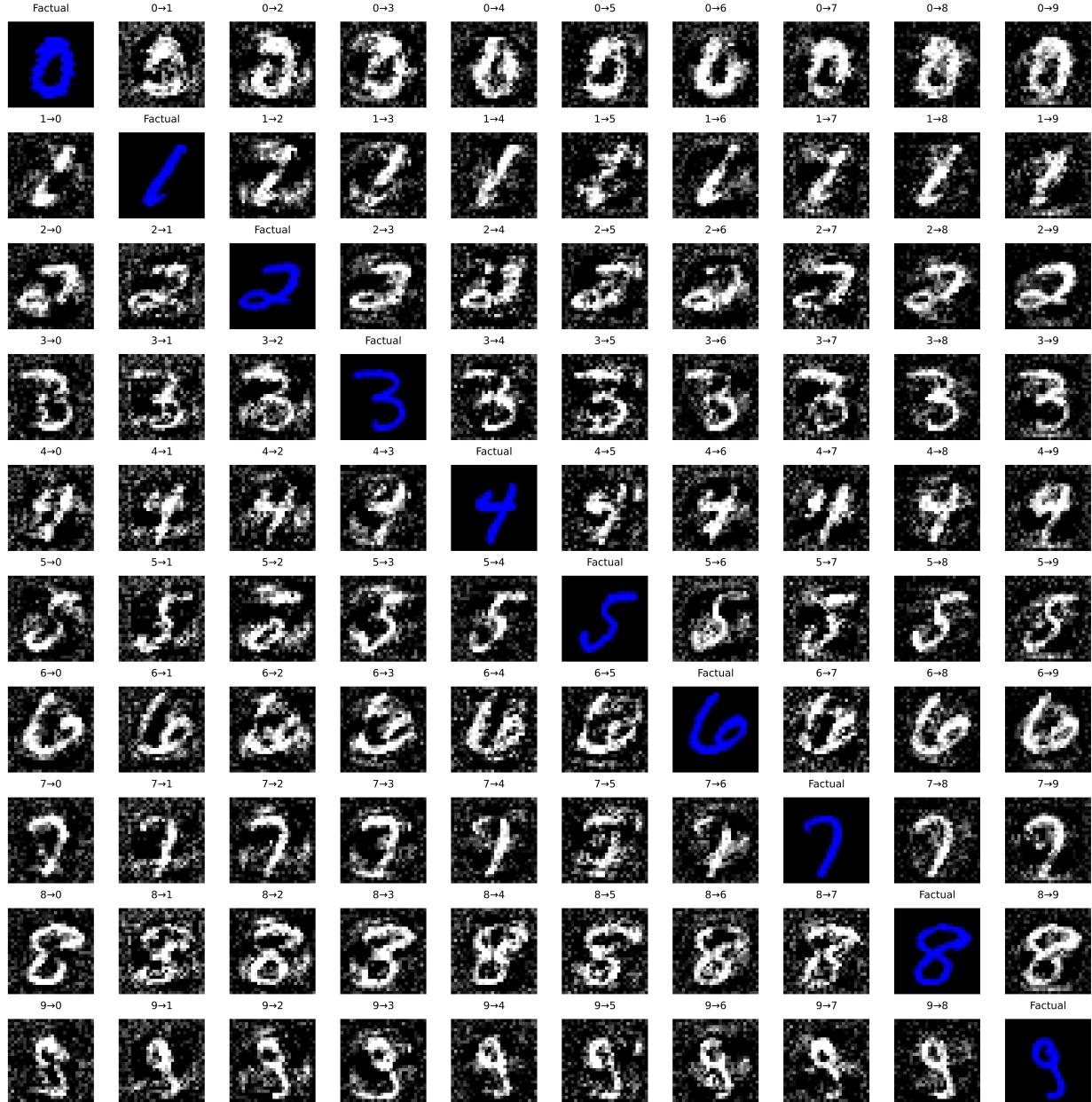


Figure A3: Counterfactual images for *MLP* with conventional training. The underlying generator, *ECCo*, aims to generate counterfactuals that are faithful to the model (Altmeyer et al. 2024).

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

609 J.1 Evaluation Details

610 To measure predictive performance, we compute the accuracy and F1-score for all models on test data (Table A3,
 611 Table A4, Table A5). With respect to explanatory performance, we report here our findings for the (im)plausibility
 612 and cost of counterfactuals at test time. Since the computation of our proposed divergence metric (Equation 5) is
 613 memory-intensive, we rely on the distance-based metric for the grid searches. For the counterfactual evaluation, we
 614 draw factual samples from the training data for the grid searches to avoid data leakage with respect to our final results
 615 reported in the body of the paper. Specifically, we want to avoid choosing our default hyperparameters based on results
 616 on the test data. Since we are optimizing for explainability, not predictive performance, we still present test accuracy
 617 and F1-scores.

618 J.1.1 Predictive Performance

619 We find that CT is associated with little to no decrease in average predictive performance for our synthetic datasets:
 620 test accuracy and F1-scores decrease by at most ~1 percentage point, but generally much less (Table A3, Table A4,
 621 Table A5). Variation across hyperparameters is negligible as indicated by small standard deviations for these metrics
 622 across the board.

623 J.1.2 Counterfactual Outcomes

624 Overall, we find that counterfactual training (CT) achieves its key objectives consistently across all hyperparameter
 625 settings and also broadly across datasets: plausibility is improved by up to ~60 percent (%) for the *Circles* data
 626 (e.g. Figure A4), ~25-30% for the *Moons* data (e.g. Figure A6) and ~10-20% for the *Linearly Separable* data (e.g.
 627 Figure A5). At the same time, the average costs of faithful counterfactuals are reduced in many cases by around
 628 ~20-25% for *Circles* (e.g. Figure A8) and up to ~50% for *Moons* (e.g. Figure A10). For the *Linearly Separable* data,
 629 costs are generally increased although typically by less than 10% (e.g. Figure A9), which reflects a common tradeoff
 630 between costs and plausibility (Altmeyer et al. 2024).

631 We do observe strong sensitivity to certain hyperparameters, with clear manageable patterns. Concerning generator
 632 parameters, we firstly find that using *REVISE* to generate counterfactuals during training typically yields the worst
 633 outcomes out of all generators, often leading to a substantial decrease in plausibility. This finding can be attributed to
 634 the fact that *REVISE* effectively assigns the task of learning plausible explanations from the model itself to a surrogate
 635 VAE. In other words, counterfactuals generated by *REVISE* are less faithful to the model than *ECCo* and *Generic*, and
 636 hence we would expect them to be a less effective and, in fact, potentially detrimental role in our training regime.
 637 Secondly, we observe that allowing for a higher number of maximum steps T for the counterfactual search generally
 638 yields better outcomes. This is intuitive, because it allows more counterfactuals to reach maturity in any given iteration.
 639 Looking in particular at the results for *Linearly Separable*, it seems that higher values for T in combination with higher
 640 decision thresholds (τ) yields the best results when using *ECCo*. But depending on the degree of class separability
 641 of the underlying data, a high decision-threshold can also affect results adversely, as evident from the results for
 642 the *Overlapping* data (Figure A7): here we find that CT generally fails to achieve its objective because only a tiny
 643 proportion of counterfactuals ever reaches maturity.

644 Regarding penalty strengths, we find that the strength of the energy regularization, λ_{reg} is a key hyperparameter, while
 645 sensitivity with respect to λ_{div} and λ_{adv} is much less evident. In particular, we observe that not regularizing energy
 646 enough or at all typically leads to poor performance in terms of decreased plausibility and increased costs, in particular
 647 for *Circles* (Figure A12), *Linearly Separable* (Figure A13) and *Overlapping* (Figure A15). High values of λ_{reg} can
 648 increase the variability in outcomes, in particular when combined with high values for λ_{div} and λ_{adv} , but this effect is
 649 less pronounced.

650 Finally, concerning other hyperparameters we observe that the effectiveness and stability of CT is positively associated
 651 with the number of counterfactuals generated during each training epoch, in particular for *Circles* (Figure A20) and
 652 *Moons* (Figure A22). We further find that a higher number of training epochs is beneficial as expected, where we
 653 tested training models for 50 and 100 epochs. Interestingly, we find that it is not necessary to employ CT during
 654 the entire training phase to achieve the desired improvements in explainability: specifically, we have tested training
 655 models conventionally during the first half of training before switching to CT after this initial burn-in period.

656 J.2 Generator Parameters

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

659 Note 1: 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`

660 Note 2: Evaluation Phase

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

661 **J.2.1 Accuracy**

Table A3: Predictive performance measures by dataset and objective averaged across training-phase parameters (Note 1) and evaluation-phase parameters (Note 2).

Dataset	Variable	Objective	Mean	Std
Circ	Accuracy	Full	0.997	0.00309
Circ	Accuracy	Vanilla	0.998	0.000557
Circ	F1-score	Full	0.997	0.00309
Circ	F1-score	Vanilla	0.998	0.000558
LS	Accuracy	Full	0.999	0.00201
LS	Accuracy	Vanilla	1	0
LS	F1-score	Full	0.999	0.00201
LS	F1-score	Vanilla	1	0
Moon	Accuracy	Full	0.999	0.000696
Moon	Accuracy	Vanilla	1	0.00111
Moon	F1-score	Full	0.999	0.000696
Moon	F1-score	Vanilla	1	0.00111
OL	Accuracy	Full	0.915	0.00477
OL	Accuracy	Vanilla	0.917	0.00123
OL	F1-score	Full	0.915	0.00478
OL	F1-score	Vanilla	0.917	0.00124

662 **J.2.2 Plausibility**

663 The results with respect to the plausibility measure are shown in Figure A4 to Figure A7.

664 **J.2.3 Cost**

665 The results with respect to the cost measure are shown in Figure A8 to Figure A11.

666 **J.3 Penalty Strengths**

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

669 Note 3: Training Phase

- Generator: `ecco`, `generic`, `revise`
- Model: `mlp`
- Training Parameters:
 - λ_{adv} : 0.1, 0.25, 1.0

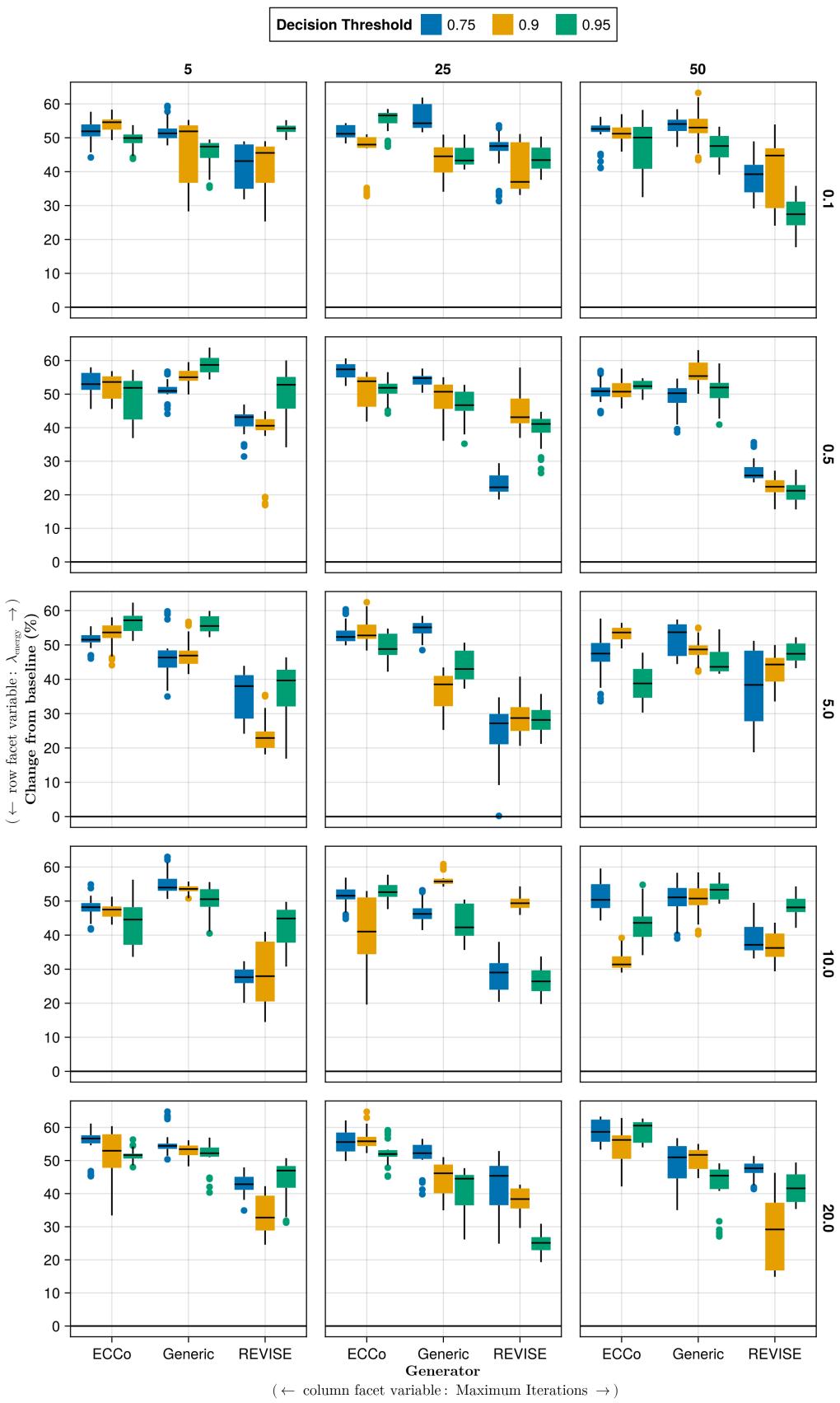


Figure A4: Average outcomes for the plausibility measure across hyperparameters. Data: Circles.

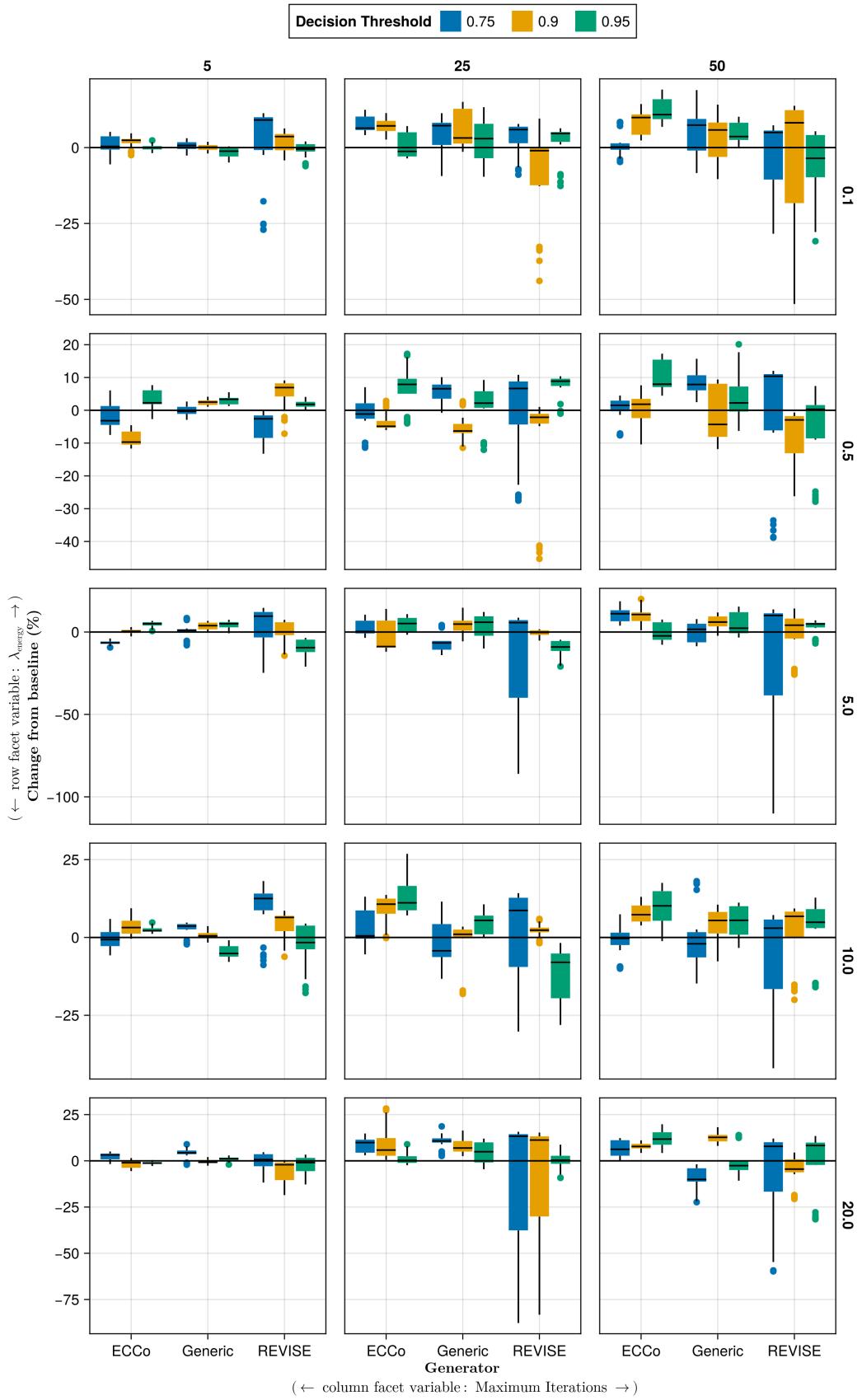


Figure A5: Average outcomes for the plausibility measure across hyperparameters. Data: Linearly Separable.

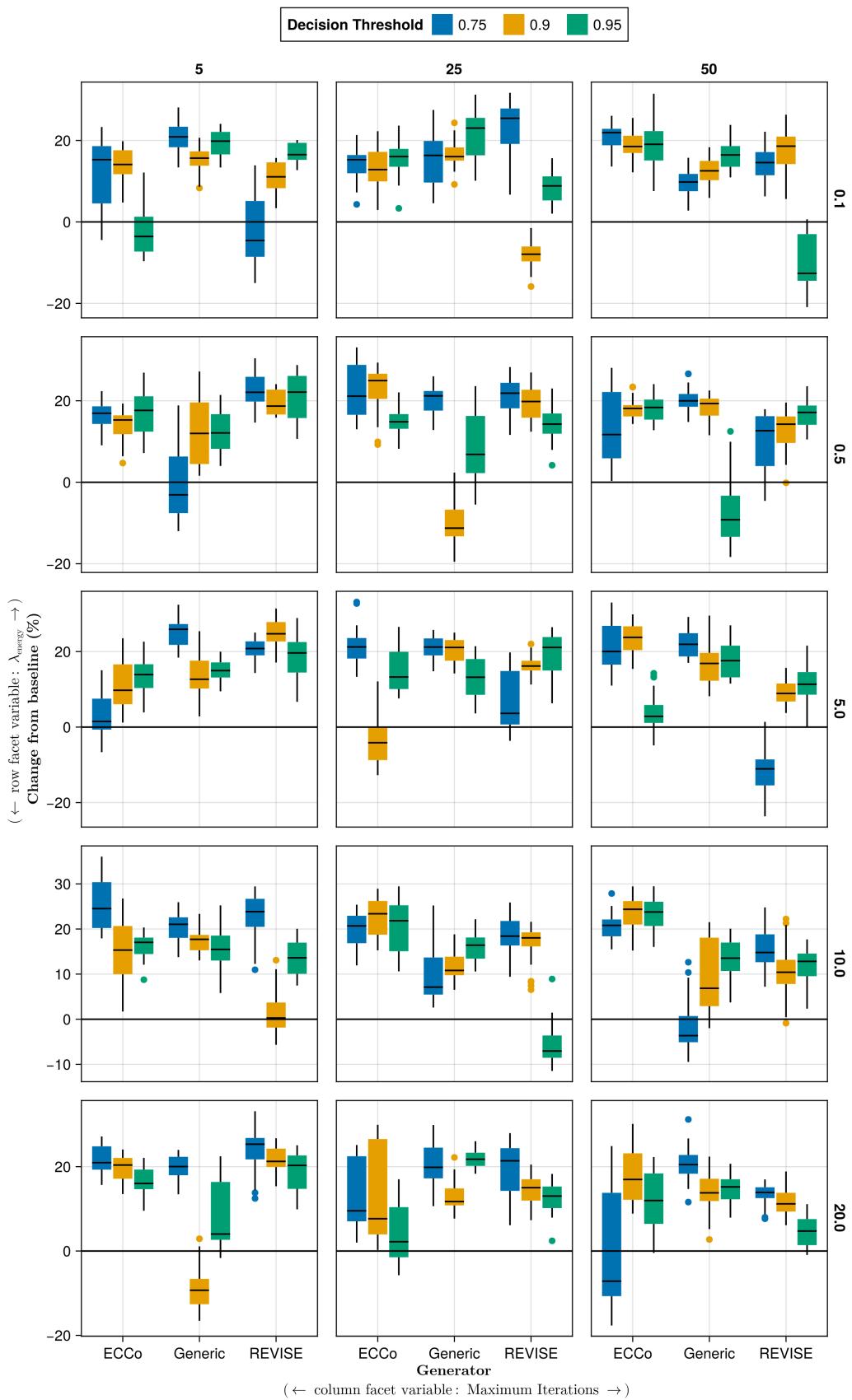


Figure A6: Average outcomes for the plausibility measure across hyperparameters. Data: Moons.

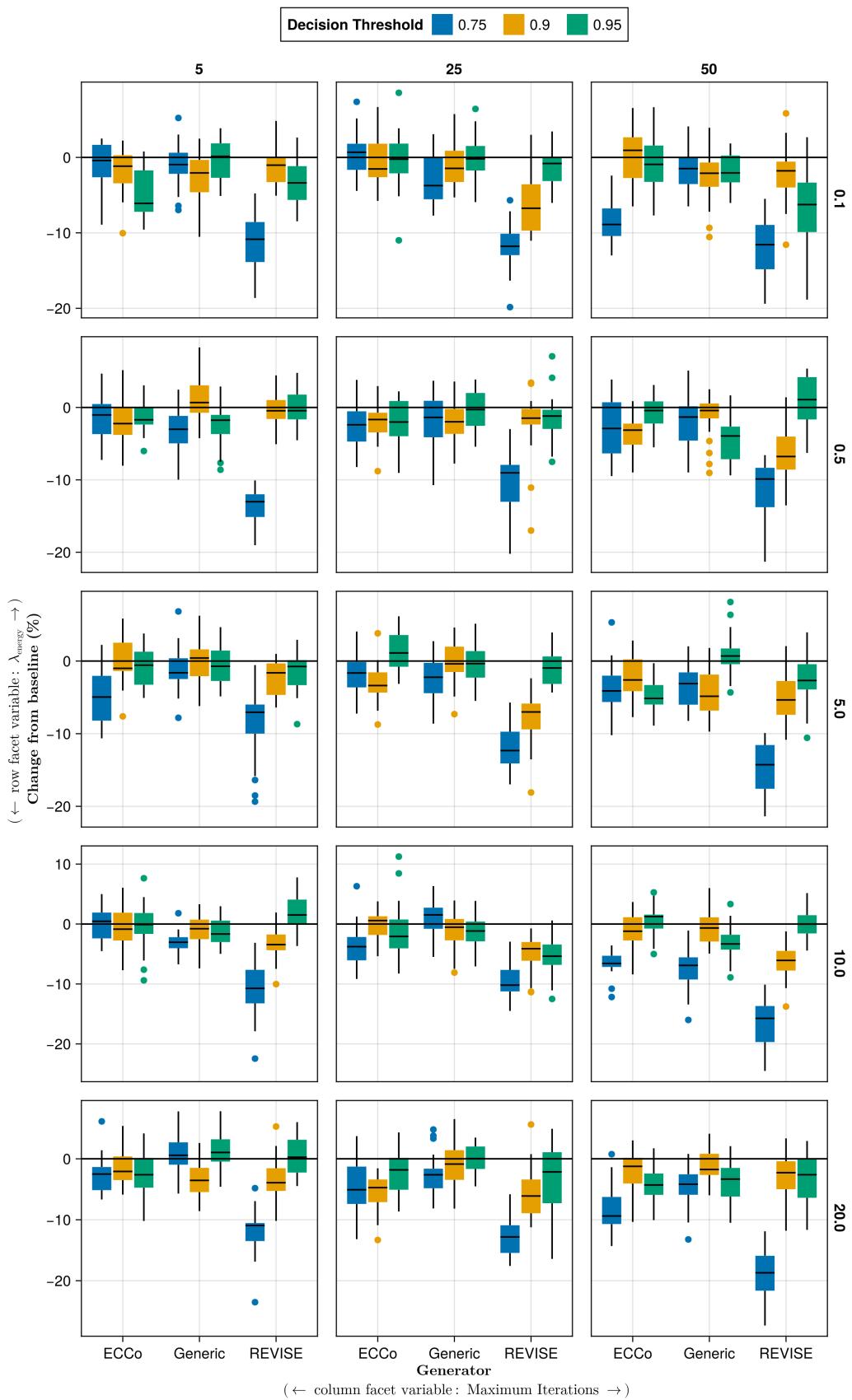


Figure A7: Average outcomes for the plausibility measure across hyperparameters. Data: Overlapping.

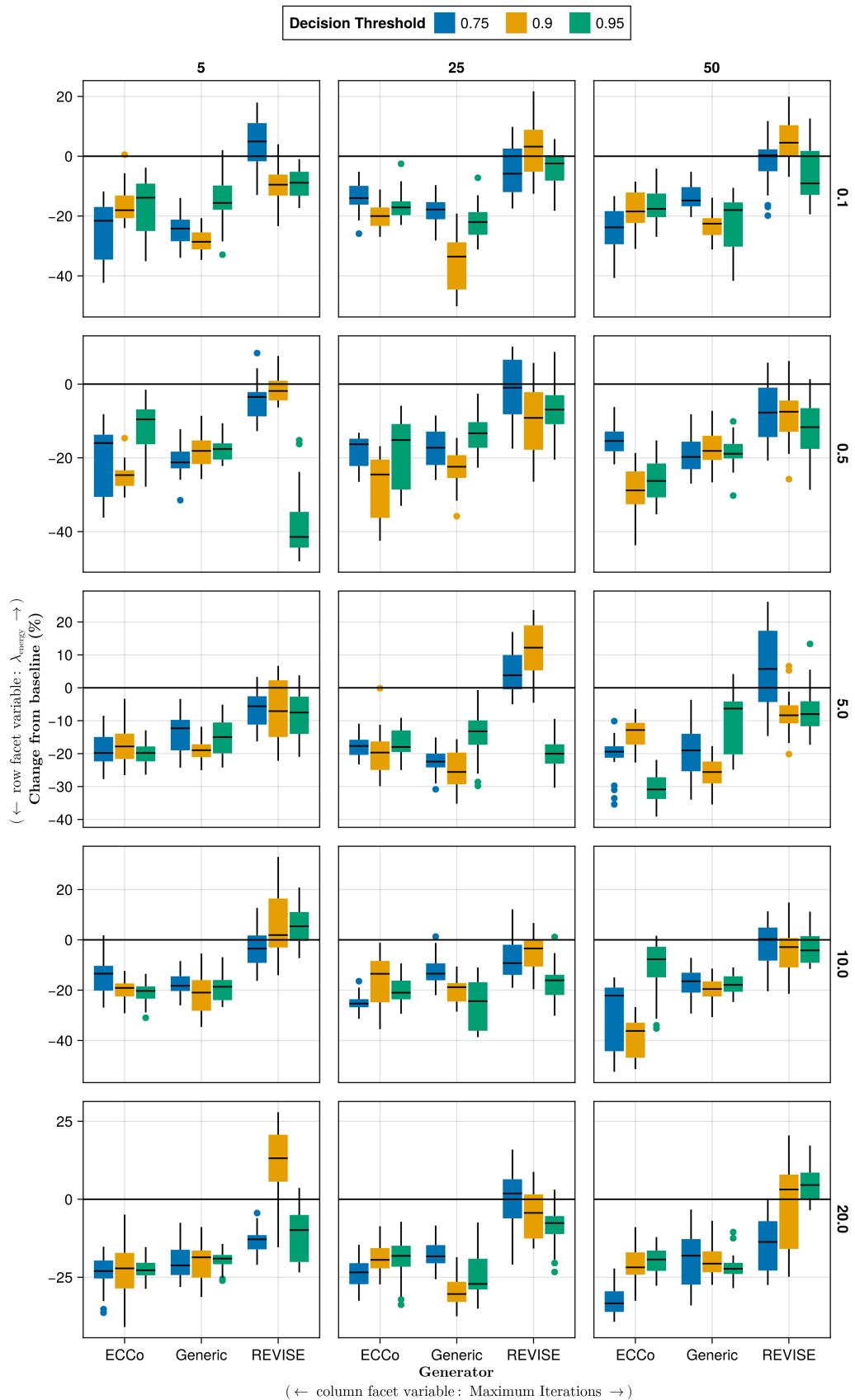


Figure A8: Average outcomes for the cost measure across hyperparameters. Data: Circles.

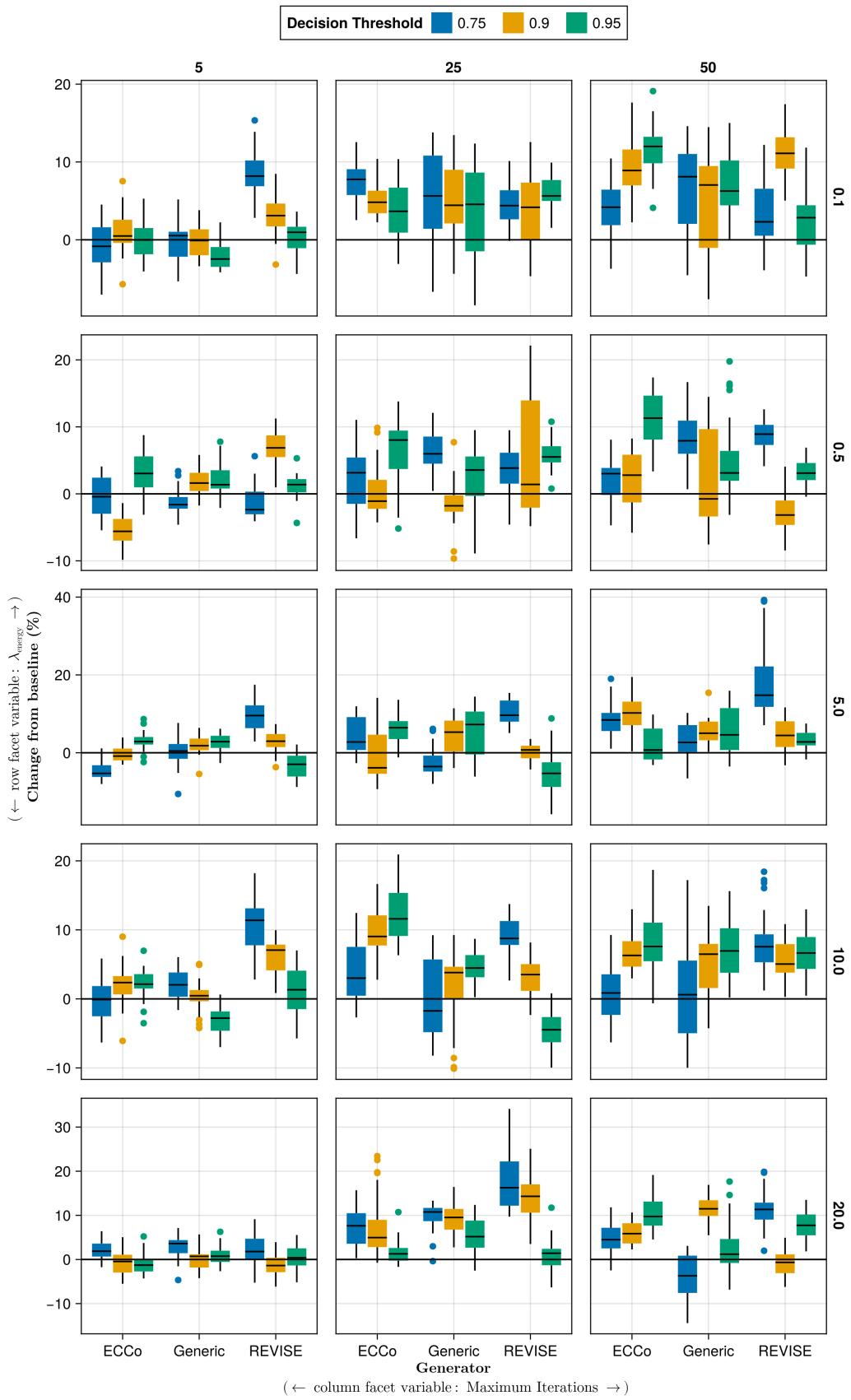


Figure A9: Average outcomes for the cost measure across hyperparameters. Data: Linearly Separable.

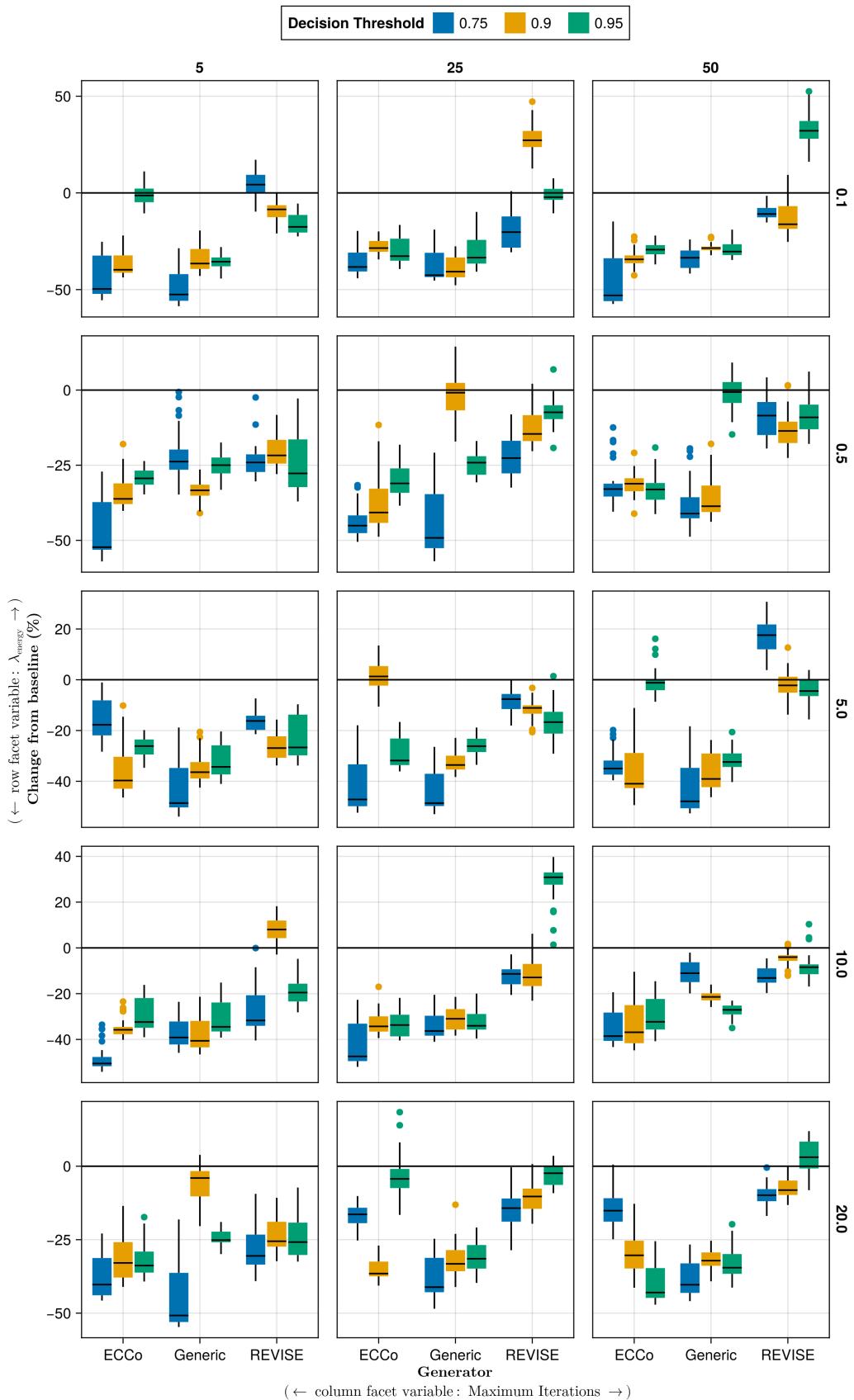


Figure A10: Average outcomes for the cost measure across hyperparameters. Data: Moons.

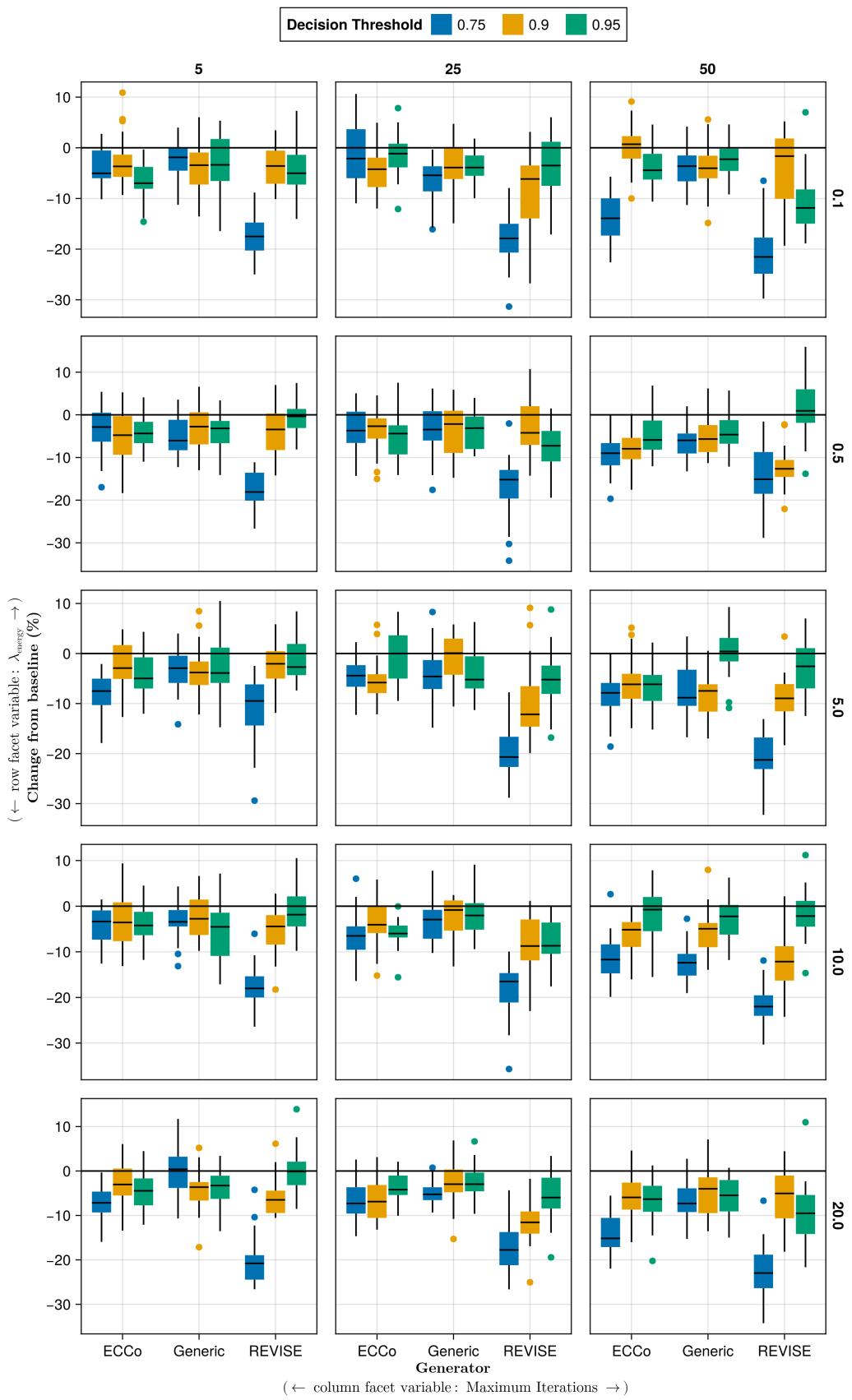


Figure A11: Average outcomes for the cost measure across hyperparameters. Data: Overlapping.

- 670
- λ_{div} : 0.01, 0.1, 1.0
 - λ_{reg} : 0.0, 0.01, 0.1, 0.25, 0.5
 - Objective: full, vanilla

671

Note 4: Evaluation Phase

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

672

J.3.1 Accuracy

Table A4: 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	0.994	0.0144
Circ	Accuracy	Vanilla	0.998	0.000875
Circ	F1-score	Full	0.994	0.0145
Circ	F1-score	Vanilla	0.998	0.000875
LS	Accuracy	Full	0.998	0.00772
LS	Accuracy	Vanilla	1	0
LS	F1-score	Full	0.998	0.00773
LS	F1-score	Vanilla	1	0
Moon	Accuracy	Full	0.987	0.0351
Moon	Accuracy	Vanilla	0.998	0.0101
Moon	F1-score	Full	0.987	0.0352
Moon	F1-score	Vanilla	0.998	0.0102
OL	Accuracy	Full	0.911	0.0217
OL	Accuracy	Vanilla	0.916	0.00236
OL	F1-score	Full	0.911	0.0219
OL	F1-score	Vanilla	0.916	0.00236

673

J.3.2 Plausibility

674 The results with respect to the plausibility measure are shown in Figure A12 to Figure A15.

675

J.3.3 Cost

676 The results with respect to the cost measure are shown in Figure A16 to Figure A19.

677

J.4 Other Parameters

678 The hyperparameter grid with other varying training parameters is shown in Note 5. The corresponding evaluation
679 grid used for these experiments is shown in Note 6.

Note 5: 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

680

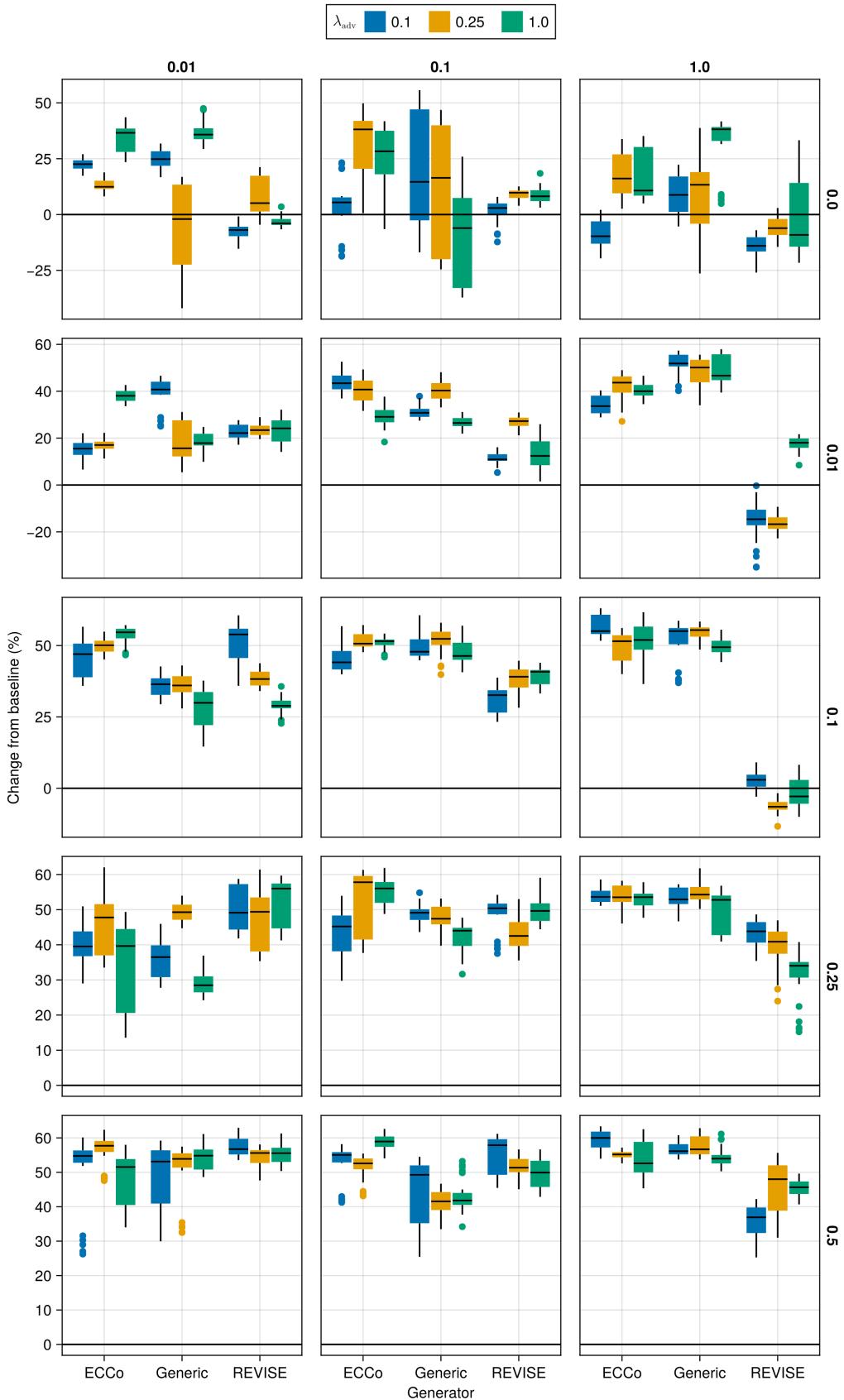


Figure A12: Average outcomes for the plausibility measure across hyperparameters. Data: Circles.

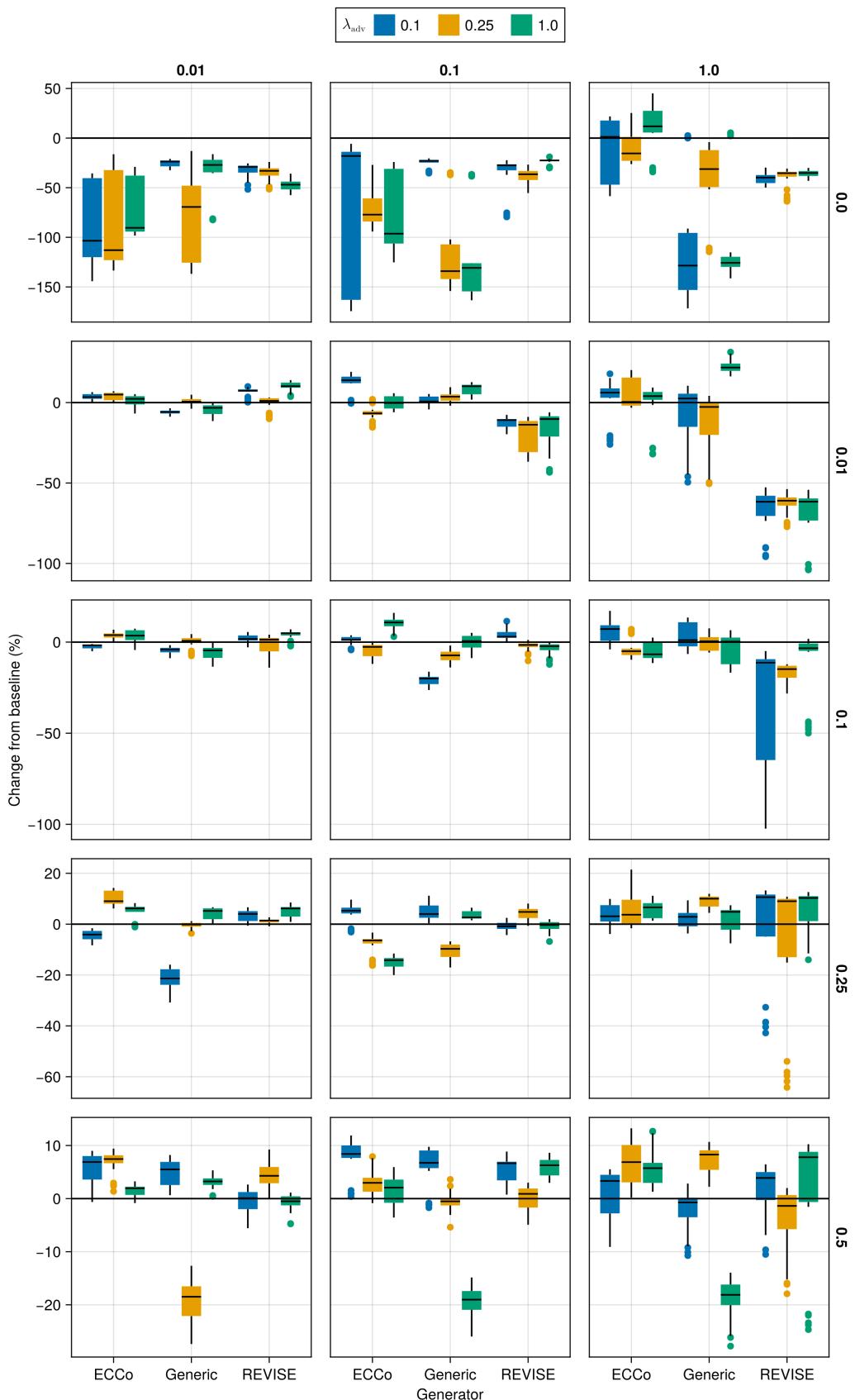


Figure A13: Average outcomes for the plausibility measure across hyperparameters. Data: Linearly Separable.

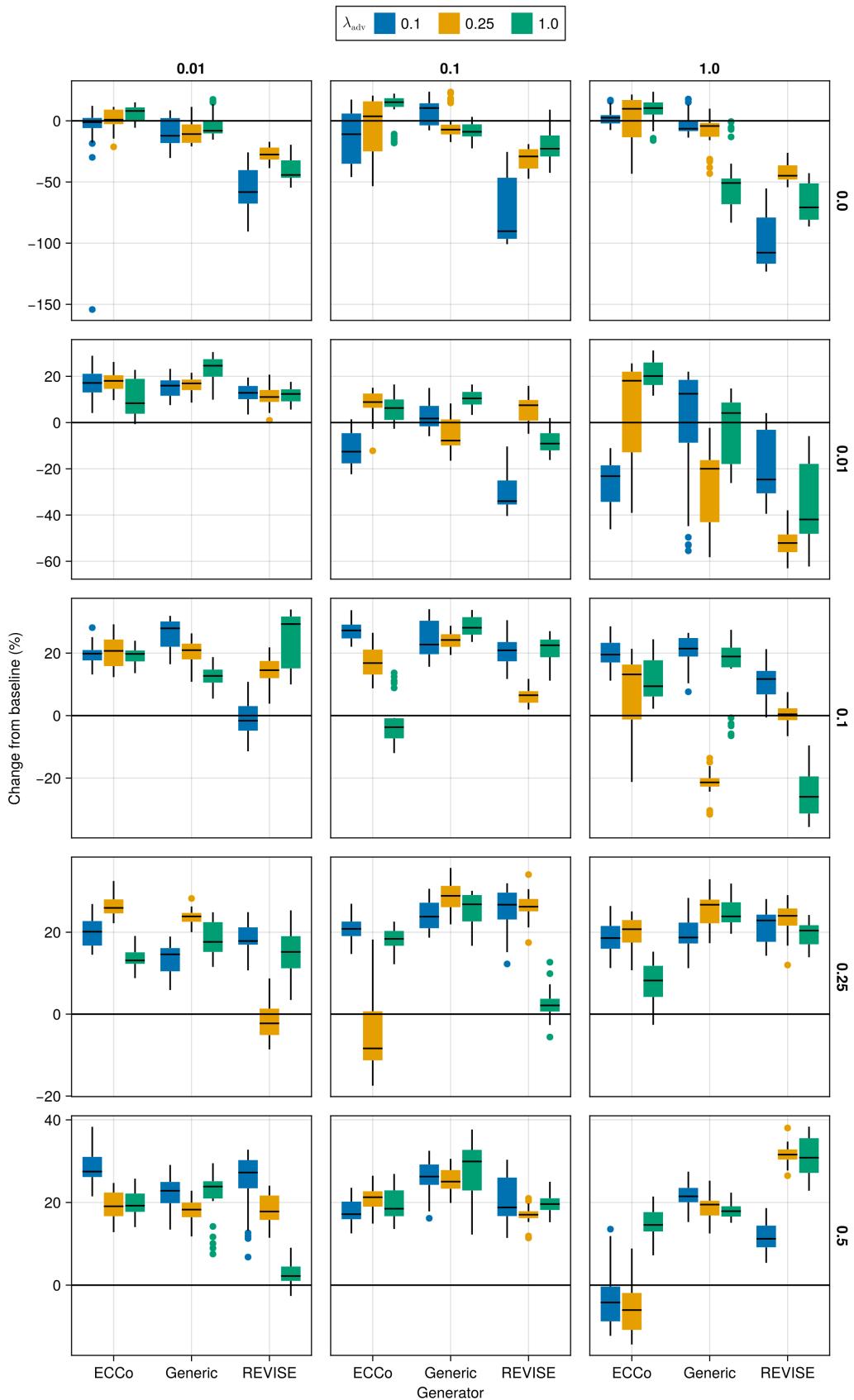


Figure A14: Average outcomes for the plausibility measure across hyperparameters. Data: Moons.

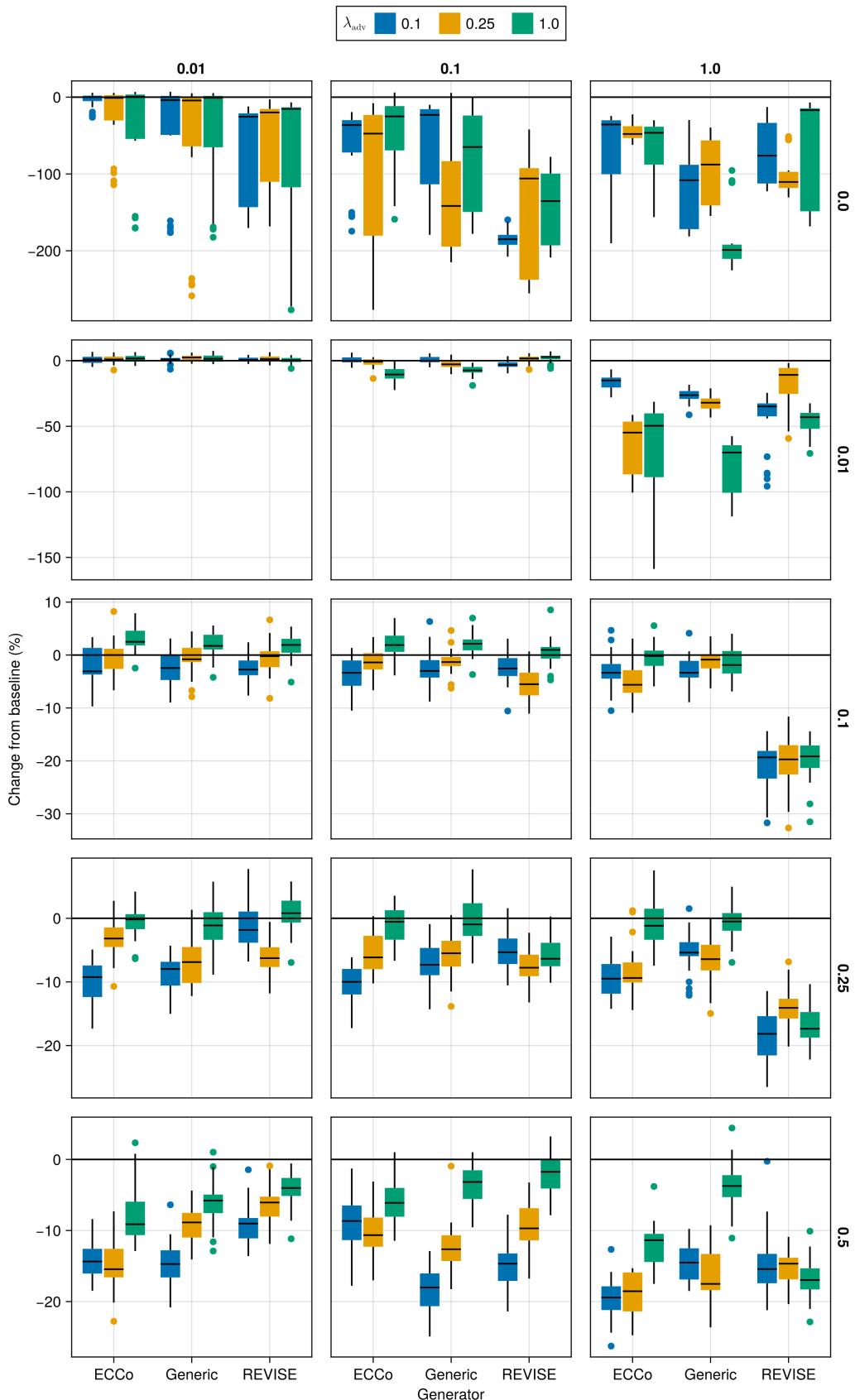


Figure A15: Average outcomes for the plausibility measure across hyperparameters. Data: Overlapping.

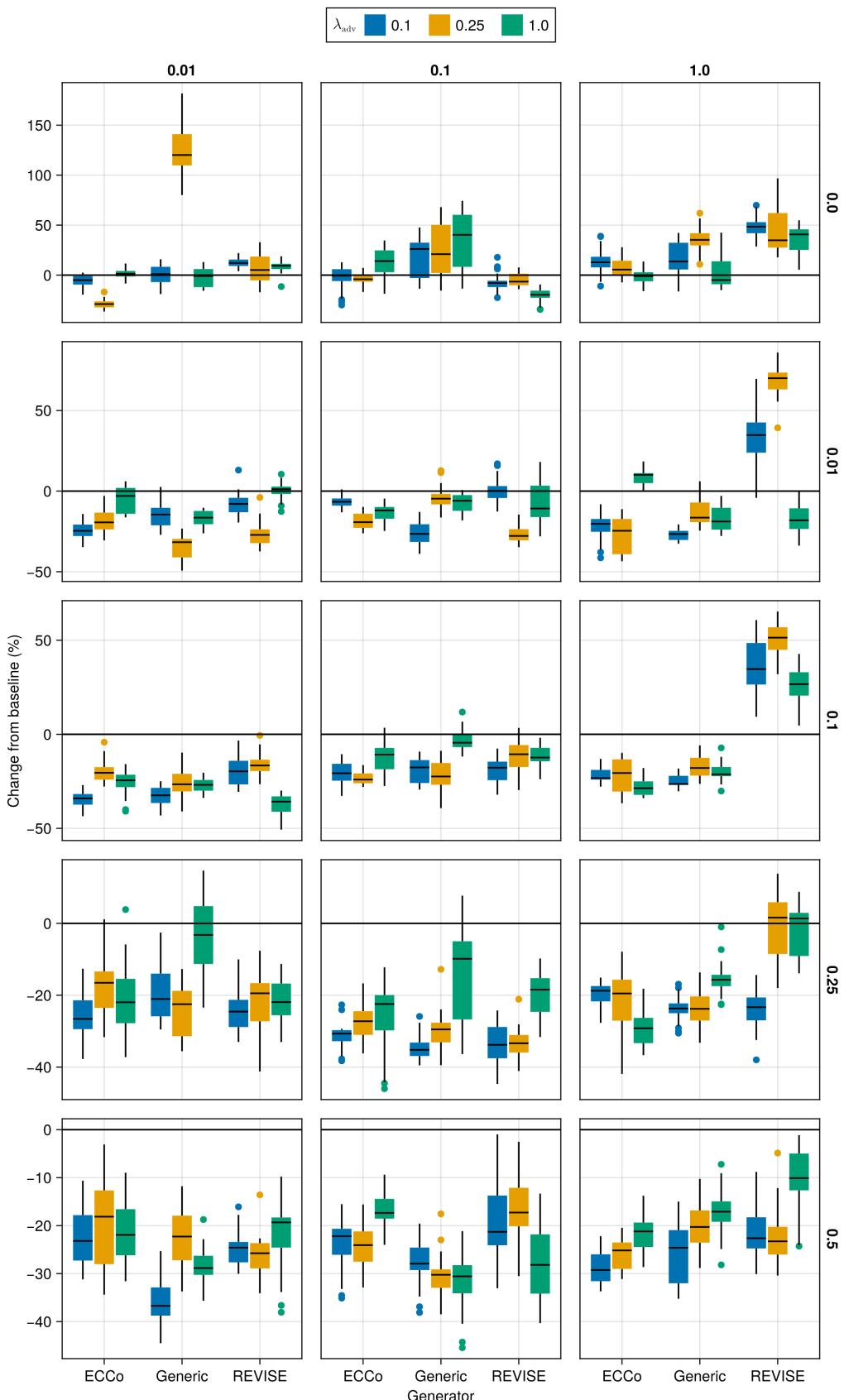


Figure A16: Average outcomes for the cost measure across hyperparameters. Data: Circles.

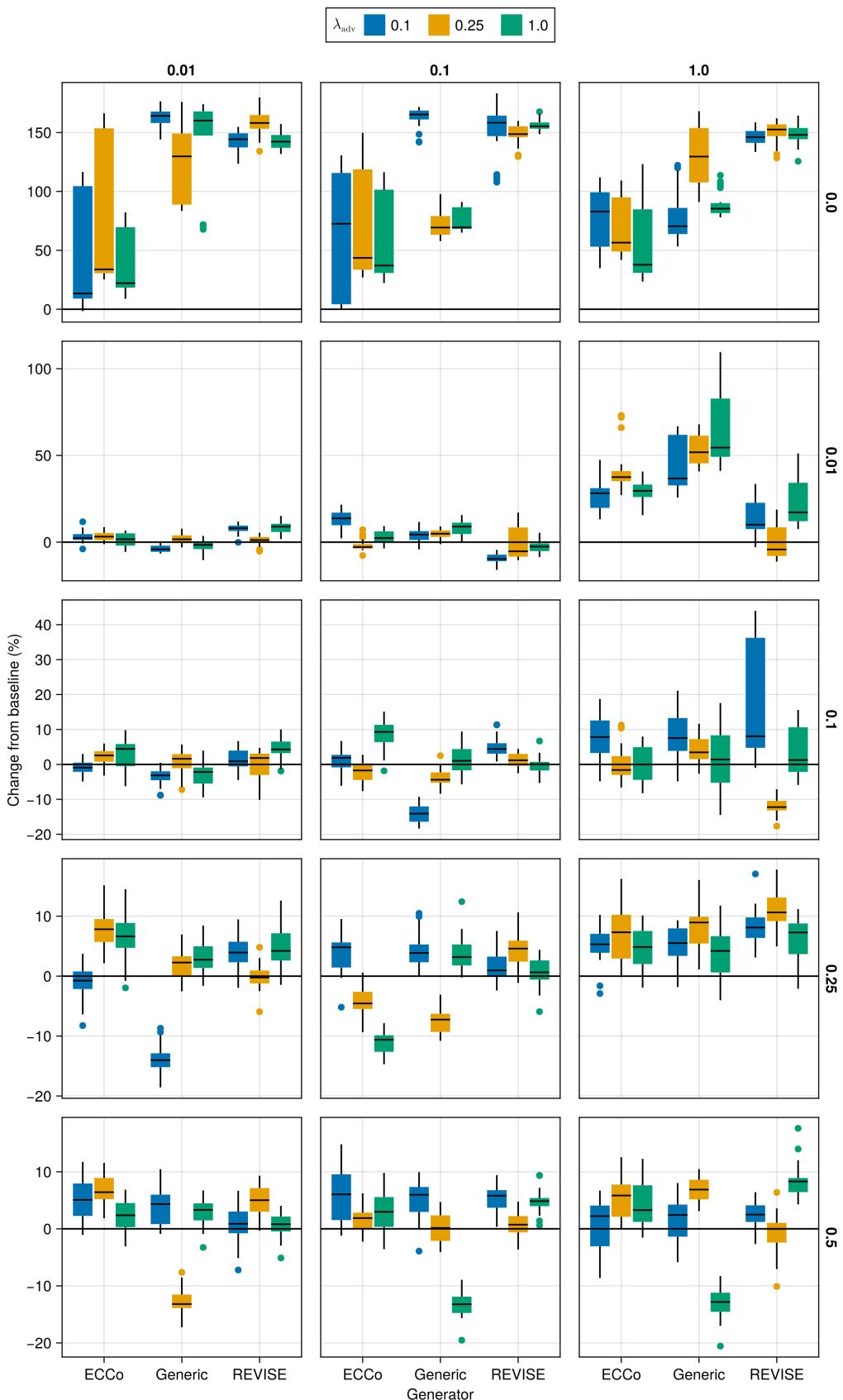


Figure A17: Average outcomes for the cost measure across hyperparameters. Data: Linearly Separable.

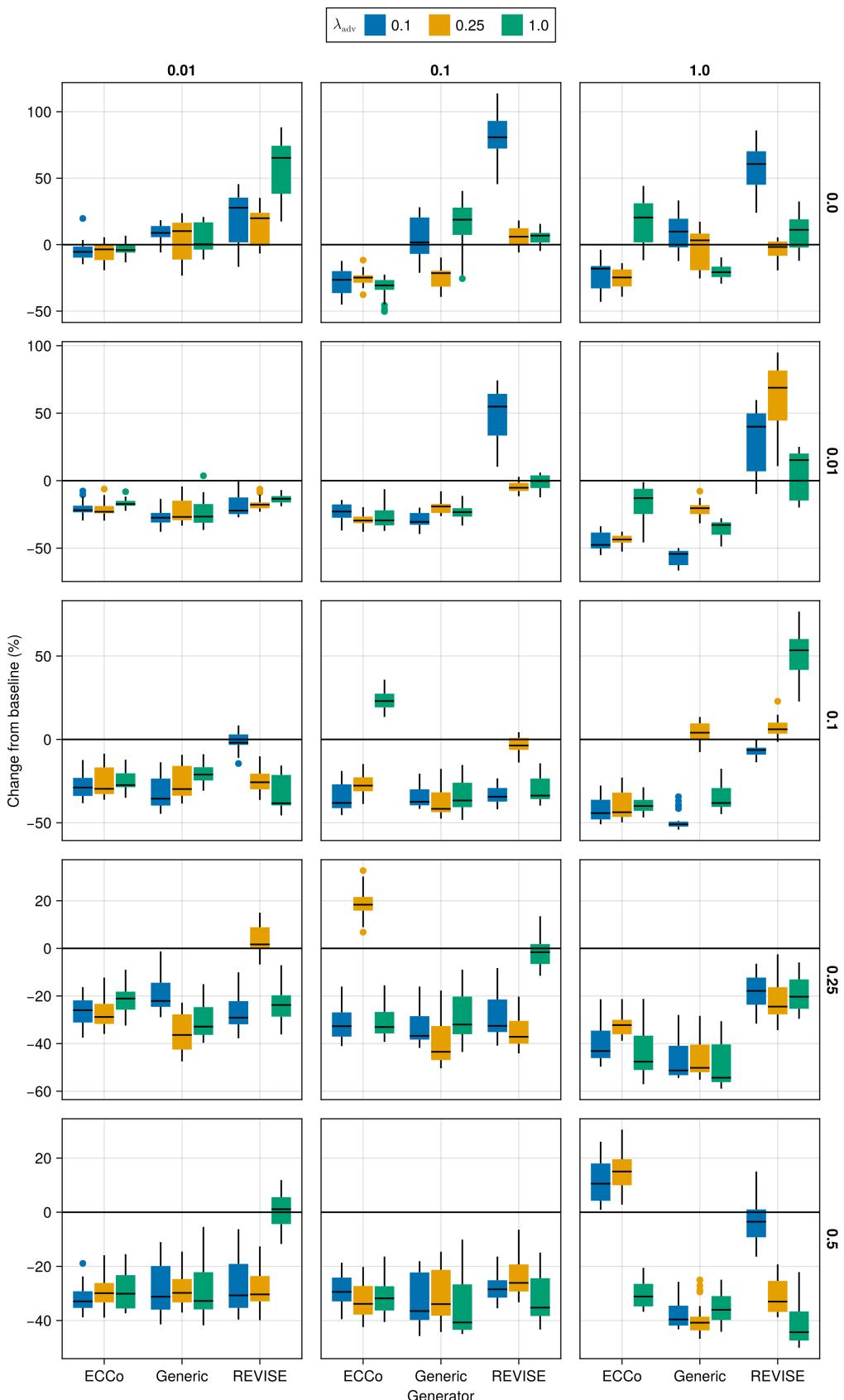


Figure A18: Average outcomes for the cost measure across hyperparameters. Data: Moons.

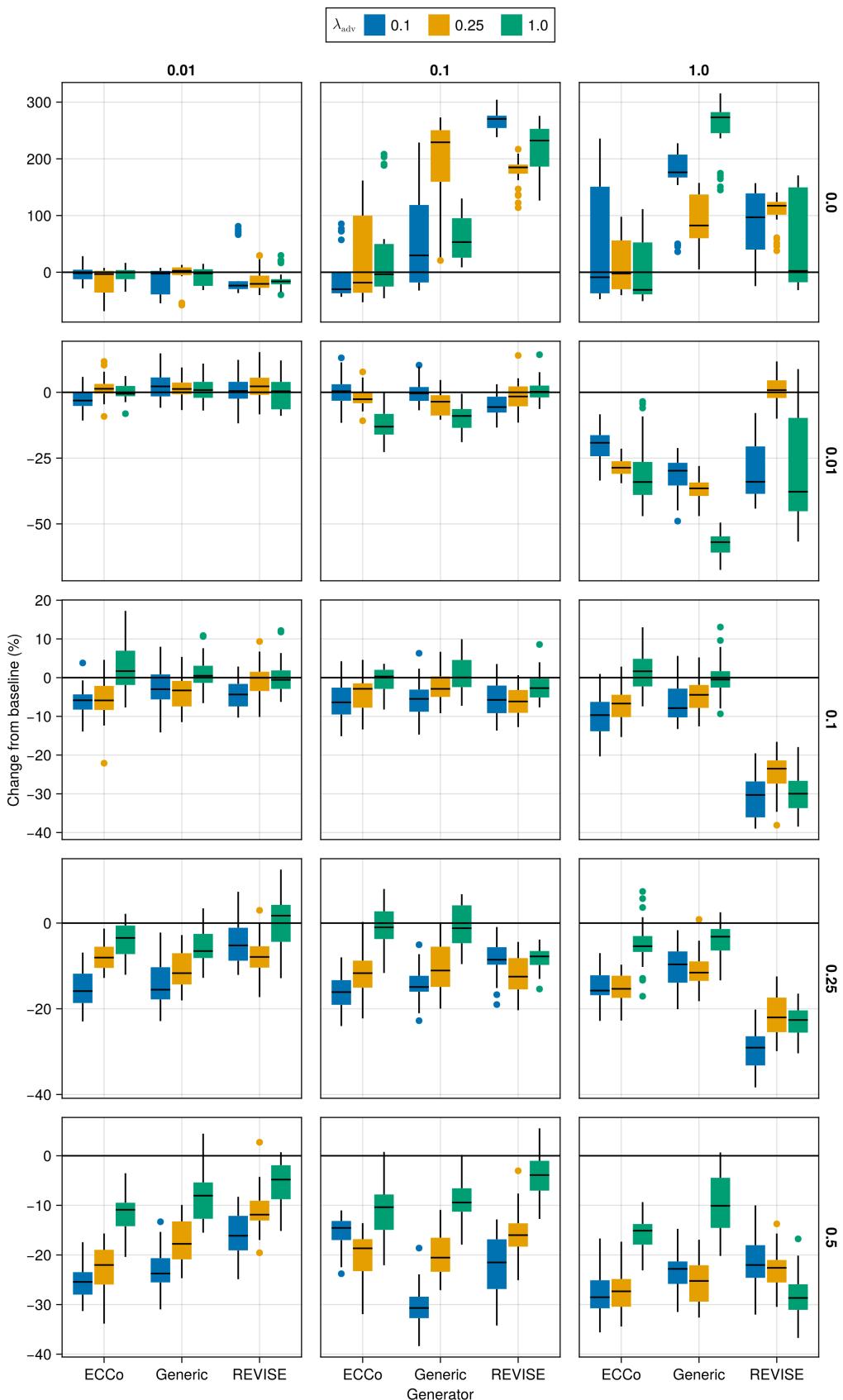


Figure A19: Average outcomes for the cost measure across hyperparameters. Data: Overlapping.

Note 6: Evaluation Phase

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

681

J.4.1 Accuracy

Table A5: 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	0.995	0.00431
Circ	Accuracy	Vanilla	0.998	0.000566
Circ	F1-score	Full	0.995	0.00432
Circ	F1-score	Vanilla	0.998	0.000566
LS	Accuracy	Full	0.999	0.00231
LS	Accuracy	Vanilla	1	0
LS	F1-score	Full	0.999	0.00231
LS	F1-score	Vanilla	1	0
Moon	Accuracy	Full	0.996	0.0136
Moon	Accuracy	Vanilla	0.988	0.022
Moon	F1-score	Full	0.996	0.0136
Moon	F1-score	Vanilla	0.988	0.022
OL	Accuracy	Full	0.914	0.00563
OL	Accuracy	Vanilla	0.918	0.00116
OL	F1-score	Full	0.914	0.0057
OL	F1-score	Vanilla	0.918	0.00116

J.4.2 Plausibility

The results with respect to the plausibility measure are shown in Figure A20 to Figure A23.

J.4.3 Cost

The results with respect to the cost measure are shown in Figure A24 to Figure A27.

K Tuning Key Parameters

Based on the findings from our initial large grid searches (Section J), we tune selected hyperparameters for all datasets: namely, the decision threshold τ and the strength of the energy regularization λ_{reg} . The final hyperparameter choices for each dataset are presented in **ADD TABLE**. Detailed results for each data set are shown in Figure A28 to Figure A45. From **ADD TABLE**, we notice that the same decision threshold of $\tau = 0.5$ is optimal for all but one dataset. We attribute this to the fact that a low decision threshold results in a higher share of mature counterfactuals and hence more opportunities for the model to learn from examples (Figure A37 to Figure A45). This has played a role in particular for our real-world tabular datasets and MNIST, which suffered from low levels of maturity for higher decision thresholds. In cases where maturity is not an issue, as for *Moons*, higher decision thresholds lead to better outcomes, which may have to do with the fact that the resulting counterfactuals are more faithful to the model. Concerning the regularization strength, we find somewhat high variation across datasets. Most notably, we find that relatively low levels of regularization are optimal for MNIST. We hypothesize that this finding may be attributed to the uniform scaling of all input features (digits).

Finally, to increase the proportion of mature counterfactuals for some datasets, we have also investigated the effect on the learning rate η for the counterfactual search and even smaller regularization strengths for a fixed decision threshold of 0.5 (Figure A46 to Figure A51). For the given low decision threshold, we find that the learning rate has no discernable impact on the proportion of mature counterfactuals (Figure A52 to Figure A57). We do notice, however, that the results for MNIST are much improved when using a low value λ_{reg} , the strength for the energy regularization: plausibility is increased by up to ~10% (Figure A50) and the proportion of mature counterfactuals reaches 100%.

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

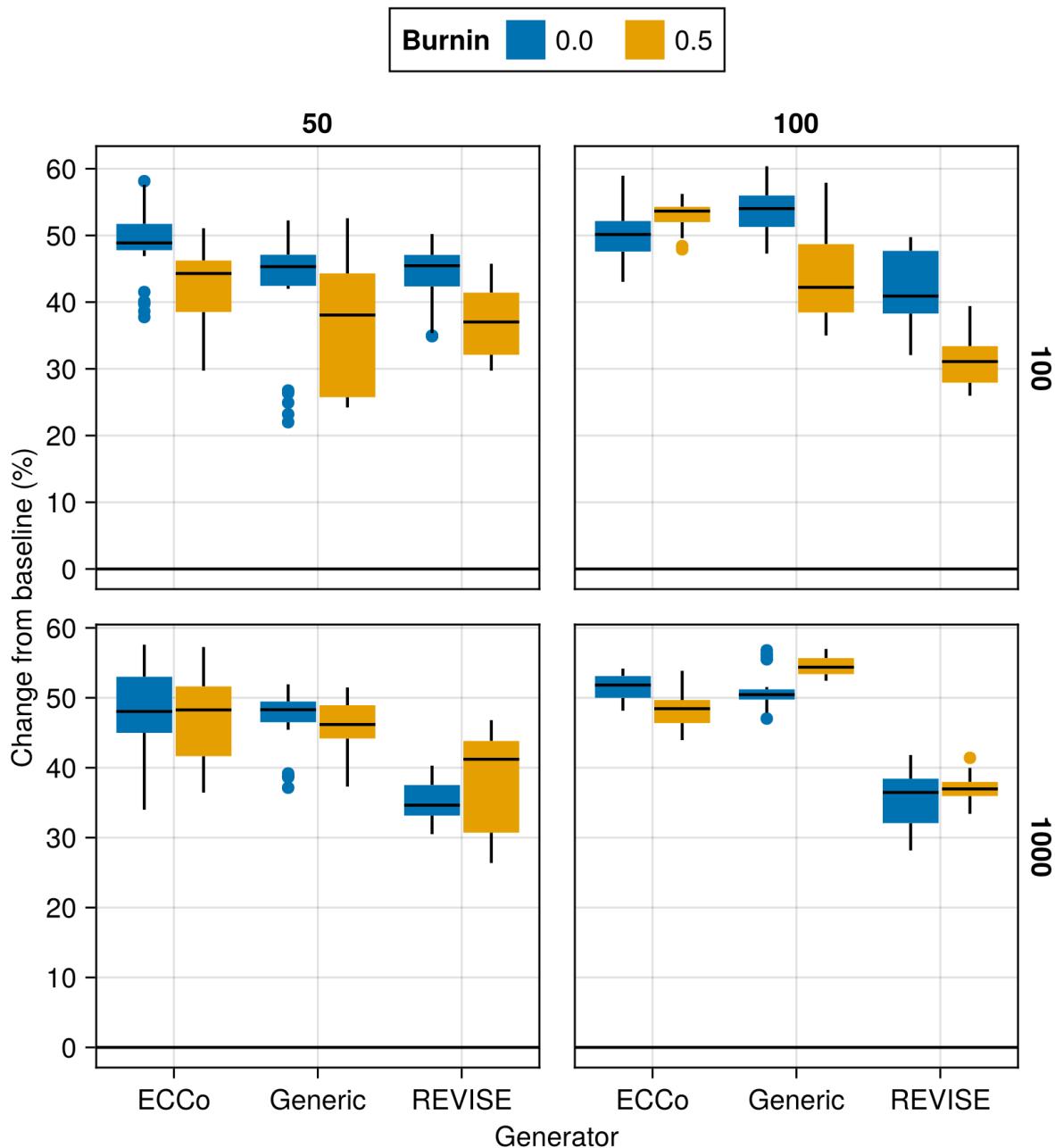


Figure A20: Average outcomes for the plausibility measure across hyperparameters. Data: Circles.

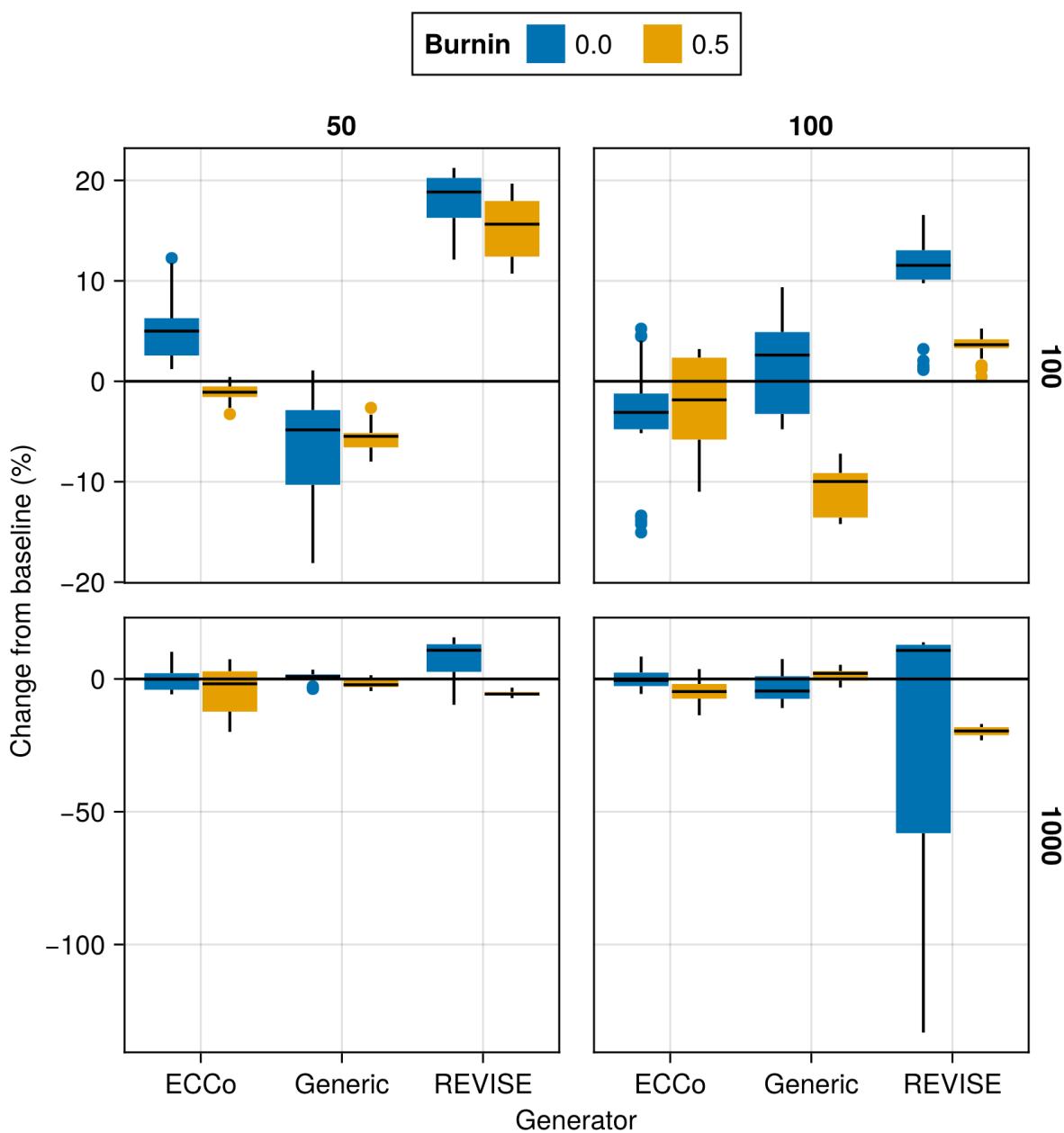


Figure A21: Average outcomes for the plausibility measure across hyperparameters. Data: Linearly Separable.

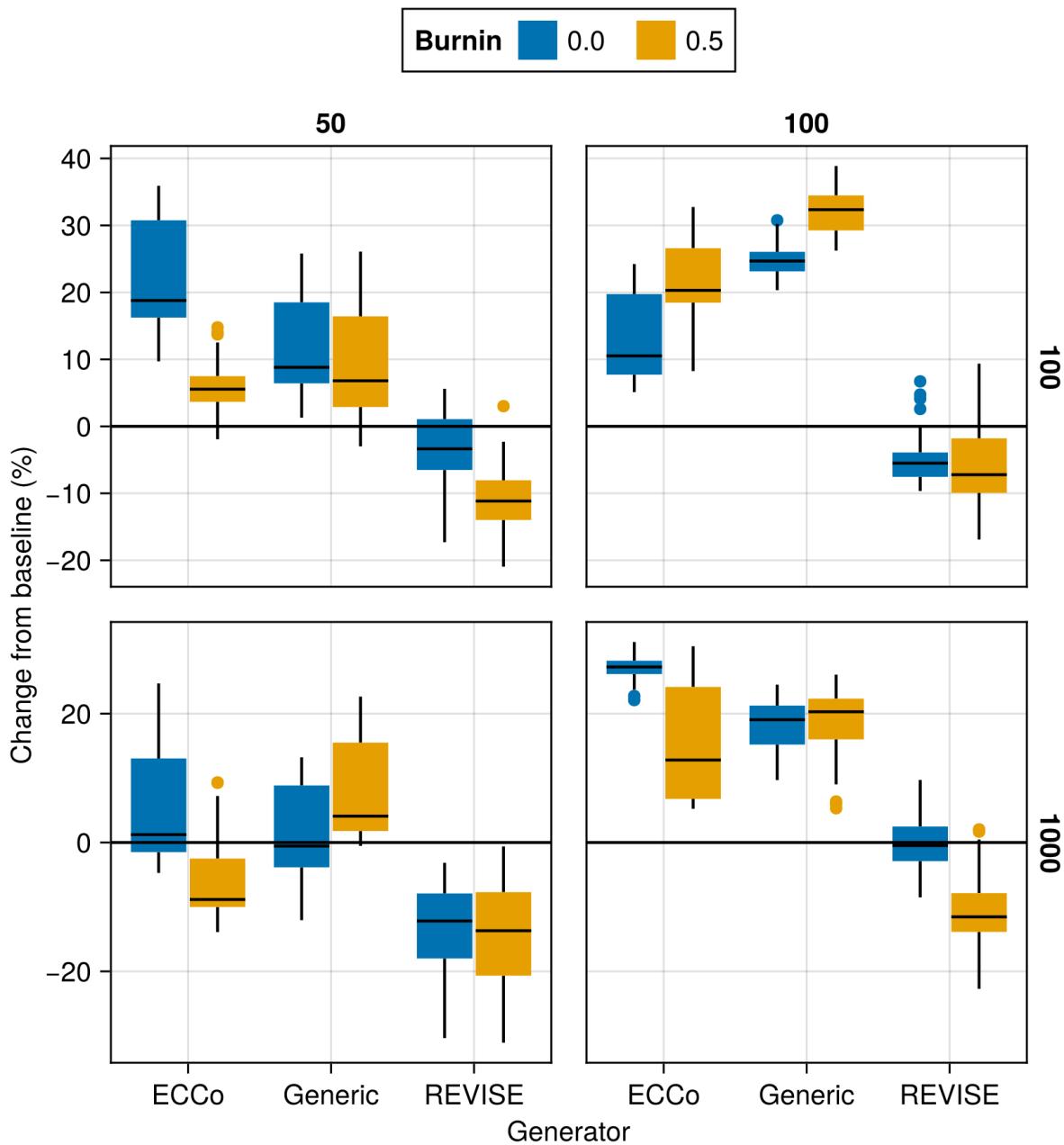


Figure A22: Average outcomes for the plausibility measure across hyperparameters. Data: Moons.

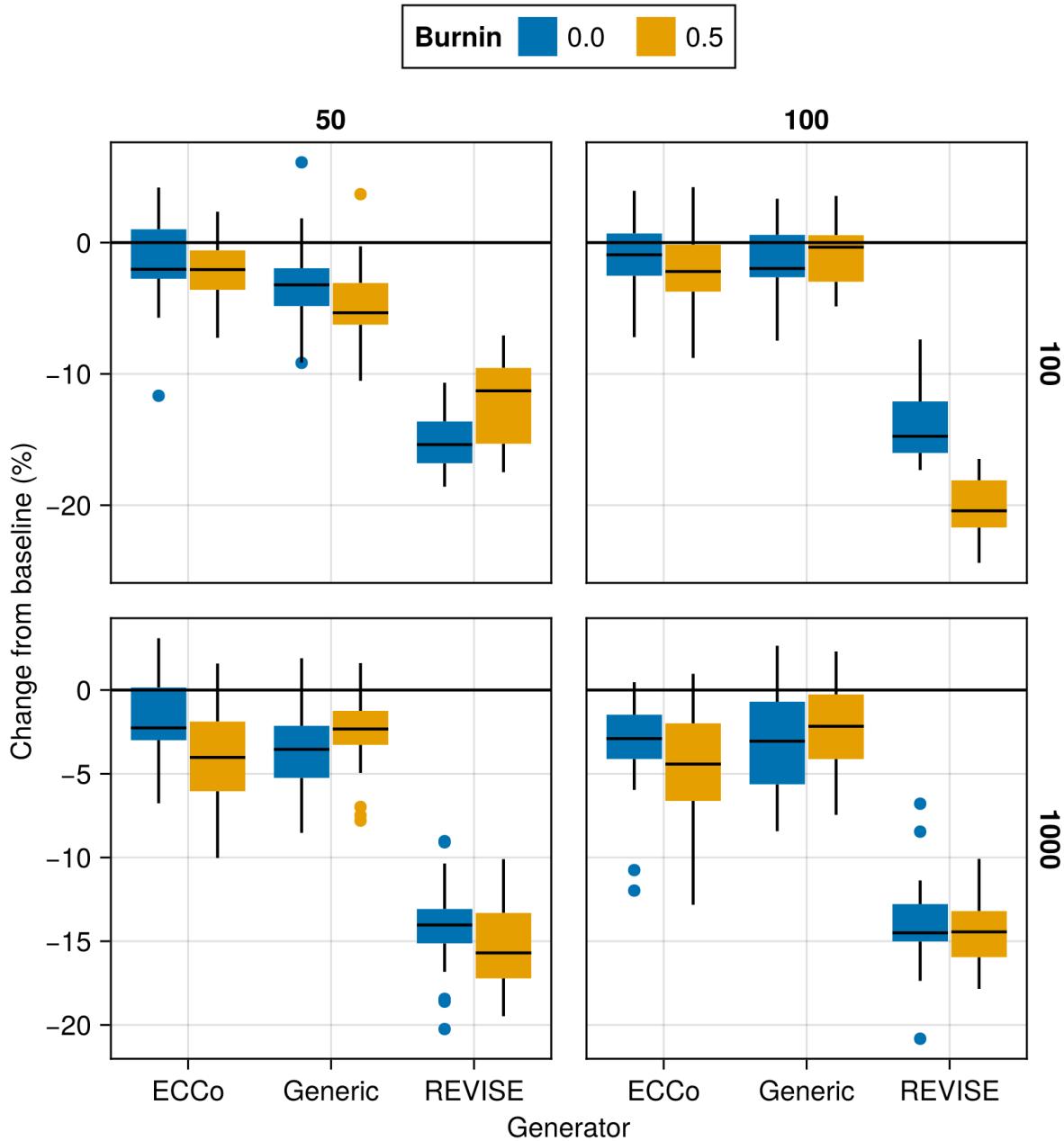


Figure A23: Average outcomes for the plausibility measure across hyperparameters. Data: Overlapping.

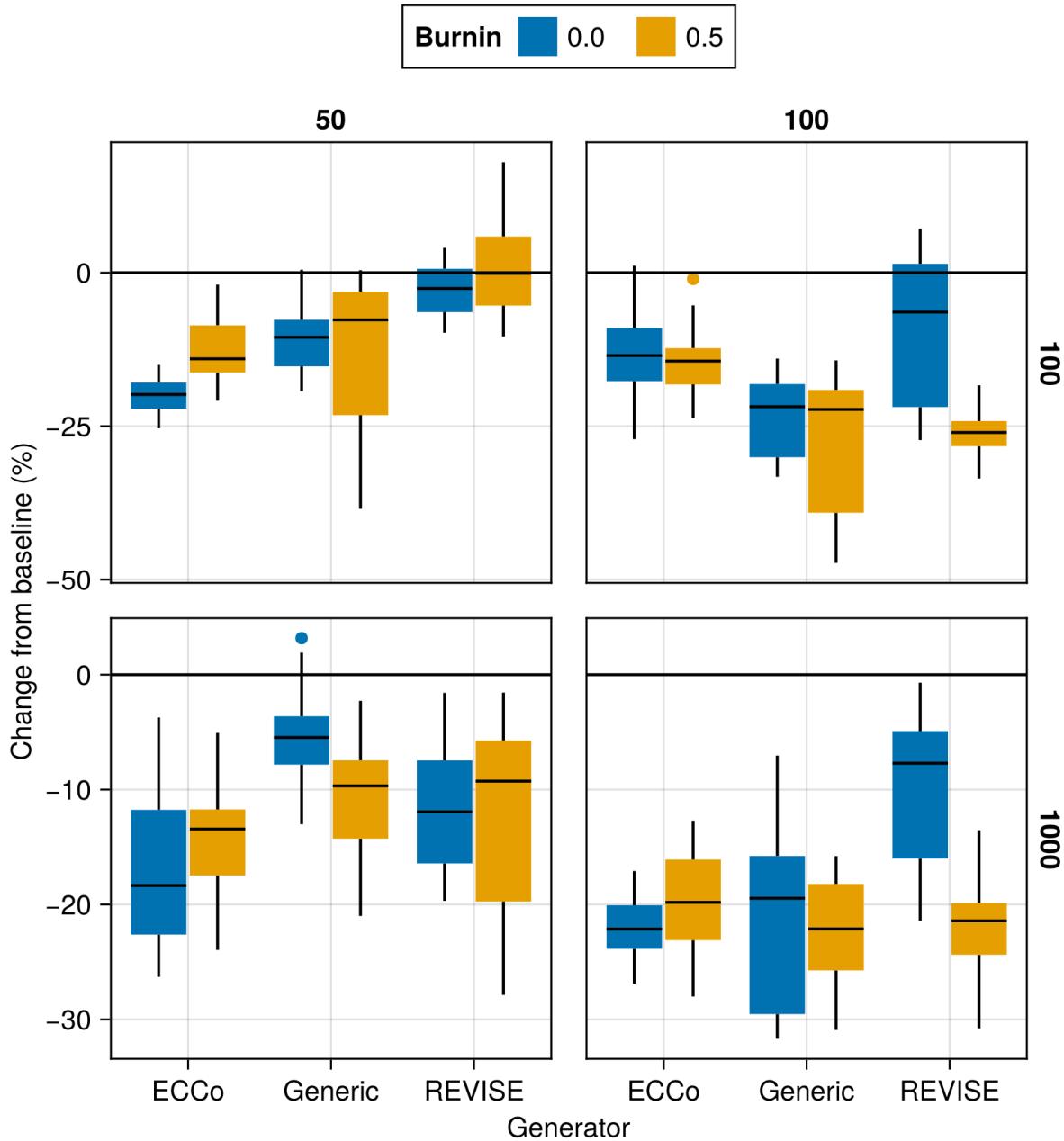


Figure A24: Average outcomes for the cost measure across hyperparameters. Data: Circles.

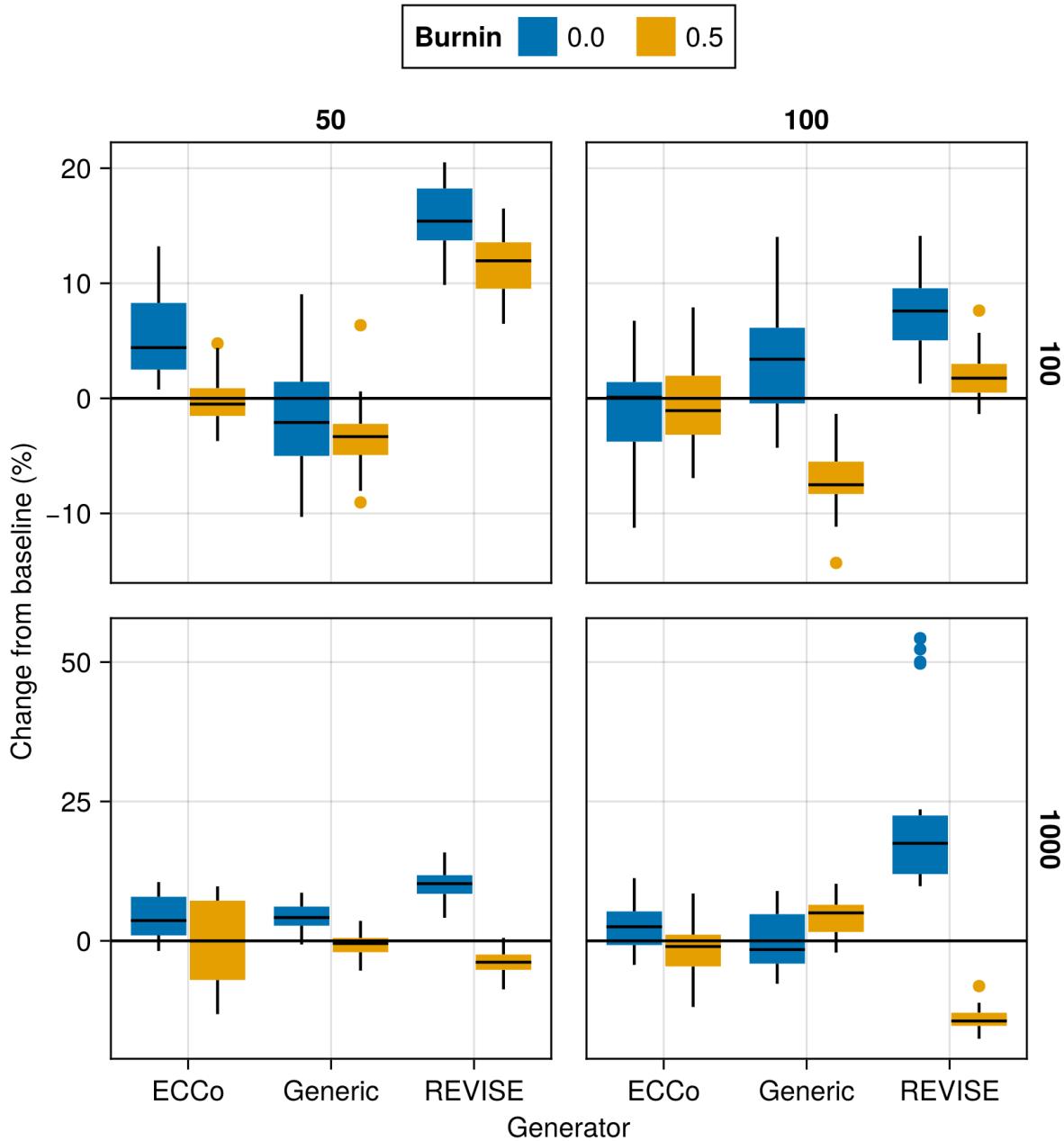


Figure A25: Average outcomes for the cost measure across hyperparameters. Data: Linearly Separable.

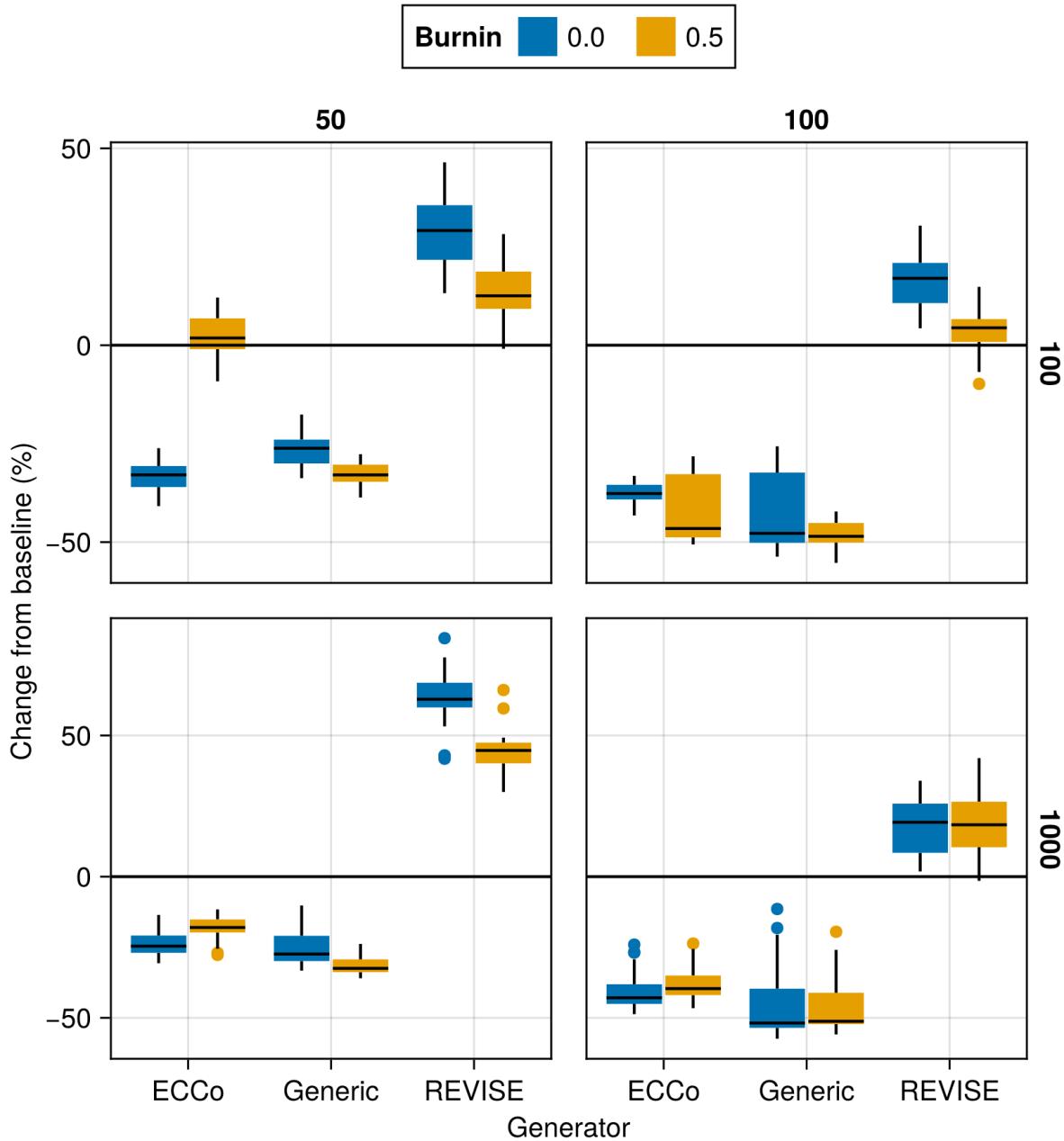


Figure A26: Average outcomes for the cost measure across hyperparameters. Data: Moons.

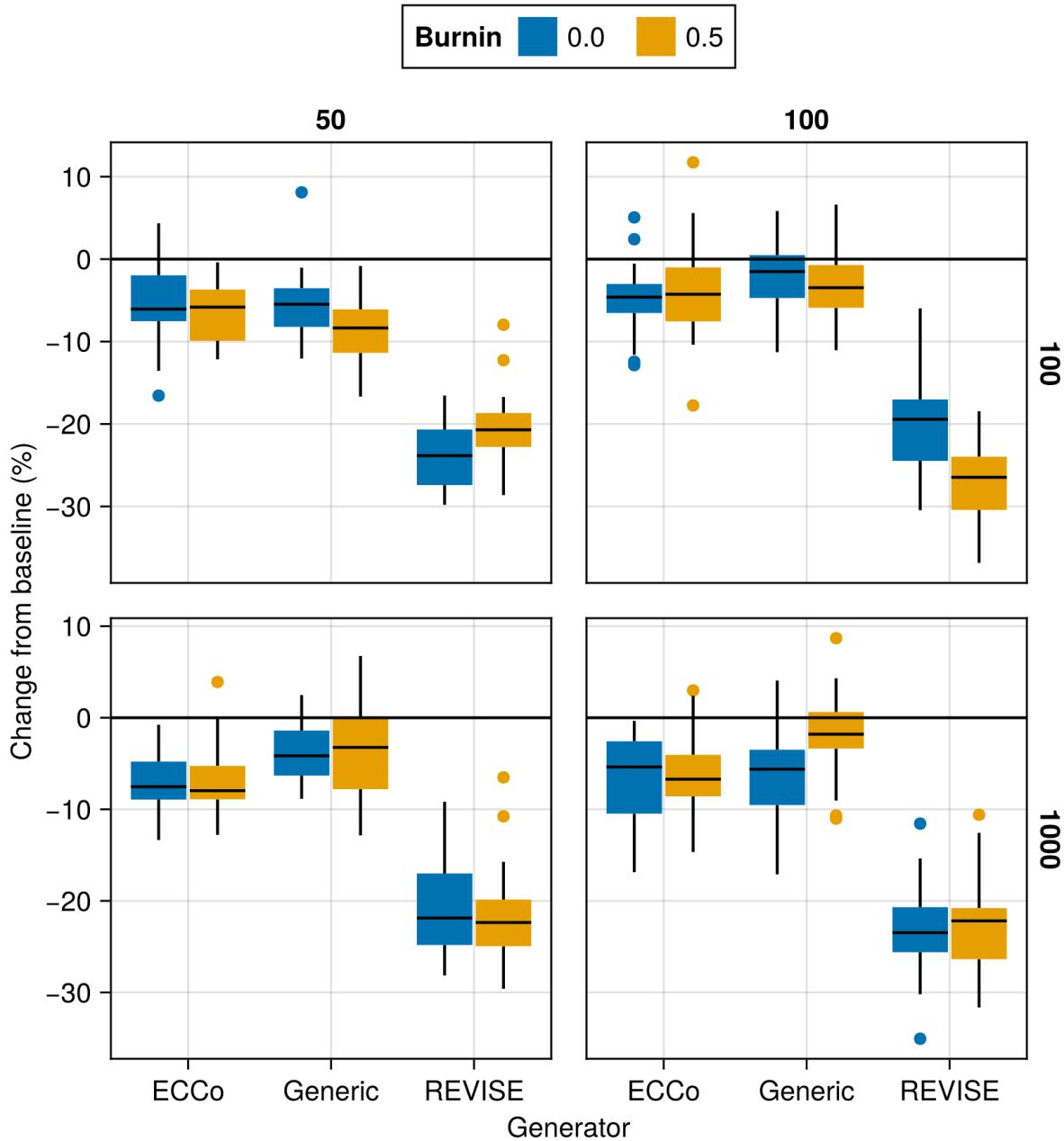


Figure A27: Average outcomes for the cost measure across hyperparameters. Data: Overlapping.

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.

708

709 K.1 Key Parameters

710 The hyperparameter grid for tuning key parameters is shown in Note 7. The corresponding evaluation grid used for
711 these experiments is shown in Note 8.

Note 7: 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`

712

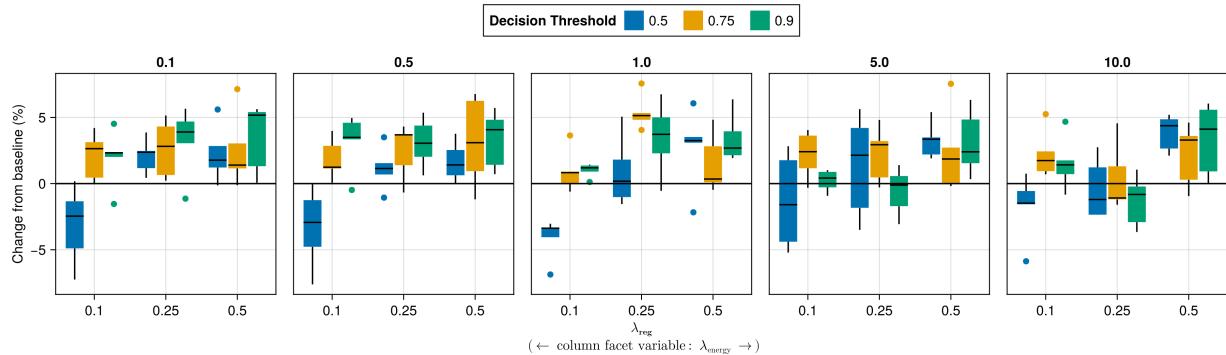
Note 8: Evaluation Phase

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

713

K.1.1 Plausibility

714 The results with respect to the plausibility measure are shown in Figure A28 to Figure A36.



715 Figure A28: Average outcomes for the plausibility measure across key hyperparameters. Data: Adult.

716 K.1.2 Proportion of Mature CE

717 The results with respect to the proportion of mature counterfactuals in each epoch are shown in Figure A37 to Figure
718 A45.

719 K.2 Learning Rate

720 The hyperparameter grid for tuning the learning rate is shown in Note 9. The corresponding evaluation grid used for
721 these experiments is shown in Note 10.

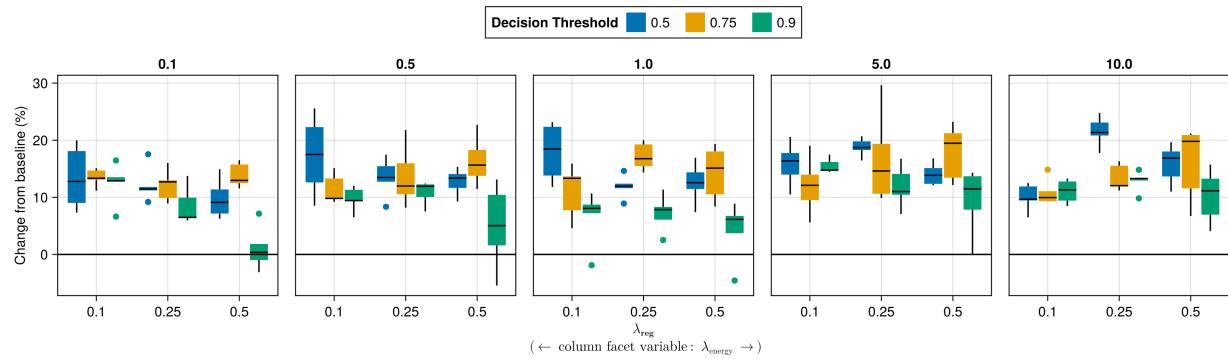


Figure A29: Average outcomes for the plausibility measure across key hyperparameters. Data: California Housing.

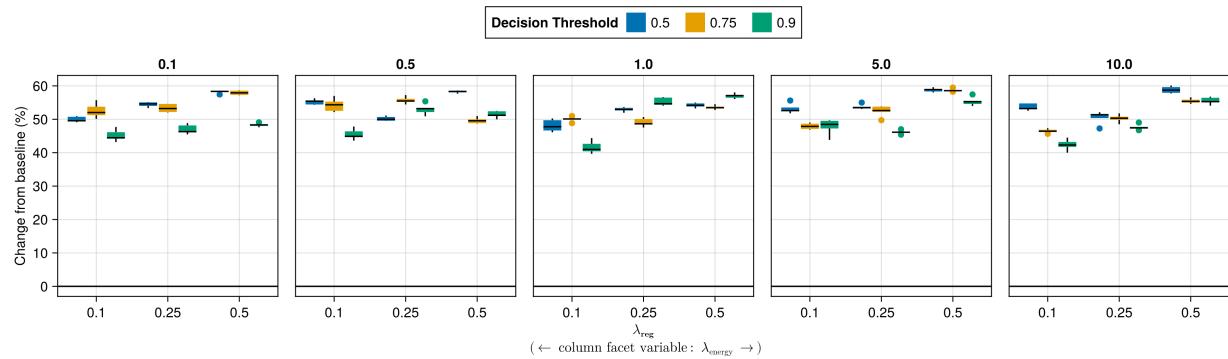


Figure A30: Average outcomes for the plausibility measure across key hyperparameters. Data: Circles.

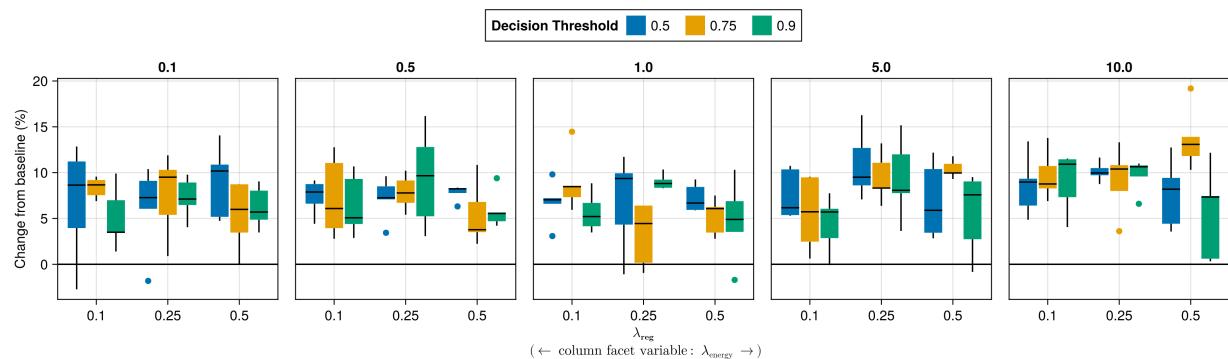


Figure A31: Average outcomes for the plausibility measure across key hyperparameters. Data: Credit.

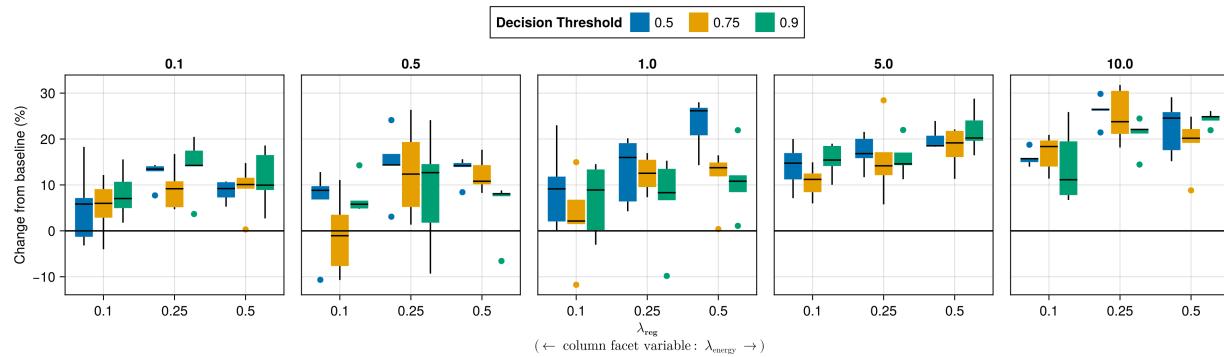


Figure A32: Average outcomes for the plausibility measure across key hyperparameters. Data: GMSC.

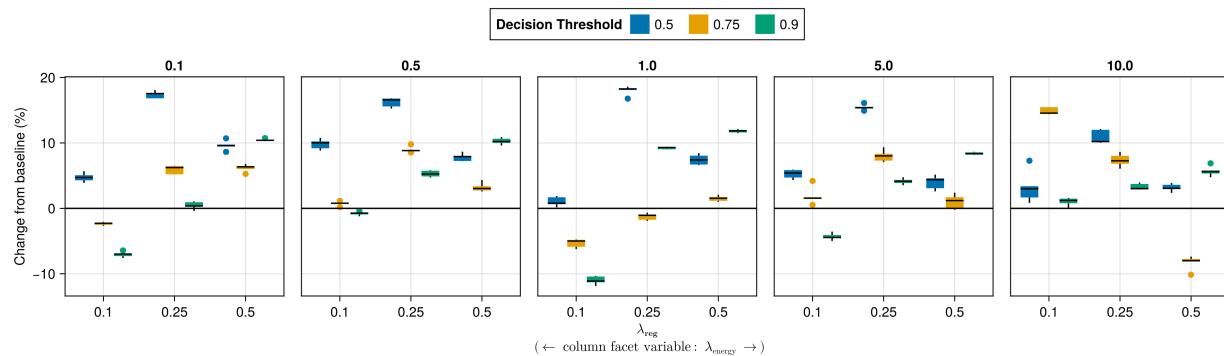


Figure A33: Average outcomes for the plausibility measure across key hyperparameters. Data: Linearly Separable.

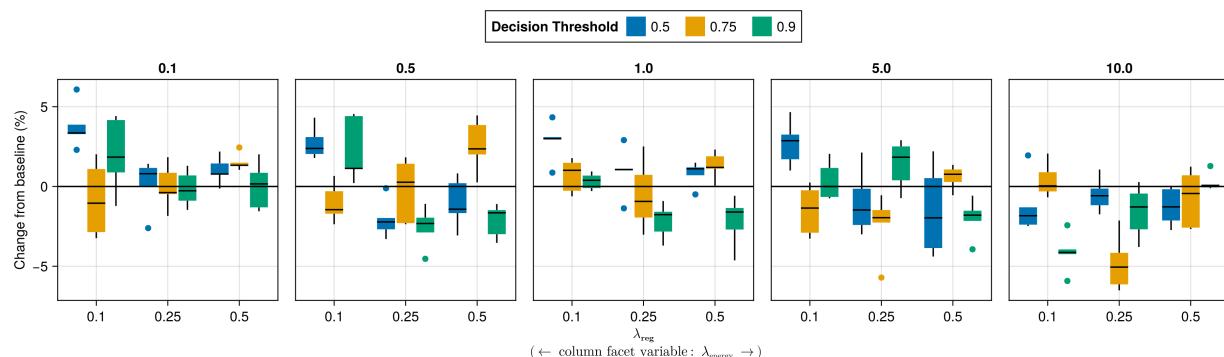


Figure A34: Average outcomes for the plausibility measure across key hyperparameters. Data: MNIST.

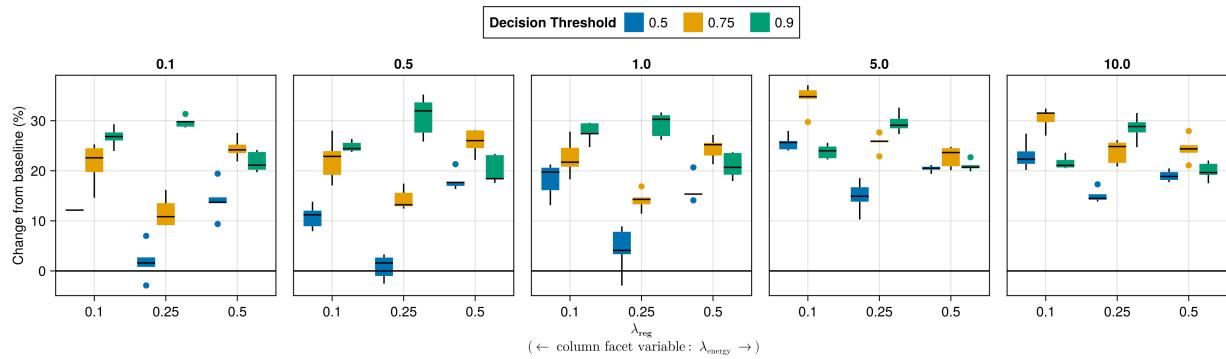


Figure A35: Average outcomes for the plausibility measure across key hyperparameters. Data: Moons.

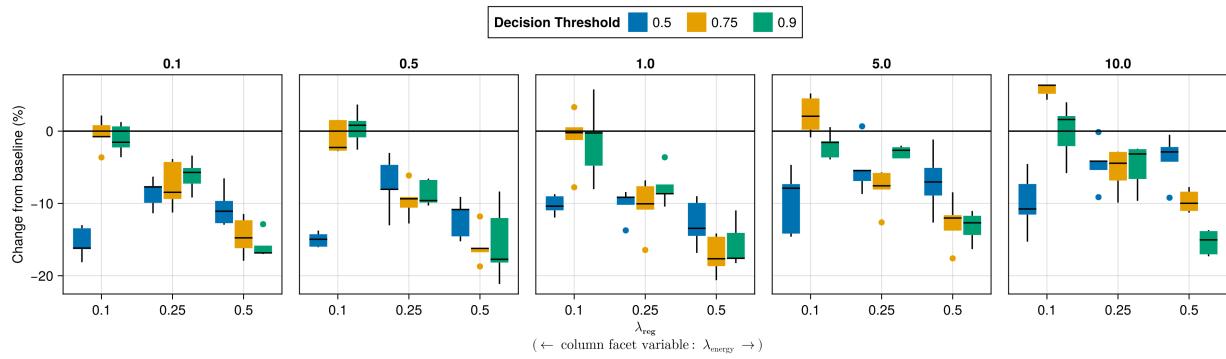


Figure A36: Average outcomes for the plausibility measure across key hyperparameters. Data: Overlapping.

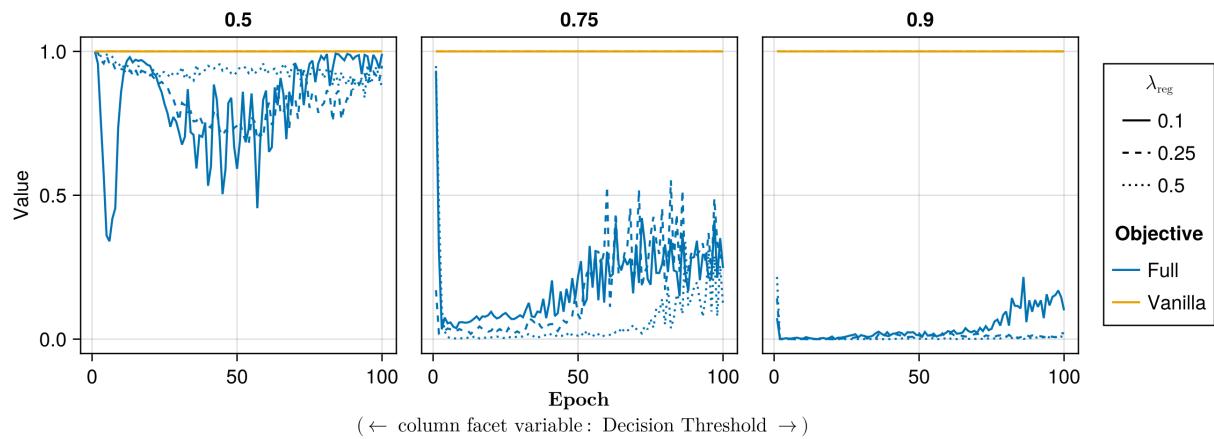


Figure A37: Proportion of mature counterfactuals in each epoch. Data: Adult.

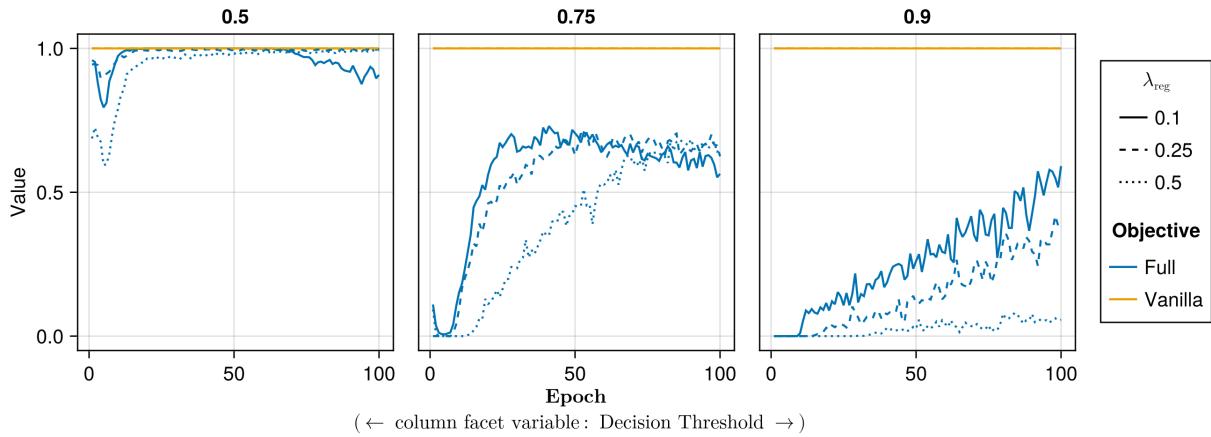


Figure A38: Proportion of mature counterfactuals in each epoch. Data: California Housing.

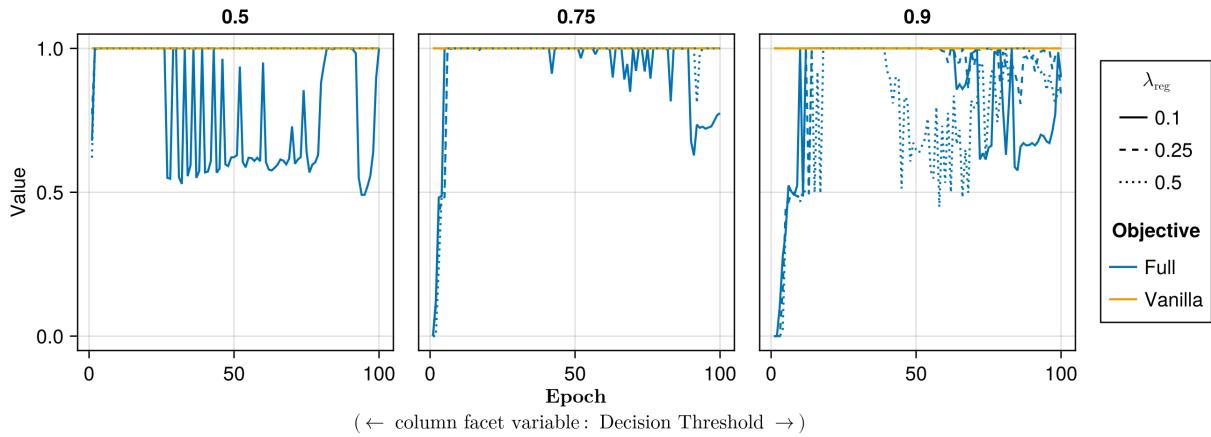


Figure A39: Proportion of mature counterfactuals in each epoch. Data: Circles.

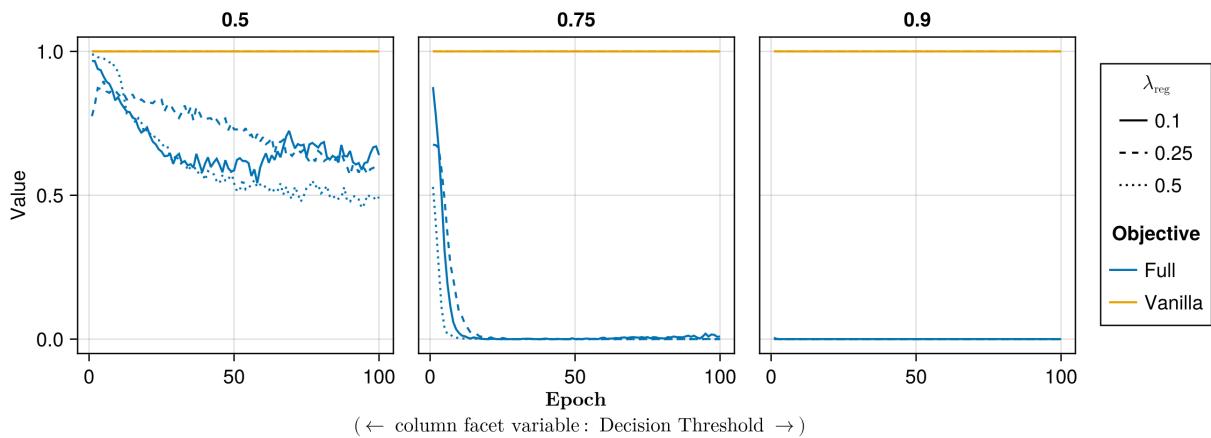


Figure A40: Proportion of mature counterfactuals in each epoch. Data: Credit.

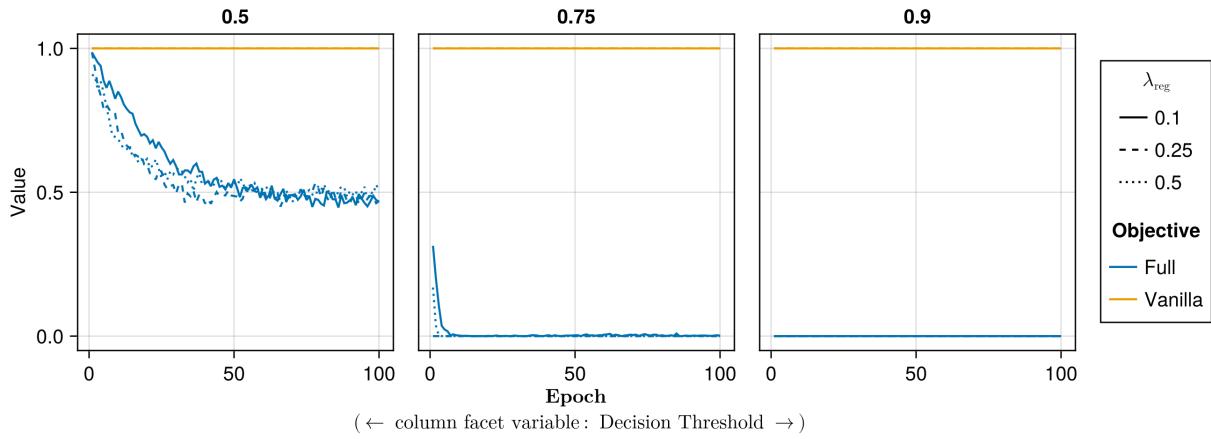


Figure A41: Proportion of mature counterfactuals in each epoch. Data: GMSC.

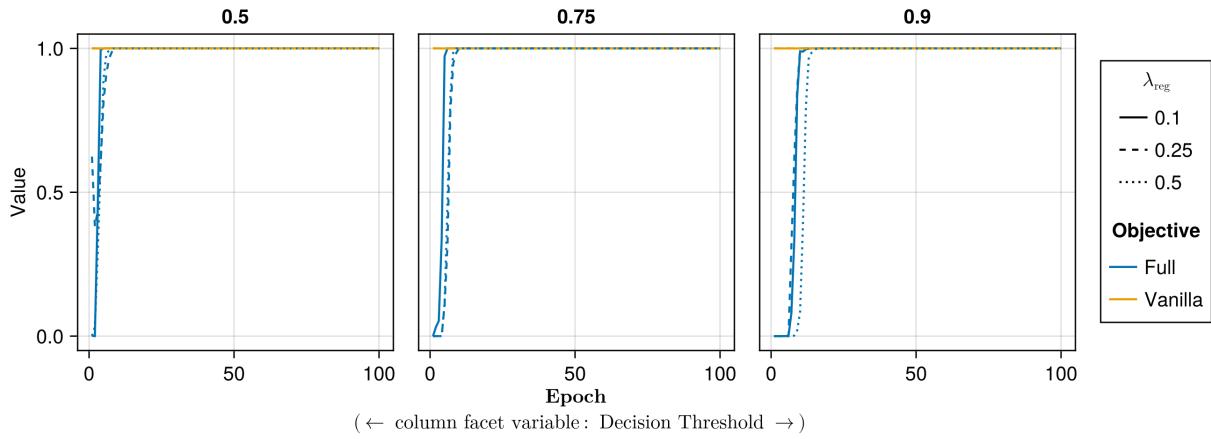


Figure A42: Proportion of mature counterfactuals in each epoch. Data: Linearly Separable.

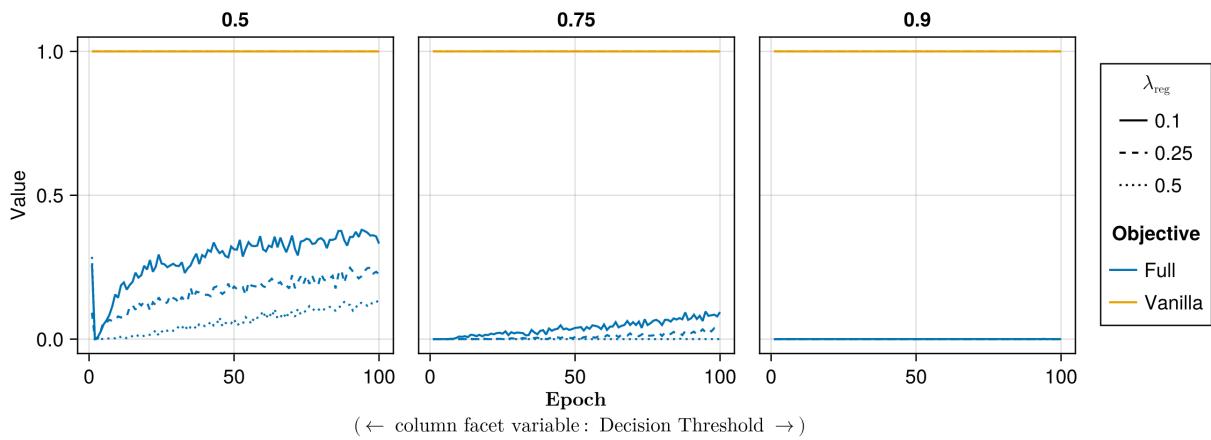


Figure A43: Proportion of mature counterfactuals in each epoch. Data: MNIST.

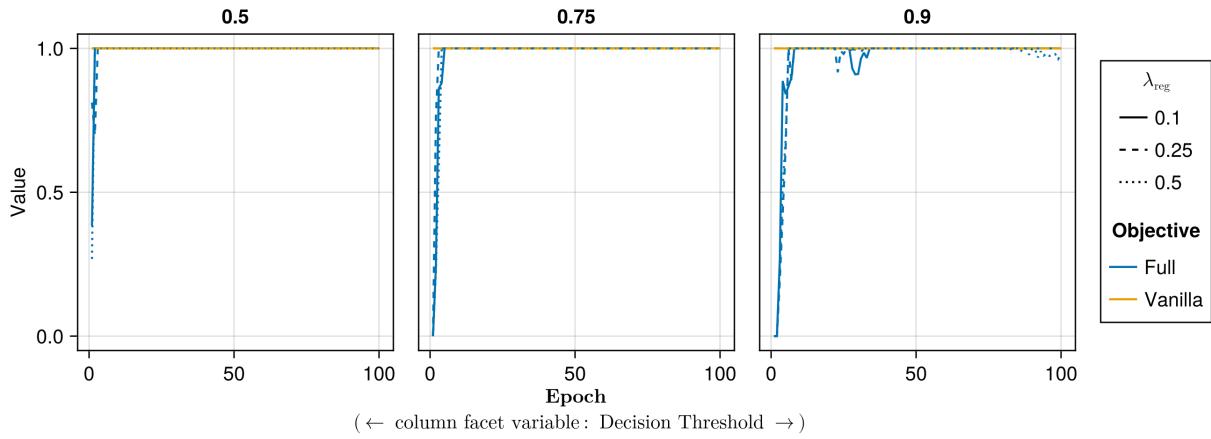


Figure A44: Proportion of mature counterfactuals in each epoch. Data: Moons.

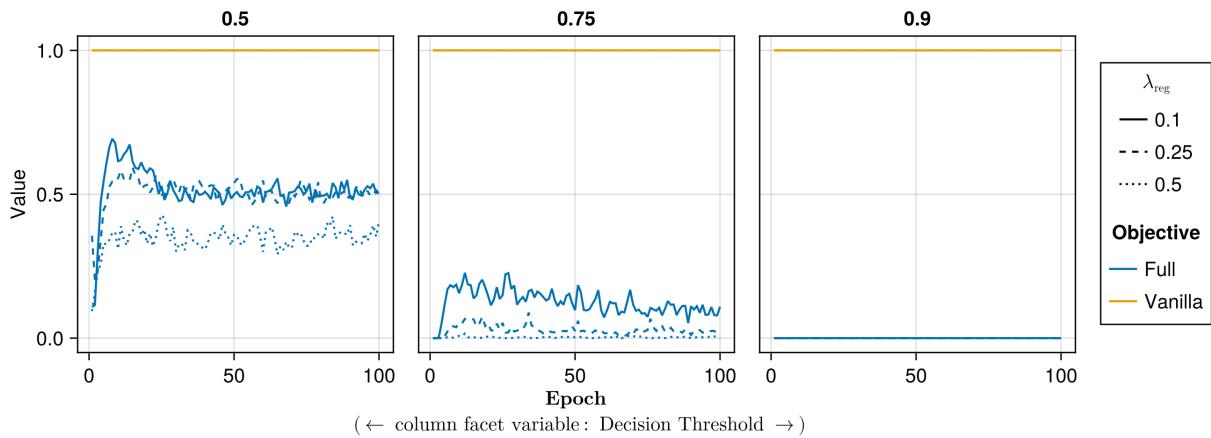


Figure A45: Proportion of mature counterfactuals in each epoch. Data: Overlapping.

Note 9: 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

722

Note 10: Evaluation Phase

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

723

724 K.2.1 Plausibility

725 The results with respect to the plausibility measure are shown in Figure A46 to Figure A51.

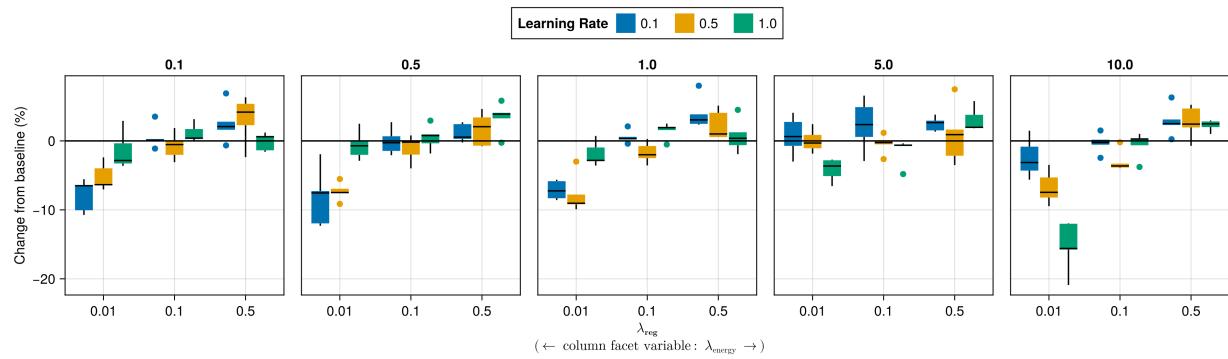


Figure A46: Average outcomes for the plausibility measure across key hyperparameters. Data: Adult.

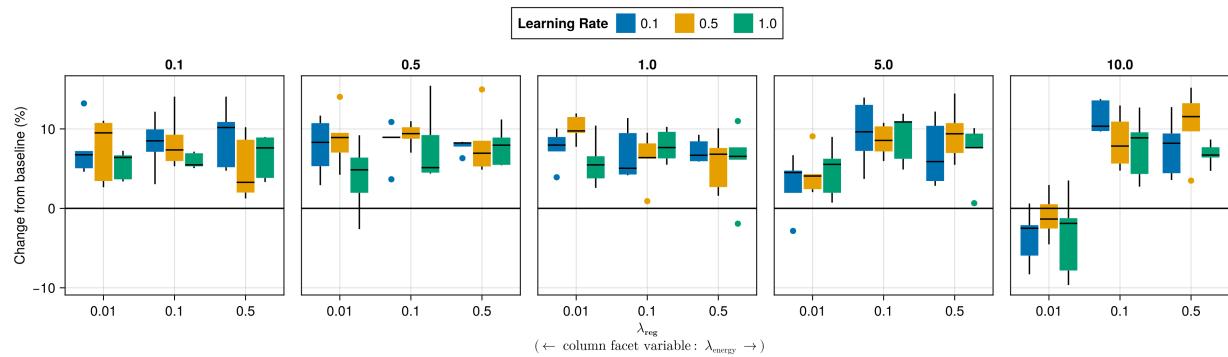


Figure A47: Average outcomes for the plausibility measure across key hyperparameters. Data: Credit.

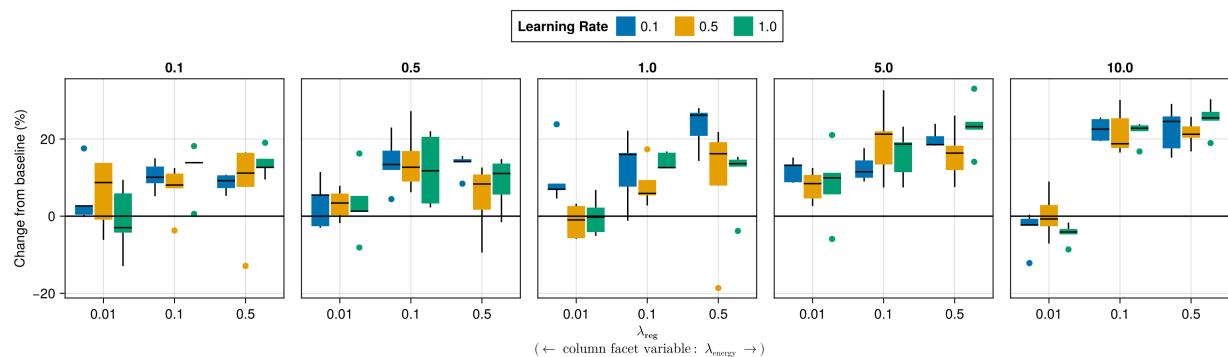


Figure A48: Average outcomes for the plausibility measure across key hyperparameters. Data: GMSC.

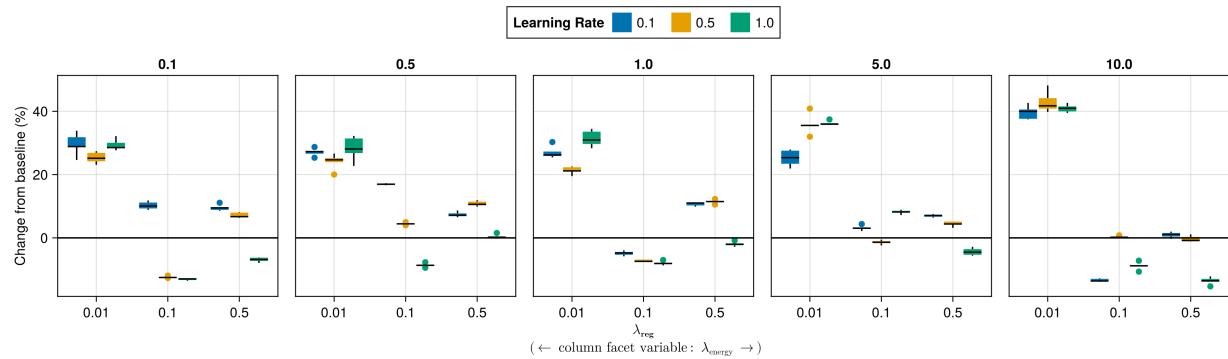


Figure A49: Average outcomes for the plausibility measure across key hyperparameters. Data: Linearly Separable.

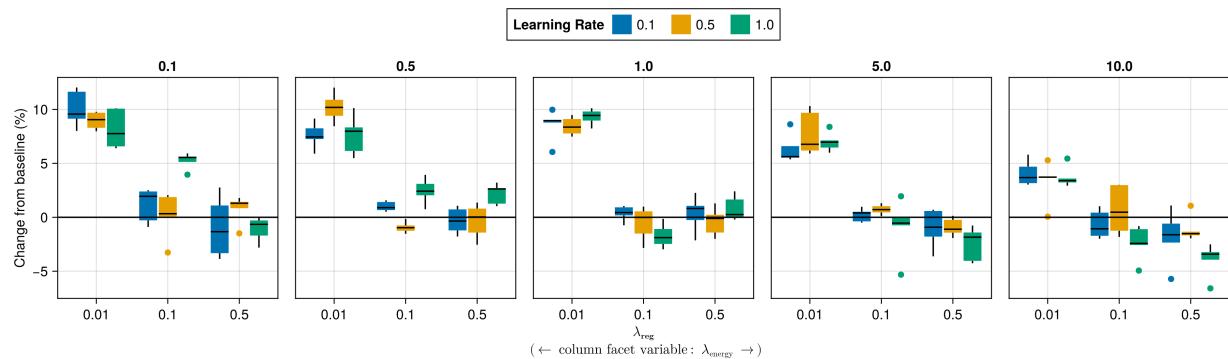


Figure A50: Average outcomes for the plausibility measure across key hyperparameters. Data: MNIST.

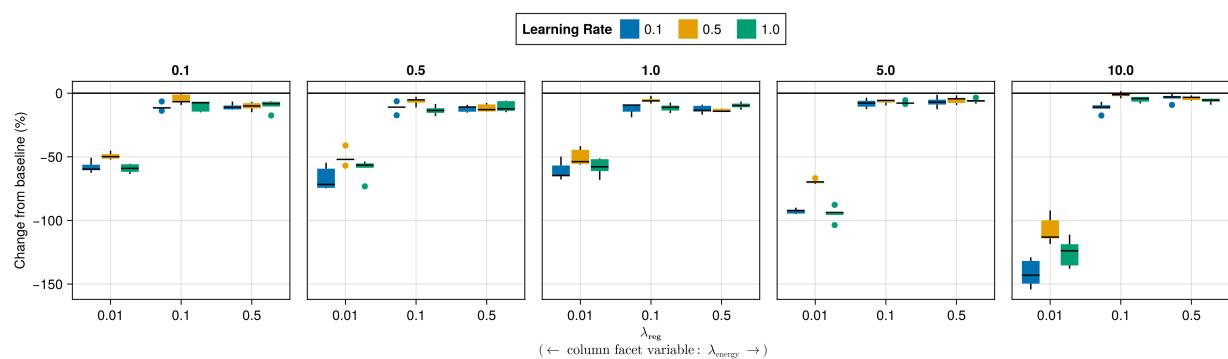


Figure A51: Average outcomes for the plausibility measure across key hyperparameters. Data: Overlapping.

726 **K.2.2 Proportion of Mature CE**

727 The results with respect to the proportion of mature counterfactuals in each epoch are shown in Figure A52 to Figure A57.

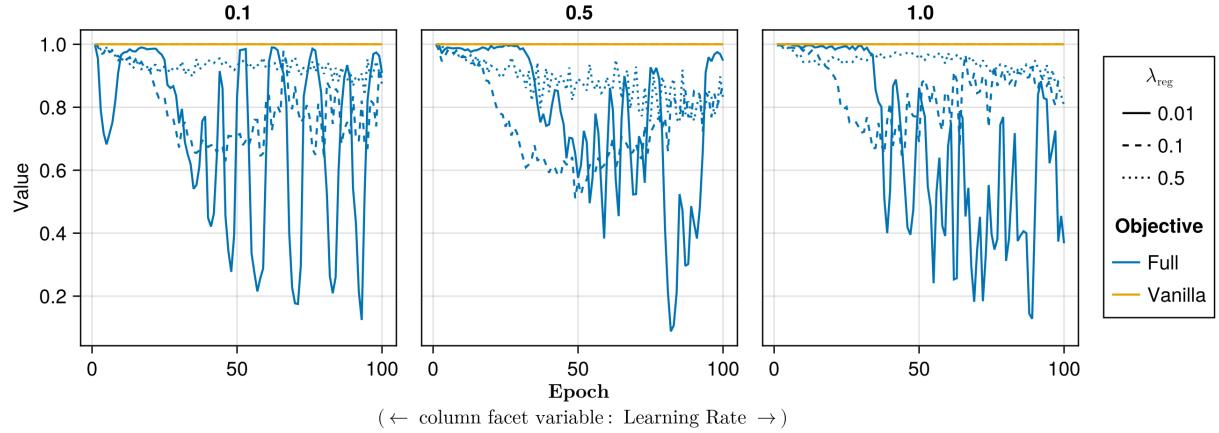


Figure A52: Proportion of mature counterfactuals in each epoch. Data: Adult.

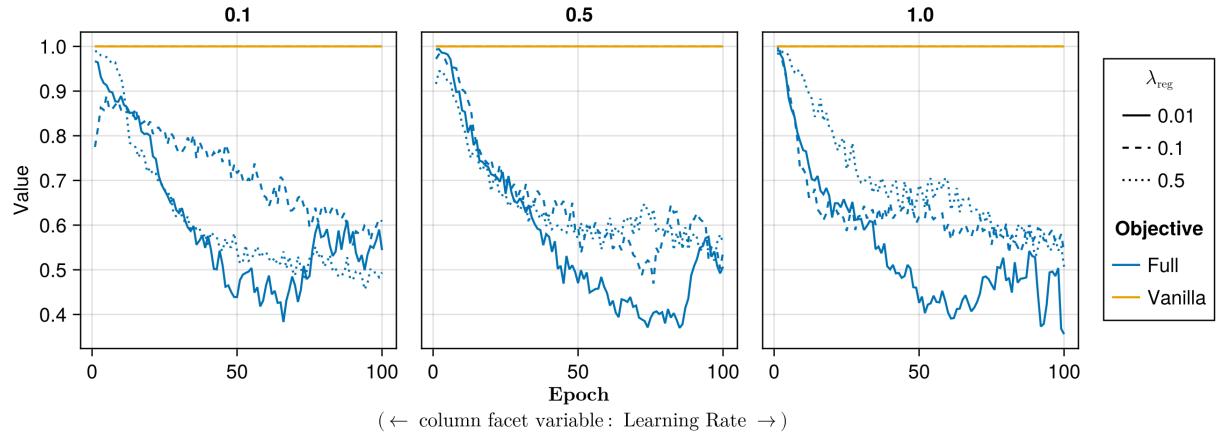


Figure A53: Proportion of mature counterfactuals in each epoch. Data: Credit.

729 **L Computation Details**

730 **L.1 Hardware**

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

735 **L.1.1 Grid Searches**

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

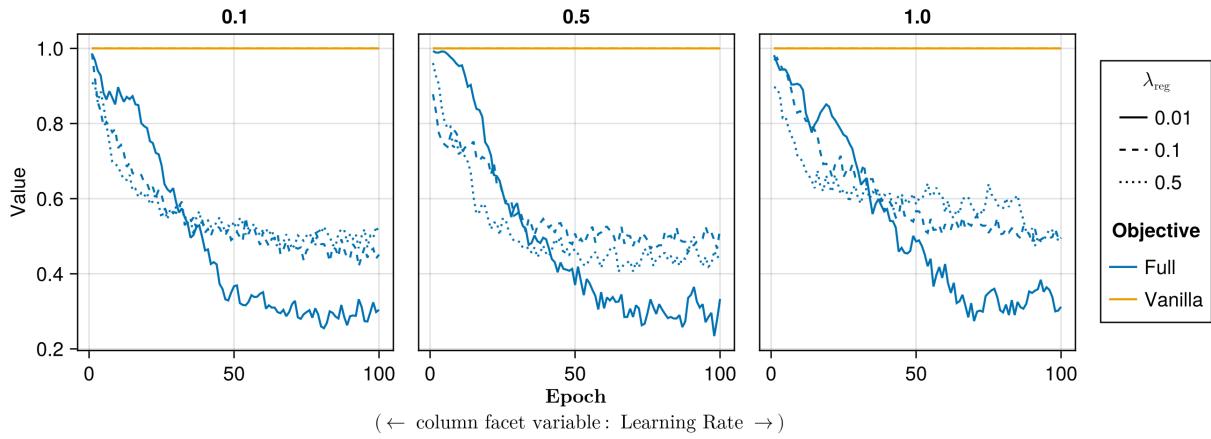


Figure A54: Proportion of mature counterfactuals in each epoch. Data: GMSC.

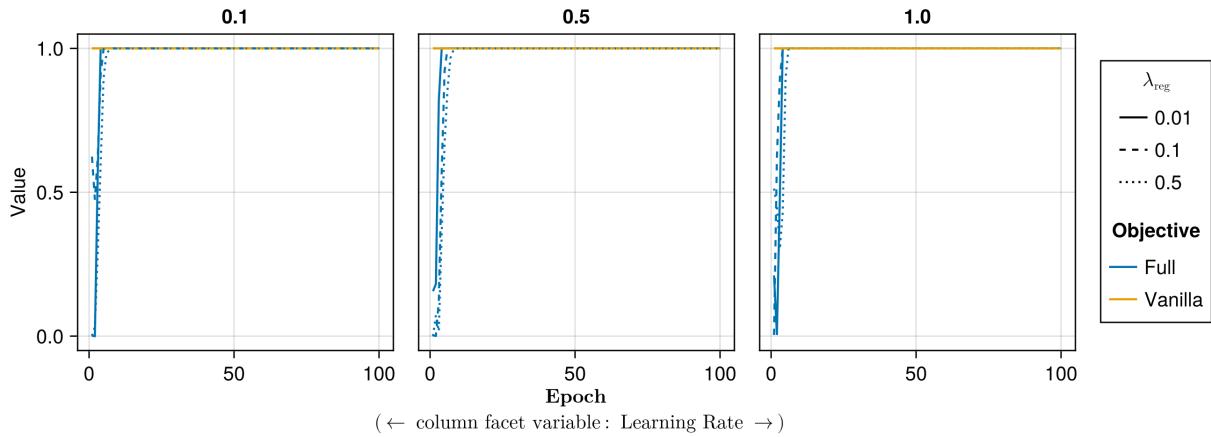


Figure A55: Proportion of mature counterfactuals in each epoch. Data: Linearly Separable.

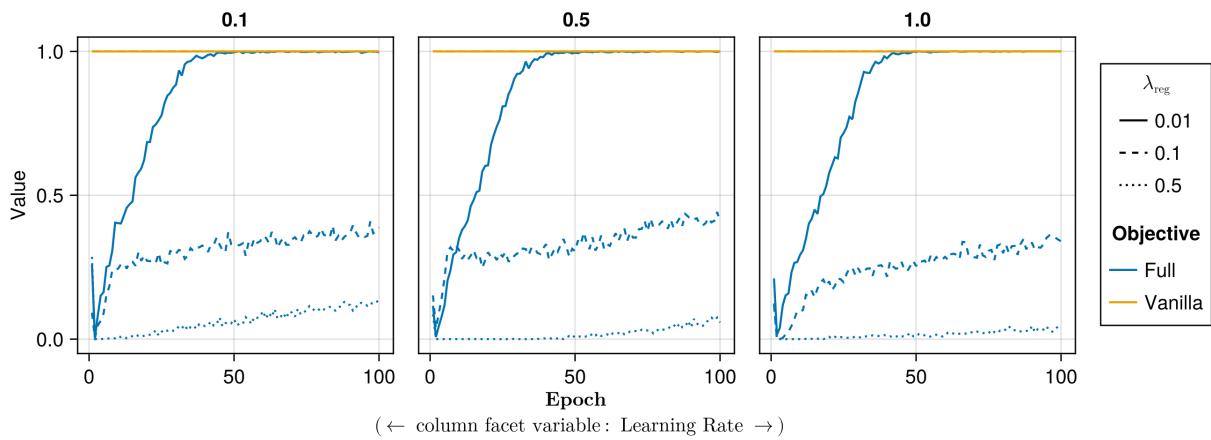


Figure A56: Proportion of mature counterfactuals in each epoch. Data: MNIST.

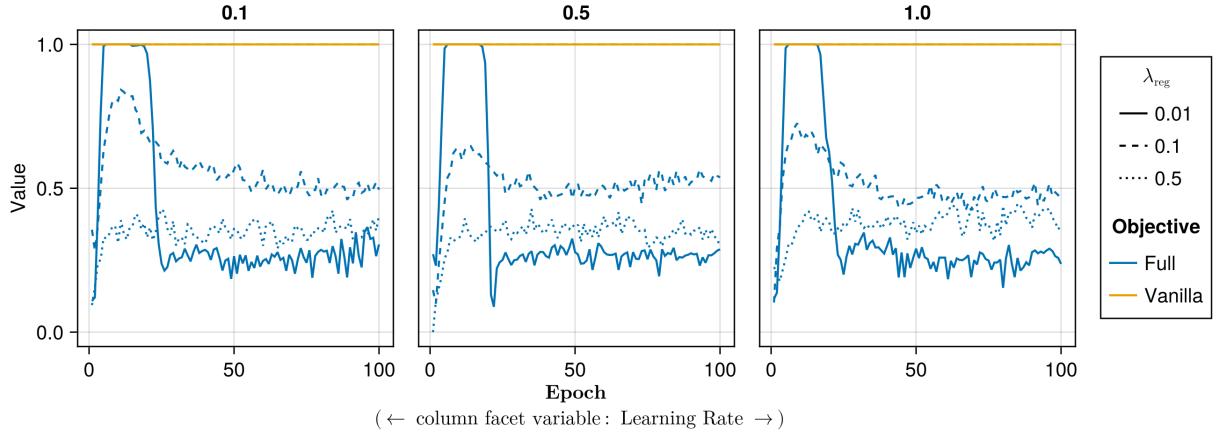


Figure A57: Proportion of mature counterfactuals in each epoch. Data: Overlapping.

741 **L.1.2 Tuning**

742 For tuning of selected hyperparameters, we distributed the task of generating counterfactuals during training across 40
 743 CPUs with 2GB memory each for all tabular datasets. Except for the *Adult* dataset, all training runs were completed
 744 in less than half an hour (<0h30m). The *Adult* dataset took around 0h35m to complete. Evaluations across 20 CPUs
 745 with 3GB memory each generally took less than 0h30m to complete. For *MNIST*, we relied on 100 CPUs with 2GB
 746 memory each. For the *MLP*, training of all models could be completed in 1h30m, while the evaluation across 20 CPUs
 747 (6GB memory) took 4h12m. For the *CNN*, training of all models took ~8h, with conventionally trained models taking
 748 ~0h15m each and model with CT taking ~0h30m-0h45m each.

749 **L.2 Software**

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