
SUPPLEMENTARY APPENDIX

COUNTERFACTUAL TRAINING: TEACHING MODELS PLAUSIBLE AND ACTIONABLE EXPLANATIONS

A PREPRINT

July 23, 2025

ABSTRACT

This is the supplementary appendix to our paper titled *Counterfactual Training: Teaching Models Plausible and Actionable Explanations*. It provides helpful details on mathematical notations and formulas, our proposed training regime, extended empirical findings, hyperparameter tuning and grid searches as well as software and computations.

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

Table of contents

A Notation	3
A.1 Other Technical Details	3
B Technical Details of Our Approach	3
B.1 Generating Counterfactuals through Gradient Descent	3
B.1.1 Background	3
B.1.2 Convergence	3
B.2 Protecting Mutability Constraints with Linear Classifiers	4
B.3 Domain Constraints	5
B.4 Training Hyperparameters	5
B.5 Evaluation Details	6
B.5.1 Robust Accuracy	6
C Details on Main Experiments	6
C.1 Final Hyperparameters	6
C.2 Final Results	7
C.2.1 Robust Performance Plots	7
C.2.2 Confidence Intervals	7
C.2.3 Qualitative Findings for Image Data	7
C.2.5 Costs and Validity	7
C.2.4 Integrated Gradients	8
D Grid Searches	9
D.1 Evaluation Details	9
D.1.1 Predictive Performance	9
D.1.2 Counterfactual Outcomes	9
D.2 Generator Parameters	10
D.2.1 Predictive Performance	10
D.2.2 Plausibility	11
D.2.3 Cost	11
D.3 Penalty Strengths	11
D.3.1 Predictive Performance	11
D.3.2 Plausibility	12
D.3.3 Cost	12
D.4 Other Parameters	12
D.4.1 Predictive Performance	12
D.4.2 Plausibility	13
D.4.3 Cost	13
E Tuning Key Parameters	38
E.1 Key Parameters	38
E.1.1 Plausibility	38
E.1.2 Proportion of Mature CE	38
E.2 Learning Rate	38
E.2.1 Plausibility	45
E.2.2 Proportion of Mature CE	45
F Computation Details	51
F.1 Hardware	51
F.1.1 Grid Searches	51
F.1.2 Tuning	51
F.2 Software	51
References	51

Appendix A Notation

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

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

A.1 Other Technical Details

Maximum mean discrepancy is defined as follows,

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

where $k(\cdot, \cdot)$ is a kernel function (Gretton et al. 2012). We make use of a Gaussian kernel with a constant length-scale parameter of 0.5. In our implementation, Equation 1 is by default applied to the entire subset of the training data for which $y = y^+$.

Appendix B Technical Details of Our Approach

B.1 Generating Counterfactuals through Gradient Descent

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

B.1.1 Background

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

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

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

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

B.1.2 Convergence

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

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

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

B.2 Protecting Mutability Constraints with Linear Classifiers

In the main paper, we explain that to avoid penalizing implausibility that arises due to mutability constraints, we impose a point mass prior on $p(\mathbf{x})$ for the corresponding feature. We argue that this approach induces models to be relatively less sensitive to immutable features, propose a theoretical result supporting this and provide empirical evidence that strengthens our argument (both in the main paper and additional findings in this appendix). Below we derive the analytical results in Prp.~??.

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

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

$$\log \frac{p(k|\mathbf{x})}{p(m|\mathbf{x})} = \mathbf{x}^\top \Sigma^{-1} (\mu_k - \mu_m) + \text{const} \quad (2)$$

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

$$\log \frac{p(k|\mathbf{x})}{p(m|\mathbf{x})} = \sum_{d=1}^D \frac{\mu_{k,d} - \mu_{m,d}}{\sigma_d^2} x_d + \text{const} \quad (3)$$

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

$$\log \frac{p(k|\mathbf{x})}{p(m|\mathbf{x})} = \sum_{d=1}^D (\theta_{k,d} - \theta_{m,d}) x_d + \text{const} \quad (4)$$

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

Based on Equation 3 and Equation 4 we can identify that $(\mu_{k,d} - \mu_{m,d}) \propto (\theta_{k,d} - \theta_{m,d})$ under the assumptions we made above. Hence, we have that $(\theta_{y^-, d_{\text{immmtbl}}} - \theta_{y^+, d_{\text{immmtbl}}}) < 0$ and $(\theta_{y^-, d_{\text{mtbl}}} - \theta_{y^+, d_{\text{mtbl}}}) > 0$

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

$$\frac{\partial}{\partial \theta_{y^+, d}} (\text{div}(\mathbf{x}^+, \mathbf{x}', \mathbf{y}; \theta)) = \frac{\partial}{\partial \theta_{y^+, d}} ((-\mathbf{M}_\theta(\mathbf{x}^+)[y^+]) - (-\mathbf{M}_\theta(\mathbf{x}') [y^+])) = x'_d - x_d^+ \quad (5)$$

and equal to zero everywhere else.

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

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

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

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

B.3 Domain Constraints

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

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

$$\begin{aligned} x_d^{\text{LB}} &= \arg \min_{x_d} \{\mu_d - n_{\sigma_d} \sigma_d, \arg \min_{x_d} x_d\} \\ x_d^{\text{UB}} &= \arg \max_{x_d} \{\mu_d + n_{\sigma_d} \sigma_d, \arg \max_{x_d} x_d\} \end{aligned} \quad (6)$$

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

B.4 Training Hyperparameters

Note 1 presents the default hyperparameters used during training.

Note 1: Training Phase

- Meta Parameters:
 - Generator: `ecco`
 - Model: `mlp`
- Model:
 - Activation: `relu`
 - No. Hidden: 32
 - No. Layers: 1
- Training Parameters:
 - Burnin: 0.0
 - Class Loss: `logitcrossentropy`
 - Convergence: `threshold`
 - Generator Parameters:
 - * Decision Threshold: 0.75
 - * λ_{cst} : 0.001
 - * λ_{egy} : 5.0
 - * Learning Rate: 0.25
 - * Maximum Iterations: 30
 - * Optimizer: `sgd`
 - * Type: ECCo
 - λ_{adv} : 0.25
 - λ_{clf} : 1.0
 - λ_{div} : 0.5
 - λ_{reg} : 0.1
 - Learning Rate: 0.001
 - No. Counterfactuals: 1000
 - No. Epochs: 100

- Objective: full
- Optimizer: adam

B.5 Evaluation Details

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

B.5.1 Robust Accuracy

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

Note 2: Evaluation Phase

- Counterfactual Parameters:
 - Convergence: threshold
 - Decision Threshold: 0.95
 - Generator Parameters:
 - * Decision Threshold: 0.75
 - * λ_{cst} : 0.001
 - * λ_{egy} : 5.0
 - * Learning Rate: 0.25
 - * Maximum Iterations: 30
 - * Optimizer: sgd
 - * Type: ECCo
 - Maximum Iterations: 50
 - No. Individuals: 100
 - No. Runs: 5

Appendix C Details on Main Experiments

C.1 Final Hyperparameters

As discussed the main paper, CT is sensitive to certain hyperparameter choices. We study the effect of many hyperparameters extensively in Section D. For the main results, we tune a small set of key hyperparameters (Section E). The final choices for the main results are presented for each data set in Table 1 along with training, test and batch sizes.

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

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

Table 2: Mean outcomes for **CT** and **BL** along with bootstrapped confidence intervals (99%) for difference in mean outcomes grouped by dataset and evaluation metric. Column **LB** and **UB** show the lower and upper bound of the intervals, respectively, and computed using the percentile method. The underlying counterfactual evaluations are the same as the ones used to produce the main table in the paper.

Variable	Data	CT	BL	LB	UB
Cost	Adult	2.26	2.2	-0.22	0.28
Cost	CH	1.46	2.46	-1.1	-0.89
Cost	Circ	0.67	1.23	-0.58	-0.53
Cost	Cred	2.68	2.29	0.16	0.63
Cost	GMSC	1.14	3.05	-2.45	-1.77
Cost	LS	3.82	4.44	-0.7	-0.56
Cost	MNIST	77.04	68.67	-3.47	18.34
Cost	Moon	1.55	1.6	-0.08	-0.01
Cost	OL	1.62	2.63	-1.15	-0.81
IP*	Adult	0.07	0.11	-0.06	-0.01
IP*	CH	0.02	0.06	-0.06	-0.04
IP*	Circ	0.0	0.0	-0.01	-0.0
IP*	Cred	0.03	0.06	-0.05	-0.01
IP*	GMSC	0.02	0.07	-0.06	-0.04
IP*	LS	0.1	0.23	-0.14	-0.12
IP*	MNIST	0.04	0.04	-0.1	0.09
IP*	Moon	0.02	0.02	-0.01	-0.0
IP*	OL	0.12	0.09	-0.01	0.05
IP	Adult	15.03	15.15	-0.68	0.26
IP	CH	6.61	7.52	-1.17	-0.63
IP	Circ	1.03	2.36	-1.37	-1.29
IP	Cred	19.31	22.03	-3.69	-1.74
IP	GMSC	6.19	8.09	-2.4	-1.49
IP	LS	2.41	3.4	-1.04	-0.94
IP	MNIST	258.83	278.54	-30.49	-7.64
IP	Moon	1.36	1.71	-0.38	-0.32
IP	OL	4.49	4.44	-0.03	0.13

C.2 Final Results

Plus/minus two standard deviations of bootstrap estimates.

C.2.1 Robust Performance Plots

C.2.2 Confidence Intervals

C.2.3 Qualitative Findings for Image Data

Figure 1 shows much more plausible (faithful) counterfactuals for a model with CT than the model with conventional training (Figure 2).

C.2.5 Costs and Validity

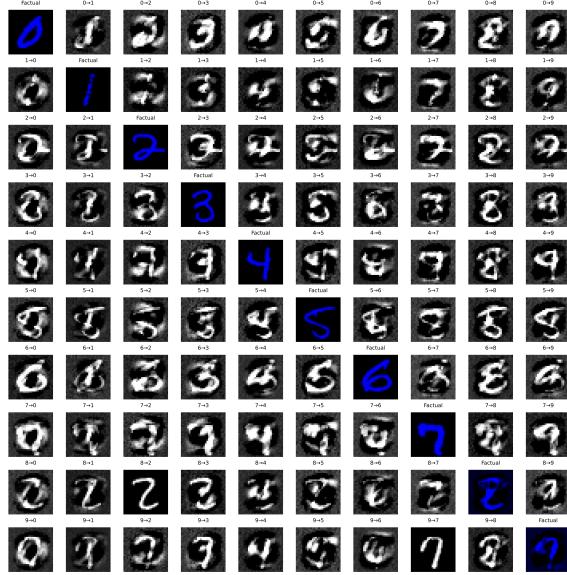


Figure 1: Counterfactual images for *MLP* with counterfactual training. Factual images are shown on the diagonal, with the corresponding counterfactual for each target class (columns) in that same row. The underlying generator, *ECCo*, aims to generate counterfactuals that are faithful to the model (Altmeyer et al. 2024).

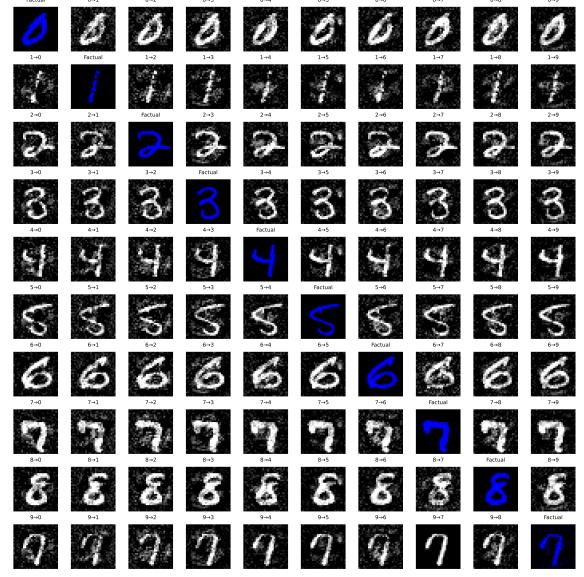


Figure 2: The same setup, factuals, model architecture and generator as in Figure 1, but the model was trained conventionally.

C.2.4 Integrated Gradients

Table 3: Integrated gradients.

Data	CT	BL
LS	0.03	10.24(240)
Circ	3.20(67)	149.76(84275)
Moon	60.84(12851)	0.55(6)
OL	0.78(12)	4.81(108)
Adult	0.43(1)	1.0
CH	0.08(1)	0.23(1)
Cred	0.0	0.43(1)
GMSC	1.0	0.21(3)
MNIST	0.18(1)	0.41(1)

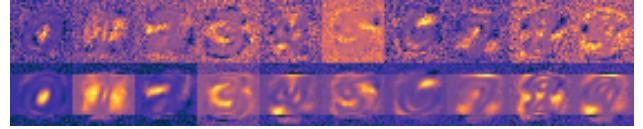


Figure 3: Class-conditional integrated gradients.

Table 4: Costs and validity.

(a) Costs		(b) Validity			(c) Validity		
Data	Cost (-%)	Data	CT	BL	Data	CT	BL
LS	-26.82(86)*	LS	1.0	1.0	LS	1.0	1.0
Circ	40.97(82)*	Circ	0.97	0.52	Circ	0.67	0.49
Moon	33.83(98)*	Moon	1.0	1.0	Moon	0.99	0.98
OL	10.35(128)*	OL	0.87	0.98	OL	0.37	0.57
Adult	1.16(353)	Adult	0.61	0.99	Adult	0.56	0.99
CH	-34.89(231)*	CH	0.96	1.0	CH	0.96	1.0
Cred	28.24(108)*	Cred	0.7	1.0	Cred	0.67	1.0
GMSC	3.54(578)	GMSC	0.63	1.0	GMSC	0.38	1.0
MNIST	-31.67(772)*	MNIST	1.0	1.0	MNIST	1.0	1.0
Avg.	2.75						

Appendix D Grid Searches

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

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

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

D.1 Evaluation Details

To measure predictive performance, we compute the accuracy and F1-score for all models on test data (Table 5, Table 6, Table 7). With respect to explanatory performance, we report here our findings for the (im)plausibility and cost of counterfactuals at test time. Since the computation of our proposed divergence-based adaption (IP*) is memory-intensive, we rely on the distance-based metric for the grid searches. For the counterfactual evaluation, we draw factual samples from the training data for the grid searches to avoid data leakage with respect to our final results reported in the body of the paper. Specifically, we want to avoid choosing our default hyperparameters based on results on the test data. Since we are optimizing for explainability, not predictive performance, we still present test accuracy and F1-scores.

D.1.1 Predictive Performance

We find that CT is associated with little to no decrease in average predictive performance for our synthetic datasets: test accuracy and F1-scores decrease by at most ~1 percentage point, but generally much less (Table 5, Table 6, Table 7). Variation across hyperparameters is negligible as indicated by small standard deviations for these metrics across the board.

D.1.2 Counterfactual Outcomes

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

We do observe strong sensitivity to certain hyperparameters, with clear and manageable patterns. Concerning generator parameters, we firstly find that using *REVISE* to generate counterfactuals during training typically yields the worst

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

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

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

D.2 Generator Parameters

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

Note 3: Training Phase

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

Note 4: Evaluation Phase

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

D.2.1 Predictive Performance

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

Table 5: Predictive performance measures by dataset and objective averaged across training-phase parameters (Note 3) and evaluation-phase parameters (Note 4).

Dataset	Variable	Objective	Mean	Se
Circ	Accuracy	Full	1.0	0.0
Circ	Accuracy	Vanilla	1.0	0.0
Circ	F1-score	Full	1.0	0.0

Continuing table below.

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

D.2.2 Plausibility

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

D.2.3 Cost

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

D.3 Penalty Strengths

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

Note 5: Training Phase

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

Note 6: Evaluation Phase

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

D.3.1 Predictive Performance

Predictive performance measures for this grid search are shown in Table 6.

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

Dataset	Variable	Objective	Mean	Se
Circ	Accuracy	Full	0.99	0.01
Circ	Accuracy	Vanilla	1.0	0.0
Circ	F1-score	Full	0.99	0.01
Circ	F1-score	Vanilla	1.0	0.0
LS	Accuracy	Full	1.0	0.01
LS	Accuracy	Vanilla	1.0	0.0

Continuing table below.

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

D.3.2 Plausibility

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

D.3.3 Cost

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

D.4 Other Parameters

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

Note 7: Training Phase

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

Note 8: Evaluation Phase

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

D.4.1 Predictive Performance

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

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

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

Continuing table below.

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

D.4.2 Plausibility

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

D.4.3 Cost

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

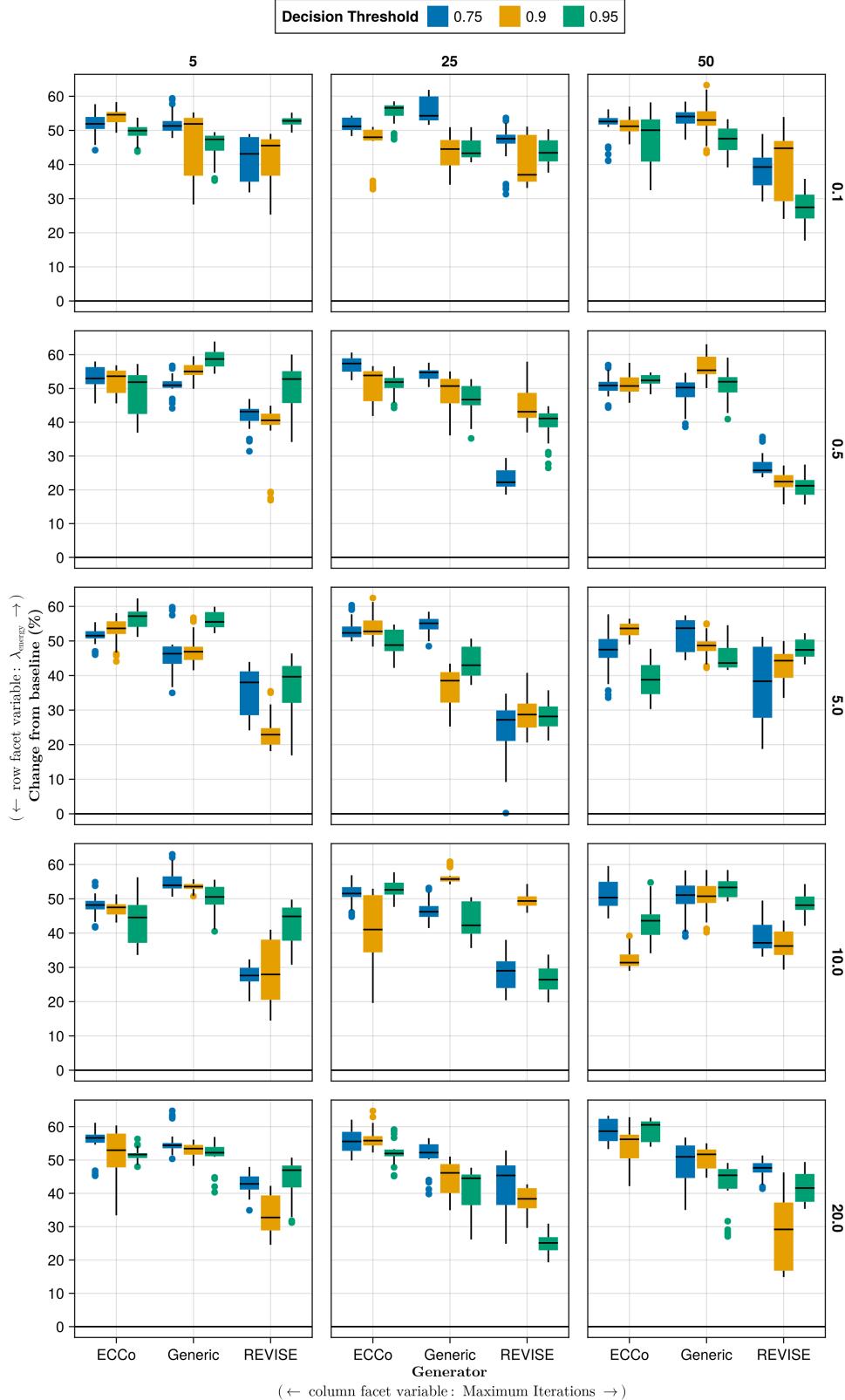


Figure 4: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\$ \text{ext}\{\text{IP}\} \$$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters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

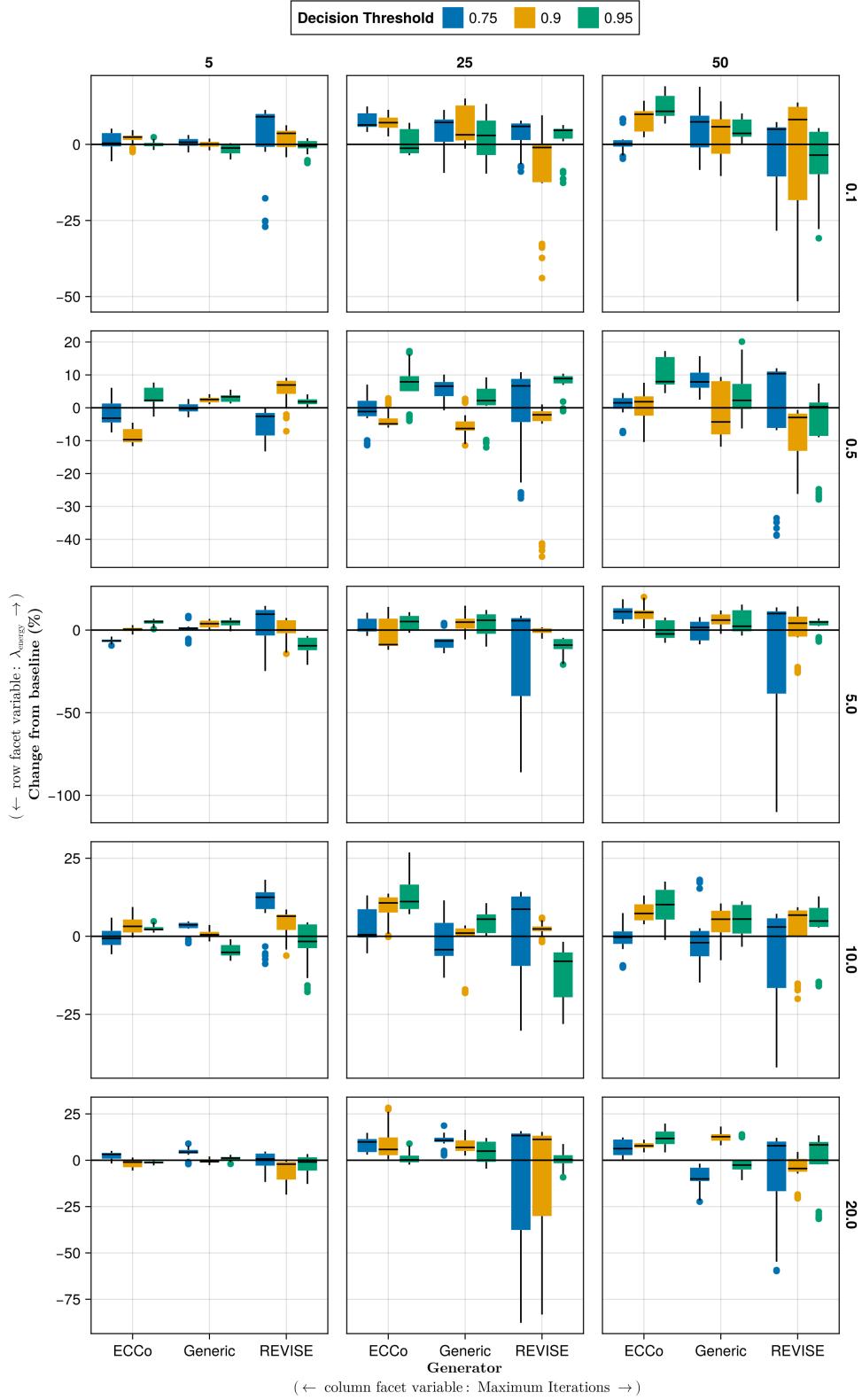


Figure 5: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\$ \text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for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

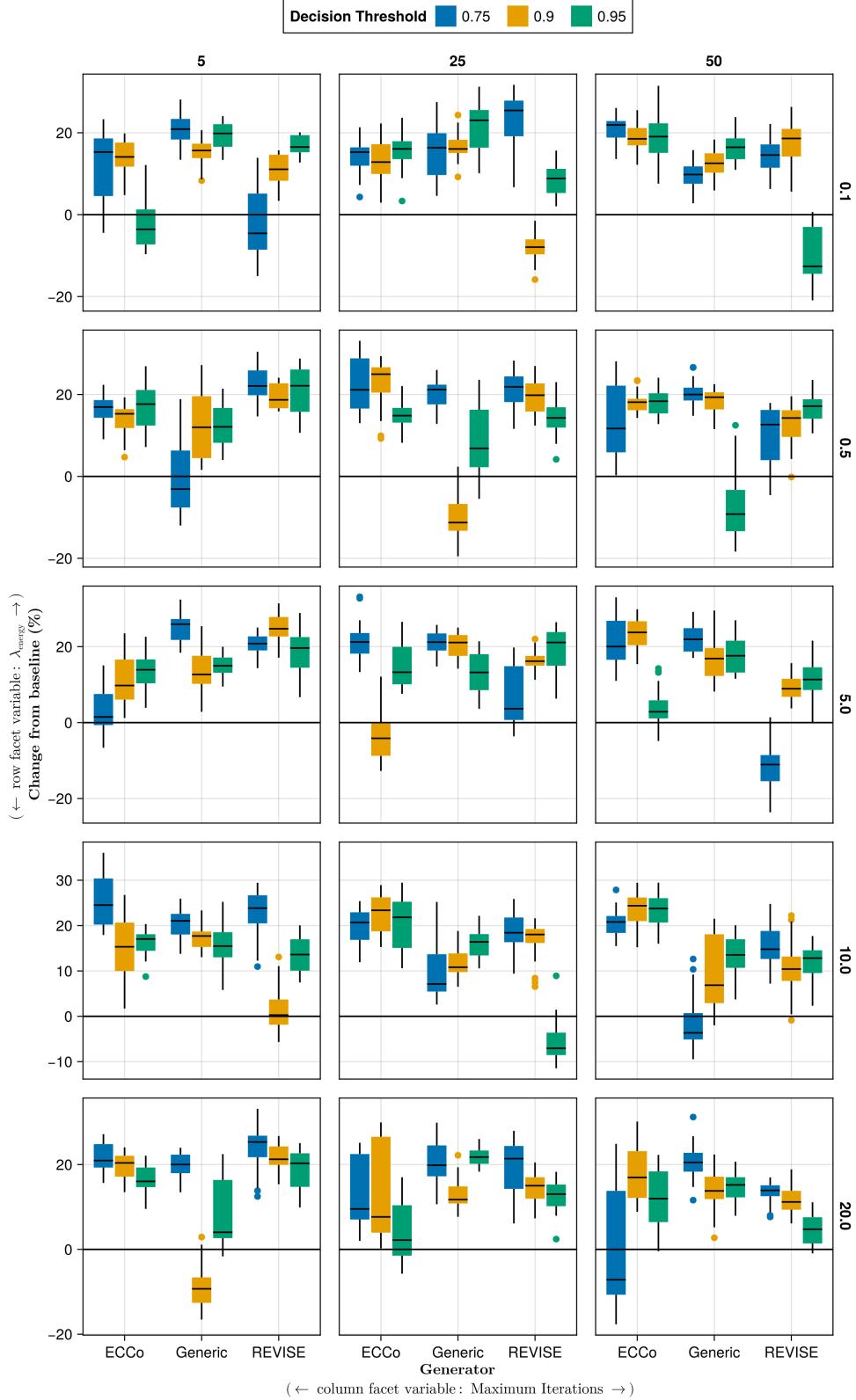


Figure 6: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

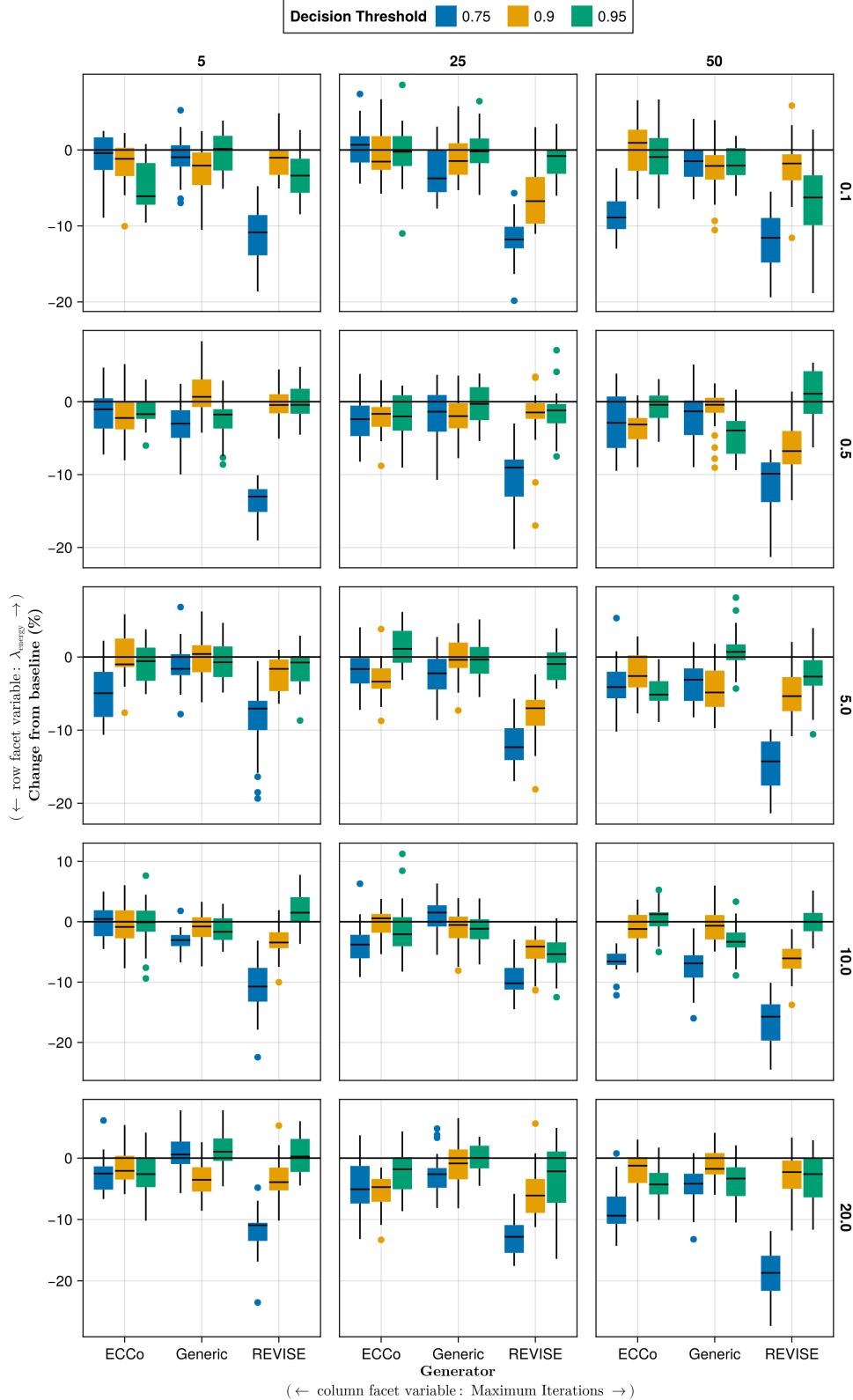


Figure 7: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric (λ_{energy}). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*). Data: Overlapping.

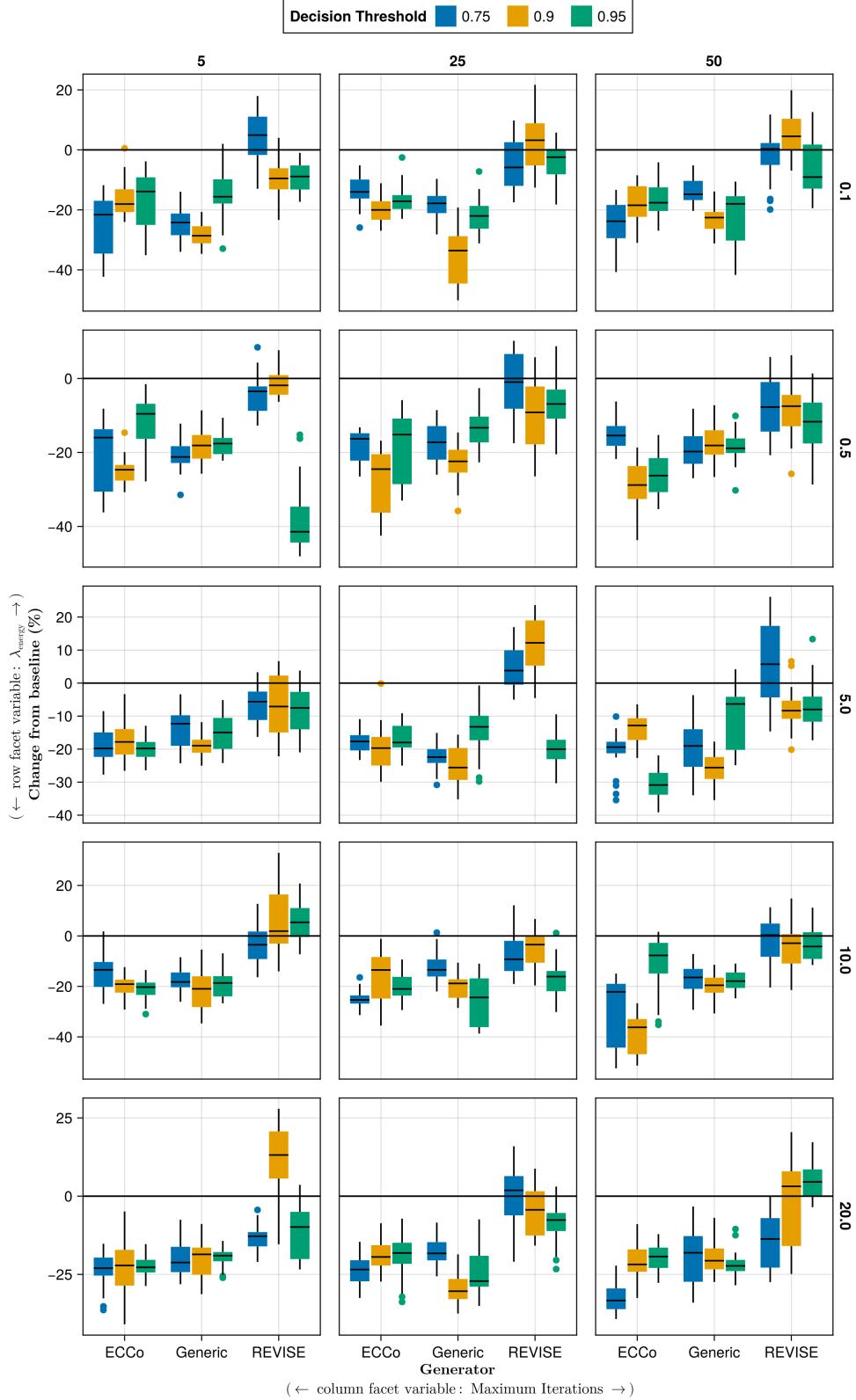


Figure 8: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

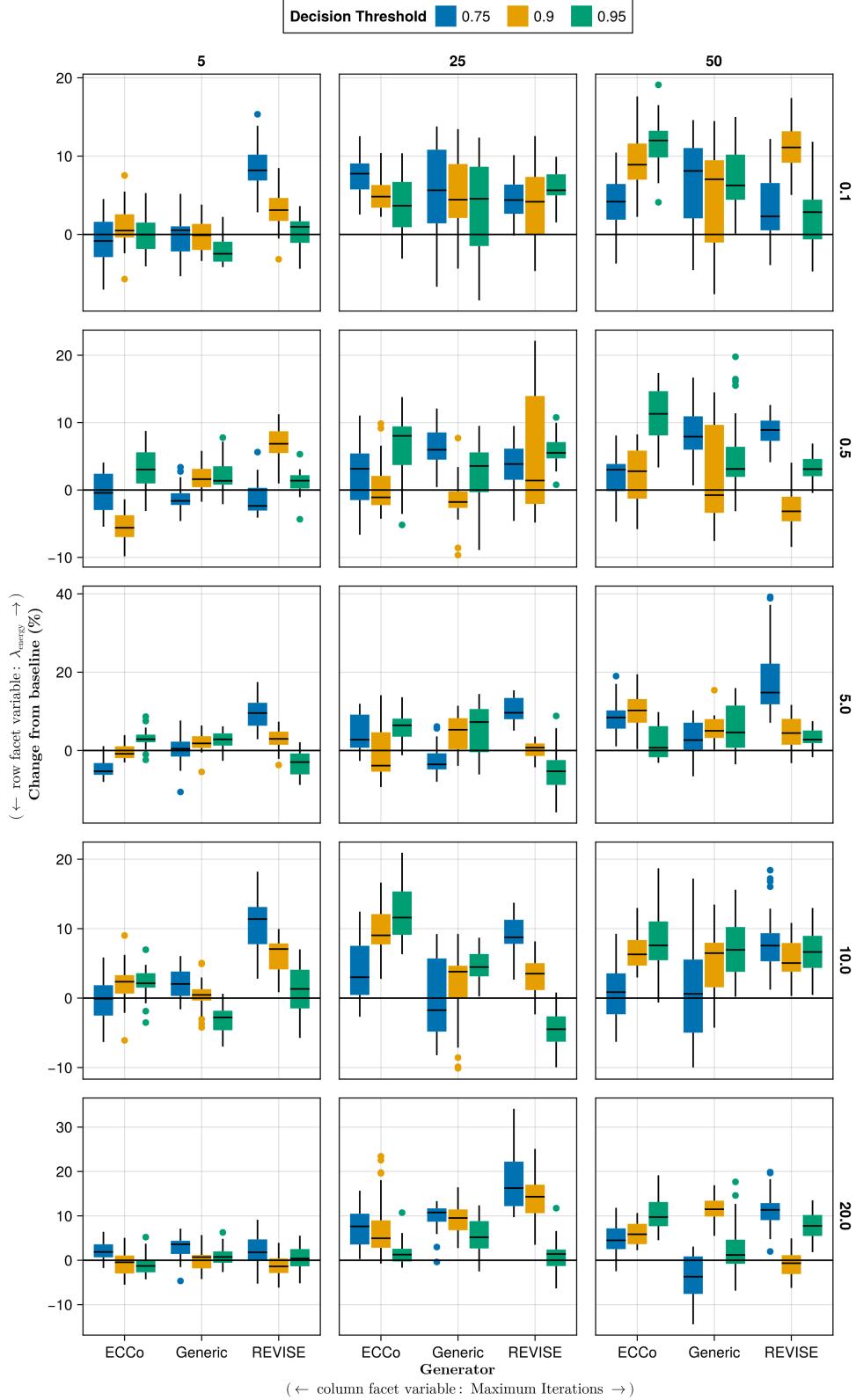


Figure 9: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*).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

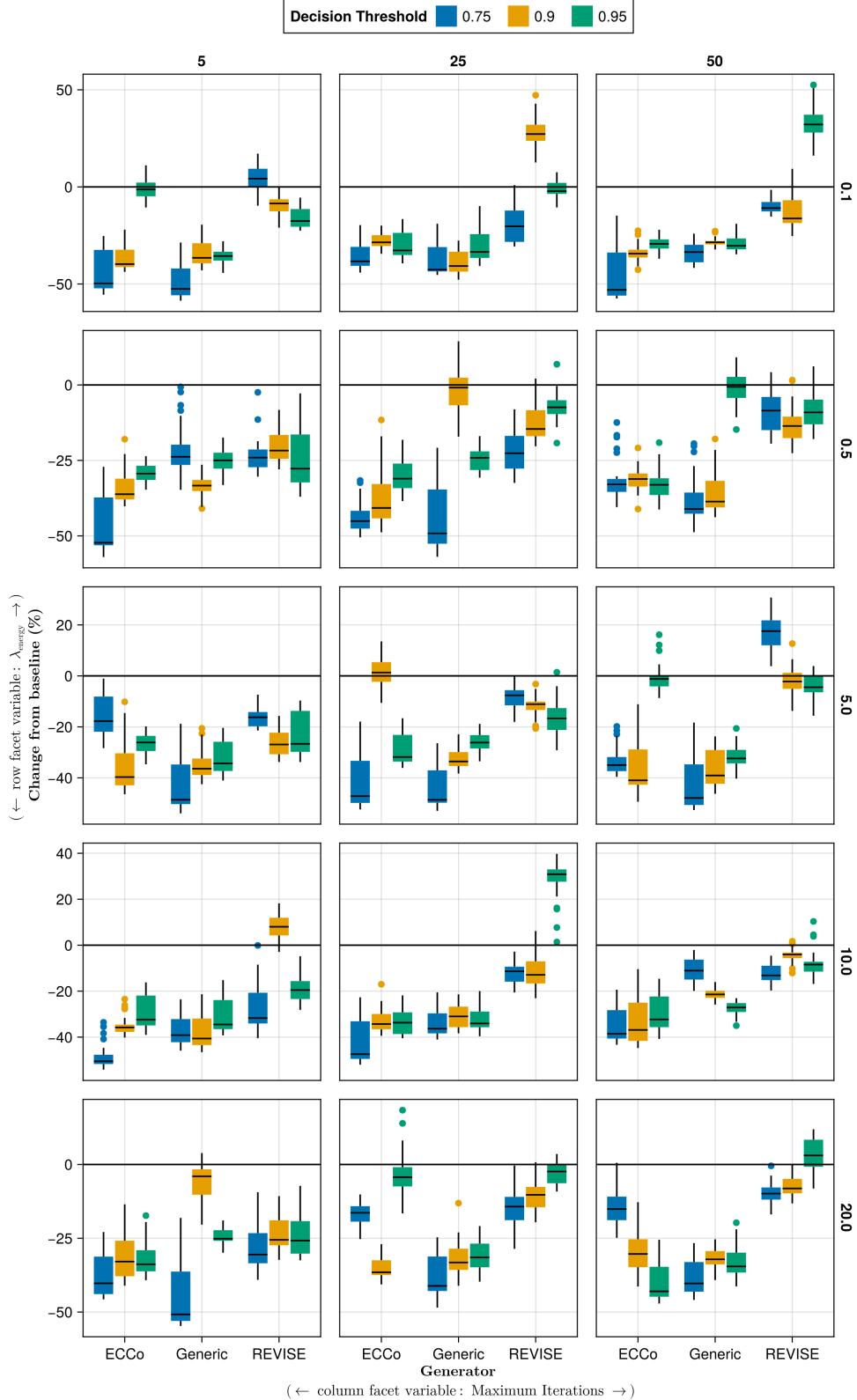


Figure 10: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*).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

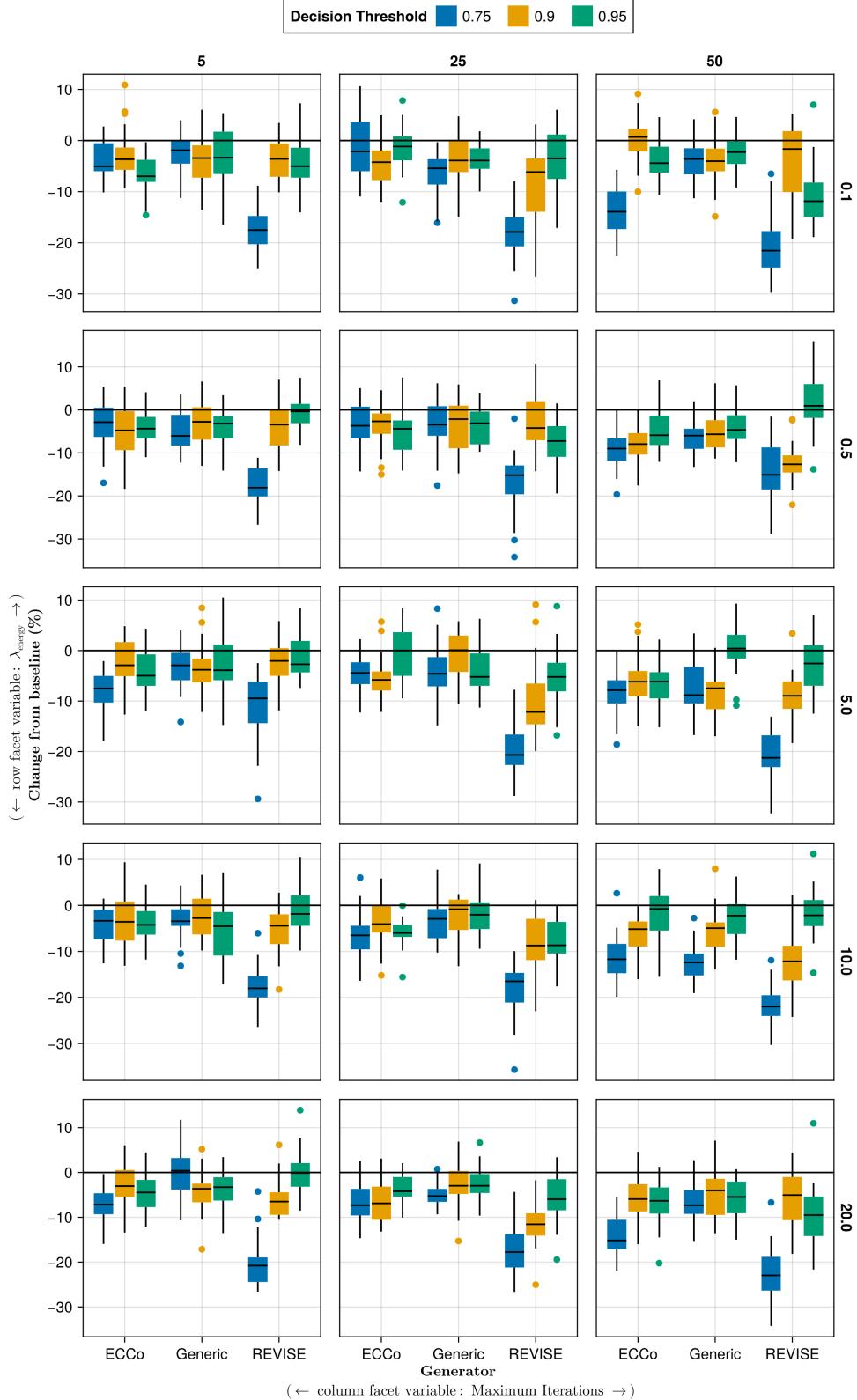


Figure 11: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric (Wachter, Mittelstadt, and Russell 2017). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

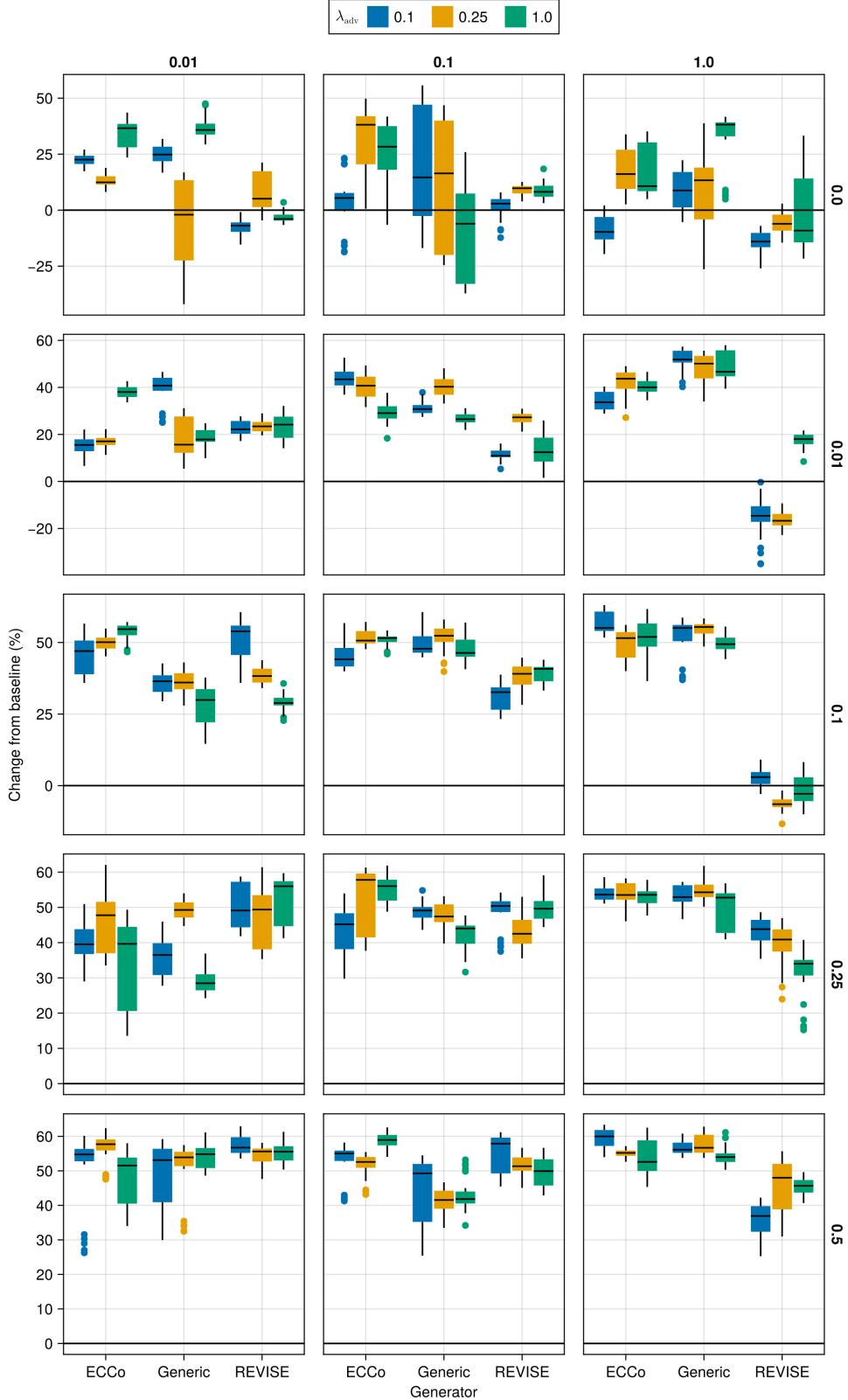


Figure 12: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

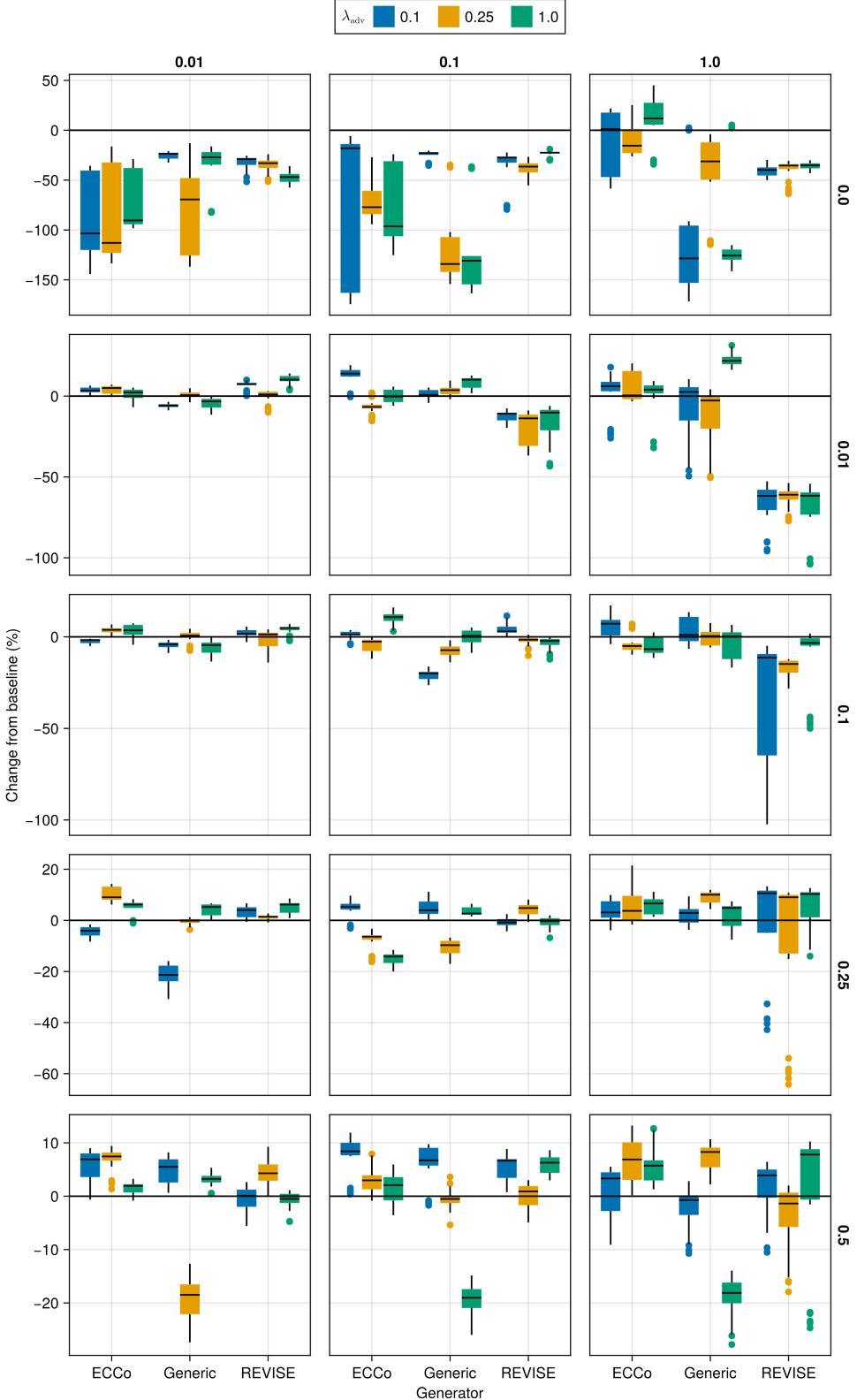


Figure 13: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

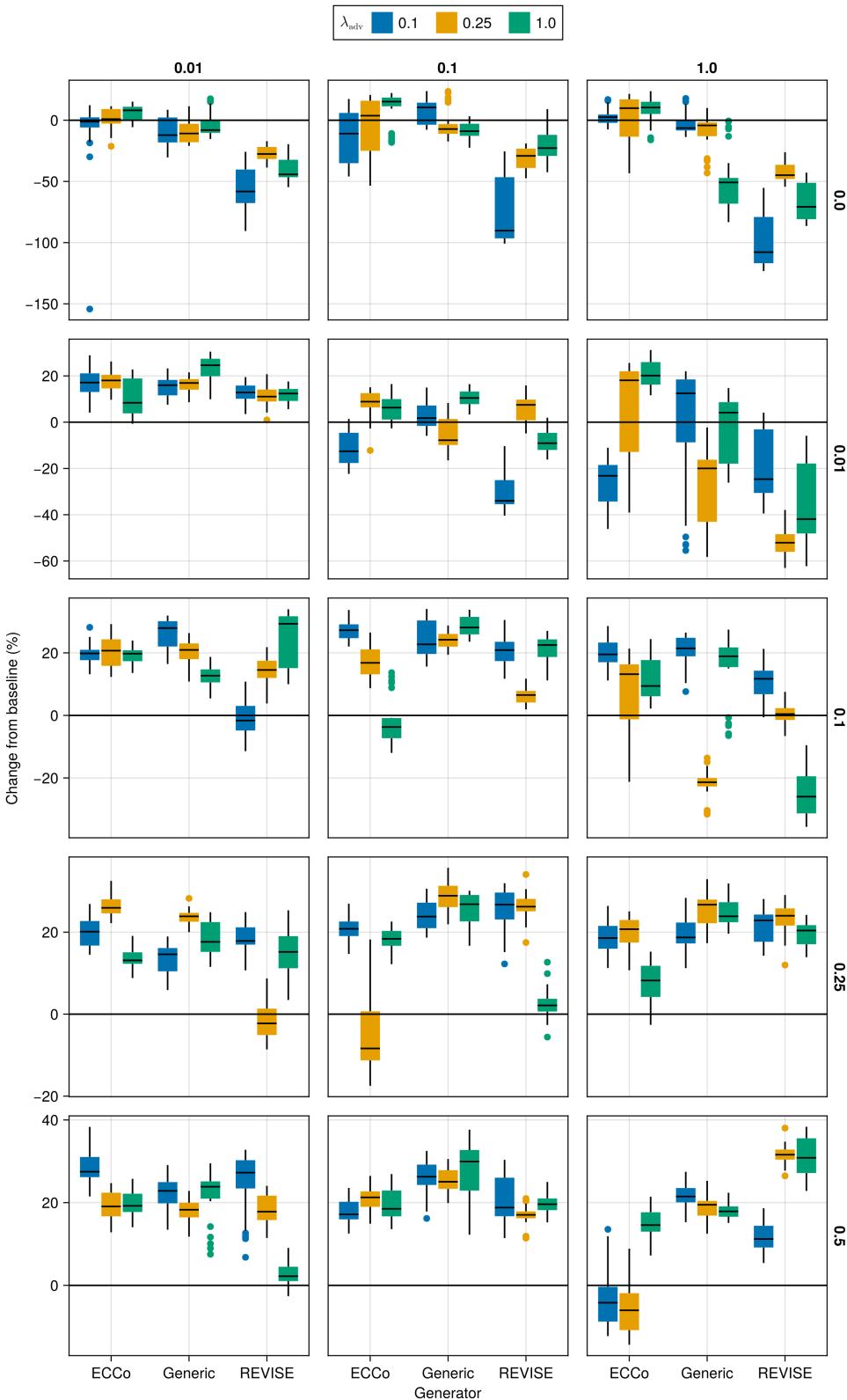


Figure 14: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

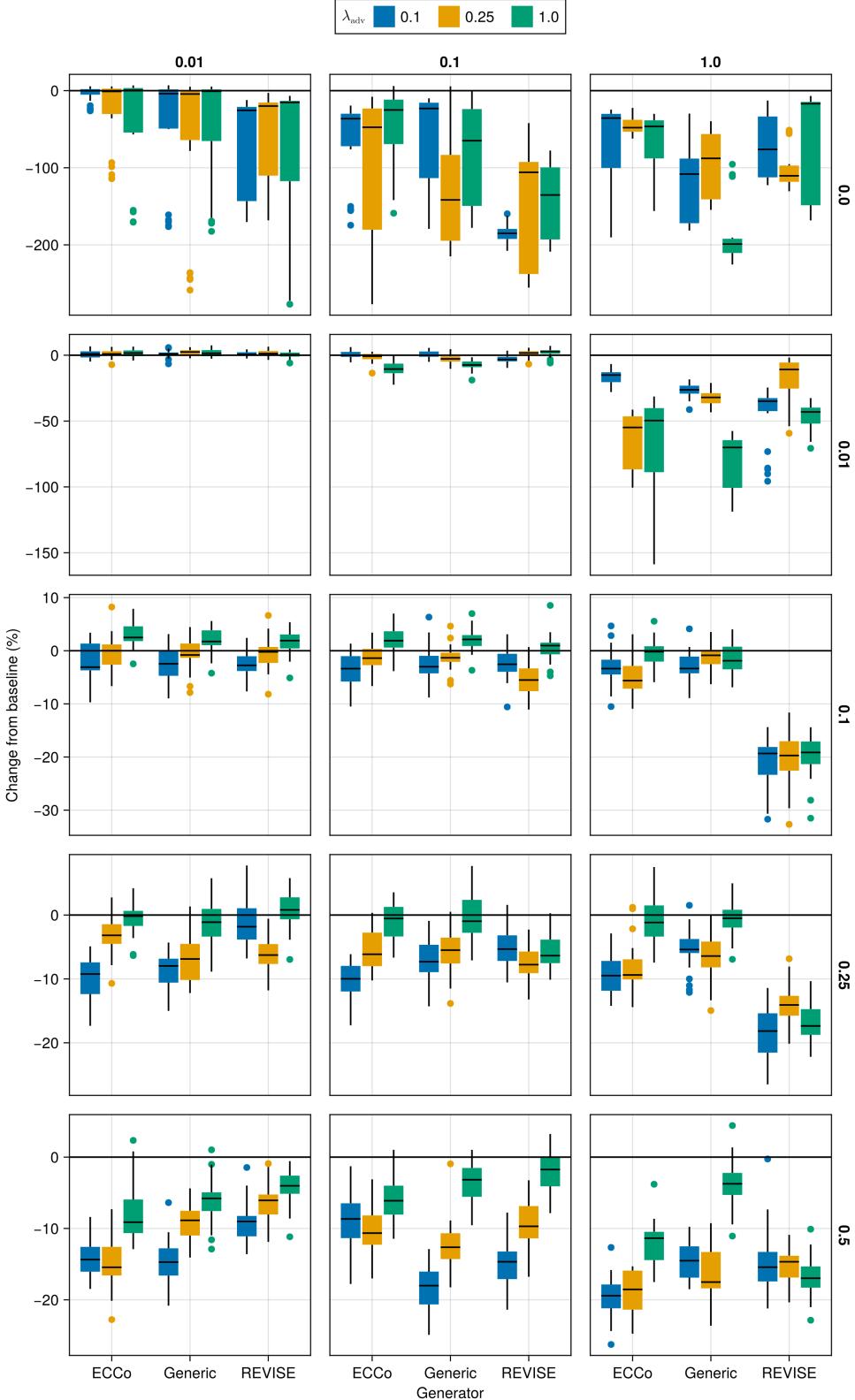


Figure 15: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

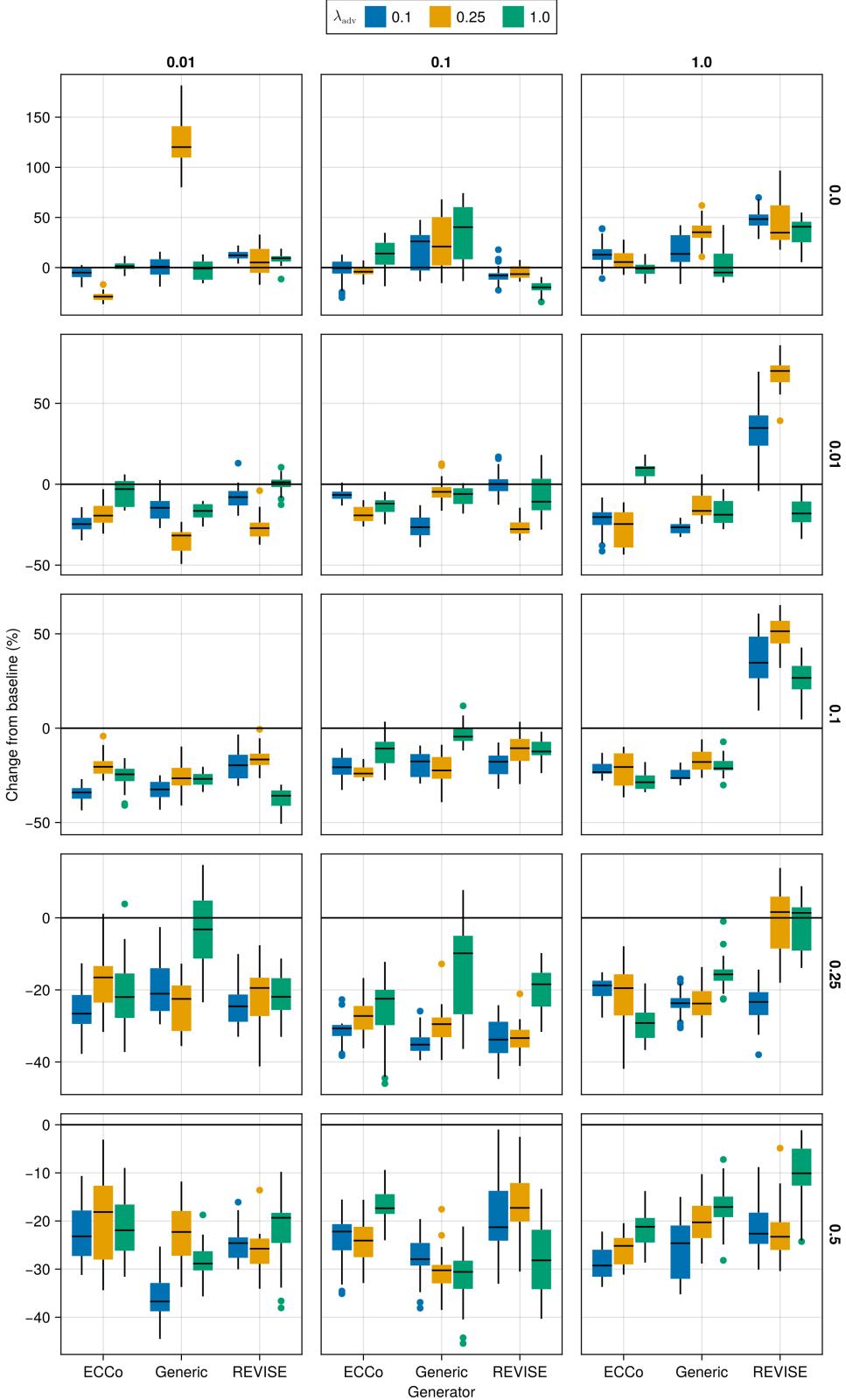


Figure 16: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

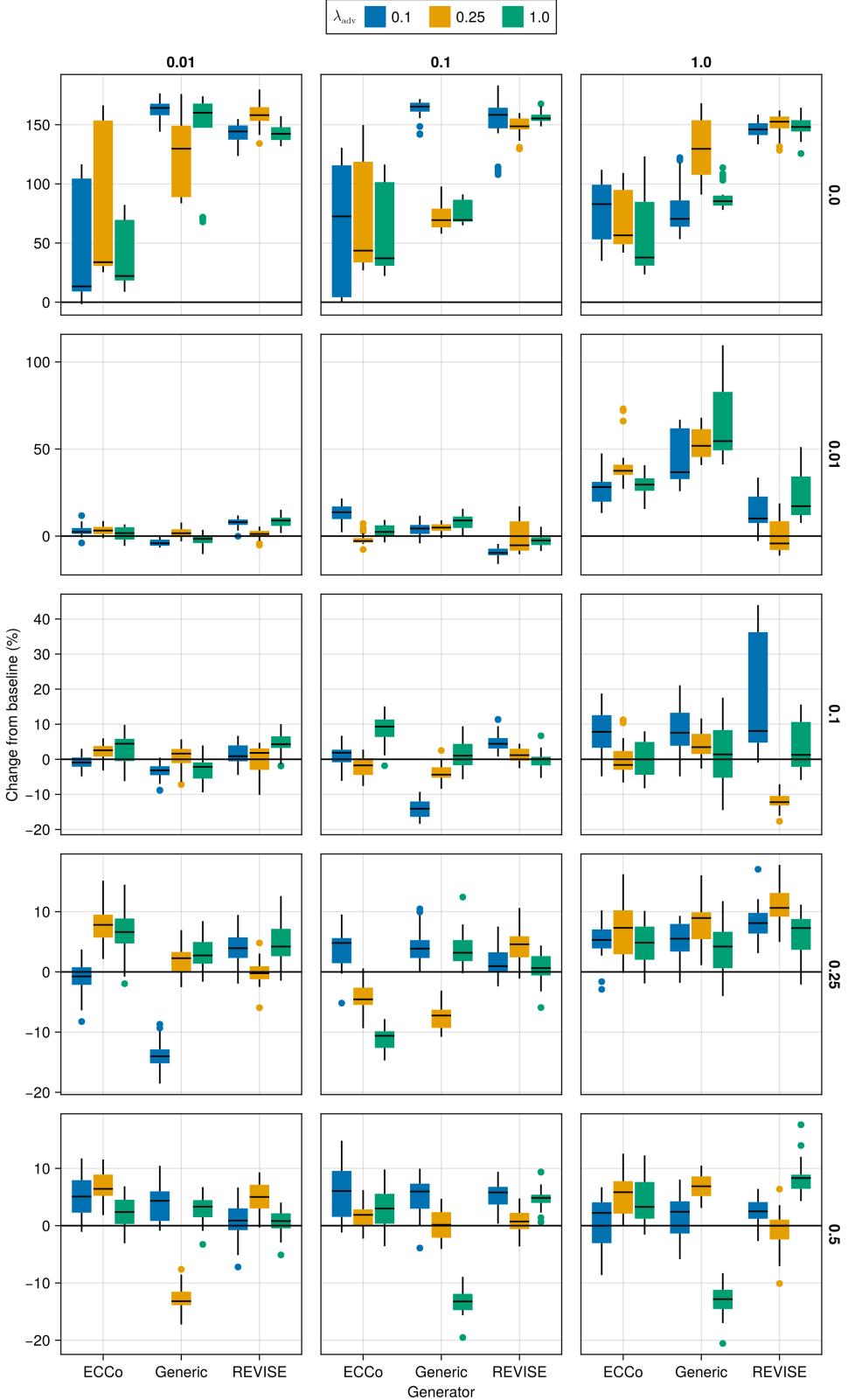


Figure 17: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

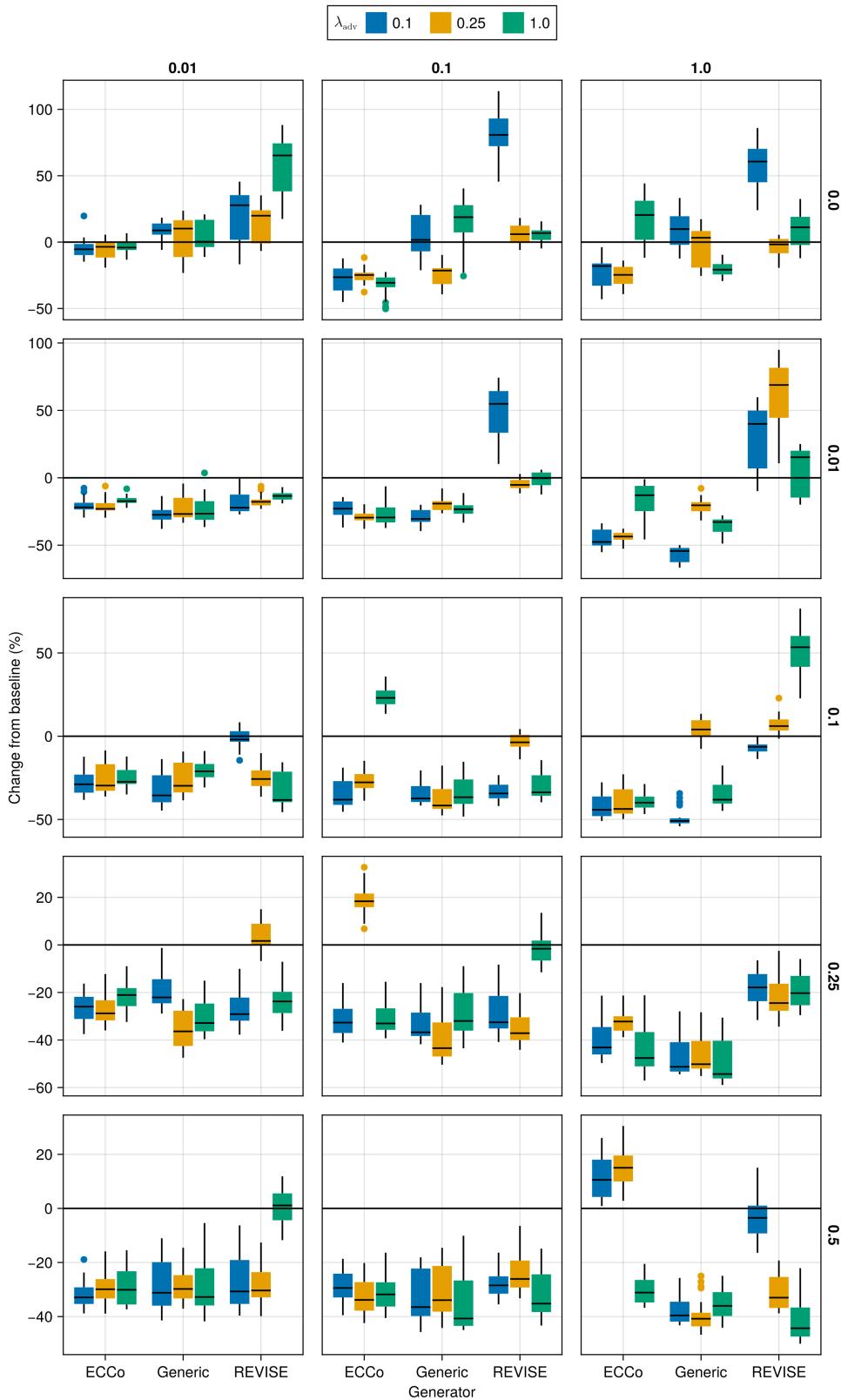


Figure 18: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

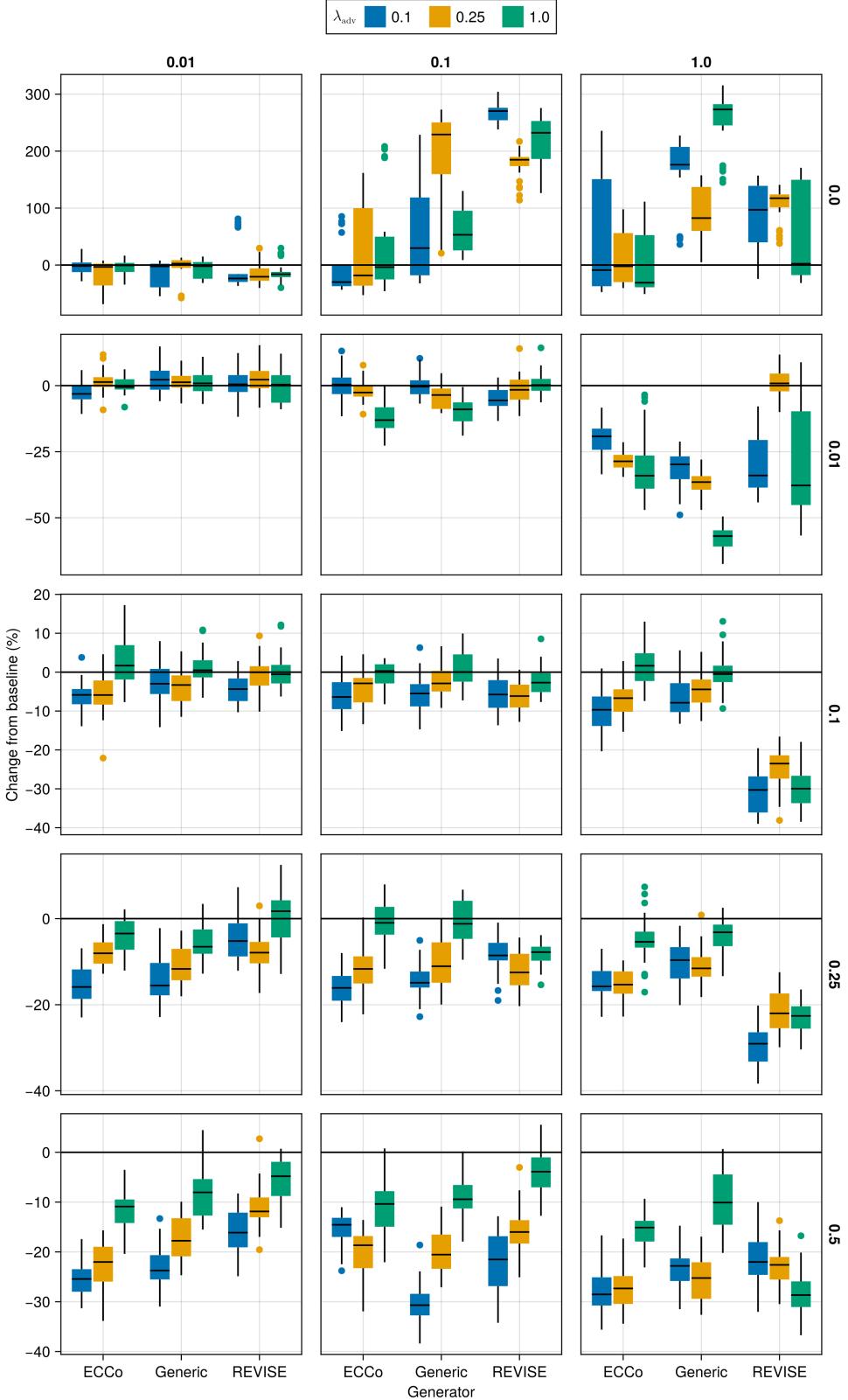


Figure 19: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

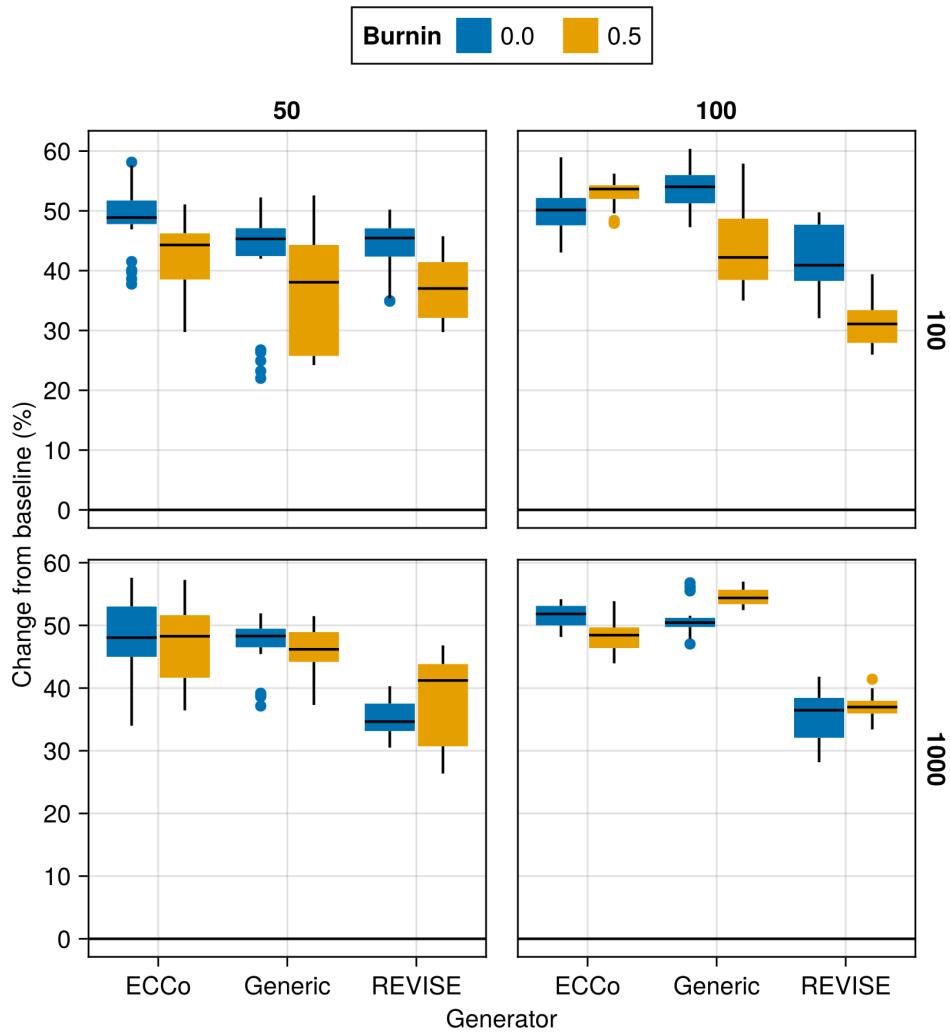


Figure 20: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Circles.

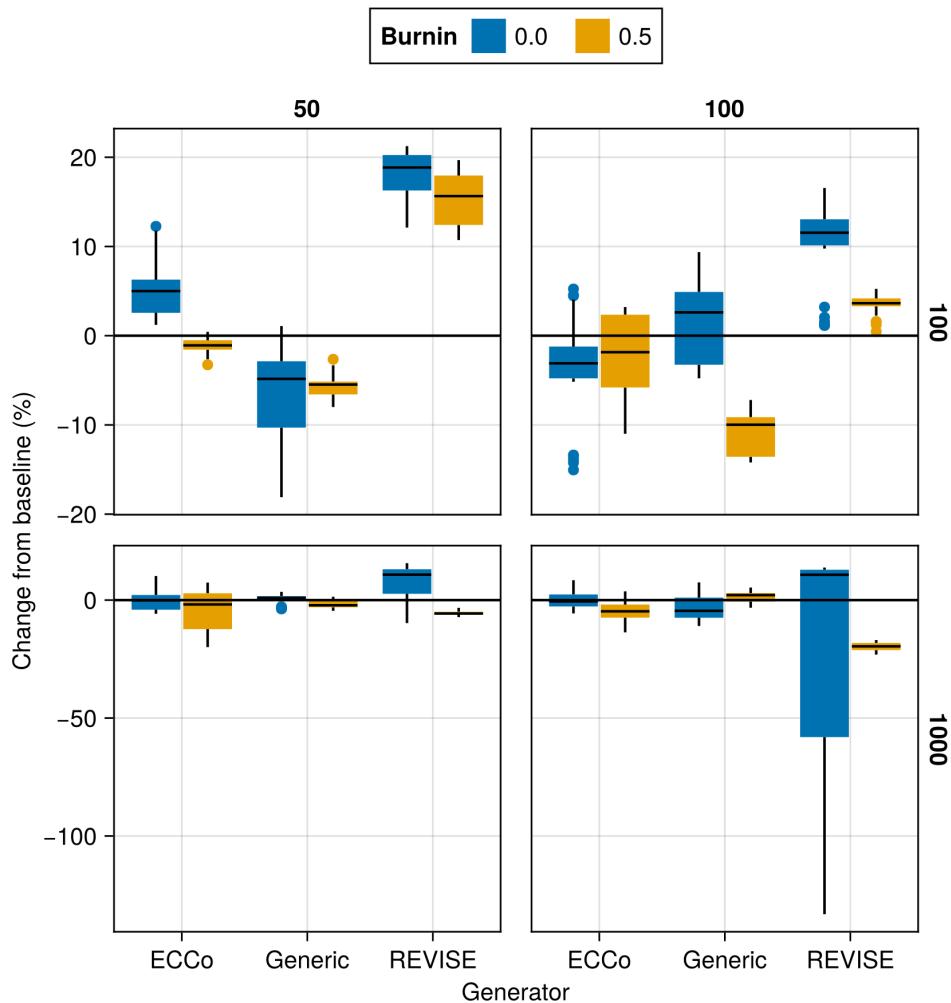


Figure 21: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

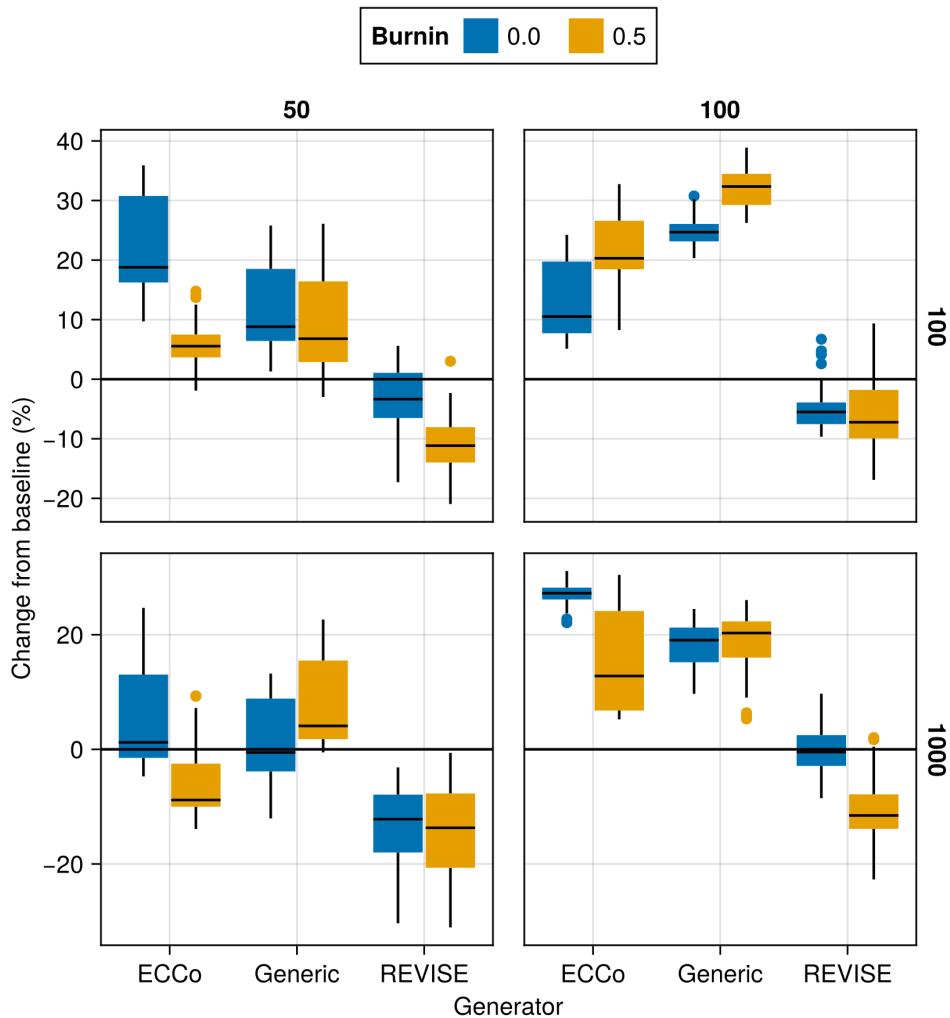


Figure 22: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

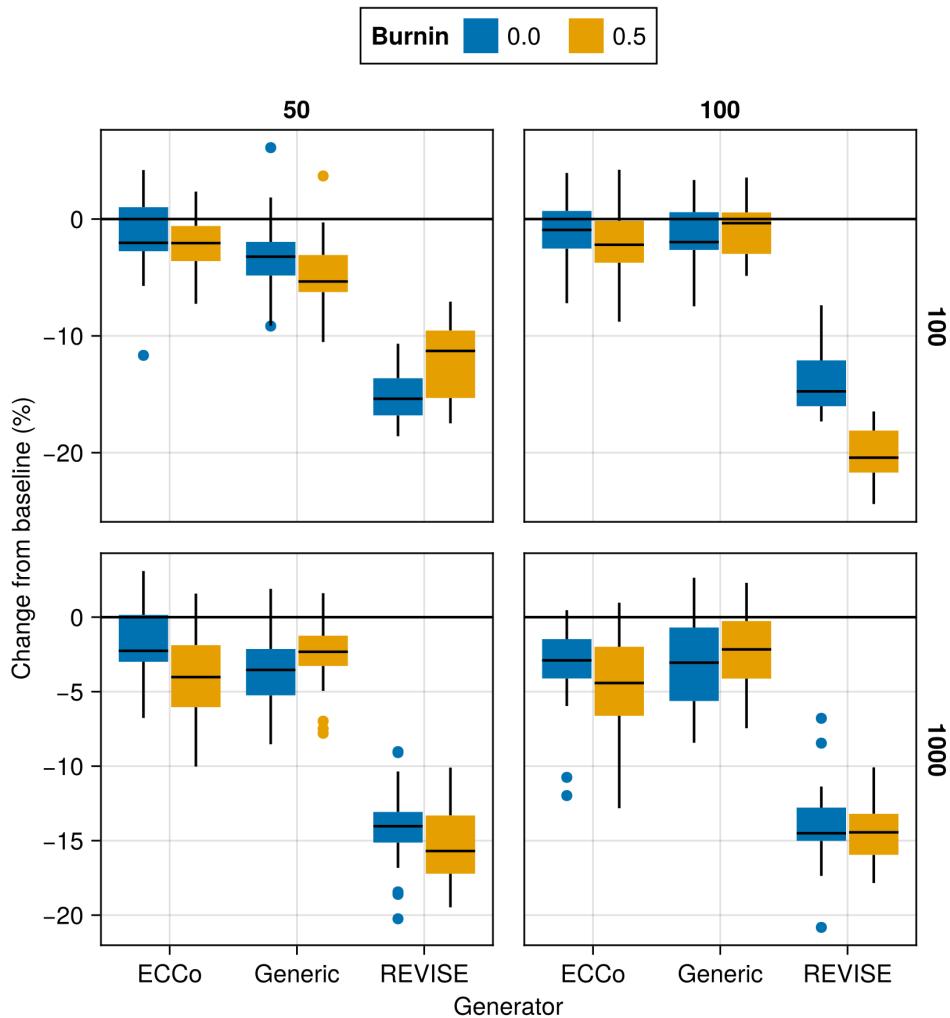


Figure 23: Average outcomes for the plausibility measure across hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

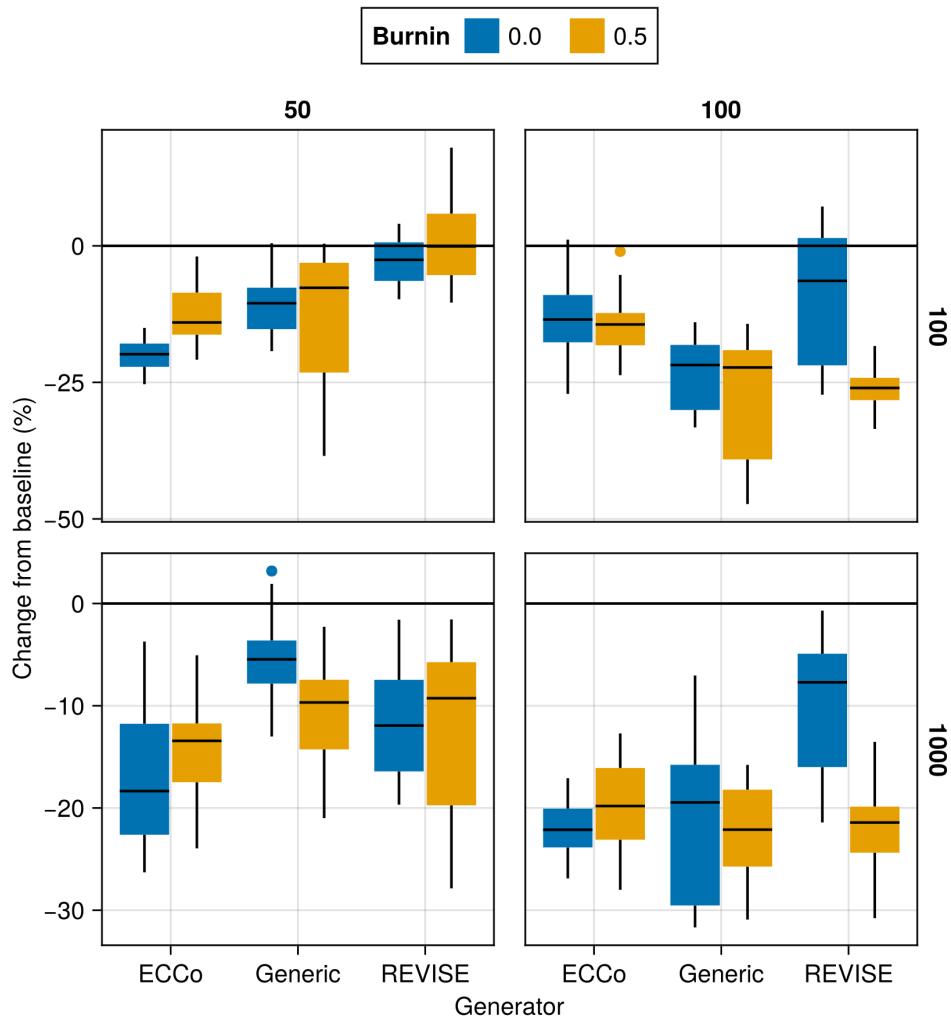


Figure 24: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

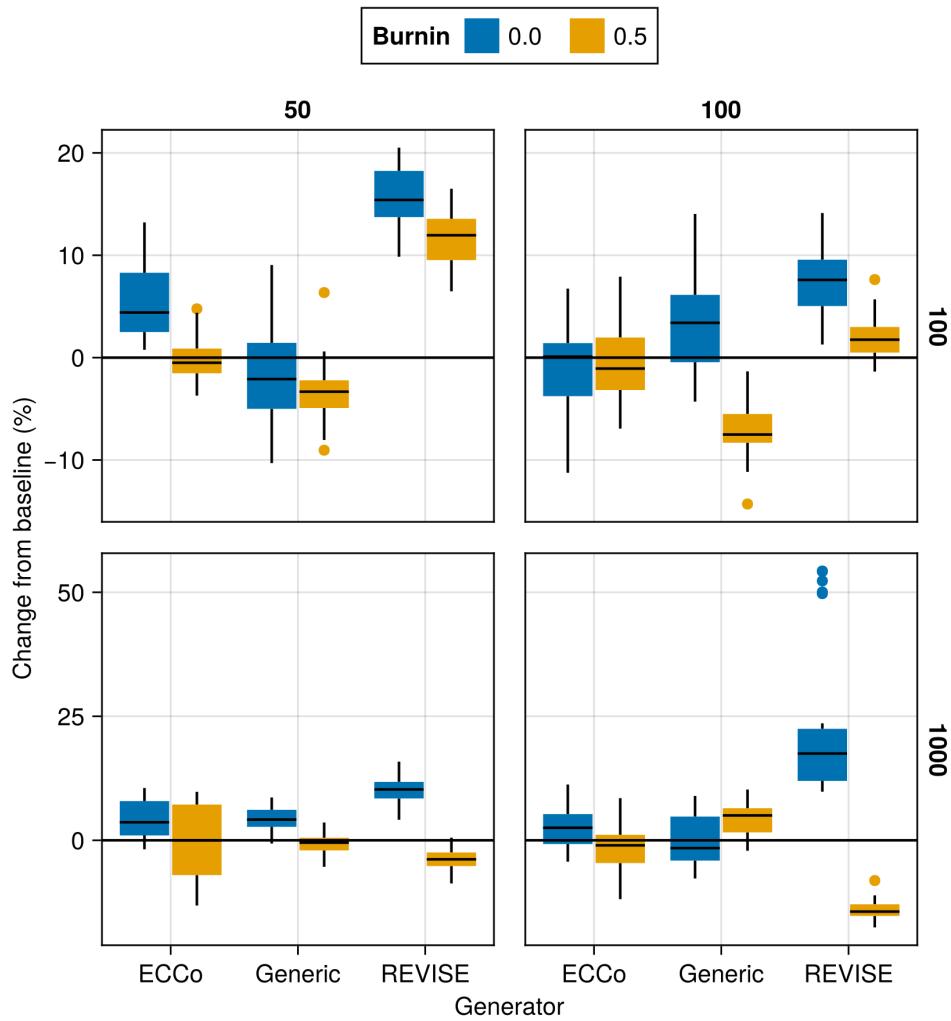


Figure 25: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

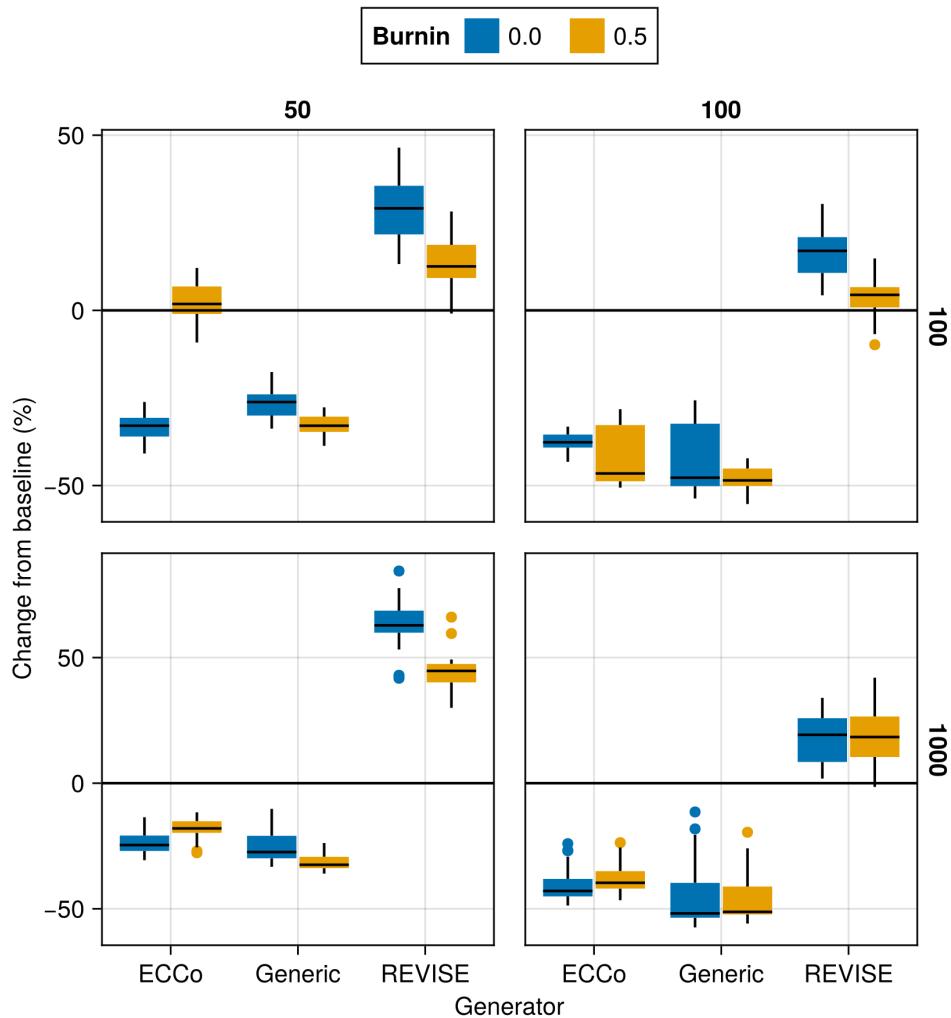


Figure 26: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

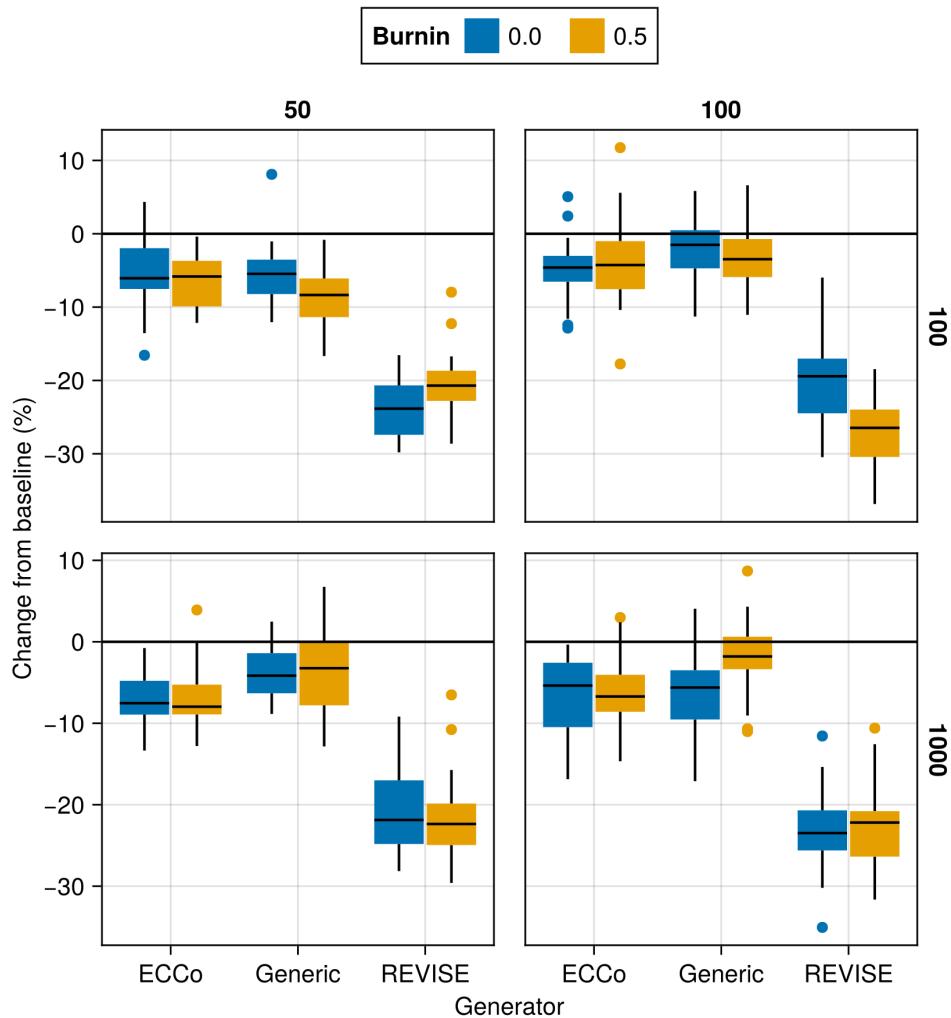


Figure 27: Average outcomes for the cost measure across hyperparameters. This shows the % change from the baseline model for the distance-based cost metric ([Wachter, Mittelstadt, and Russell 2017](#)). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCo*).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

Appendix E Tuning Key Parameters

Based on the findings from our initial large grid searches (Section D), we tune selected hyperparameters for all datasets: namely, the decision threshold τ and the strength of the energy regularization λ_{reg} . The final hyperparameter choices for each dataset are presented in Table 1 in Section C. Detailed results for each data set are shown in Figure 28 to Figure 45. From Table 1, 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 37 to Figure 45). 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 46 to Figure 54). For the given low decision threshold, we find that the learning rate has no discernable impact on the proportion of mature counterfactuals (Figure 55 to Figure 63). We do notice, however, that the results for MNIST are much improved when using a low value λ_{reg} , the strength for the energy regularization: plausibility is increased by up to $\sim 10\%$ (Figure 52) and the proportion of mature counterfactuals reaches 100%.

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

E.1 Key Parameters

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

Note 9: Training Phase

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

Note 10: Evaluation Phase

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

E.1.1 Plausibility

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

E.1.2 Proportion of Mature CE

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

E.2 Learning Rate

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

Note 11: Training Phase

- Generator Parameters:
 - Learning Rate: 0.1, 0.5, 1.0

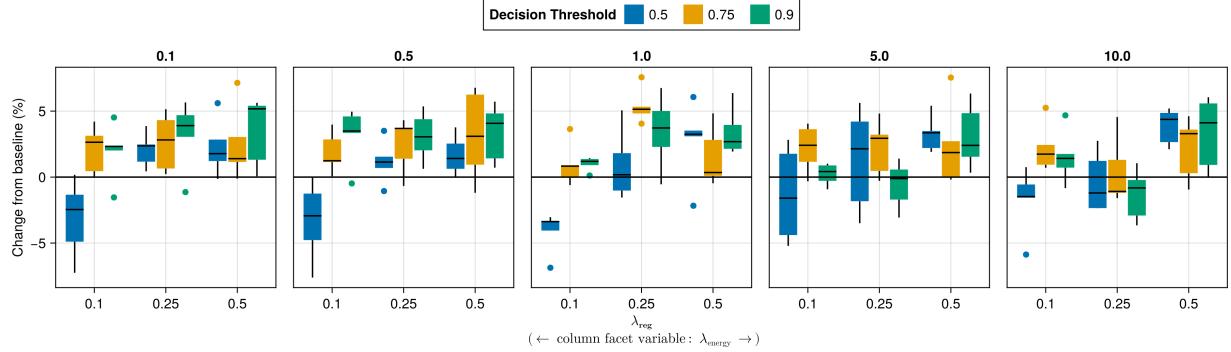


Figure 28: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for $ECCCo$). Data: Adult.

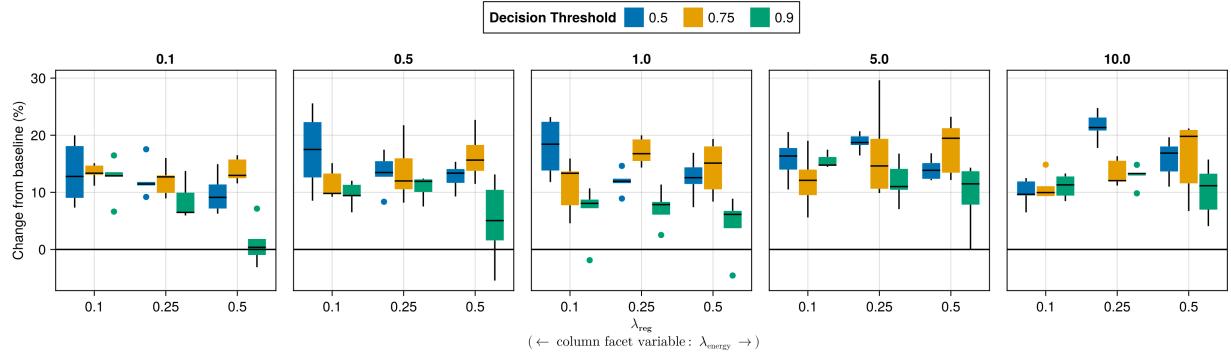


Figure 29: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for $ECCCo$). Data: California Housing.

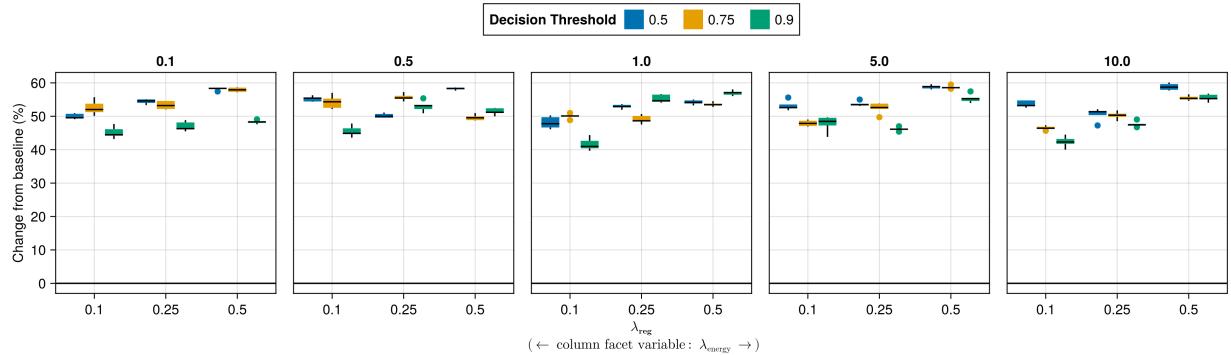


Figure 30: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for $ECCCo$). Data: Circles.

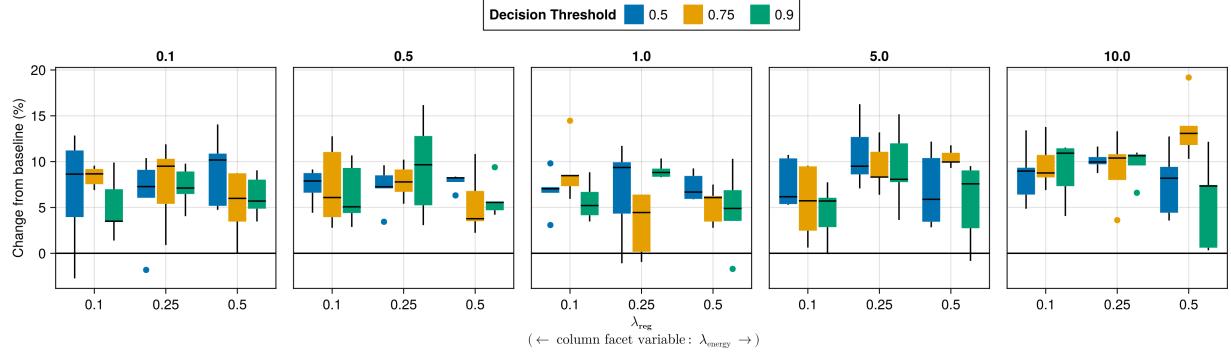


Figure 31: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Credit.

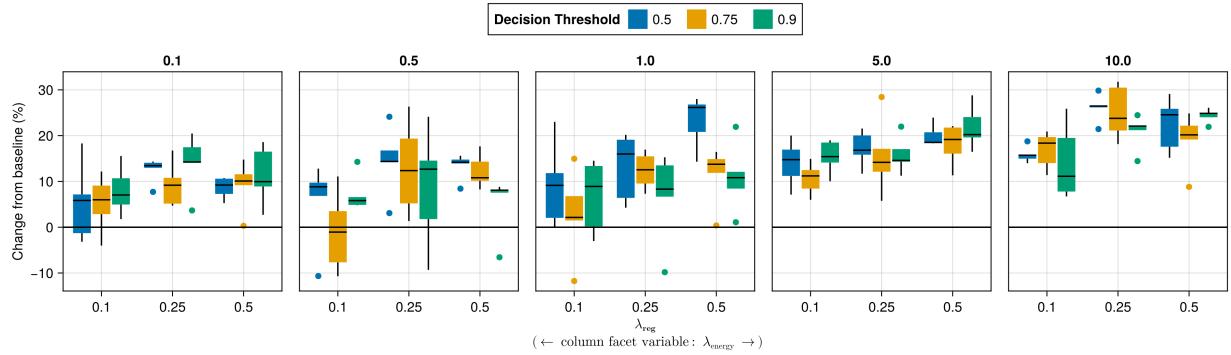


Figure 32: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: GMSC.

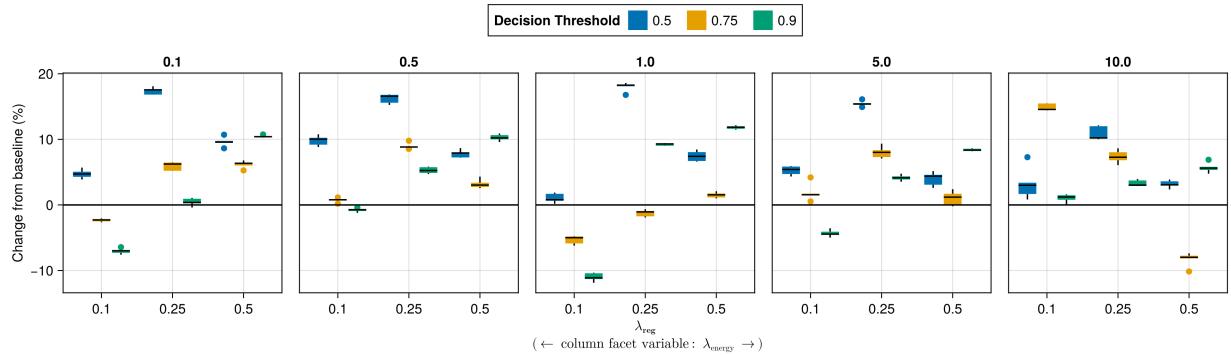


Figure 33: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for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

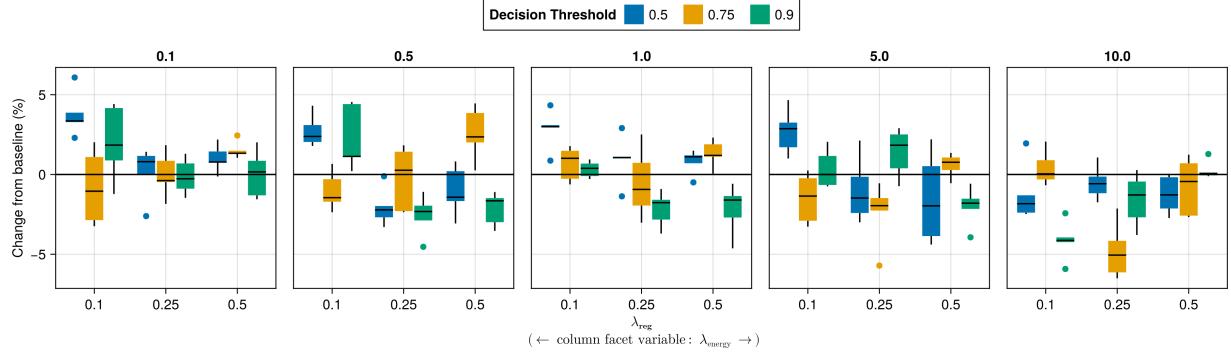


Figure 34: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for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

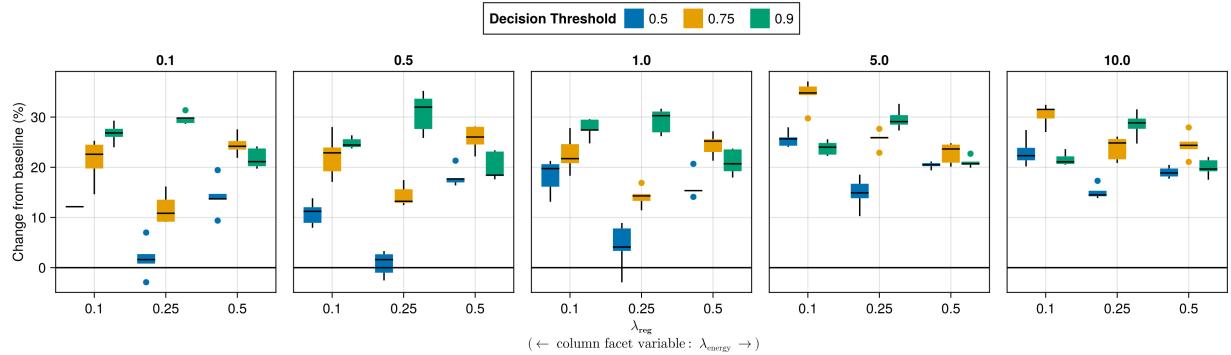


Figure 35: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for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

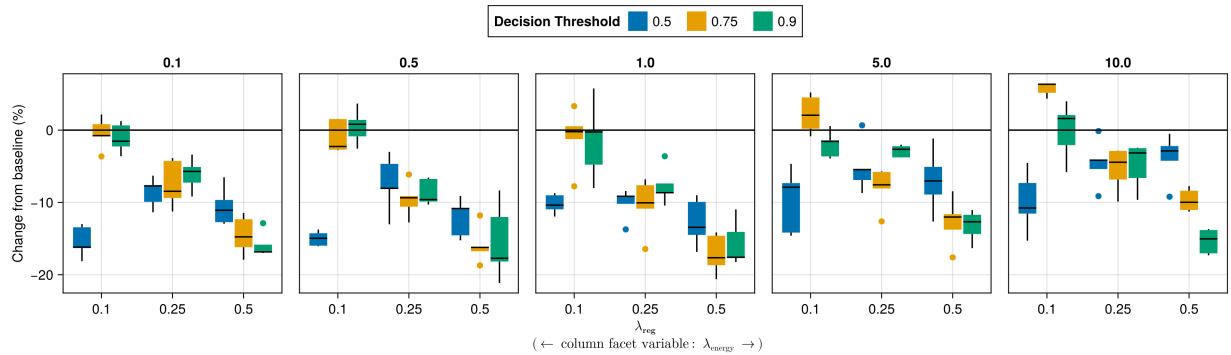


Figure 36: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for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

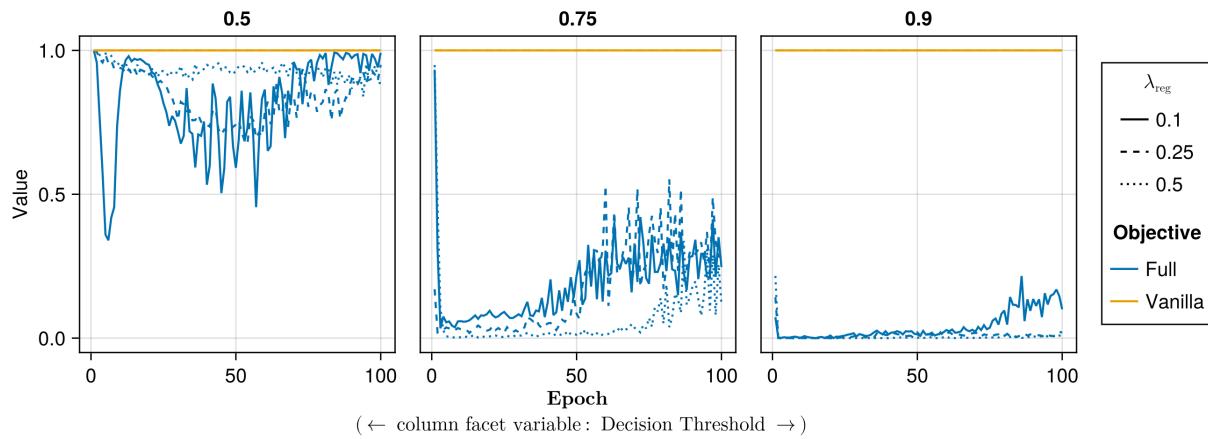


Figure 37: Proportion of mature counterfactuals in each epoch. Data: Adult.

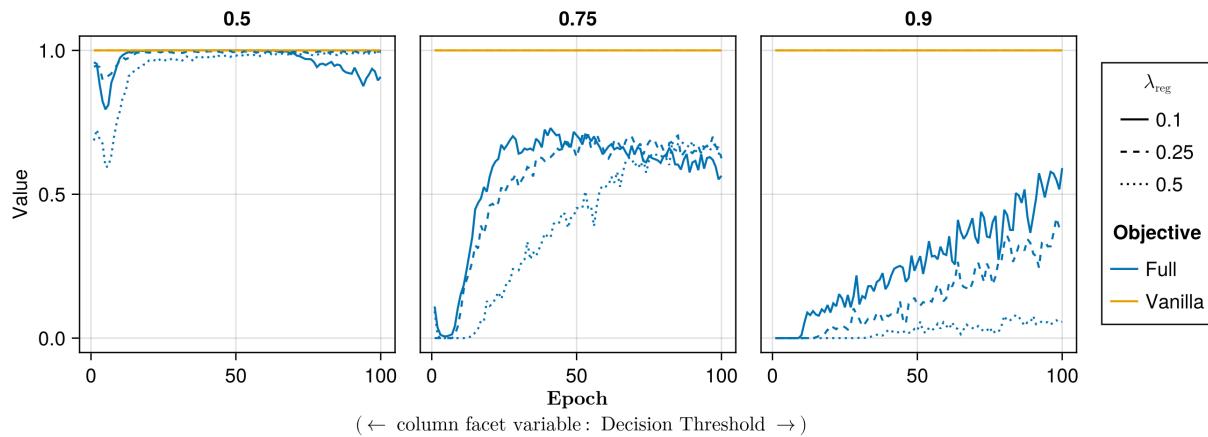


Figure 38: Proportion of mature counterfactuals in each epoch. Data: California Housing.

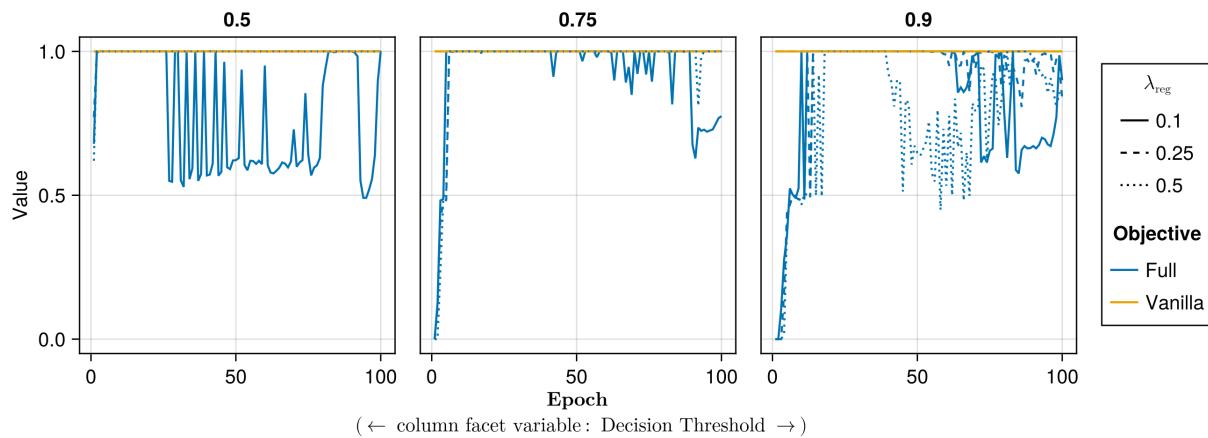


Figure 39: Proportion of mature counterfactuals in each epoch. Data: Circles.

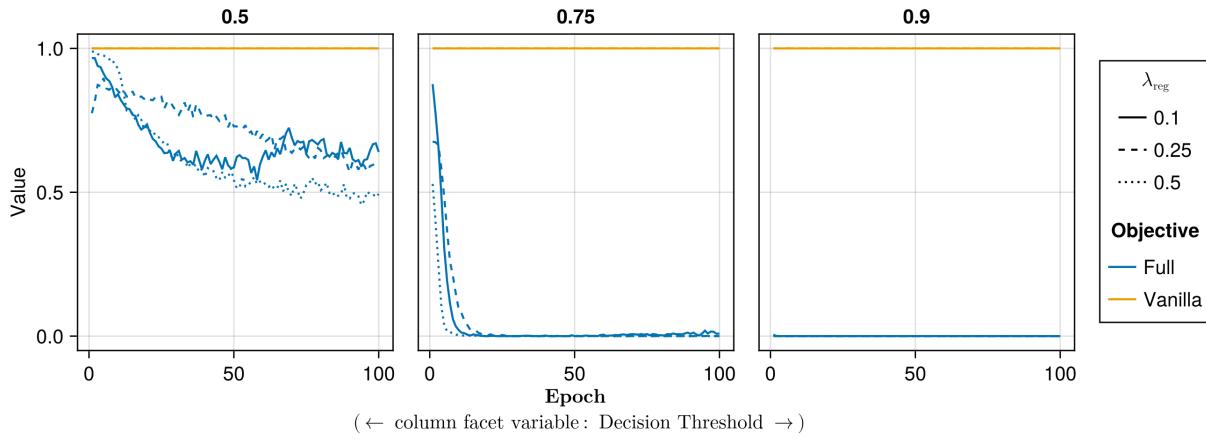


Figure 40: Proportion of mature counterfactuals in each epoch. Data: Credit.

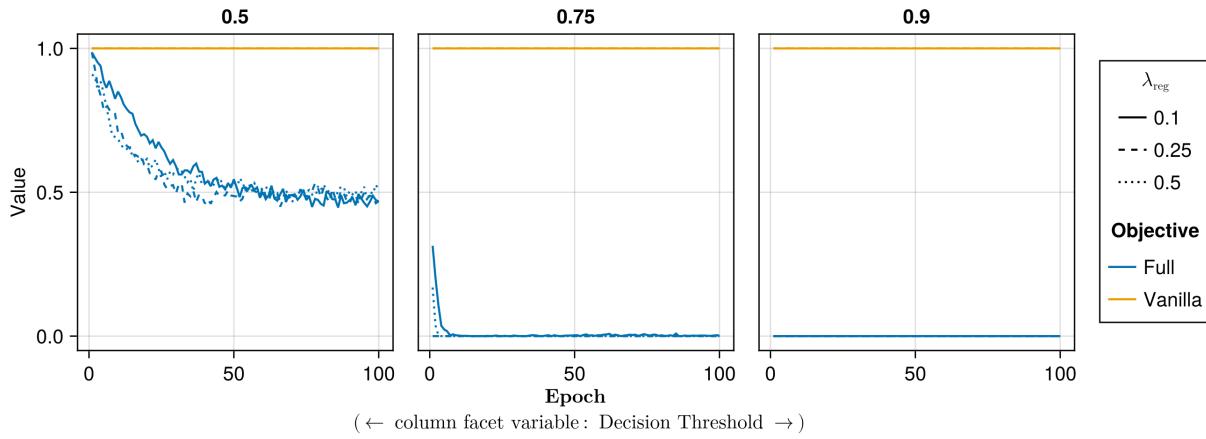


Figure 41: Proportion of mature counterfactuals in each epoch. Data: GMSC.

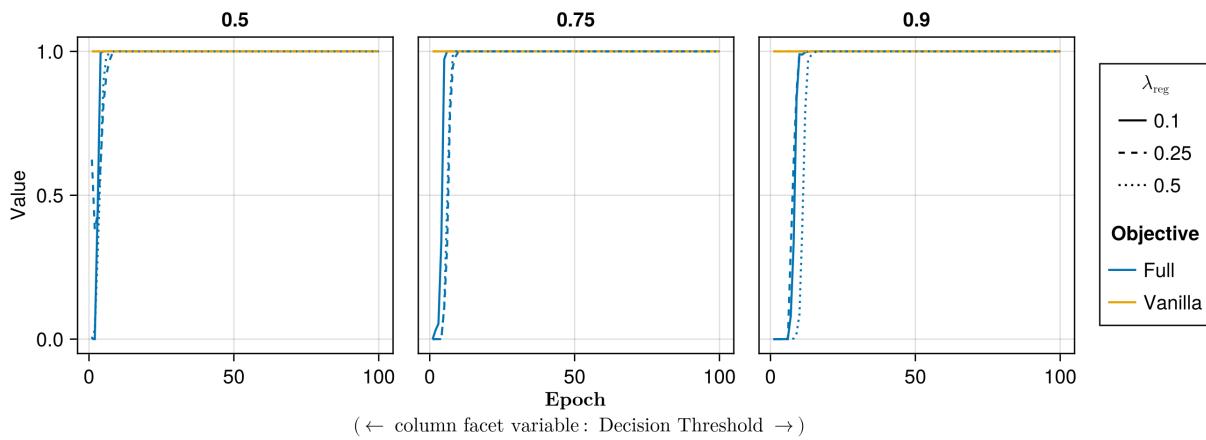


Figure 42: Proportion of mature counterfactuals in each epoch. Data: Linearly Separable.

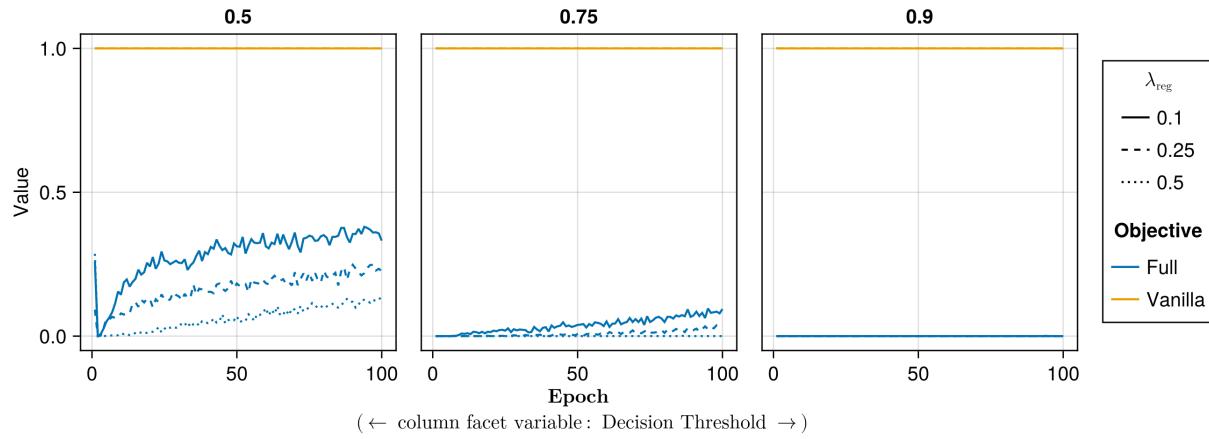


Figure 43: Proportion of mature counterfactuals in each epoch. Data: MNIST.

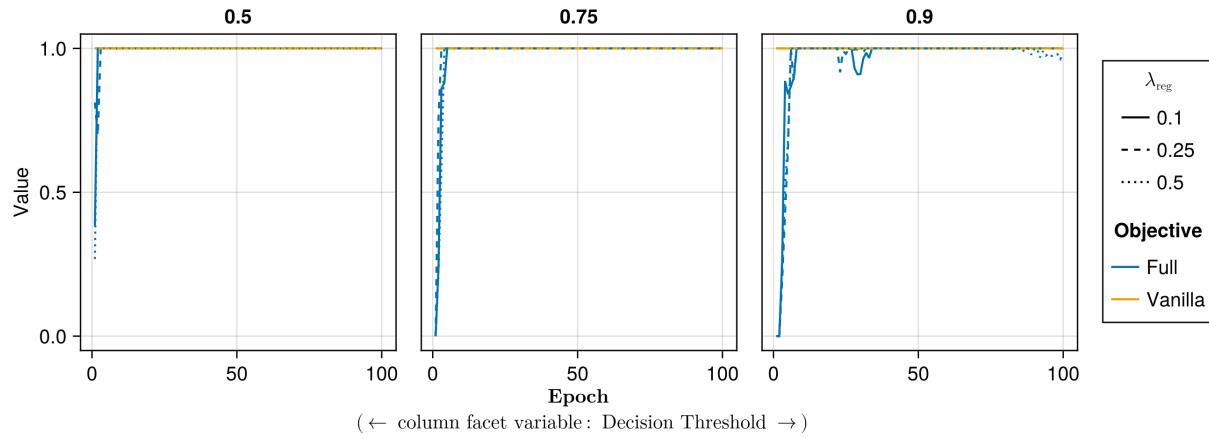


Figure 44: Proportion of mature counterfactuals in each epoch. Data: Moons.

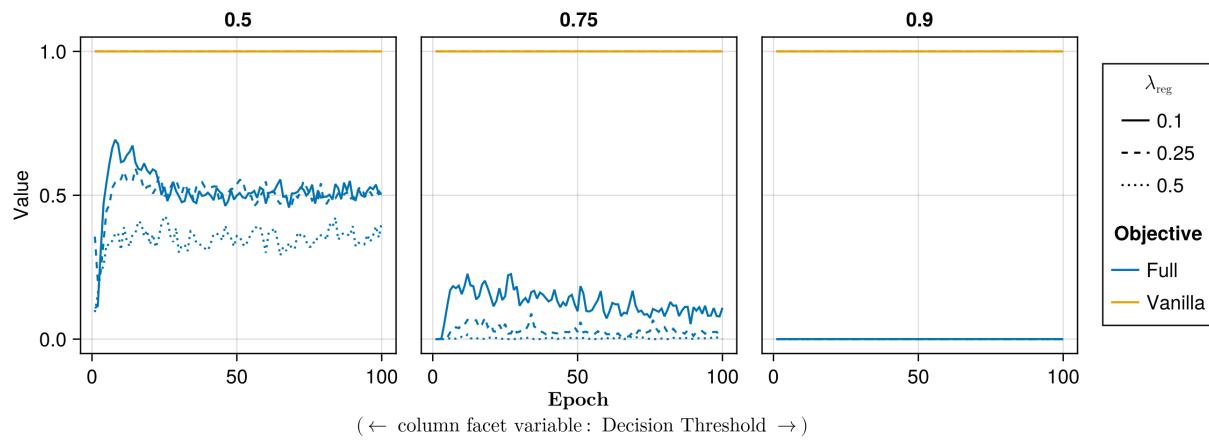


Figure 45: Proportion of mature counterfactuals in each epoch. Data: Overlapping.

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

Note 12: Evaluation Phase

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

E.2.1 Plausibility

The results with respect to the plausibility measure are shown in Figure 46 to Figure 54.

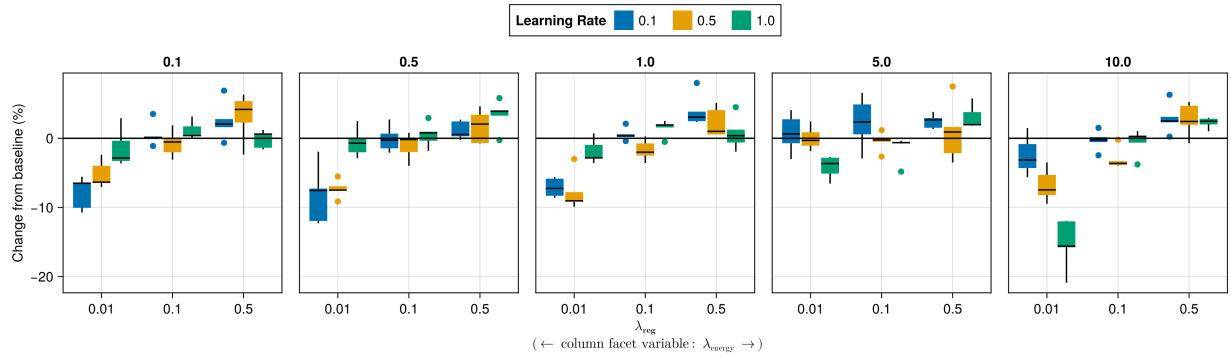


Figure 46: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: Adult.

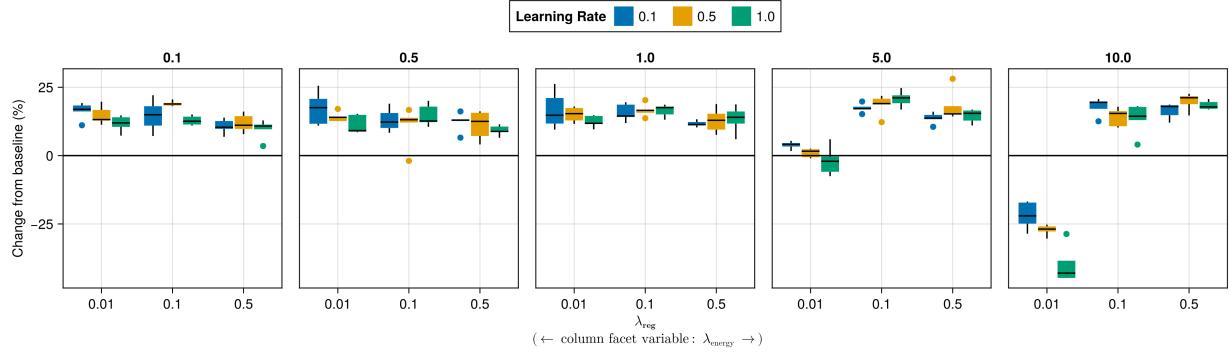


Figure 47: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for *ECCCo*). Data: California Housing.

E.2.2 Proportion of Mature CE

The results with respect to the proportion of mature counterfactuals in each epoch are shown in Figure 55 to Figure 63.

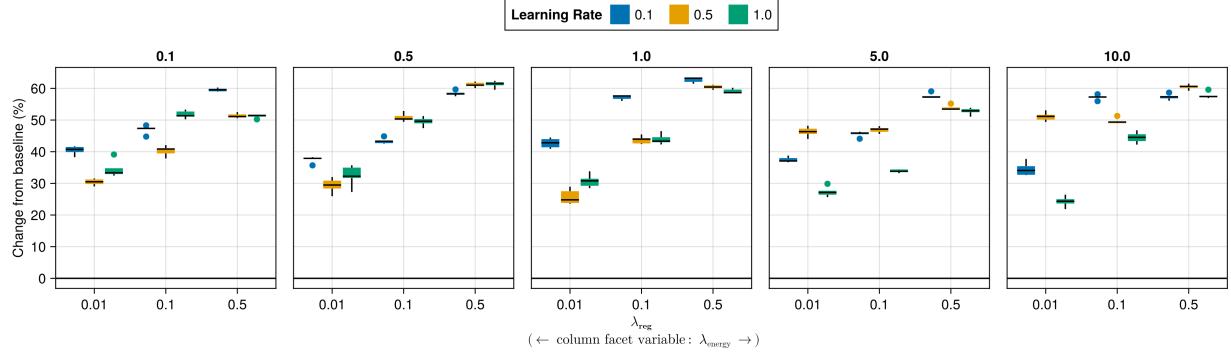


Figure 48: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for $ECCCo$). Data: Circles.

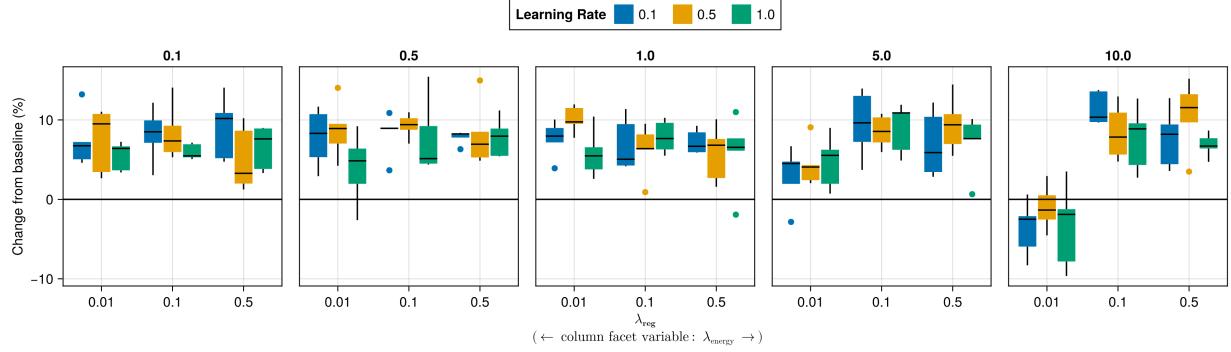


Figure 49: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for $ECCCo$). Data: Credit.

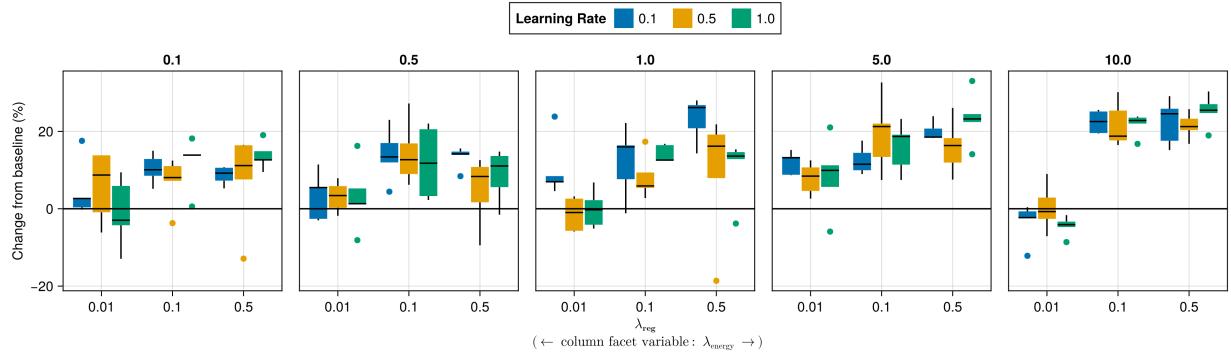


Figure 50: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for $ECCCo$). Data: GMSC.

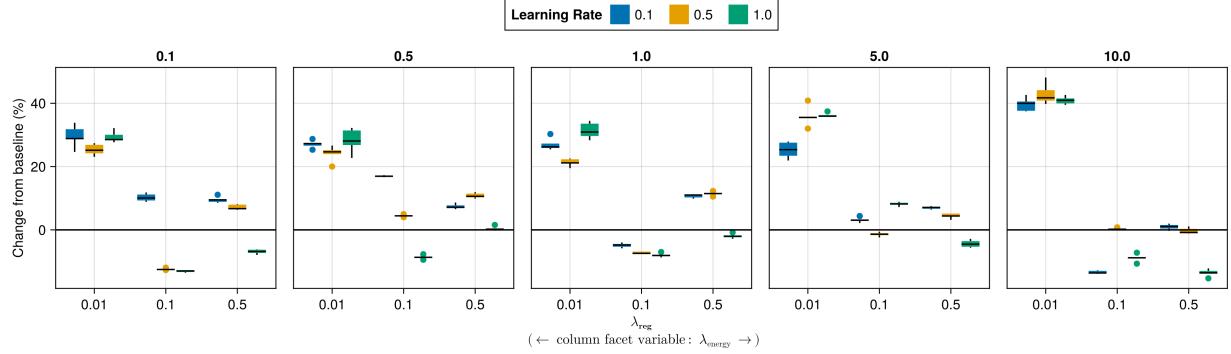


Figure 51: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for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

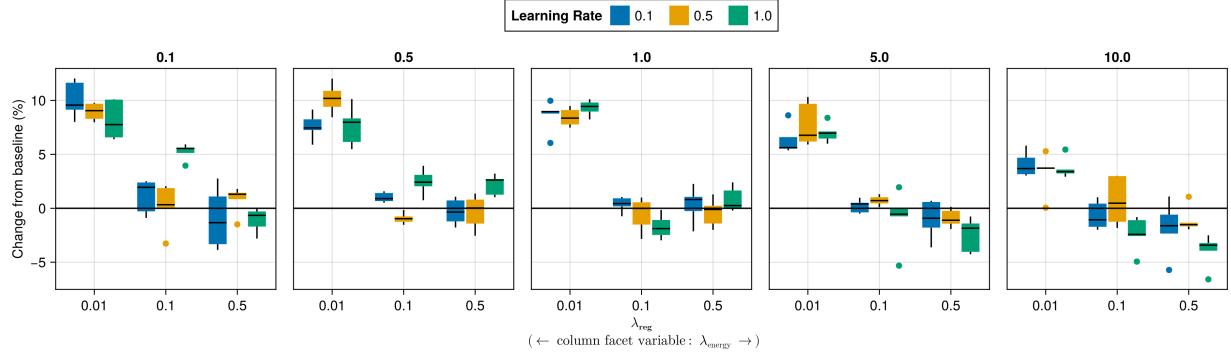


Figure 52: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for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

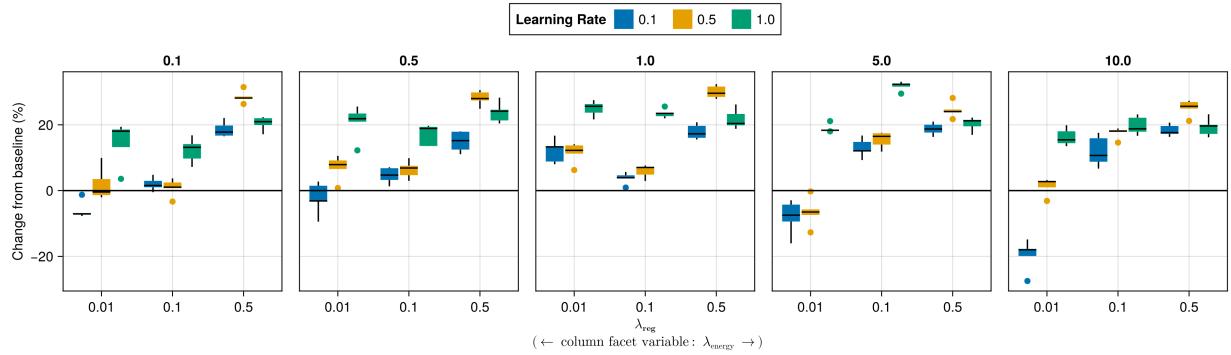


Figure 53: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for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

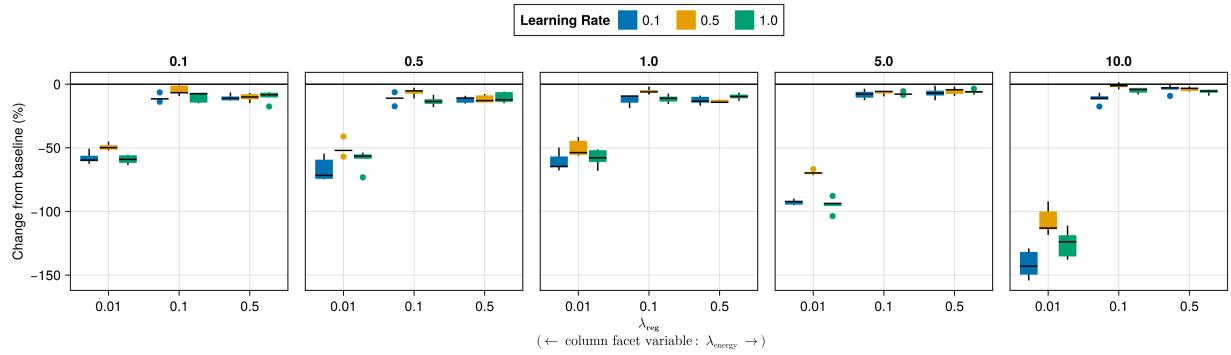


Figure 54: Average outcomes for the plausibility measure across key hyperparameters. This shows the % change from the baseline model for the distance-based implausibility metric ($\text{ext}\{\text{IP}\}$). Boxplots indicate the variation across evaluation runs and test settings (varying parameters for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

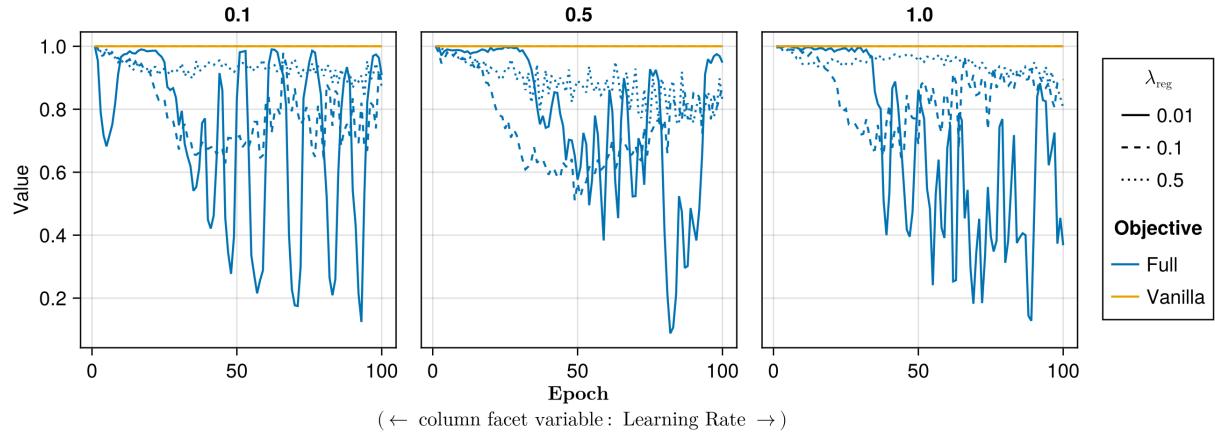


Figure 55: Proportion of mature counterfactuals in each epoch. Data: Adult.

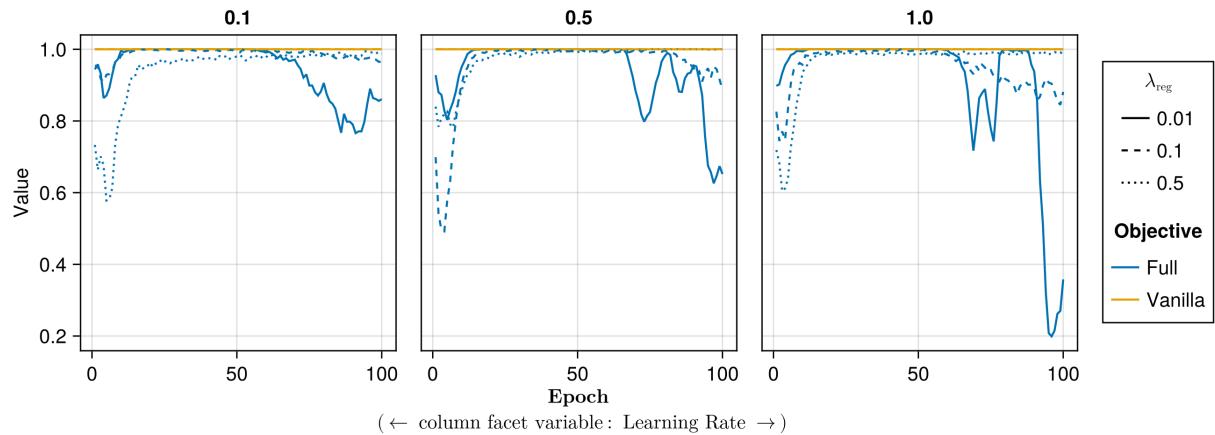


Figure 56: Proportion of mature counterfactuals in each epoch. Data: California Housing.

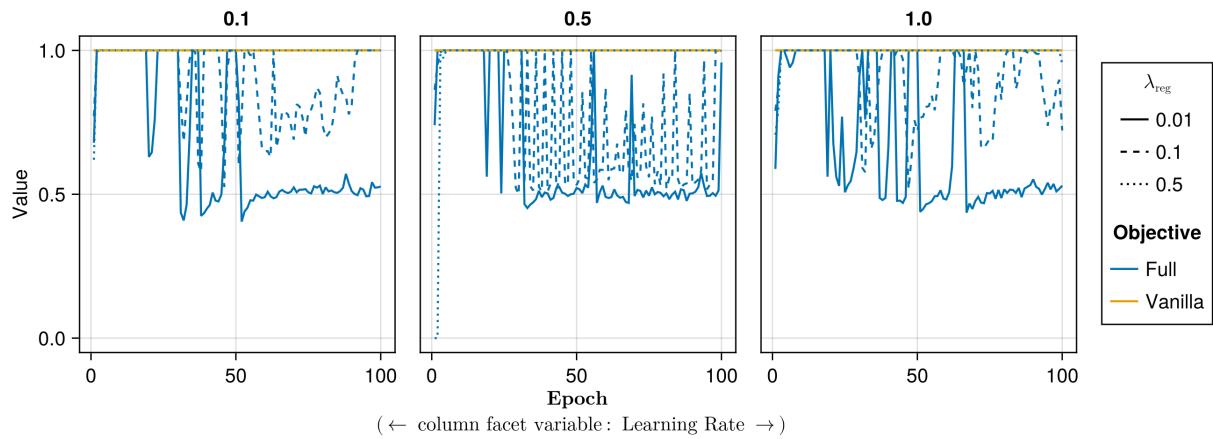


Figure 57: Proportion of mature counterfactuals in each epoch. Data: Circles.

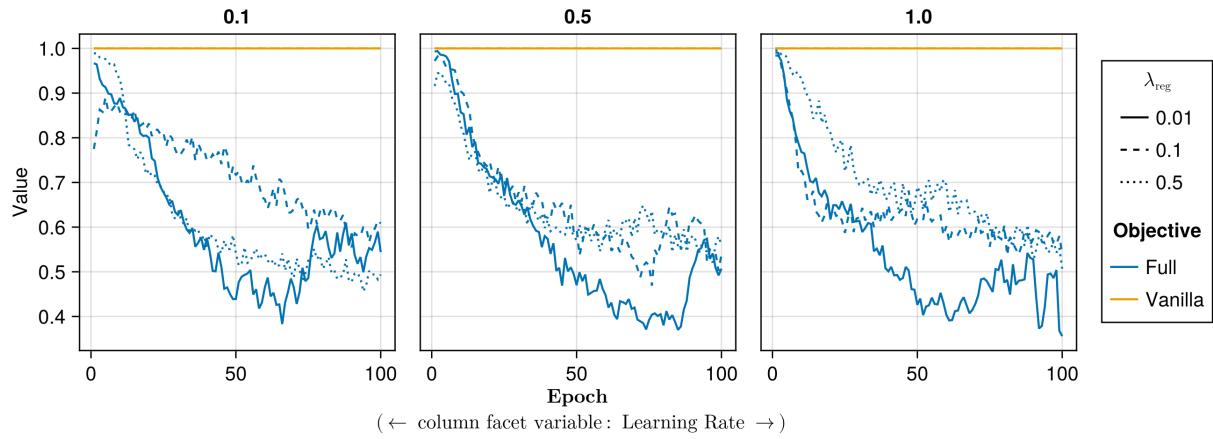


Figure 58: Proportion of mature counterfactuals in each epoch. Data: Credit.

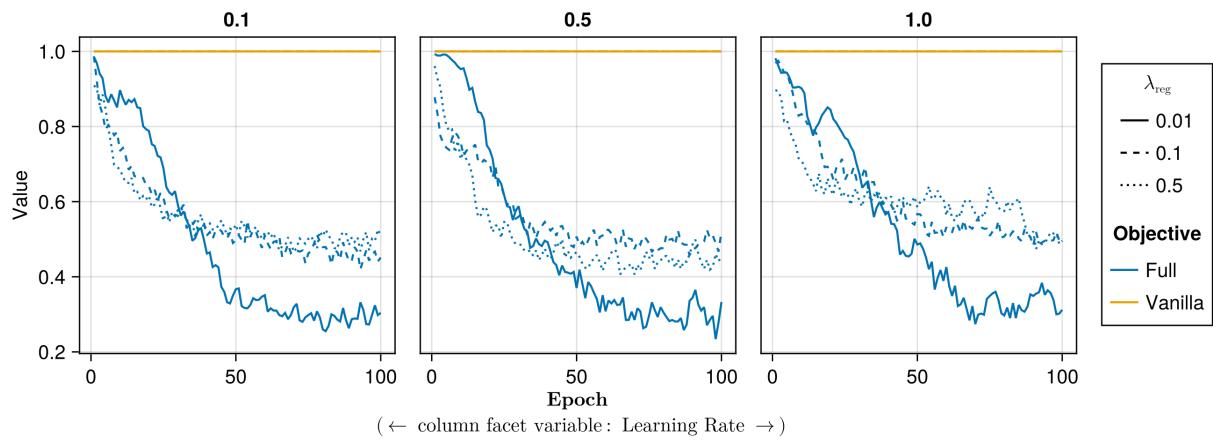


Figure 59: Proportion of mature counterfactuals in each epoch. Data: GMSC.

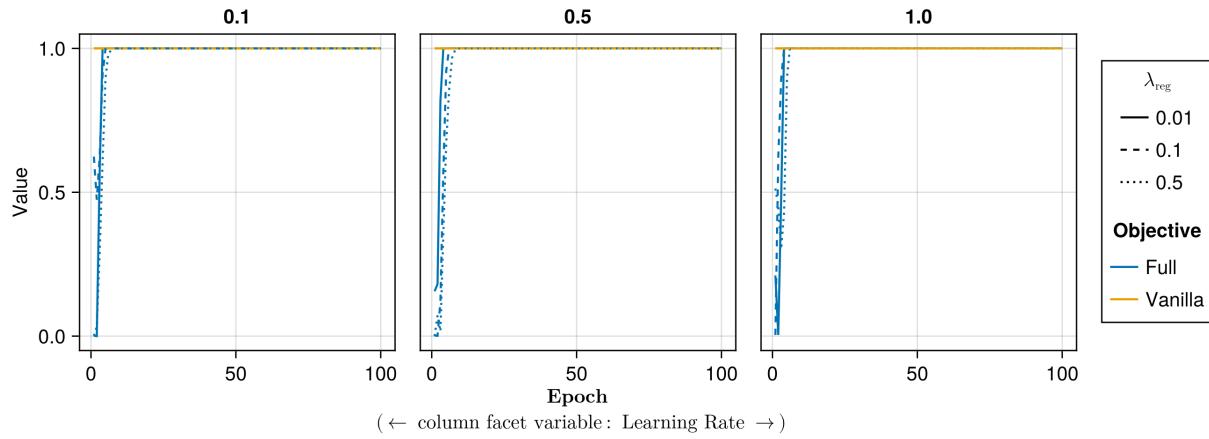


Figure 60: Proportion of mature counterfactuals in each epoch. Data: Linearly Separable.

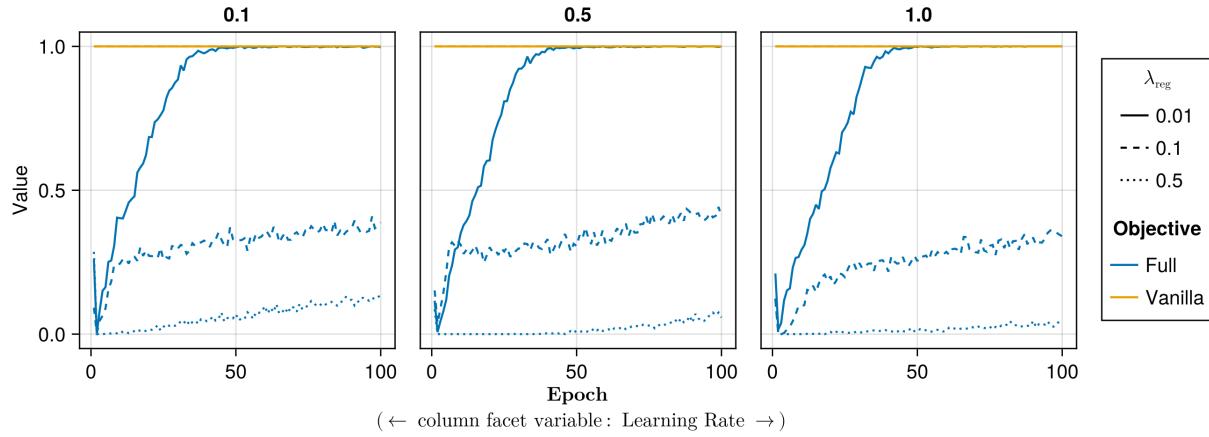


Figure 61: Proportion of mature counterfactuals in each epoch. Data: MNIST.

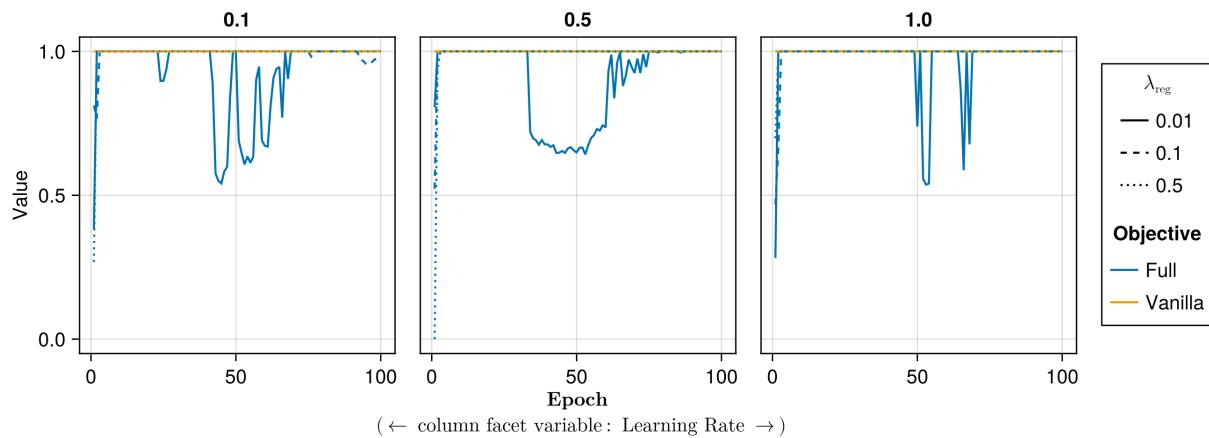


Figure 62: Proportion of mature counterfactuals in each epoch. Data: Moons.

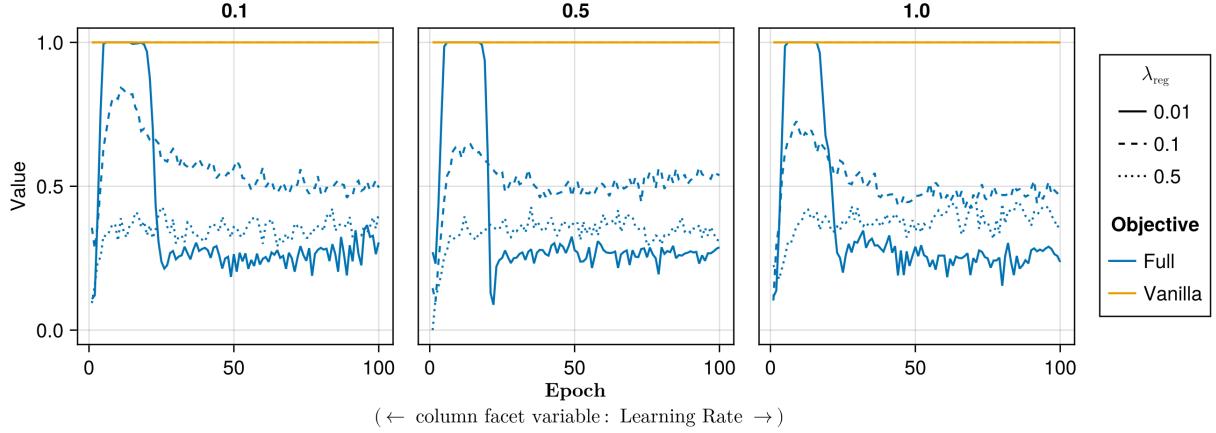


Figure 63: Proportion of mature counterfactuals in each epoch. Data: Overlapping.

Appendix F Computation Details

F.1 Hardware

We performed our experiments on a high-performance cluster. 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.

F.1.1 Grid Searches

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

F.1.2 Tuning

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

F.2 Software

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

References

- Altmeyer, Patrick, Arie van Deursen, and Cynthia C. S. Liem. 2023. “Explaining Black-Box Models through Counterfactuals.” In *Proceedings of the JuliaCon Conferences*, 1:130.
- Altmeyer, Patrick, Mojtaba Farmanbar, Arie van Deursen, and Cynthia C. S. Liem. 2024. “Faithful Model Explanations through Energy-Constrained Conformal Counterfactuals.” In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence*, 38:10829–37. 10. <https://doi.org/10.1609/aaai.v38i10.28956>.

- Bezanson, Jeff, Alan Edelman, Stefan Karpinski, and Viral B Shah. 2017. “Julia: A Fresh Approach to Numerical Computing.” *SIAM Review* 59 (1): 65–98. <https://doi.org/10.1137/141000671>.
- Bouchet-Valat, Milan, and Bogumi Kamiski. 2023. “DataFrames.jl: Flexible and Fast Tabular Data in Julia.” *Journal of Statistical Software* 107 (4): 1–32. <https://doi.org/10.18637/jss.v107.i04>.
- Byrne, Simon, Lucas C. Wilcox, and Valentin Churavy. 2021. “MPI.jl: Julia Bindings for the Message Passing Interface.” *Proceedings of the JuliaCon Conferences* 1 (1): 68. <https://doi.org/10.21105/jcon.00068>.
- Chagas, Ronan Arraes Jardim, Ben Baumgold, Glen Hertz, Hendrik Ranocha, Mark Wells, Nathan Boyer, Nicholas Ritchie, et al. 2024. “Ronisbr/PrettyTables.jl: V2.4.0.” Zenodo. <https://doi.org/10.5281/zenodo.1383553>.
- Christ, Simon, Daniel Schwabeneder, Christopher Rackauckas, Michael Krabbe Borregaard, and Thomas Breloff. 2023. “Plots.jl – a User Extendable Plotting API for the Julia Programming Language.” <https://doi.org/https://doi.org/10.5334/jors.431>.
- Danisch, Simon, and Julius Krumbiegel. 2021. “Makie.jl: Flexible High-Performance Data Visualization for Julia.” *Journal of Open Source Software* 6 (65): 3349. <https://doi.org/10.21105/joss.03349>.
- Goodfellow, Ian, Jonathon Shlens, and Christian Szegedy. 2015. “Explaining and Harnessing Adversarial Examples.” <https://arxiv.org/abs/1412.6572>.
- Grathwohl, Will, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. 2020. “Your Classifier Is Secretly an Energy Based Model and You Should Treat It Like One.” In *International Conference on Learning Representations*.
- Gretton, Arthur, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. 2012. “A Kernel Two-Sample Test.” *The Journal of Machine Learning Research* 13 (1): 723–73.
- Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2009. *The Elements of Statistical Learning*. Springer New York. <https://doi.org/10.1007/978-0-387-84858-7>.
- Innes, Michael, Elliot Saba, Keno Fischer, Dhairyा Gandhi, Marco Conchetto Rudilosso, Neethu Mariya Joy, Tejan Karmali, Avik Pal, and Viral Shah. 2018. “Fashionable Modelling with Flux.” <https://arxiv.org/abs/1811.01457>.
- Innes, Mike. 2018. “Flux: Elegant Machine Learning with Julia.” *Journal of Open Source Software* 3 (25): 602. <https://doi.org/10.21105/joss.00602>.
- Joshi, Shalmali, Oluwasanmi Koyejo, Warut Vigitbenjaronk, Been Kim, and Joydeep Ghosh. 2019. “Towards Realistic Individual Recourse and Actionable Explanations in Black-Box Decision Making Systems.” <https://arxiv.org/abs/1907.09615>.
- Wachter, Sandra, Brent Mittelstadt, and Chris Russell. 2017. “Counterfactual Explanations Without Opening the Black Box: Automated Decisions and the GDPR.” *Harv. JL & Tech.* 31: 841. <https://doi.org/10.2139/ssrn.3063289>.