

---

# ECCCOs from the Black Box: Faithful Explanations through Energy-Constrained Conformal Counterfactuals

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Counterfactual Explanations offer an intuitive and straightforward way to explain black-box models and offer Algorithmic Recourse to individuals. To address the need for plausible explanations, existing work has primarily relied on surrogate models to learn how the input data is distributed. This effectively reallocates the task of learning realistic explanations for the data from the model itself to the surrogate. Consequently, the generated explanations may seem plausible to humans but need not necessarily describe the behaviour of the black-box model faithfully. We formalise this notion of faithfulness through the introduction of a tailored evaluation metric and propose a novel algorithmic framework for generating Energy-Constrained Conformal Counterfactuals (ECCCOs) that are only as plausible as the model permits. Through extensive empirical studies involving multiple synthetic and real-world datasets, we demonstrate that ECCCOs reconcile the need for plausibility and faithfulness. In particular, we show that it is possible to achieve state-of-the-art plausibility for models with gradient access without the need for surrogate models. To do so, our framework relies solely on properties defining the black-box model itself by leveraging recent advances in energy-based modelling and conformal prediction. To our knowledge, this is the first venture in this direction for generating faithful Counterfactual Explanations. Thus, we anticipate that ECCCOs can serve as a baseline for future research. We believe that our work opens avenues for researchers and practitioners seeking tools to better distinguish trustworthy from unreliable models.

## 1 Introduction

Counterfactual Explanations (CE) provide a powerful, flexible and intuitive way to not only explain black-box models but also help affected individuals through the means of Algorithmic Recourse. Instead of opening the Black Box, CE works under the premise of strategically perturbing model inputs to understand model behaviour [29]. Intuitively speaking, we generate explanations in this context by asking what-if questions of the following nature: ‘Our credit risk model currently predicts that this individual is not credit-worthy. What if they reduced their monthly expenditures by 10%?’

This is typically implemented by defining a target outcome  $\mathbf{y}^+ \in \mathcal{Y}$  for some individual  $\mathbf{x} \in \mathcal{X} = \mathbb{R}^D$  described by  $D$  attributes, for which the model  $M_\theta : \mathcal{X} \mapsto \mathcal{Y}$  initially predicts a different outcome:  $M_\theta(\mathbf{x}) \neq \mathbf{y}^+$ . Counterfactuals are then searched by minimizing a loss function that compares the predicted model output to the target outcome:  $y_{\text{loss}}(M_\theta(\mathbf{x}), \mathbf{y}^+)$ . Since Counterfactual Explanations work directly with the black-box model, valid counterfactuals always have full local fidelity by construction where fidelity is defined as the degree to which explanations approximate the predictions of a black-box model [17, 16].

In situations where full fidelity is a requirement, CE offers a more appropriate solution to Explainable Artificial Intelligence (XAI) than other popular approaches like LIME [22] and SHAP [13], which involve local surrogate models. But even full fidelity is not a sufficient condition for ensuring that an explanation faithfully describes the behaviour of a model. That is because multiple very distinct explanations can all lead to the same model prediction, especially when dealing with heavily parameterized models like deep neural networks, which are typically underspecified by the data [31].

In the context of CE, the idea that no two explanations are the same arises almost naturally. A key focus in the literature has therefore been to identify those explanations and algorithmic recourses that are most appropriate based on a myriad of desiderata such as sparsity, actionability and plausibility. In this work, we draw closer attention to model faithfulness rather than fidelity as a desideratum for counterfactuals. Our key contributions are as follows:

- We show that fidelity is an insufficient evaluation metric for counterfactuals (Section 3) and propose a definition of faithfulness that gives rise to more suitable metrics (Section 4).
- We introduce a novel algorithmic approach for generating Energy-Constrained Conformal Counterfactuals (ECCCos) in Section 5.
- We provide extensive empirical evidence demonstrating that ECCCos faithfully explain model behaviour without sacrificing plausibility (Section 6).

Thus, we believe that our work opens avenues for researchers and practitioners seeking tools to better distinguish trustworthy from unreliable models.

## 2 Background

While Counterfactual Explanations can be generated for arbitrary regression models [24], existing work has primarily focused on classification problems. Let  $\mathcal{Y} = (0, 1)^K$  denote the one-hot-encoded output domain with  $K$  classes. Then most counterfactual generators rely on gradient descent to optimize different flavours of the following counterfactual search objective:

$$\mathbf{Z}' = \arg \min_{\mathbf{Z}' \in \mathcal{Z}^L} \{ \text{yloss}(M_\theta(f(\mathbf{Z}')), \mathbf{y}^+) + \lambda \text{cost}(f(\mathbf{Z}')) \} \quad (1)$$

Here  $\text{yloss}$  denotes the primary loss function,  $f(\cdot)$  is a function that maps from the counterfactual state space to the feature space and  $\text{cost}$  is either a single penalty or a collection of penalties that are used to impose constraints through regularization. Equation 1 restates the baseline approach to gradient-based counterfactual search proposed by Wachter et al. [29] in general form as introduced by Altmeyer et al. [2]. To explicitly account for the multiplicity of explanations  $\mathbf{Z}' = \{\mathbf{z}_l\}_L$  denotes an  $L$ -dimensional array of counterfactual states.

The baseline approach, which we will simply refer to as **Wachter** [29], searches a single counterfactual directly in the feature space and penalises its distance to the original factual. In this case,  $f(\cdot)$  is simply the identity function and  $\mathcal{Z}$  corresponds to the feature space itself. Many derivative works of Wachter et al. [29] have proposed new flavours of Equation 1, each of them designed to address specific *desiderata* that counterfactuals ought to meet in order to properly serve both AI practitioners and individuals affected by algorithmic decision-making systems. The list of desiderata includes but is not limited to the following: sparsity, proximity [29], actionability [27], diversity [17], plausibility [8, 21, 23], robustness [26, 20, 2] and causality [11]. Different counterfactual generators addressing these needs have been extensively surveyed and evaluated in various studies [28, 10, 19, 4, 7].

Perhaps unsurprisingly, the different desiderata are often positively correlated. For example, Artelt et al. [4] find that plausibility typically also leads to improved robustness. Similarly, plausibility has also been connected to causality in the sense that plausible counterfactuals respect causal relationships [14]. Consequently, the plausibility of counterfactuals has been among the primary concerns for researchers. Achieving plausibility is equivalent to ensuring that the generated counterfactuals comply with the true and unobserved data-generating process (DGP). We define plausibility formally in this work as follows:

**Definition 2.1** (Plausible Counterfactuals). *Let  $\mathcal{X}|\mathbf{y}^+$  denote the true conditional distribution of samples in the target class  $\mathbf{y}^+$ . Then for  $\mathbf{x}'$  to be considered a plausible counterfactual, we need:  $\mathbf{x}' \sim \mathcal{X}|\mathbf{y}^+$ .*

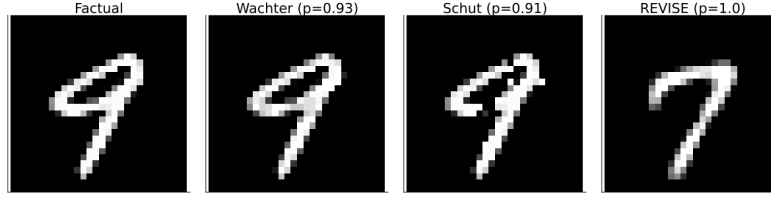


Figure 1: Counterfactuals for turning a 9 (nine) into a 7 (seven): original image (left); then from left to right the counterfactuals generated using Wachter, Schut and REVISE

To generate plausible counterfactuals, we need to be able to quantify the DGP:  $\mathcal{X}|\mathbf{y}^+$ . One straightforward way to do this is to use surrogate models for the task. Joshi et al. [8], for example, suggest that instead of searching counterfactuals in the feature space  $\mathcal{X}$ , we can instead traverse a latent embedding  $\mathcal{Z}$  (Equation 1) that implicitly codifies the DGP. To learn the latent embedding, they propose using a generative model such as a Variational Autoencoder (VAE). Provided the surrogate model is well-trained, their proposed approach called **REVISE** can yield plausible explanations. Others have proposed similar approaches: Dombrowski et al. [5] traverse the base space of a normalizing flow to solve Equation 1; Poyiadzi et al. [21] use density estimators ( $\hat{p} : \mathcal{X} \mapsto [0, 1]$ ) to constrain the counterfactuals to dense regions in the feature space; and, finally, Karimi et al. [11] assume knowledge about the structural causal model that generates the data.

A competing approach towards plausibility that is also closely related to this work instead relies on the black-box model itself. Schut et al. [23] show that to meet the plausibility objective we need not explicitly model the input distribution. Pointing to the undesirable engineering overhead induced by surrogate models, they propose that we rely on the implicit minimisation of predictive uncertainty instead. Their proposed methodology, which we will refer to as **Schut**, solves Equation 1 by greedily applying JSMA in the feature space with standard cross-entropy loss and no penalty at all. The authors demonstrate theoretically and empirically that their approach yields counterfactuals for which the model  $M_\theta$  predicts the target label  $\mathbf{y}^+$  with high confidence. Provided the model is well-specified, these counterfactuals are plausible. This idea hinges on the assumption that the black-box model provides well-calibrated predictive uncertainty estimates.

### 3 Why Fidelity is not Enough

As discussed in the introduction, any valid Counterfactual Explanation also has full fidelity by construction: solutions to Equation 1 are considered valid as soon as the label predicted by the model matches the target class. So while fidelity always applies, counterfactuals that address the various desiderata introduced above can look vastly different from each other. The following motivating example illustrates this point and demonstrates why fidelity is an insufficient evaluation metric to assess the faithfulness of Counterfactual Explanations.

We have trained a simple image classifier  $M_\theta$  on the well-known MNIST dataset [12]: a Multi-Layer Perceptron (MLP) with above 90 percent test accuracy. No measures have been taken to improve the model’s adversarial robustness or its capacity for predictive uncertainty quantification. The far left panel of Figure 1 shows a random sample drawn from the dataset. The underlying classifier correctly predicts the label ‘nine’ for this image. For the given factual image and model, we have used Wachter, Schut and REVISE to generate one counterfactual each in the target class ‘seven’. The perturbed images are shown next to the factual image from left to right in Figure 1. Captions on top of the individual images indicate the generator along with the predicted probability that the image belongs to the target class. In all three cases that probability is above 90 percent and yet the counterfactuals look very different from each other.

Since Wachter is only concerned with proximity, the generated counterfactual is almost indistinguishable from the factual. The approach by Schut expects a well-calibrated model that can generate predictive uncertainty estimates. Since this is not the case, the generated counterfactual looks like an adversarial example. Finally, the counterfactual generated by REVISE looks much more plausible than the other two. But is it also more faithful to the behaviour of our MNIST classifier? That is much less clear because the surrogate used by REVISE introduces friction: the generated explanations no longer depend exclusively on the black-box model itself.

129 So which of the counterfactuals most faithfully explains the behaviour of our image classifier? Fidelity  
 130 cannot help us to make that judgement, because all of these counterfactuals have full fidelity. To  
 131 bridge this gap, we introduce a new notion of faithfulness in the following section.

## 132 4 A new Notion of Faithfulness

133 Analogous to Definition 2.1, we propose to define faithfulness in the context of Counterfactual  
 134 Explanations as follows:

135 **Definition 4.1** (Faithful Counterfactuals). *Let  $\mathcal{X}_\theta|\mathbf{y}^+ = p_\theta(\mathbf{X}_{\mathbf{y}^+})$  denote the conditional distribution*  
 136 *of  $\mathbf{x}$  in the target class  $\mathbf{y}^+$ , where  $\theta$  denotes the parameters of model  $M_\theta$ . Then for  $\mathbf{x}'$  to be considered*  
 137 *a conformal counterfactual, we need:  $\mathbf{x}' \sim \mathcal{X}_\theta|\mathbf{y}^+$ .*

138 In doing this, we merge in and nuance the concept of plausibility (Definition 2.1) where the notion of  
 139 ‘consistent with the data’ becomes ‘consistent with what the model has learned about the data’.

### 140 4.1 Quantifying the Model’s Generative Property

141 To assess counterfactuals with respect to Definition 4.1, we need a way to quantify the posterior  
 142 conditional distribution  $p_\theta(\mathbf{x}|\mathbf{y}^+)$ . To this end, we draw on recent advances in Energy-Based  
 143 Modelling (EBM), a subdomain of machine learning that is concerned with generative or hybrid  
 144 modelling [6? ]. In particular, note that if we fix  $\mathbf{y}$  to our target value  $\mathbf{y}^+$ , we can conditionally draw  
 145 from  $p_\theta(\mathbf{x}|\mathbf{y}^+)$  using Stochastic Gradient Langevin Dynamics (SGLD) as follows,

$$\mathbf{x}_{j+1} \leftarrow \mathbf{x}_j - \frac{\epsilon^2}{2} \mathcal{E}(\mathbf{x}_j|\mathbf{y}^+) + \epsilon \mathbf{r}_j, \quad j = 1, \dots, J \quad (2)$$

146 where  $\mathbf{r}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  is the stochastic term and the step-size  $\epsilon$  is typically polynomially decayed [30].  
 147 The term  $\mathcal{E}(\mathbf{x}_j|\mathbf{y}^+)$  denotes the model energy conditioned on the target class label  $\mathbf{y}^+$  which we  
 148 specify as the negative logit corresponding to the target class label  $\mathbf{y}^*$ . To allow for faster sampling,  
 149 we follow the common practice of choosing the step-size  $\epsilon$  and the standard deviation of  $\mathbf{r}_j$  separately.  
 150 While  $\mathbf{x}_J$  is only guaranteed to distribute as  $p_\theta(\mathbf{x}|\mathbf{y}^*)$  if  $\epsilon \rightarrow 0$  and  $J \rightarrow \infty$ , the bias introduced for a  
 151 small finite  $\epsilon$  is negligible in practice [18, 6]. Appendix A provides additional implementation details  
 152 for any tasks related to energy-based modelling.

153 Generating multiple samples using SGLD thus yields an empirical distribution  $\hat{\mathbf{X}}_{\theta, \mathbf{y}^+}$  that approx-  
 154 imates what the model has learned about the input data. While in the context of Energy-Based  
 155 Modelling, this is usually done during training, we propose to repurpose this approach during  
 156 inference in order to evaluate and generate faithful model explanations.

### 157 4.2 Evaluating Plausibility and Faithfulness

158 The parallels between our definitions of plausibility and faithfulness imply that we can also use  
 159 similar evaluation metrics in both cases. Since existing work has focused heavily on plausibility,  
 160 it offers a useful starting point. In particular, Guidotti [7] have proposed an implausibility metric  
 161 that measures the distance of the counterfactual from its nearest neighbour in the target class. As  
 162 this distance is reduced, counterfactuals get more plausible under the assumption that the nearest  
 163 neighbour itself is plausible in the sense of Definition 2.1. In this work, we use the following adapted  
 164 implausibility metric that relaxes this assumption,

$$\text{impl} = \frac{1}{|\mathbf{x} \in \mathbf{X}_{\mathbf{y}^+}|} \sum_{\mathbf{x} \in \mathbf{X}_{\mathbf{y}^+}} \text{dist}(\mathbf{x}', \mathbf{x}) \quad (3)$$

165 where  $\mathbf{X}_{\mathbf{y}^+}$  is a subsample of the training data in the target class  $\mathbf{y}^+$ .

166 This gives rise to a very similar evaluation metric for unfaithfulness. We merely swap out the  
 167 subsample of individuals in the target class for a subset  $\hat{\mathbf{X}}_{\theta, \mathbf{y}^+}^{n_E}$  of the generated conditional samples:

$$\text{unfaith} = \frac{1}{|\mathbf{x} \in \hat{\mathbf{X}}_{\theta, \mathbf{y}^+}^{n_E}|} \sum_{\mathbf{x} \in \hat{\mathbf{X}}_{\theta, \mathbf{y}^+}^{n_E}} \text{dist}(\mathbf{x}', \mathbf{x}) \quad (4)$$

Specifically, we form this subset based on the  $n_E$  generated samples with the lowest energy.

## 5 Energy-Constrained Conformal Counterfactuals (ECCCo)

In this section, we describe our proposed framework for generating Energy-Constrained Conformal Counterfactuals (ECCCos). It is based on the premise that counterfactuals should be faithful, first and foremost. Plausibility, as a secondary concern, is then still attainable but only to the degree that the black-box model itself has learned plausible explanations for the underlying data.

We begin by stating our proposed objective function, which involves tailored loss and penalty functions that we will explain in the following. In particular, we extend Equation 1 as follows:

$$\begin{aligned} \mathbf{Z}' = \arg \min_{\mathbf{Z}' \in \mathcal{Z}^M} \{ & \text{yloss}(M_\theta(f(\mathbf{Z}')), \mathbf{y}^+) + \lambda_1 \text{dist}(f(\mathbf{Z}'), \mathbf{x}) \\ & + \lambda_2 \text{dist}(f(\mathbf{Z}'), \hat{\mathbf{x}}_\theta) + \lambda_3 \Omega(C_\theta(f(\mathbf{Z}'); \alpha)) \} \end{aligned} \quad (5)$$

The first penalty term involving  $\lambda_1$  induces proximity like in Wachter et al. [29]. Our default choice for  $\text{dist}(\cdot)$  is the L1 Norm due to its sparsity-inducing properties. The second penalty term involving  $\lambda_2$  constrains the energy of the generated counterfactual by penalising its distance from the lowest-energy conditional samples as defined in Equation 4. Intuitively, this component induces faithfulness which coincides with plausibility to the extent that the model  $M_\theta$  has learned the true posterior conditional distribution of inputs:  $p_\theta(\mathbf{X}_{\mathbf{y}^+}) \rightarrow p(\mathbf{X}_{\mathbf{y}^+})$ .

The third and final penalty term involving  $\lambda_3$  introduces a new but familiar concept: it ensures that the generated counterfactual is associated with low predictive uncertainty. As mentioned above, Schut et al. [23] have shown that plausible counterfactuals can be generated implicitly through predictive uncertainty minimization. Unfortunately, this relies on the assumption that the model itself can provide predictive uncertainty estimates, which may be too restrictive in practice.

To relax this assumption, we leverage recent advances in Conformal Prediction (CP), an approach to predictive uncertainty quantification that has recently gained popularity [3, 15]. Crucially for our intended application, CP is model-agnostic and can be applied during inference without placing any restrictions on model training. Intuitively, CP works under the premise of turning heuristic notions of uncertainty into rigorous uncertainty estimates by repeatedly sifting through the training data or a dedicated calibration dataset. Conformal classifiers produce prediction sets for individual inputs that include all output labels that can be reasonably attributed to the input. These sets tend to be larger for inputs that do not conform with the training data and are therefore characterized by high predictive uncertainty.

In order to generate counterfactuals that are associated with low predictive uncertainty, we use a smooth set size penalty introduced by Stutz et al. [25] in the context of conformal training:

$$\Omega(C_\theta(\mathbf{x}; \alpha)) = \max \left( 0, \sum_{\mathbf{y} \in \mathcal{Y}} C_{\theta, \mathbf{y}}(\mathbf{x}_i; \alpha) - \kappa \right) \quad (6)$$

Here,  $\kappa \in \{0, 1\}$  is a hyper-parameter and  $C_{\theta, \mathbf{y}}(\mathbf{x}_i; \alpha)$  can be interpreted as the probability of label  $\mathbf{y}$  being included in the prediction set.

In order to compute this penalty for any black-box model we merely need to perform a single calibration pass through a holdout set  $\mathcal{D}_{\text{cal}}$ . Arguably, data is typically abundant and in most applications, practitioners tend to hold out a test data set anyway. Consequently, CP removes the restriction on the family of predictive models, at the small cost of reserving a subset of the available data for calibration. This particular case of conformal prediction is referred to as Split Conformal Prediction (SCP) as it involves splitting the training data into a proper training dataset and a calibration dataset.

207 In addition to the smooth set size penalty, we have also experimented with the use of a tailored  
 208 function for  $y_{\text{loss}}(\cdot)$  that enforces that only the target label  $y^+$  is included in the prediction set Stutz  
 209 et al. [25]. Further details are described Appendix B.

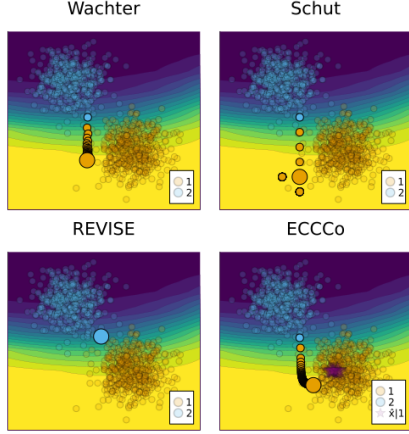


Figure 2: An example involving linearly separable synthetic data and illustrating how ECCCo compares to Wachter, Schut and REVISE. The factual class is 2 and the target class is 1. Contours indicate the predicted probability given by a Joint Energy Model that the counterfactual belongs to the target class.

Algorithm 1: Generating ECCCos (For more details, see Appendix C)

**Input:**

$\mathbf{x}, \mathbf{y}^+, M_\theta, f, \Lambda, \alpha, \mathcal{D}, T, \eta, n_B, N_B$   
 where  $M_\theta(\mathbf{x}) \neq \mathbf{y}^+$

**Output:**  $\mathbf{x}'$

- 1: Initialize  $\mathbf{z}' \leftarrow f^{-1}(\mathbf{x})$
- 2: Generate buffer  $\mathcal{B}$  of  $N_B$  conditional samples  $\hat{\mathbf{x}}_\theta | \mathbf{y}^+$  using SGLD (Equation 2)
- 3: Run SCP for  $M_\theta$  using  $\mathcal{D}$
- 4: Initialize  $t \leftarrow 0$
- 5: **while** not converged or  $t < T$  **do**
- 6:    $\hat{\mathbf{x}}_{\theta,t} \leftarrow \text{rand}(\mathcal{B}, n_B)$
- 7:    $\mathbf{z}' \leftarrow \mathbf{z}' - \eta \nabla_{\mathbf{z}'} \mathcal{L}(\mathbf{z}', \mathbf{y}^+, \hat{\mathbf{x}}_{\theta,t}; \Lambda, \alpha)$
- 8:    $t \leftarrow t + 1$
- 9: **end while**
- 10:  $\mathbf{x}' \leftarrow f(\mathbf{z}')$

211 The entire procedure for generating ECCCos is described in Algorithm 1. For the sake of simplicity  
 212 and without loss of generality, we limit our attention to generating a single counterfactual  $\mathbf{x}' = f(\mathbf{z}')$   
 213 where in contrast to Equation 5  $\mathbf{z}'$  denotes a 1-dimensional array containing a single counterfactual  
 214 state. That state is initialized by passing the factual  $\mathbf{x}$  through the encoder  $f^{-1}$  which in our case cor-  
 215 responds to a simple feature transformer, rather than the encoder part of VAE as in REVISE [8]. Next,  
 216 we generate a buffer of  $N_B$  conditional samples  $\hat{\mathbf{x}}_\theta | \mathbf{y}^+$  using SGLD (Equation 2) and conformalise  
 217 the model  $M_\theta$  through Split Conformal Prediction on training data  $\mathcal{D}$ .

218 Finally, we search counterfactuals through gradient descent. Let  $\mathcal{L}(\mathbf{z}', \mathbf{y}^+, \hat{\mathbf{x}}_{\theta,t}; \Lambda, \alpha)$  denote our loss  
 219 function defined in Equation 5. Then in each iteration, we first randomly draw  $n_B$  samples from  
 220 the buffer  $\mathcal{B}$  before updating the counterfactual state  $\mathbf{z}'$  by moving in the negative direction of that  
 221 loss function. The search terminates once the convergence criterium is met or the maximum number  
 222 of iterations  $T$  has been exhausted. Note that the choice of convergence criterium has important  
 223 implications on the final counterfactual (for more detail on this see Appendix C).

224 Figure 2 illustrates how ECCCos compare to counterfactuals generated using Wachter, Schut and  
 225 REVISE. The example involves synthetically generated linearly separable data that belong to one  
 226 of two classes. Contours indicate the predicted probabilities of a Joint Energy Model that has been  
 227 jointly trained to predict the output class and generate inputs Grathwohl et al. [6]. We have drawn a  
 228 random sample from the factual class 1 and used each generator to produce a counterfactual in the  
 229 target class 2. Both Wachter and Schut yield valid counterfactuals but fail to achieve plausibility in the  
 230 sense that the generated counterfactuals are far away from the densely populated region in the target  
 231 class. Conversely, ECCCo yields a faithful and plausible counterfactual in the neighbourhood of the  
 232 generated conditional samples. REVISE fails to yield a valid counterfactual because the underlying  
 233 surrogate has failed to learn the DGP.

## 234 6 Empirical Analysis

235 Our goal in this section is to shed light on the following research questions:

236 **Research Question 6.1 (Faithfulness).** *Are ECCCos more faithful to the black-box model than*  
 237 *counterfactuals produced by benchmark generators?*

238 **Research Question 6.2** (Plausibility). *How do ECCCo compare to state-of-the-art generators with*  
239 *respect to plausibility?*

240 We first briefly describe our experimental setup, before presenting our main results.

## 241 6.1 Key Evaluation Metrics

242 While we focus on these key evaluation metrics in the body of this paper, we also sporadically discuss  
243 outcomes with respect to other common measures used to evaluate the validity, proximity and sparsity  
244 of counterfactuals. Details can be found in Appendix E.

## 245 6.2 Experimental Setup

246 To assess and benchmark the performance of ECCCo against the state of the art, we generate multiple  
247 counterfactuals for different black-box models and datasets. In particular, we compare ECCCo to  
248 the following counterfactual generators that were introduced above: firstly; Schut [23], which works  
249 under the premise of minimizing predictive uncertainty; secondly, REVISE [8], which uses a VAE as  
250 its surrogate model; and, finally, Wachter [29], which serves as our baseline.

251 We use both synthetic and real-world datasets from different domains, all of which are publically  
252 available and commonly used to train and benchmark classification algorithms. The synthetic datasets  
253 include: a dataset containing two **Linearly Separable** Gaussian clusters ( $n = 1000$ ), as well as  
254 the well-known **Circles** ( $n = 1000$ ) and **Moons** ( $n = 2500$ ) data. As for real-world data, we  
255 follow Schut et al. [23] and use the **MNIST** [12] dataset containing images of handwritten digits such  
256 as the examples shown above. From the social sciences domain, we include Give Me Some Credit  
257 (**GMSC**) [9]: a tabular dataset that has been studied extensively in the literature on Algorithmic  
258 Recourse [19]. It consists of 11 numeric features that can be used to predict the binary outcome  
259 variable indicating whether or not retail borrowers experience financial distress.

260 As with the example in Section 5, we use simple neural networks (**MLP**) and Joint Energy Models  
261 (**JEM**). For the more complex real-world datasets we also use ensembling in each case. To account  
262 for stochasticity, we generate multiple counterfactuals for each possible target class, generator, model  
263 and dataset. Specifically, we randomly sample  $n^-$  times from the subset of individuals for which  
264 the given model predicts the non-target class  $y^-$  given the current target. We set  $n^- = 25$  for all  
265 of our synthetic datasets,  $n^- = 10$  for GMSC and  $n^- = 5$  for MNIST. Full details concerning our  
266 parameter choices, training procedures and model performance can be found in Appendix D.

## 267 6.3 Results

268 Table 1 shows the key results for the synthetic datasets separated by model (first columns) and  
269 generator (second column). The numerical columns show the average values of our key evaluation  
270 metrics computed across all counterfactuals. Standard deviations are shown in parentheses. In bold  
271 we have highlighted the best outcome for each model and metric. To provide some sense of the  
272 statistical significance of our findings, we have added asterisks to indicate that a given value is at  
273 least one (\*) or two (\*\*) standard deviations lower than the baseline (Wachter).

274 Starting with the high-level results for our Linearly Separable data, we find that ECCCo produces  
275 the most faithful counterfactuals for both black-box models. This is not surprising, since ECCCo  
276 directly enforces faithfulness through regularization. Crucially though, ECCCo also produces the  
277 most plausible counterfactuals for the Joint Energy Model, which was explicitly trained to learn  
278 plausible representations of the input data. This high-level pattern is broadly consistent across all  
279 datasets and supportive of our narrative, so it is worth highlighting: ECCCos consistently achieve  
280 high faithfulness, which—subject to the quality of the model itself—coincides with high plausibility.

281 Zooming in on the granular details for the Linearly Separable data, note that the list of generators  
282 in Table 1 includes ‘ECCCo (no CP)’ and ‘ECCCo (no EBM)’ in addition to ‘ECCCo’ and our  
283 benchmark generators. These have been added to gain some sense of the degree to which the two  
284 components underlying ECCCo—namely energy-based modelling (EBM) and conformal prediction  
285 (CP)—drive the results. Specifically, ‘ECCCo (no CP)’ involves no set size penalty ( $\lambda_3 = 0$  in  
286 Equation 5), while ‘ECCCo (no EBM)’ does not penalise the distance to samples generated through  
287 SGLD ( $\lambda_2 = 0$  in Equation 5). The corresponding results indicate that the positive results are

Table 1: Results for synthetic datasets. Standard deviations across samples are shown in parentheses. Best outcomes are highlighted in bold. Asterisks indicate that the given value is more than one (\*) or two (\*\*) standard deviations away from the baseline (Wachter).

Model	Generator	Linearly Separable		Moons		Circles	
		Unfaithfulness ↓	Implausibility ↓	Unfaithfulness ↓	Implausibility ↓	Unfaithfulness ↓	Implausibility ↓
JEM	ECCCo	<b>0.03 (0.06)**</b>	<b>0.20 (0.08)**</b>	<b>0.31 (0.30)*</b>	<b>1.20 (0.15)**</b>	0.52 (0.36)	1.22 (0.46)
	ECCCo (no CP)	0.03 (0.06)**	0.20 (0.08)**	0.37 (0.30)*	1.21 (0.17)**	0.54 (0.39)	1.21 (0.46)
	ECCCo (no EBM)	0.16 (0.11)	0.34 (0.19)	0.91 (0.32)	1.71 (0.25)	0.70 (0.33)	1.30 (0.37)
	REVISE	0.19 (0.03)	0.41 (0.01)**	0.78 (0.23)	1.57 (0.26)	<b>0.48 (0.16)*</b>	<b>0.95 (0.32)*</b>
	Schut	0.39 (0.07)	0.73 (0.17)	0.67 (0.27)	1.50 (0.22)*	0.54 (0.43)	1.28 (0.53)
	Wachter	0.18 (0.10)	0.44 (0.17)	0.80 (0.27)	1.78 (0.24)	0.68 (0.34)	1.33 (0.32)
MLP	ECCCo	<b>0.29 (0.05)**</b>	0.23 (0.06)**	0.80 (0.62)	1.69 (0.40)	0.65 (0.53)	1.17 (0.41)
	ECCCo (no CP)	0.29 (0.05)**	<b>0.23 (0.07)**</b>	<b>0.79 (0.62)</b>	1.68 (0.42)	<b>0.49 (0.35)</b>	1.19 (0.44)
	ECCCo (no EBM)	0.46 (0.05)	0.28 (0.04)**	1.34 (0.47)	1.68 (0.47)	0.84 (0.51)	1.23 (0.31)
	REVISE	0.56 (0.05)	0.41 (0.01)	1.45 (0.44)	<b>1.64 (0.31)</b>	0.58 (0.52)	<b>0.95 (0.32)</b>
	Schut	0.43 (0.06)*	0.47 (0.36)	1.45 (0.55)	1.73 (0.48)	0.58 (0.37)	1.23 (0.43)
	Wachter	0.51 (0.04)	0.40 (0.08)	1.32 (0.41)	1.69 (0.32)	0.83 (0.50)	1.24 (0.29)

dominated by the effect of quantifying and leveraging the model’s generative property (EBM) in our search for counterfactuals. Conformal Prediction alone only leads to marginally improved faithfulness and plausibility relative to the benchmark generators for our JEM. As a final observation for the Linearly Separable data we note that for the MLP, increased faithfulness comes at the cost of reduced plausibility. Specifically, this means that counterfactuals generated through ECCCo end up further away from individuals in the target class than those produced by our benchmark generators.

The findings for the Moons dataset are broadly in line with the findings so far: for the JEM, ECCCo yields significantly more faithful and plausible counterfactuals than all other generators. For the MLP, faithfulness is maintained but counterfactuals are not plausible. By comparison, REVISE yields fairly plausible counterfactuals in both cases, but it does so at the cost of faithfulness. We also observe that the best results for ECCCo are achieved when using both penalties. Once again though, the generative component (EBM) has a stronger impact on the positive results for the JEM.

For the Circles data, the most faithful counterfactuals are generated by ECCCo. While it appears that REVISE generates the most plausible counterfactuals in this case, we note that they are valid only half of the time (see Appendix E for a complete overview of all evaluation metrics). It turns out that in this case, the underlying VAE with default parameters has not adequately learned the data-generating process. Of course, it is possible to achieve better generative performance through hyperparameter tuning. But this example serves to illustrate that REVISE depends strictly on the quality of the surrogate model. Independent of the outcome for REVISE, however, the results do not seem to indicate that ECCCo significantly improves our plausibility metric for the Circles data.

Moving on to our real-world datasets, the results are shown in Table 2. Once again the findings indicate that the plausibility of ECCCos is positively correlated with the capacity of the black-box model to distinguish plausible from implausible inputs. The case is very clear for MNIST: ECCCos are consistently more faithful than the corresponding counterfactuals produced by any of the benchmark generators and their plausibility gradually improves through ensembling and joint-energy modelling. For the JEM Ensemble, ECCCo is essentially on par with REVISE and does significantly better than the baseline generator. We also note that ECCCo is the only generator that consistently achieves full validity for all models (Appendix E). Interestingly, ECCCo also yields lower-cost outcomes than the baseline generator for the JEMs.

For the tabular credit dataset (GMSC) we have struggled to get good generative and discriminative performance for our JEMs. Consequently, it is not surprising to find that ECCCo never achieves state-of-the-art plausibility, although it does improve outcomes compared to the baseline (Wachter). Concerning faithfulness, ECCCo once again consistently outperforms all other generators.

To conclude this section, we summarize our findings with reference to the opening questions. Concerning the feasibility of our proposed methodology (Research Question 6.1), our findings demonstrate that it is indeed possible to generate plausible counterfactuals without the need for surrogate models. A related important finding is that ECCCo never sacrifices faithfulness for plausibility: any plausible ECCCo also faithfully describes model behaviour. This mitigates the risk of generating plausible explanations for models that are, in fact, highly susceptible to implausible counterfactuals as well.



Table 2: Results for real-world datasets. Standard deviations across samples are shown in parentheses. Best outcomes are highlighted in bold. Asterisks indicate that the given value is more than one (\*) or two (\*\*) standard deviations away from the baseline (Wachter).

Model	Generator	MNIST		GMSC	
		Unfaithfulness ↓	Implausibility ↓	Unfaithfulness ↓	Implausibility ↓
<b>JEM</b>	ECCCo	<b>81.78 (17.49)**</b>	299.40 (29.48)**	199.40 (38.02)	17.26 (5.64)**
	REVISE	190.01 (28.89)**	<b>263.01 (46.46)**</b>	206.57 (41.88)	<b>4.86 (0.90)**</b>
	Schut	210.14 (27.35)**	286.50 (40.67)**	197.85 (37.95)	6.46 (2.11)**
	Wachter	280.70 (26.17)	499.25 (38.25)	<b>195.02 (32.35)</b>	68.48 (60.80)
<b>JEM Ensemble</b>	ECCCo	<b>72.46 (11.11)**</b>	276.48 (26.75)**	<b>182.04 (26.62)</b>	16.85 (4.49)**
	REVISE	173.81 (22.22)**	<b>248.50 (41.54)**</b>	206.02 (41.79)	<b>4.76 (0.63)**</b>
	Schut	202.69 (22.90)**	282.77 (39.72)**	204.53 (24.20)	6.53 (1.55)**
	Wachter	272.32 (23.03)	494.77 (37.74)	185.59 (33.79)	59.26 (49.56)
<b>MLP</b>	ECCCo	<b>155.25 (22.13)**</b>	519.58 (33.92)	<b>177.98 (39.03)</b>	19.09 (5.00)**
	REVISE	367.93 (14.90)**	<b>256.16 (44.23)**</b>	201.61 (30.74)	<b>5.33 (1.74)**</b>
	Schut	382.40 (16.67)*	286.14 (41.39)**	199.35 (32.06)	6.84 (1.96)**
	Wachter	406.24 (17.34)	488.30 (39.64)	195.51 (23.99)	81.62 (54.15)
<b>MLP Ensemble</b>	ECCCo	<b>144.74 (20.08)**</b>	484.56 (31.26)	<b>196.45 (34.80)*</b>	20.18 (5.20)**
	REVISE	340.33 (13.32)**	<b>251.30 (42.13)**</b>	202.67 (27.80)*	<b>4.82 (0.40)**</b>
	Schut	358.83 (13.17)*	283.12 (43.27)**	199.64 (42.29)*	6.35 (1.66)**
	Wachter	375.22 (18.91)	456.68 (47.21)	244.65 (44.55)	63.00 (53.77)

Our findings here indicate that ECCCo achieves this result primarily by leveraging the model’s generative property. We think that further work is needed, however, to definitively answer Research Question ??, on which we elaborate in the following section.

## 7 Limitations

Even though we have taken considerable measures to study our proposed methodology carefully, this work is limited in scope, which caveats our findings. In particular, we have found that the performance of ECCCo is sensitive to hyperparameter choices. In order to achieve faithfulness, we generally had to penalise the distance from generated samples slightly more than the distance from factual values. This choice is associated with relatively higher costs to individuals since the proposed recourses typically involve more substantial feature changes than for our benchmark generators.

Conversely, we have not found that penalising prediction set sizes disproportionately strongly had any discernable effect on our results. As discussed above, we also struggled to achieve good results by relying on conformal prediction alone. We want to caveat this finding by acknowledging that the role of CP in this context needs to be investigated more thoroughly through future work. Our suggested approach involving a smooth set size penalty may be insufficient in this context.

The fact that our findings are primarily driven by applying ideas from energy-based modelling presents a challenge in itself: while our approach is readily applicable to models with gradient access like deep neural networks, more work is needed to generalise our methodology to other popular machine learning models such as gradient-boosted trees. Relatedly, we have encountered common challenges associated with energy-based modelling during our experiments including sensitivity to scale, training instabilities and sensitivity to hyperparameters. We have also struggled to apply our proposed approach to low-dimensional tabular data.

## 8 Conclusion

This work leverages recent advances in energy-based modelling and conformal prediction in the context of Explainable Artificial Intelligence. We have proposed a new way to generate Counterfactual Explanations that are maximally faithful to the black-model they aim to explain. Our proposed counterfactual generator, ECCCo, produces plausible counterfactual if and only if the black-model itself has learned realistic representations of the data. This should enable researchers and practitioners

to use counterfactuals in order to discern trustworthy models from unreliable ones. While the scope of this work limits its generalizability, we believe that ECCCo offers a solid baseline for future work on faithful Counterfactual Explanations.

## References

- [1] Patrick Altmeyer. Conformal Prediction in Julia. URL <https://www.paltmeyer.com/blog/posts/conformal-prediction/>.
- [2] Patrick Altmeyer, Giovan Angela, Aleksander Buszydlík, Karol Dobiczek, Arie van Deursen, and Cynthia Liem. Endogenous Macrodynamics in Algorithmic Recourse. In *First IEEE Conference on Secure and Trustworthy Machine Learning*, 2023.
- [3] Anastasios N. Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. 2021.
- [4] André Artelt, Valerie Vaquet, Riza Velioglu, Fabian Hinder, Johannes Brinkrolf, Malte Schilling, and Barbara Hammer. Evaluating Robustness of Counterfactual Explanations. Technical report, arXiv. URL <http://arxiv.org/abs/2103.02354>. arXiv:2103.02354 [cs] type: article.
- [5] Ann-Kathrin Dombrowski, Jan E Gerken, and Pan Kessel. Diffeomorphic explanations with normalizing flows. In *ICML Workshop on Invertible Neural Networks, Normalizing Flows, and Explicit Likelihood Models*, 2021.
- [6] Will Grathwohl, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. March 2020. URL <https://openreview.net/forum?id=HkxzxONtDB>.
- [7] Riccardo Guidotti. Counterfactual explanations and how to find them: literature review and benchmarking. ISSN 1573-756X. doi: 10.1007/s10618-022-00831-6. URL <https://doi.org/10.1007/s10618-022-00831-6>.
- [8] Shalmali Joshi, Oluwasanmi Koyejo, Warut Vijitbenjaronk, Been Kim, and Joydeep Ghosh. Towards realistic individual recourse and actionable explanations in black-box decision making systems. 2019.
- [9] Kaggle. Give me some credit, Improve on the state of the art in credit scoring by predicting the probability that somebody will experience financial distress in the next two years., 2011. URL <https://www.kaggle.com/c/GiveMeSomeCredit>.
- [10] Amir-Hossein Karimi, Gilles Barthe, Bernhard Schölkopf, and Isabel Valera. A survey of algorithmic recourse: Definitions, formulations, solutions, and prospects. 2020.
- [11] Amir-Hossein Karimi, Bernhard Schölkopf, and Isabel Valera. Algorithmic recourse: From counterfactual explanations to interventions. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 353–362, 2021.
- [12] Yann LeCun. The MNIST database of handwritten digits. 1998.
- [13] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 4768–4777, 2017.
- [14] Divyat Mahajan, Chenhao Tan, and Amit Sharma. Preserving Causal Constraints in Counterfactual Explanations for Machine Learning Classifiers. Technical report, arXiv. URL <http://arxiv.org/abs/1912.03277>. arXiv:1912.03277 [cs, stat] type: article.
- [15] Valery Manokhin. Awesome conformal prediction.
- [16] Christoph Molnar. *Interpretable Machine Learning*. Lulu. com, 2020.
- [17] Ramaravind K Mothilal, Amit Sharma, and Chenhao Tan. Explaining machine learning classifiers through diverse counterfactual explanations. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 607–617, 2020.

- [18] Kevin P. Murphy. *Probabilistic machine learning: Advanced topics*. MIT Press.
- [19] Martin Pawelczyk, Sascha Bielowski, Johannes van den Heuvel, Tobias Richter, and Gjergji Kasneci. Carla: A python library to benchmark algorithmic recourse and counterfactual explanation algorithms. 2021.
- [20] Martin Pawelczyk, Teresa Datta, Johannes van-den Heuvel, Gjergji Kasneci, and Himabindu Lakkaraju. Probabilistically Robust Recourse: Navigating the Trade-offs between Costs and Robustness in Algorithmic Recourse. *arXiv preprint arXiv:2203.06768*, 2022.
- [21] Rafael Poyiadzi, Kacper Sokol, Raul Santos-Rodriguez, Tijl De Bie, and Peter Flach. FACE: Feasible and actionable counterfactual explanations. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 344–350, 2020.
- [22] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144, 2016.
- [23] Lisa Schut, Oscar Key, Rory Mc Grath, Luca Costabello, Bogdan Sacaleanu, Yarin Gal, et al. Generating Interpretable Counterfactual Explanations By Implicit Minimisation of Epistemic and Aleatoric Uncertainties. In *International Conference on Artificial Intelligence and Statistics*, pages 1756–1764. PMLR, 2021.
- [24] Thomas Spooner, Danial Dervovic, Jason Long, Jon Shepard, Jiahao Chen, and Daniele Magazzeni. Counterfactual Explanations for Arbitrary Regression Models. 2021.
- [25] David Stutz, Krishnamurthy Dj Dvijotham, Ali Taylan Cemgil, and Arnaud Doucet. Learning Optimal Conformal Classifiers. May 2022. URL <https://openreview.net/forum?id=t80-4LKfVx>.
- [26] Sohini Upadhyay, Shalmali Joshi, and Himabindu Lakkaraju. Towards Robust and Reliable Algorithmic Recourse. 2021.
- [27] Berk Ustun, Alexander Spangher, and Yang Liu. Actionable recourse in linear classification. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 10–19, 2019.
- [28] Sahil Verma, John Dickerson, and Keegan Hines. Counterfactual explanations for machine learning: A review. 2020.
- [29] Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harv. JL & Tech.*, 31:841, 2017.
- [30] M. Welling and Y. Teh. Bayesian Learning via Stochastic Gradient Langevin Dynamics. URL <https://www.semanticscholar.org/paper/Bayesian-Learning-via-Stochastic-Gradient-Langevin-Welling-Teh/aed631d6a84100b5e9a021ec1914095c66de415>.
- [31] Andrew Gordon Wilson. The case for Bayesian deep learning. 2020.

## Appendices

### A JEM

While  $\mathbf{x}_J$  is only guaranteed to distribute as  $p_\theta(\mathbf{x}|\mathbf{y}^+)$  if  $\epsilon \rightarrow 0$  and  $J \rightarrow \infty$ , the bias introduced for a small finite  $\epsilon$  is negligible in practice [18, 6]. While Grathwohl et al. [6] use Equation 2 during training, we are interested in applying the conditional sampling procedure in a post-hoc fashion to any standard discriminative model.

## 443 B Conformal Prediction

444 The fact that conformal classifiers produce set-valued predictions introduces a challenge: it is not  
 445 immediately obvious how to use such classifiers in the context of gradient-based counterfactual  
 446 search. Put differently, it is not clear how to use prediction sets in Equation 1. Fortunately, Stutz et al.  
 447 [25] have recently proposed a framework for Conformal Training that also hinges on differentiability.  
 448 Specifically, they show how Stochastic Gradient Descent can be used to train classifiers not only  
 449 for the discriminative task but also for additional objectives related to Conformal Prediction. One  
 450 such objective is *efficiency*: for a given target error rate  $\alpha$ , the efficiency of a conformal classifier  
 451 improves as its average prediction set size decreases. To this end, the authors introduce a smooth set  
 452 size penalty defined in Equation 6 in the body of this paper

453 Formally, it is defined as  $C_{\theta, \mathbf{y}}(\mathbf{x}_i; \alpha) := \sigma((s(\mathbf{x}_i, \mathbf{y}) - \alpha)T^{-1})$  for  $\mathbf{y} \in \mathcal{Y}$ , where  $\sigma$  is the sigmoid  
 454 function and  $T$  is a hyper-parameter used for temperature scaling [25].

455 Intuitively, CP works under the premise of turning heuristic notions of uncertainty into rigorous  
 456 uncertainty estimates by repeatedly sifting through the data. It can be used to generate prediction  
 457 intervals for regression models and prediction sets for classification models [1]. Since the literature  
 458 on CE and AR is typically concerned with classification problems, we focus on the latter. A particular  
 459 variant of CP called Split Conformal Prediction (SCP) is well-suited for our purposes, because it  
 460 imposes only minimal restrictions on model training.

461 Specifically, SCP involves splitting the data  $\mathcal{D}_n = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1, \dots, n}$  into a proper training set  $\mathcal{D}_{\text{train}}$   
 462 and a calibration set  $\mathcal{D}_{\text{cal}}$ . The former is used to train the classifier in any conventional fashion.  
 463 The latter is then used to compute so-called nonconformity scores:  $\mathcal{S} = \{s(\mathbf{x}_i, \mathbf{y}_i)\}_{i \in \mathcal{D}_{\text{cal}}}$  where  
 464  $s : (\mathcal{X}, \mathcal{Y}) \mapsto \mathbb{R}$  is referred to as *score function*. In the context of classification, a common choice for  
 465 the score function is just  $s_i = 1 - M_{\theta}(\mathbf{x}_i)[\mathbf{y}_i]$ , that is one minus the softmax output corresponding  
 466 to the observed label  $\mathbf{y}_i$  [3].

467 Finally, classification sets are formed as follows,

$$C_{\theta}(\mathbf{x}_i; \alpha) = \{\mathbf{y} : s(\mathbf{x}_i, \mathbf{y}) \leq \hat{q}\} \quad (7)$$

468 where  $\hat{q}$  denotes the  $(1 - \alpha)$ -quantile of  $\mathcal{S}$  and  $\alpha$  is a predetermined error rate. As the size of the  
 469 calibration set increases, the probability that the classification set  $C(\mathbf{x}_{\text{test}})$  for a newly arrived sample  
 470  $\mathbf{x}_{\text{test}}$  does not cover the true test label  $\mathbf{y}_{\text{test}}$  approaches  $\alpha$  [3].

471 Observe from Equation 7 that Conformal Prediction works on an instance-level basis, much like  
 472 Counterfactual Explanations are local. The prediction set for an individual instance  $\mathbf{x}_i$  depends only  
 473 on the characteristics of that sample and the specified error rate. Intuitively, the set is more likely  
 474 to include multiple labels for samples that are difficult to classify, so the set size is indicative of  
 475 predictive uncertainty. To see why this effect is exacerbated by small choices for  $\alpha$  consider the case  
 476 of  $\alpha = 0$ , which requires that the true label is covered by the prediction set with probability equal to  
 477 1.

## 478 C Conformal Prediction

## 479 D Experimental Setup

## 480 E Results

Table 3: All results for all datasets. Standard deviations across samples are shown in parentheses. Best outcomes are highlighted in bold. Asterisks indicate that the given value is more than one (\*) or two (\*\*) standard deviations away from the baseline (Wachter).

Model	Data	Generator	Cost ↓	Unfaithfulness ↓	Implausibility ↓	Redundancy ↑	Uncertainty ↓	Validity ↑
California Housing	JEM	ECCCo	39.14 (3.71)	<b>236.79 (51.16)</b>	39.78 (3.18)	0.00 (0.00)	2.00 (0.00)	1.00 (0.00)
		REVISE	4.39 (2.08)	284.51 (52.74)	<b>5.58 (0.81)**</b>	0.01 (0.03)	<b>1.85 (0.32)</b>	1.00 (0.00)
		Schut	4.17 (1.84)	263.55 (60.56)	8.00 (2.03)	<b>0.25 (0.24)*</b>	1.88 (0.31)	1.00 (0.00)
		Wachter	<b>2.03 (1.01)</b>	274.55 (51.17)	7.32 (1.80)	0.00 (0.00)	1.90 (0.31)	1.00 (0.00)
	JEM Ensemble	ECCCo	34.85 (4.67)	<b>249.44 (58.53)</b>	35.09 (5.56)	0.00 (0.00)	2.00 (0.00)	1.00 (0.00)
		REVISE	4.53 (1.97)	268.45 (66.87)	<b>5.44 (0.74)**</b>	0.00 (0.00)	1.95 (0.21)	1.00 (0.00)
		Schut	<b>0.98 (0.38)**</b>	279.38 (63.23)	7.64 (1.47)	<b>0.84 (0.06)**</b>	2.00 (0.00)	1.00 (0.00)
		Wachter	2.00 (0.59)	268.59 (68.66)	7.16 (1.46)	0.00 (0.00)	<b>1.90 (0.31)</b>	1.00 (0.00)
	MLP	ECCCo	37.47 (4.59)	<b>230.92 (48.86)</b>	37.53 (5.40)	0.00 (0.00)	1.00 (0.00)**	1.00 (0.00)
		REVISE	3.38 (2.06)	281.10 (53.01)	<b>5.34 (0.67)**</b>	0.00 (0.00)	1.10 (0.31)	1.00 (0.00)
		Schut	<b>0.88 (0.51)**</b>	285.12 (56.00)	6.48 (1.18)**	<b>0.72 (0.22)**</b>	<b>1.00 (0.00)**</b>	1.00 (0.00)
		Wachter	5.35 (10.88)	262.50 (56.87)	9.21 (10.41)	0.00 (0.00)	1.05 (0.22)	1.00 (0.00)
	MLP Ensemble	ECCCo	38.33 (4.99)	<b>212.47 (59.27)*</b>	38.17 (6.18)	0.00 (0.00)	1.00 (0.00)**	1.00 (0.00)
		REVISE	3.41 (1.79)	284.65 (49.52)	<b>5.64 (1.13)*</b>	0.00 (0.00)	1.05 (0.22)	1.00 (0.00)
		Schut	<b>0.84 (0.56)**</b>	269.19 (46.08)	7.30 (1.94)	<b>0.81 (0.11)**</b>	<b>1.00 (0.00)**</b>	1.00 (0.00)
		Wachter	2.00 (1.39)	278.09 (73.65)	7.32 (1.75)	0.00 (0.00)	1.07 (0.23)	1.00 (0.00)
Circles	JEM	ECCCo	1.34 (1.48)	<b>0.63 (1.58)</b>	1.44 (1.37)	0.00 (0.00)	0.98 (0.14)	0.98 (0.14)
		ECCCo (no CP)	1.33 (1.49)	0.64 (1.61)	1.45 (1.38)	0.00 (0.00)	0.98 (0.14)	0.98 (0.14)
		ECCCo (no EBM)	0.85 (1.49)	1.41 (1.51)	1.50 (1.38)	0.00 (0.00)	1.04 (0.28)	0.98 (0.14)
		REVISE	0.99 (0.35)	0.96 (0.32)*	<b>0.95 (0.32)*</b>	0.00 (0.00)	<b>0.50 (0.51)</b>	0.50 (0.51)
		Schut	1.00 (0.43)	0.99 (0.80)	1.28 (0.53)	<b>0.25 (0.25)</b>	1.11 (0.38)	<b>1.00 (0.00)**</b>
		Wachter	<b>0.74 (1.50)</b>	1.41 (1.50)	1.51 (1.35)	0.00 (0.00)	0.98 (0.14)	0.98 (0.14)
	MLP	ECCCo	1.39 (0.23)	<b>0.37 (0.65)**</b>	1.30 (0.68)	0.00 (0.00)	1.00 (0.00)**	<b>1.00 (0.00)</b>
		ECCCo (no CP)	1.33 (0.28)	0.50 (0.85)*	1.28 (0.66)	0.00 (0.00)	1.04 (0.20)*	<b>1.00 (0.00)</b>
		ECCCo (no EBM)	1.15 (0.69)	2.00 (1.46)	1.83 (1.00)	0.00 (0.00)	0.97 (0.10)**	<b>1.00 (0.00)</b>
		REVISE	0.98 (0.36)	1.16 (1.05)	<b>0.95 (0.32)*</b>	0.00 (0.00)	<b>0.50 (0.51)*</b>	0.50 (0.51)
		Schut	0.61 (0.11)	1.60 (1.15)	1.24 (0.44)	<b>0.34 (0.24)*</b>	1.00 (0.00)**	<b>1.00 (0.00)</b>
		Wachter	<b>0.53 (0.15)</b>	1.67 (1.05)	1.31 (0.43)	0.00 (0.00)	1.28 (0.46)	<b>1.00 (0.00)</b>
FashionMNIST	JEM	ECCCo	859.68 (91.05)	<b>40.65 (5.67)**</b>	605.67 (19.56)	0.00 (0.00)	3.00 (0.00)**	<b>1.00 (0.00)</b>
		REVISE	500.28 (86.07)	693.81 (118.47)*	<b>467.88 (132.24)</b>	0.00 (0.00)	3.20 (2.28)**	0.80 (0.45)
		Schut	<b>10.00 (0.00)**</b>	871.82 (64.75)	561.81 (94.76)	<b>0.99 (0.00)**</b>	<b>0.00 (0.00)**</b>	0.00 (0.00)
		Wachter	100.86 (13.85)	902.84 (88.79)	586.49 (97.17)	0.00 (0.00)	10.00 (0.00)	<b>1.00 (0.00)</b>
	JEM Ensemble	ECCCo	679.19 (66.95)	<b>59.61 (32.93)**</b>	500.50 (27.51)	0.00 (0.00)	4.00 (0.00)**	<b>1.00 (0.00)</b>
		REVISE	476.47 (147.09)	533.64 (102.81)*	<b>356.60 (79.57)*</b>	0.00 (0.00)	4.80 (1.30)**	<b>1.00 (0.00)</b>
		Schut	<b>10.00 (0.00)**</b>	688.61 (86.83)	445.55 (99.03)	<b>0.99 (0.00)**</b>	<b>0.00 (0.00)**</b>	0.00 (0.00)
		Wachter	92.50 (9.31)	714.63 (54.58)	470.54 (96.18)	0.00 (0.00)	10.00 (0.00)	<b>1.00 (0.00)</b>
	MLP	ECCCo	885.97 (29.70)	<b>65.36 (20.64)**</b>	791.07 (14.51)	0.00 (0.00)	2.00 (0.00)**	<b>1.00 (0.00)**</b>
		REVISE	323.10 (102.63)	856.08 (73.66)	<b>394.73 (252.67)</b>	0.00 (0.00)	1.00 (1.00)**	0.60 (0.55)
		Schut	<b>10.00 (0.00)**</b>	928.77 (42.27)	518.98 (143.30)	<b>0.99 (0.00)**</b>	<b>0.00 (0.00)**</b>	0.00 (0.00)
		Wachter	94.57 (10.26)	916.45 (50.09)	546.35 (145.24)	0.00 (0.00)	3.61 (4.01)	0.80 (0.45)
	MLP Ensemble	ECCCo	869.65 (67.92)	<b>47.37 (7.72)**</b>	751.83 (11.87)	0.00 (0.00)	1.00 (0.00)**	<b>1.00 (0.00)</b>
		REVISE	267.88 (69.67)	822.34 (57.55)	<b>307.50 (105.09)*</b>	0.00 (0.00)	3.00 (4.00)	0.80 (0.45)
		Schut	<b>10.00 (0.00)**</b>	891.57 (70.10)	449.79 (149.32)	<b>0.99 (0.00)**</b>	<b>0.00 (0.00)**</b>	0.00 (0.00)
		Wachter	91.50 (16.35)	874.21 (59.36)	476.59 (150.76)	0.00 (0.00)	4.60 (4.93)	<b>1.00 (0.00)</b>
GMSC	JEM	ECCCo	40.78 (8.79)**	<b>41.65 (17.24)**</b>	40.57 (8.74)**	0.00 (0.00)	1.50 (0.51)	<b>1.00 (0.00)**</b>
		REVISE	5.10 (6.48)**	74.89 (15.82)**	<b>6.01 (5.75)**</b>	0.00 (0.00)	1.81 (0.40)	<b>1.00 (0.00)**</b>
		Schut	<b>1.10 (0.39)**</b>	76.23 (15.54)**	6.02 (0.72)**	<b>0.77 (0.09)**</b>	1.55 (0.51)	<b>1.00 (0.00)**</b>
		Wachter	127.26 (75.11)	146.02 (64.48)	128.93 (74.00)	0.00 (0.00)	<b>1.00 (1.03)</b>	0.50 (0.51)
	JEM Ensemble	ECCCo	33.87 (8.25)**	<b>26.55 (12.94)**</b>	33.65 (8.33)**	0.00 (0.00)	2.00 (0.00)	<b>1.00 (0.00)**</b>
		REVISE	6.00 (4.92)**	52.47 (14.12)**	6.69 (3.37)**	0.00 (0.00)	1.80 (0.52)	0.95 (0.22)**
		Schut	<b>1.29 (0.92)**</b>	56.34 (15.00)**	<b>6.27 (1.06)**</b>	<b>0.74 (0.16)**</b>	1.62 (0.52)	<b>1.00 (0.00)**</b>
		Wachter	124.35 (95.08)	125.72 (70.80)	126.55 (93.75)	0.00 (0.00)	<b>1.00 (1.03)</b>	0.50 (0.51)
	MLP	ECCCo	38.91 (7.68)**	<b>46.90 (15.80)**</b>	37.78 (8.40)**	0.00 (0.00)	1.00 (0.00)	1.00 (0.00)
		REVISE	4.16 (2.35)**	81.08 (19.53)**	<b>4.60 (0.72)**</b>	0.00 (0.00)	1.23 (0.40)	1.00 (0.00)
		Schut	<b>0.72 (0.32)**</b>	90.67 (20.80)**	5.56 (0.81)**	<b>0.87 (0.06)**</b>	<b>1.00 (0.00)</b>	1.00 (0.00)
		Wachter	199.28 (14.78)	191.68 (30.86)	200.23 (15.05)	0.00 (0.00)	<b>1.00 (0.00)</b>	1.00 (0.00)
	MLP Ensemble	ECCCo	72.42 (145.72)	<b>74.65 (144.69)*</b>	71.87 (145.19)	0.00 (0.00)	1.00 (0.00)	1.00 (0.00)
		REVISE	4.75 (2.94)**	80.90 (14.59)**	<b>5.20 (1.52)**</b>	0.00 (0.00)	1.07 (0.12)	1.00 (0.00)
		Schut	<b>0.65 (0.24)**</b>	85.63 (19.15)**	6.00 (0.99)**	<b>0.88 (0.04)**</b>	<b>1.00 (0.00)**</b>	1.00 (0.00)
		Wachter	202.64 (14.71)	220.05 (17.41)	203.65 (14.77)	0.00 (0.00)	1.00 (0.00)	1.00 (0.00)
Linearly Separable	JEM	ECCCo	0.91 (0.14)	0.10 (0.06)**	0.19 (0.03)**	0.00 (0.00)	0.97 (0.03)**	<b>1.00 (0.00)</b>
		ECCCo (no CP)	0.91 (0.14)	<b>0.10 (0.07)**</b>	<b>0.19 (0.03)**</b>	0.00 (0.00)	0.98 (0.03)**	<b>1.00 (0.00)</b>
		ECCCo (no EBM)	0.90 (0.17)	0.37 (0.28)	0.38 (0.26)	0.00 (0.00)	1.23 (0.49)	<b>1.00 (0.00)</b>
		REVISE	<b>0.42 (0.14)*</b>	0.41 (0.02)**	0.41 (0.01)**	0.00 (0.00)	<b>0.81 (0.82)</b>	0.50 (0.51)
		Schut	1.14 (0.27)	0.66 (0.23)	0.66 (0.22)	<b>0.21 (0.25)</b>	1.74 (0.43)	<b>1.00 (0.00)</b>
		Wachter	0.61 (0.12)	0.44 (0.16)	0.44 (0.15)	0.00 (0.00)	1.50 (0.50)	<b>1.00 (0.00)</b>
	MLP	ECCCo	1.52 (0.16)	<b>0.03 (0.02)**</b>	0.69 (0.10)	0.00 (0.00)	1.00 (0.00)**	<b>1.00 (0.00)</b>
		ECCCo (no CP)	1.52 (0.16)	<b>0.03 (0.02)**</b>	0.68 (0.10)	0.00 (0.00)	1.00 (0.00)**	<b>1.00 (0.00)</b>
		ECCCo (no EBM)	2.66 (1.10)	1.25 (0.87)	1.84 (1.10)	0.00 (0.00)	1.00 (0.00)**	<b>1.00 (0.00)</b>
		REVISE	<b>0.44 (0.13)*</b>	1.10 (0.10)	<b>0.40 (0.01)**</b>	0.00 (0.00)	1.64 (0.78)	0.82 (0.39)
		Schut	0.76 (0.14)	0.81 (0.10)*	0.47 (0.24)	<b>0.26 (0.25)*</b>	<b>1.00 (0.00)**</b>	<b>1.00 (0.00)</b>
		Wachter	0.60 (0.14)	0.94 (0.11)	0.44 (0.15)	0.00 (0.00)	1.54 (0.50)	<b>1.00 (0.00)</b>
MNIST	JEM	ECCCo	269.99 (57.02)**	<b>116.09 (30.70)**</b>	281.33 (41.51)**	0.00 (0.00)	NA	<b>1.00 (0.00)**</b>
		REVISE	143.79 (43.43)**	348.74 (65.65)**	<b>246.69 (36.69)*</b>	0.00 (0.01)	NA	0.80 (0.40)
		Schut	<b>9.90 (0.55)**</b>	355.58 (64.84)**	270.06 (40.41)**	<b>0.99 (0.00)**</b>	NA	0.15 (0.36)
		Wachter	453.86 (16.96)	694.08 (50.86)	630.99 (33.01)	0.00 (0.00)	NA	0.90 (0.30)
	JEM Ensemble	ECCCo	260.94 (52.14)**	<b>89.89 (27.26)**</b>	240.59 (37.41)**	0.00 (0.00)	NA	<b>1.00 (0.00)**</b>
		REVISE	138.82 (33.99)**	292.52 (53.13)**	<b>240.50 (35.73)*</b>	0.00 (0.01)	NA	0.81 (0.39)
		Schut	<b>9.97 (0.28)**</b>	319.45 (59.02)**	266.80 (40.46)**	<b>0.99 (0.00)**</b>	NA	0.05 (0.22)
		Wachter	365.46 (35.14)	582.52 (58.46)	543.90 (44.24)	0.00 (0.00)	NA	0.96 (0.20)
	MLP	ECCCo	658.48 (65.03)	<b>212.45 (36.70)**</b>	649.63 (58.80)	0.00 (0.00)	NA	<b>1.00 (0.00)</b>
		REVISE	150.41 (51.81)**	839.79 (77.14)*	<b>244.33 (38.69)**</b>	0.00 (0.00)	NA	0.95 (0.22)
		Schut	<b>9.95 (0.41)**</b>	842.80 (82.01)*	264.94 (42.18)**	<b>0.99 (0.00)**</b>	NA	0.06 (0.25)
		Wachter	400.08 (34.33)	982.32 (61.81)	561.23 (45.08)	0.00 (0.00)	NA	<b>1.00 (0.00)</b>
	MLP Ensemble	ECCCo	616.12 (102.01)	<b>162.21 (36.21)**</b>	587.65 (95.01)	0.00 (0.00)	NA	<b>1.00 (0.00)**</b>
		REVISE	149.48 (47.90)**	741.30 (125.98)*	<b>242.76 (41.16)**</b>	0.00 (0.01)	NA	0.92 (0.27)
		Schut	<b>9.98 (0.23)**</b>	754.35 (132.26)	266.94 (42.55)**	<b>0.99 (0.00)**</b>	NA	0.03 (0.18)
		Wachter	374.37 (41.37)	871.09 (92.36)	536.24 (48.73)	0.00 (0.00)	NA	1.00 (0.05)
Moons	JEM	ECCCo	1.87 (0.79)	<b>0.57 (0.58)**</b>	<b>1.29 (0.21)*</b>	0.00 (0.00)	0.99 (0.18)**	1.00 (0.00)
		ECCCo (no CP)	1.83 (0.80)	0.63 (0.64)*	1.30 (0.21)*	0.00 (0.00)	1.13 (0.35)	1.00 (0.00)
		ECCCo (no EBM)	1.30 (1.72)	1.73 (1.34)	1.73 (1.42)	0.00 (0.00)	<b>0.94 (0.27)*</b>	1.00 (0.00)
		REVISE	1.07 (0.26)	1.59 (0.55)	1.55 (0.20)	0.00 (0.00)	1.30 (0.40)	1.00 (0.00)
		Schut	1.36 (0.35)	1.55 (0.61)	1.42 (0.16)*	<b>0.03 (0.12)</b>	1.11 (0.30)*	1.00 (0.00)
		Wachter	<b>0.89 (0.21)</b>	1.77 (0.48)	1.67 (0.15)	0.00 (0.00)	1.45 (0.47)	1.00 (0.00)
	MLP	ECCCo	2.53 (1.24)	1.68 (1.74)	2.02 (0.86)	0.00 (0.00)	1.11 (0.31)	<b>1.00 (0.00)</b>
		ECCCo (no CP)	2.45 (1.36)	<b>1.34 (1.66)</b>	2.11 (0.88)	0.00 (0.00)	1.24 (0.41)	<b>1.00 (0.00)</b>
		ECCCo (no EBM)	2.53 (2.03)	2.98 (1.89)	2.29 (1.75)	0.00 (0.00)	0.99 (0.07)**	<b>1.00 (0.00)</b>
		REVISE	0.98 (0.33)*	2.46 (1.05)	<b>1.54 (0.27)*</b>	0.00 (0.00)	1.40 (0.49)	<b>1.00 (0.00)</b>
		Schut	<b>0.75 (0.23)**</b>	2.71 (1.15)	1.62 (0.42)	<b>0.31 (0.27)*</b>	<b>0.94 (0.24)*</b>	0.94 (0.24)
		Wachter	1.49 (1.76)	2.95 (1.42)	1.84 (1.33)	0.00 (0.00)	1.33 (0.48)	<b>1.00 (0.00)</b>