UNCERTAINTY-AWARE COUNTERFACTUAL EXPLANA TIONS USING BAYESIAN NEURAL NETS

Anonymous authors

Paper under double-blind review

ABSTRACT

A counterfactual explanation describes the smallest input change required to alter the prediction of an AI model towards a desired outcome. When using neural networks, counterfactuals are obtained using variants of projected gradient descent. Such counterfactuals have been shown to be brittle and implausible, potentially jeopardising the explanatory aspects of counterfactuals. Numerous approaches for obtaining better counterfactuals have been put forward. Even though these solutions address some of the shortcomings, they often fall short of providing an all-around solution for robust and plausible counterfactuals. We hypothesise this is due to the deterministic nature and limitations of neural networks, which fail to capture the uncertainty of the training data. Bayesian Neural Networks (BNNs) are a well-known class of probabilistic models that could be used to overcome these issues; unfortunately, there is currently no framework for developing counterfactuals for them. In this paper, we fill this gap by proposing a formal framework to define counterfactuals for BNNs and develop algorithmic solutions for computing them. We evaluate our framework on a set of commonly used benchmarks and observe that BNNs produce counterfactuals that are more robust, plausible, and less costly than deterministic baselines.¹

004

010 011

012

013

014

015

016

017

018

019

021

023

025

1 INTRODUCTION

As Artificial Intelligence (AI) and Machine Learning (ML) increasingly influence critical decisions 031 in areas such as finance (Cao, 2022) and healthcare (Shaheen, 2021), the need for reliable explana-032 tions of the decisions made by AI is becoming increasingly important. *Counterfactual Explanations* 033 have emerged as a powerful tool for interpreting the decision-making processes of ML models, of-034 fering actionable insights into how the input to an ML model needs to be changed for the model to produce a different, and often desirable, outcome (CFXs) (see (Guidotti, 2024) for a recent survey). This is particularly useful in the context of algorithmic recourse (Karimi et al., 2023), where 037 CFXs are used to generate recourse recommendations for users that have been negatively affected 038 by the decisions of an ML model. CFXs are particularly suited for this task given their intelligibility (Byrne, 2019), appeal to users (Barocas et al., 2020), information capacity (Kenny & Keane, 040 2021) and alignment with human reasoning (Miller, 2019).

041 To see what makes CFXs useful, consider a (fictional) loan application where a customer applies 042 for a loan with a bank which uses an ML model to process the application and predict whether the 043 customer will be able to repain the loan or not. For illustration, assume the application is modelled 044 by an input x with features 32 years of age, \$10,000 loan amount and \$25,000 salary. Assume that the application is initially rejected, based on the prediction made by the AI that the customer 046 will not be able to repay the loan back. A possible CFX for this rejection could be an altered input x', where a salary of \$30,000 (with the other features unchanged) would result in the loan being 047 accepted, thus pointing the user to what they would need to change in their application for the loan 048 to be accepted. 049

Despite their potential, current approaches to generating CFXs often fall short in terms of satisfying two key properties: *plausibility* (Laugel et al., 2019) and *robustness* (Jiang et al., 2024). The former requires that CFXs adhere as much as possible to the data manifold, to avoid suggesting

⁰⁵³

¹The code for reproducing the results is provided in the supplementary materials.

unrealistic input changes. The latter instead requires that similar CFXs be generated for similar inputs (Artelt et al., 2021), to ensure fairness in applications such as algorithmic recourse (Slack et al., 2021). These properties are more than just metrics characterising the utility of CFXs; they are core desiderata without which CFXs may erode trust in the model they are trying to explain, rather than engendering it.

We posit that these limitations stem from the deterministic nature of traditional neural networks, which fail to capture the inherent uncertainty in the data. To address this fundamental issue, we propose a novel framework for generating counterfactual explanations using Bayesian Neural Networks (BNNs). Our approach leverages the uncertainty quantification capabilities of BNNs to produce CFXs that are more plausible and robust than those generated by deterministic models. Specifically, our contributions are as follows:

- Defining counterfactual explanations for BNNs. We first introduce a formal definition of counterfactual explanations in the context of Bayesian Neural Networks. This definition extends the concept of CFXs to this class of probabilistic models, accounting for the distribution over model parameters, which in turn enables a more nuanced understanding of the decision boundary.
- Demonstrating enhanced plausibility. Through extensive experiments on both vision and tabular datasets, we show that CFXs generated using our proposed BNN-based approach consistently lie closer to the data manifold than those produced by deterministic MLPs or ensembles. In this way, our explanations are more realistic and usable in practice.
- Demonstrating improved robustness. We demonstrate that our BNN-based CFXs exhibit superior robustness, meaning that similar inputs map to similar counterfactual explanations. This property is crucial for building trust in the explanations provided, as it ensures consistency across meaningful perturbations on the data manifold.

To validate our approach, we conduct a comprehensive empirical evaluation across multiple datasets, including *MNIST* (LeCun et al., 1998) for vision tasks and several tabular datasets, including *German Credit Risk* (Dua & Graff, 2017), *Diabetes* (Smith et al., 1988), *News Categorisation* (Fernandes et al., 2015), and *Spam Base* (Hopkins et al., 1999), covering various domains such as finance and healthcare. Our results consistently show that BNN-based CFXs outperform their deterministic counterparts across various metrics, including plausibility and robustness. Notably, this result holds when comparing against previously-proposed uncertainty-aware models.

The remainder of this paper is organised as follows. We provide the essential background for this paper in Section 2. We then present our key contribution in Section 3, where we formally define counterfactual explanations for BNNs and show how they can be computed. We validate our proposal in Section 4 and present an extensive experimental evaluation using common datasets from the literature on CFXs. Finally, we discuss related work in Section 5 and discuss the broader implications of our work for the field of explainable AI and the practical deployment of machine learning models.

092

065

2 BACKGROUND

093 094

103 104

Counterfactual explanations. Counterfactual explanations (CFXs) provide a way to interpret the decisions of ML models by showing how changes to the input of a model would lead to different outcomes. Mainstream approaches to compute CFXs characterise these explanations in terms of the solutions of an optimisation problem (Wachter et al., 2017; Mohammadi et al., 2021), which we present next for a binary classification setting without loss of generality. Let \mathcal{M} be a machine learning model mapping an input $x \in \mathcal{X}$ to label $\ell \in \{0, 1\}$. For ease of exposition, we refer to $\mathcal{M}(x) = 0$ as the *negative outcome* and to $\mathcal{M}(x) = 1$ as the *positive outcome*. Assuming \mathcal{M} initially produces a negative outcome for an input x, a CFX x_c for this decision can be obtained as:

$$\underset{\boldsymbol{x}_{c} \in \mathcal{X}}{\arg\min} d(\boldsymbol{x}, \boldsymbol{x}_{c}) \text{ s.t. } \mathcal{M}(\boldsymbol{x}_{c}) = 1,$$
(1)

where $d: \mathcal{X} \times \mathcal{X} \to \mathbb{R}^+$ is a distance metric defined over the input space from which x and x_c are drawn. Since computing an exact solution for the problem presented in Equation (1) may be viable only for certain types of machine learning models, the following relaxation is typically considered for more general classes of differentiable models: $\underset{\boldsymbol{x}_{c} \in \mathcal{X}}{\arg\min} \mathcal{L}(\mathcal{M}(\boldsymbol{x}_{c}), 1) + \lambda \cdot d(\boldsymbol{x}, \boldsymbol{x}_{c})$ (2)

110 111

108

where \mathcal{L} is a differentiable loss function that guides the search towards an input \boldsymbol{x}_c for which \mathcal{M} yield a positive outcome with high confidence, and λ is a parameter controlling the trade-off between the first term and a distance loss d defined as in Equation (1).

116 Several metrics have been proposed to assess the quality of CFXs (Karimi et al., 2023). For example, validity captures the basic requirement that a CFX should change the output of a model, 117 turning a negative outcome into a positive one. Validity is typically considered in tandem with prox-118 mity (Wachter et al., 2017), which gives a higher preference to CFXs that are closer to the original 119 input. Additionally, CFXs are typically required to be actionable (Ustun et al., 2019) and only alter 120 features that can be realistically modified by the user (e.g. users cannot modify their age but they 121 can act on credit score). Sparsity (Wachter et al., 2017) is also deemed important in many cases, 122 whereby CFXs requiring changes on fewer features are to be preferred to avoid overloading users 123 with too much information. Another important requirement is *plausibility* (Dhurandhar et al., 2018; 124 Altmeyer et al., 2024), which requires that counterfactual explanations adhere as much as possible to 125 the data manifold, to avoid causing unrealistic changes to input features. Finally, robustness (Artelt 126 et al., 2021; Slack et al., 2021; Leofante & Potyka, 2024), advocates for the generation of similar 127 CFXs for similar inputs, to ensure CFXs are not perceived as potentially malicious or discriminatory. Validity, plausibility, and robustness will be the focus of the experimental analysis presented 128 in this paper. 129

130

131 **Bayesian Neural Networks (BNNs)** A BNN \mathcal{B} is a probabilistic model based on a Neural Net-132 work (NN) architecture, where for each layer l = 1, ..., L, the parameters w are sampled from a 133 posterior distribution P(w).

134 135 136 137 138 Definition 1 (BNN). A BNN \mathcal{B} is a pair $(f_{w \sim P(w)}(x), P(w))$, where $f_w(x)$ defines the architecture and operations of the network and P(w) is the posterior distribution over the parameters of the BNN. Thus, the output of a BNN, denoted by $\mathcal{B}(x)$ for simplicity, is the expected value of the forward pass over $f_{w \sim P(w)}(x)$ with respect to the distribution of weights. Formally,

139

142 143

 $\mathcal{B}(\boldsymbol{x}) = \mathbb{E}_{\boldsymbol{w} \sim P(\boldsymbol{w})}[f_{\boldsymbol{w}}(\boldsymbol{x})] = \int_{\boldsymbol{w}} f_{\boldsymbol{w}}(\boldsymbol{x}) P(\boldsymbol{w}) \, d\boldsymbol{w}.$

In practice, computing the output of a BNN as defined in Equation (3) is intractable. Thus, we approximate Equation (3) using Monte Carlo sampling of the posterior distribution P(w). The approximate BNN output is given by

When considering classification models we denote the *l*-th output unit of a BNN as $\mathcal{B}(\boldsymbol{x})_l$.

147

148

149 150

$$\tilde{\mathcal{B}}(\boldsymbol{x}) = \frac{1}{N} \sum_{i=1}^{N} f_{\boldsymbol{w}}(\boldsymbol{x}), \tag{4}$$

(3)

where $w_1, \ldots, w_N \sim P(w)$ are iid samples from the posterior.

153 While deterministic NNs are trained via maximum likelihood estimation (MLE), training a BNN 154 corresponds to performing Bayesian inference on P(w). For this, we begin with a prior over the 155 BNN parameters, $\Pi(w)$, and update this prior using observations \mathcal{D} , $P(w) = \Pi(w|\mathcal{D})$. Various 156 BNN inference approaches exist. Bayes-by-backprop updates the parameters of the prior iteratively 157 in a process that mirrors MLE (Blundell et al., 2015). Markov chain Monte Carlo (MCMC) methods 158 directly sample from the posterior using accept-reject-style algorithms such as Metropolis-Hastings 159 (Borkar, 1953), or Hamiltonian Monte Carlo (HMC) (Duane et al., 1987). In this paper, we focus on the definition and procedure for obtaining counterfactuals on BNNs trained using HMC, as it is 160 the most precise inference algorithm. More discussion on how we train the BNNs is available in 161 Section 4.

162 3 COUNTERFACTUAL EXPLANATIONS FOR BNNs

Counterfactuals do not have a commonly accepted definition in probabilistic models and, to the best of our knowledge, they have never been formally defined for BNNs. Here, we propose a framework for defining and computing counterfactuals specific to BNNs. In contrast to deterministic networks, the parameters of a BNN are modelled as distributions rather than fixed values, complicating the definition of counterfactuals. Specifically, while counterfactuals for deterministic networks usually require the computation of model gradients, the gradient of a BNN's output with respect to its input is distributional, and its expected value is difficult to compute exactly. Moreover, to compute the true gradient of a BNN with respect to its input, we differentiate Equation (3) with respect to *x*:

$$\partial_{\boldsymbol{x}} \mathcal{B}(\boldsymbol{x}) = \partial_{\boldsymbol{x}} \left(\int_{\boldsymbol{w}} f_{\boldsymbol{w}}(\boldsymbol{x}) P(\boldsymbol{w}) \, d\boldsymbol{w} \right) = \int_{\boldsymbol{w}} \partial_{\boldsymbol{x}} f_{\boldsymbol{w}}(\boldsymbol{x}) P(\boldsymbol{w}) \, d\boldsymbol{w}.$$
(5)

However, both Equation (3) and its gradient in Equation (5) are intractable to compute directly.
Similarly to Equation (4) we can approximate the expected gradient of a BNN through Monte Carlo sampling,

177 178

188 189

172 173

$$\partial_{\boldsymbol{x}}\tilde{\mathcal{B}}(\boldsymbol{x}) = \frac{1}{N} \sum_{i=1}^{N} \partial_{\boldsymbol{x}} f_{\boldsymbol{w}_i}(\boldsymbol{x}).$$
(6)

179 i=1180 In this setting, the gradient $\partial_{\boldsymbol{x}} f_{\boldsymbol{w}_i}(\boldsymbol{x})$ can be computed in the same way as a standard, deterministic MLP.

Having established the mathematical framework, we can now formally define probabilistic counterfactuals for Bayesian Neural Networks (BNNs). This definition not only computes counterfactuals with minimal distance from the original input but also incorporates the model's inherent uncertainty. **Definition 2** (Probabilistic Counterfactual). Given a BNN \mathcal{B} , an input x, with observed negative outcome, $\mathcal{B}(x) = 0$, a probabilistic counterfactual is an input x_c such that the output achieves the desired outcome i.e., $\mathcal{B}(x_c) = \mathbb{E}_w[f_{w \sim P(w)}(x)] = 1$. Formally,

$$\boldsymbol{x}_{c} = \operatorname*{arg\,min}_{\boldsymbol{x}_{c}} d(\boldsymbol{x}, \boldsymbol{x}_{c}) \text{ s.t. } \operatorname*{arg\,max}_{l} \mathcal{B}(\boldsymbol{x}_{c})_{l} = 1. \tag{7}$$

Equation (7) describes the output constraint for the classification setting on which we focus. For regression tasks we can replace the constraint with a bound on the output units of the network, $l_i \leq \mathcal{B}(\boldsymbol{x}_c)_i \leq u_i, i = 1, ..., n$ where \mathcal{B} has n output units. Moreover, where we have focused on the binary classification setting, this definition can be extended to the multi-class case by replacing the negative outcome with the original class y, and the positive outcome with some target class, t. We continue with this more generalised notation for the classification setting.

To compute counterfactuals, we parallel the optimisation formulation given in Equation (2). Concretely, and focusing on the classification case, we use a linear loss, \mathcal{L}_{lin} , for a specified target class, *t*, and write our objective function as

$$\underset{\boldsymbol{x}_{c}}{\arg\min} \mathcal{L}_{\text{lin}}(\tilde{\mathcal{B}}(\boldsymbol{x}_{c}), t) + \lambda \cdot d(\boldsymbol{x}, \boldsymbol{x}_{c}), \tag{8}$$

201 where $\mathcal{L}_{\text{lin}}(\hat{\mathcal{B}}(\boldsymbol{x}_c), t) = \hat{\mathcal{B}}(\boldsymbol{x}_c)_y - \hat{\mathcal{B}}(\boldsymbol{x}_c)_t$ with $\hat{\mathcal{B}}(\boldsymbol{x}')_y$ being the average value of the output unit 202 corresponding to the observed class y and $\hat{\mathcal{B}}(\boldsymbol{x}_c)_t$ that for the target class. We have selected linear 203 loss due to its computational efficiency and its frequent application in the robustness literature, where 204 it is known for prompting rapid changes in model outputs (Carlini & Wagner, 2017). However, 205 alternative loss functions, such as cross-entropy, may also be employed, as discussed in Section 4. 206 Echoing Equation (2), the first term in the objective of Equation (8) accounts for the validity of 207 candidate CFXs, while the second term in Equation (8) promotes CFXs that are closer to the original input x. 208

Based on this objective function we outline our algorithm for computing probabilistic counterfactuals in Algorithm 1. The procedure begins by initialising the counterfactual with the original input vector and proceeding to the main loop. Within the main loop, we alternate between computing approximate gradients using Equation (6), and stepping the counterfactual according to the gradient. In our experiments, we set L and U as the upper and lower bounds on the input, though it is possible to limit this to an l_p ball if there is a pre-defined budget for the counterfactuals. We note that, as for the deterministic setting, this algorithm does not guarantee a valid counterfactual and that the choice of λ, ε, and N will dictate this as tunable parameters.

216 **Algorithm 1:** Generating CFX for BNNs 217 **Input** : BNN \mathcal{B} , input sample x, target class y, stepsize ϵ , distance weight λ , number of 218 iterations N, lower and upper bounds on input L and U. 219 **Output:** Counterfactual x_c . 220 ▷ Initialise the counterfactual 1 $x_c \leftarrow x$ 221 ² for $n \leftarrow 1, \ldots, N$ do 222 $\delta \leftarrow \partial_{\boldsymbol{x}_c} [\mathcal{L}_{\text{lin},t}(\mathcal{B}(\boldsymbol{x}_c)) + \lambda(\|\boldsymbol{x} - \boldsymbol{x}_c\|_p)]$ ▷ Compute loss' gradient w.r.t. to the input 3 223 $oldsymbol{x}_{c} \leftarrow oldsymbol{x}_{c} + \epsilon \cdot \delta$ > Update counterfactual using the gradient 4 224 ▷ Clip the adversarial example to ensure it is within bounds $\boldsymbol{x}_c \leftarrow \operatorname{clip}(\boldsymbol{x}_c, L, U)$ 5 225 6 end 226 7 return x_c 227

4 EVALUATION

230 231

228 229

In this section, we evaluate the properties of counterfactuals produced on BNNs. We focus on three 232 main properties, i.e. validity, robustness, and plausibility, and show that CFXs obtained for BNNs 233 outperform those produced for traditional Multi-Layer Perceptrons (MLPs). We also test other meth-234 ods for uncertainty quantification, namely ensemble methods, and show that BNNs produce better 235 CFXs in most instances. We conducted experiments on various popular datasets to cover different 236 data types and classification tasks. They include one vision dataset, MNIST (LeCun et al., 1998), 237 and four tabular datasets: credit (Dua & Graff, 2017), diabetes (Smith et al., 1988), news (Fernan-238 des et al., 2015), and spambase (Hopkins et al., 1999). For each dataset, we trained the following 239 models: a single standard deterministic multi-layer perceptron (MLP), an ensemble of 50 randomly 240 initialised MLPs (Ensemble), and a Bayesian Neural Network (BNN). We keep the architectures of 241 the three models consistent with 2 hidden layers, 150 nodes each, to aid comparison. In training 242 our BNNs we use an adaptive variant of the HMC algorithm called NUTS provided as part of the numpyro package. 243

244 We have chosen these benchmarks as they represent BNN's closest deterministic counterparts. An 245 MLP is the least complex form of deep neural network and is also used exhaustively in the CFX 246 literature as a case study making it a key benchmark. We also compare with ensembles of MLPs as 247 these have previously been studied in the context of uncertainty-aware CFX by Schut et al. (2021). 248 MLP ensembles also have functional similarities to BNNs; as we use a sampling-based Bayesian inference algorithm, our BNNs can be considered as a finite ensemble of samples in the same way 249 as an ensemble of MLPs. In these comparisons, we hypothesise that the BNN will greater be able 250 to capture the underlying data manifold, leading to more robust and plausible counterfactuals in 251 practice. 252

253 Counterfactual explanations are generated using gradient-based optimization methods tailored to each model type. For MLPs and ensembles, we used standard projected gradient descent to find min-254 imal input perturbations that change the model's prediction according to Equation (2). For BNNs, 255 we utilized the probabilistic counterfactual framework defined in Section 3, leveraging Monte Carlo 256 sampling to approximate gradients and defined in Equation (6). For every dataset, we compute 257 counterfactuals for 50 random samples from the test set. All experiments are performed on an RTX 258 3080 GPU and AMD Ryzen Threadripper 3960x 24-core CPU with 256GB of RAM running Ubuntu 259 22.04. 260

In the rest of this section, we first look at a visual example from the MNIST dataset in Section 4.1, before defining our metrics and discussing numeric results in Sections 4.2 and 4.3. Finally, we compare against previous work on uncertainty-aware CFX in Section 4.4.

264 265

266

4.1 VISUAL INTERPRETATION

We begin with an example from the MNIST image classification dataset where we randomly select a target class for each image. Figure 1 shows snapshots of computing a counterfactual explanation on the three model types. Each row is run with the same hyperparameters, original image, and target class.



Figure 1: Generation of a counterfactual for an MLP, an Ensemble, and a BNN. a) shows the original image and then a counterfactual with target '6' is progressively generated, with b)-d) showing snapshots of the image as the number of iterations increases.

292 293

289

291

270

We observe that as the number of iterations increases from left to right, the images increasingly 295 resemble the number 6. However, the MLP suffers some erroneous fragments on either side of the 296 number in b) and c). By e) the MLP's CFX has begun to degrade, particularly on the right-hand side, 297 to a point where it is nearly unrecognisable as a 6 to the human eye. The ensemble row also shows 298 some noise around the number, but the degradation of the 'key' pixels in the number are less affected and the final image is more recognisable as a 6. The BNN also suffers from noise above and to the 299 sides of the number at the b) and c) stages; however, the noise is less pronounced in these stages 300 than we observe for the Ensemble and MLP. At stage e) the BNN's CFX is much more complete 301 than for the Ensemble or MLP, even though the noise has become quite pronounced, the key pixels 302 are largely preserved and the number six is evident. 303

We attribute the improved preservation of 'key' pixels in the BNN to the better representation of the data distribution captured by this model. Similarly, we propose that the model averaging in the Ensemble prevents the major deterioration of the key pixels that we observe for the MLP. Specifically, for pixels to significantly change in the Ensemble model, they must have a significant impact in the output across all models of the ensemble. This helps to mitigate any local minima we might observe in any single MLP. We emphasise that these counterfactuals are produced with no CFX-specific regularisation scheme in either training or the CFX algorithm, with the intention of examining the explanations produced by these models in an unmodified state.

312 313 4.2 METRICS

314 We use three metrics for evaluating the counterfactuals produced on each model: Local Outlier 315 Factor (Breunig et al., 2000) (LOF) provides a measure of how closely a counterfactual lies to the 316 manifold of training data. It is frequently used as a measure of *plausibility* in the CFX literature. We 317 also use the *Implausibility* measure from (Altmeyer et al., 2024) as a secondary measure of plausi-318 bility. As outlined in (Altmeyer et al., 2024), this metric considers the sample-averaged Euclidean distance between a counterfactual and any in-class sample from the training set. Finally, we define 319 a novel metric, the *Robustness Ratio*, to measure the *robustness* of counterfactuals. This metric is 320 inspired by experimental protocols used to evaluate the robustness of CFXs to input changes (Artelt 321 et al., 2021; Leofante & Potyka, 2024) and is formally defined as follows. 322

Definition 3 (Robustness Ratio). Given an original input, x, and a counterfactual explanation, x_c , based on x. We compute a second counterfactual, x'_c , on a point x' where x' is sampled uniformly

324 Table 1: Numeric results for counterfactuals produced on the MNIST, credit, diabetes, news, and 325 spambase datasets, and the MLP, Ensemble, and BNN model types. For each dataset/model pair we 326 report three metrics covering the plausibility and robustness of the counterfactuals. We also report the clean accuracy, percentage of *valid* counterfactuals found for each pair, and the mean l_2 cost. 327 Arrows indicate for each metric whether high is better (\uparrow) or lower is better (\downarrow). 328

220	Dataset	Model	Clean Accuracy (%)	Valid CFX (%)		etric	l2 Cost	
330					$\text{LOF}\uparrow$	Implausibility \downarrow	Robustness Ratio $(10^{-3})\downarrow$	· · 2 · · · · · · · · · · · · · · · · ·
331		MLP	95.6	88.0	0.721	87.4	41.4	18.7
332	MNIST	Ensemble	97.3	82.0	0.463	91.0	3.94	54.4
333		BNN	98.2	74.0	0.838	87.1	44.1	7.57
334		MLP	75.5	100.0	1.0	3.45	36.5	0.730
00-1	credit	Ensemble	76.5	96.0	1.0	3.35	36.6	0.422
335		BNN	71.0	88.0	1.0	3.30	28.0	0.850
336		MLP	77.9	100.0	0.949	0.428	15.0	0.302
227	diabetes	Ensemble	76.0	100.0	0.897	0.422	15.5	0.301
557		BNN	72.7	100.0	1.0	0.423	25.4	0.254
338		MLP	65.0	92.0	1.0	65.7	629	22.2
339	news	Ensemble	65.9	74.0	1.0	68.0	636	18.8
340		BNN	65.9	80.0	1.0	64.2	460	22.8
0.4.4		MLP	92.9	80.0	0.886	55.0	497	45.6
341	spambase	Ensemble	92.9	92.0	0.902	64.6	399	44.1
342		BNN	92.9	92.0	0.854	58.4	369	44.0

from the neighborhood of $x: ||x - x'||_{\infty} < b$ with some budget b. We term the distribution governing these samples $U_b(x)$. The Robustness Ratio is then defined as the ratio of the distance between the two counterfactuals, x and x', to the cost of the initial counterfactual. Formally,

Robustness Ratio =
$$\mathbb{E}_{\boldsymbol{x}' \sim U_b(\boldsymbol{x})} \left[\frac{\|\boldsymbol{x}'_c - \boldsymbol{x}_c\|_p}{\|\boldsymbol{x}_c - \boldsymbol{x}\|_p} \right].$$
 (9)

350 We compute the expected value in Equation (9) using Monte Carlo sample averaging of $U_b(x)$. For 351 all our experiments we instantiate Definition 3 using l_2 norms, although this metric can be applied 352 with any l_p norm in general. For the budget vector, we use 5% of the element-wise input domain. 353 Where a feature has no obvious predefined input domain we use the range of that feature in the 354 training set. Any instance where we are unable to find a valid counterfactual is discarded when 355 computing the Robustness Ratio. 356

357 4.3 NUMERIC RESULTS 358

343 344

345

346

347 348

349

In Table 1 we report numeric results across three model types and five datasets. In addition to the 359 metrics, we report the clean accuracy of the models and the percentage of valid counterfactuals 360 found. Valid Counterfactuals is defined as the percentage of counterfactuals that our algorithm ob-361 tained that successfully change the model output to the intended target class. A validity of 100% 362 implies that we found a valid counterfactual for all 50 sampled test inputs. For these results we performed hyperparameter tuning of ϵ , λ , and N in Algorithm 1 independently for each dataset/model 364 run, and excluded any run which obtained a Valid Counterfactuals score of less than 70%.

Focusing initially on LOF, the results show that the BNN achieved the largest or equivalent LOF 366 score across all datasets except spambase. This indicates that on the BNN model, the counterfactuals 367 we find better represent the training data distribution. We observe similarly that the Implausibility 368 score was better or the same across all datasets except for spambase. Regarding CFX robustness, 369 we note that in the two benchmarks where LOF was equal, the Robustness Ratio was lower for the 370 BNN model than either baseline. We noted in our experiments that there appeared to be a trade-off 371 between LOF and Robustness Ratio. In Table 1, we have prioritised LOF in selecting the best runs 372 as this is a highly recognised metric in the CFX literature. For each dataset/model pair, we have 373 also reported the mean l_2 cost of counterfactuals found. Although this is not a metric pertaining to 374 plausibility or robustness, the literature on counterfactuals generally considers cheap counterfactuals 375 to be desirable. Our results show that in three of the five datasets, the counterfactuals produced on the BNN were the cheapest, as well as maintaining high scores over our metric suite. In two cases, 376 providing both the cheapest counterfactuals and the highest LOF score; these results support our 377 hypothesis. It is important to note that this is achieved with no hyperparameter tuning for low cost.

Metrics	Model	Datasets					
	110001	MNIST	credit	diabetes	news	spambase	
LOF ↑	MLP	0.721	1.0	0.949	1.0	0.886	
	Ensemble	0.463	1.0	0.897	1.0	0.902	
	BNN	1.0	1.0	1.0	1.0	0.853	
Implausibility ↓	MLP	87.4	3.34	0.428	64.0	55.0	
	Ensemble	91.0	3.35	0.422	63.9	58.3	
	BNN	86.9	3.30	0.423	64.2	57.0	
Robustness Ratio $(10^{-3})\downarrow$	MLP	10.5	34.3	6.48	513	370	
	Ensemble	2.38	27.6	6.46	405	239	
	BNN	5.38	9.79	25.4	289	226	

Table 2: Best performance results by metric across all hyperparameter tuning runs. Arrows indicate for each metric whether high is better (\uparrow) or lower is better (\downarrow).

In Table 2 we provide the best scores for each metric across all runs. Here we see a clear divide between the models with some component of averaging (Ensemble and BNN) and the MLP, which only achieves a best result in *Implausibility* on *spambase*. When comparing the Ensemble and BNN models the metrics are similar with the BNN yielding improved scores in seven of the twelve headto-heads with the Ensemble model. We note that these scores are to give a full picture only and that prioritising a single metric is usually at heavy detriment to other metrics or the l_2 cost of the counterfactual.

4.4 COMPARISON WITH (SCHUT ET AL., 2021)

In (Schut et al., 2021) the authors propose producing counterfactuals that consider aleatoric and epistemic uncertainty. They use softmax output and an ensemble of MLPs to capture the aleatoric and epistemic uncertainties respectively. Differently to us, their counterfactuals are generated using cross-entropy loss rather than our linear loss function \mathcal{L}_{lin} . Moreover, in (Schut et al., 2021) the au-thors apply a variation on the Jacobian-based saliency map (JSMA) originally applied in adversarial attacks (Papernot et al., 2016). Specifically, the JSMA limits updates to the input to only consider the input dimension with the largest partial gradients.

Here we apply these modifications to align our CFX algorithm with that from (Schut et al., 2021) and compare the CFXs produced by the Ensemble models and our BNNs. As in Section 4.3 we compare each model/dataset pair using three established metrics and show the results in Table 3.

The results in Table 3 show that the BNN outperforms or ties the Ensemble in LOF for every benchmark, and where there is a tie the BNN maintains a lower Robustness Ratio. The performance is more consistently in favour of the BNN than we see in Table 1. However, we note that the cost of the BNN's counterfactuals are often higher than the Ensemble, in contrast to what we observe in Table 1.

 Table 3: Numeric results for MNIST, credit, diabetes, news, and spambase datasets on the (Schut et al., 2021) MLP Ensembles and our BNNs. For each metric arrows indicate whether higher is better(\uparrow) or lower is better (\downarrow).

400	Dataset	Model	Clean Accuracy (%)	Valid CFX (%)		la Cost		
423	Dutaber	moder	clouir recuracy (,c)	(<i>i</i> , <i>i</i>)	LOF ↑	Implausibility \downarrow	Robustness Ratio $(10^{-3})\downarrow$	• •2 0000 ¢
424) O HOT	Ensemble	97.3	100.0	0.560	98.3	58.0	14.0
425	MNIST	BNN	98.2	74.0	0.946	94.5	124	3.40
426	credit	Ensemble	76.5	100.0	1.0	3.70	19.7	1.58
427	cicuit	BNN	71.0	100.0	1.0	4.03	3.54	2.38
100	diabatas	Ensemble	100.0	100.0	0.949	0.417	22.3	0.202
420	ulabeles	BNN	72.7	100.0	1.0	0.435	23.7	0.277
429		Ensemble	65.9	74.0	1.0	66.8	561	19.0
430	news	BNN	65.9	80.0	1.0	64.1	517	22.5
431	spambase	Ensemble	92.9	76.0	0.879	52.7	463	47.8
		BNN	92.9	78.0	0.882	54.3	530	47.0

In Figure 2 we compare counterfactuals produced by the Ensemble and BNN models. Here we observe a much clearer counterfactual for both models than in Figure 1, this is due to the JSMA filtering applied. However, we note that the Ensemble model continues to produce more prominent erroneous artifacts than the BNN.



Figure 2: Visual comparison of two counterfactuals produced under the same setting as in (Schut et al., 2021) for (b) an Ensemble model, and (c) a BNN. Original inputs are shown in (a). The top counterfactual is for an original input of 5, with a target of 6 and the lower counterfactual is for an original input of 9 with target class 4.

5 RELATED WORK

451

452

453

454

455 456 457

458

Various methods for generating CFXs have been proposed for a wide range of machine learning classifiers; see, e.g. (Guidotti, 2024; Karimi et al., 2023) for recent surveys. These include approaches targeting tree-based classifiers (Tolomei et al., 2017), linear classifiers (Ustun et al., 2019) as well as non-linear ones implemented by means of deep neural networks (Wachter et al., 2017). These algorithms typically cast the problem of finding explanations as an optimisation problem aimed at generating explanations that satisfy properties of interest, including validity, actionability, sparsity, and robustness. We refer the reader to Section 2 for a more detailed discussion on this.

466 Closely related to this work are approaches that generate CFXs for probabilistic models. For exam-467 ple, Bayesian classifiers are considered in (Albini et al., 2020), where CFXs are given in the form of 468 influence relations between features. This is different from the type of counterfactuals that we aim 469 to generate, in that our explanations are built from feature-wise modifications as commonly studied 470 in the literature (Guidotti, 2024). Bayesian Neural Networks are considered in (Antorán et al., 2021; Ley et al., 2022), where counterfactual explanations are defined in terms of minimal modifications 471 to input features that would result in an increase in confidence for the prediction produced by the 472 BNN. Our objective is different, as our counterfactuals are designed to change the prediction of the 473 classifier, in line with common definitions encountered in the literature on CFX (Guidotti, 2024; 474 Karimi et al., 2023). Other approaches have considered techniques for uncertainty quantification 475 to improve the quality of CFX in the presence of uncertainty. For instance, conformal prediction 476 sets and deep ensembling techniques were used in (Altmeyer et al., 2024) to generate CFXs that lie 477 closer to the data manifold. While these techniques are effective at improving the plausibility of 478 counterfactuals, they differ from our approach in that we aim to generate explanations by reasoning 479 directly on a model trained to incorporate uncertainty in its decision-making process. 480

We would like to note that while in this paper we will focus mostly on tabular data and images, explanation algorithms for other data types have been proposed, including graph data (Bajaj et al., 2021), vision tasks (Augustin et al., 2022) and time series classification tasks (Delaney et al., 2021).

 Unlike deterministic NNs, Bayesian NNs provide a natural, yet powerful, way of quantifying uncertainty in Deep Learning Models (Gal, 2016). BNNs treat model parameters as probability distributions, allowing for the computation of predictive uncertainty (MacKay, 1992). In this paper, we use

9

486 this powerful feature to synthesise counterfactuals that are more plausible than those generated by 487 deterministic or even ensemble models. 488

Although, to our knowledge, counterfactual explanations have not been defined for BNNs, there is 489 work that considers related concepts. For instance, Ali et al. (2023) study counterfactual explana-490 tions of Bayesian model uncertainty. Unlike the work presented here (Ali et al., 2023) adapts existing 491 counterfactual generation techniques to work with BNNs; thus, not fully utilising the BNNs. Raman 492 et al. (2023) on the other hand, consider deterministic NNs but treat the feature perturbations as 493 random variables endowed with prior distribution functions to provide several alternative explana-494 tions rather than a single point solution. Schut et al. (2021) use aletoric and epistemic uncertainties 495 obtained from an ensemble of models to generate more interpretable counterfactuals.

- 496 497 498
- CONCLUSION

499 In this paper we have presented the first formal study of counterfactual explanations for Bayesian 500 Neural Networks. We proposed a definition for CFX within the context of BNNs as well as a 501 framework for computing them in practice. We reported results on five commonly used datasets 502 from the CFX literature and compared the performance of our method against two baselines: MLPs 503 and MLP ensembles. We have shown that BNNs often produce cheaper, more robust, and more 504 plausible explanations. We observe that some state-of-the-art metrics appear to exhibit trade-off 505 behaviour in all models. Notably, obtaining highly robust explanations is observed to be more 506 costly, confirming previous observations made in the context of deterministic models (Jiang et al., 507 2024).

508

509 LIMITATIONS AND FUTURE WORK 510

511 In this work, we have explored CFX produced on BNNs in a straightforward setting, namely, by 512 solving Equation (8) without applying additional regularisations to promote robustness or plausi-513 bility as sometimes used in the literature (Karimi et al., 2023; Jiang et al., 2024). It would be 514 interesting to see how our results hold up in such a setting and we leave this investigation for fu-515 ture work. Furthermore, we consider only HMC-based Bayesian inference algorithms and a single 516 network architecture for ease of comparison. We leave exploring CFX on BNNs produced by less 517 precise inference algorithms and varying architectures as points for future work.

518

523 524

525 526

527

528 529

531

519 **Reproducibility Statement.** We have provided the models and code for reproducing all exper-520 iments in the supplementary materials. All datasets used in this study are open source and freely 521 available on the Internet. Additionally, we have included relevant citations to facilitate locating these 522 resources online for result reproduction.

REFERENCES

- Emanuele Albini, Antonio Rago, Pietro Baroni, and Francesca Toni. Relation-based counterfactual explanations for bayesian network classifiers. In International Joint Conference on Artificial Intelligence, IJCAI, pp. 451–457. ijcai.org, 2020. 9
- Gohar Ali, Feras Al-Obeidat, Abdallah Tubaishat, Tehseen Zia, Muhammad Ilyas, and Alvaro 530 Rocha. Counterfactual explanation of bayesian model uncertainty. Neural Computing and Applications, pp. 1-8, 2023. 10 532
- 533 Patrick Altmeyer, Mojtaba Farmanbar, Arie van Deursen, and Cynthia C. S. Liem. Faithful model 534 explanations through energy-constrained conformal counterfactuals. In Michael J. Wooldridge, 535 Jennifer G. Dy, and Sriraam Natarajan (eds.), AAAI Conference on Artificial Intelligence, pp. 536 10829–10837. AAAI Press, 2024. 3, 6, 9
- 537
- Javier Antorán, Umang Bhatt, Tameem Adel, Adrian Weller, and José Miguel Hernández-Lobato. 538 Getting a CLUE: A method for explaining uncertainty estimates. In International Conference on 539 Learning Representations, ICLR. OpenReview.net, 2021. 9

540 541 542	André Artelt, Valerie Vaquet, Riza Velioglu, Fabian Hinder, Johannes Brinkrolf, Malte Schilling, and Barbara Hammer. Evaluating robustness of counterfactual explanations. In <i>IEEE Symposium Series on Computational Intelligence, SSCI</i> , pp. 1–9. IEEE, 2021. 2, 3, 6
543 544 545 546	Maximilian Augustin, Valentyn Boreiko, Francesco Croce, and Matthias Hein. Diffusion visual counterfactual explanations. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), <i>Neural Information Processing Systems, NeurIPS</i>, 2022. 9
547 548 549	Mohit Bajaj, Lingyang Chu, Zi Yu Xue, Jian Pei, Lanjun Wang, Peter Cho-Ho Lam, and Yong Zhang. Robust counterfactual explanations on graph neural networks. In <i>Advances in Neural Information Processing Systems, NeurIPS</i> , pp. 5644–5655, 2021. 9
550 551 552 553	Solon Barocas, Andrew D. Selbst, and Manish Raghavan. The hidden assumptions behind coun- terfactual explanations and principal reasons. In <i>Conference on Fairness, Accountability, and</i> <i>Transparency, FAT*</i> , pp. 80–89, 2020. 1
554 555 556	Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. In <i>International Conference on International Conference on Machine Learning, ICML</i> , ICML'15, pp. 1613–1622. JMLR.org, 2015. 3
557 558 559	Vivek S. Borkar. Equation of state calculations by fast computing machines. <i>Resonance</i> , 27:1263 – 1269, 1953. 3
560 561 562	Markus M. Breunig, Hans-Peter Kriegel, Raymond T. Ng, and Jörg Sander. Lof: identifying density- based local outliers. In <i>Proceedings of the 2000 ACM SIGMOD International Conference on</i> <i>Management of Data</i> , pp. 93–104. Association for Computing Machinery, 2000. 6
563 564 565	Ruth M. J. Byrne. Counterfactuals in explainable artificial intelligence (XAI): evidence from human reasoning. In <i>International Joint Conference on Artificial Intelligence, IJCAI</i> , pp. 6276–6282, 2019. 1
567 568	Longbing Cao. AI in finance: challenges, techniques, and opportunities. <i>ACM Computing Surveys</i> (<i>CSUR</i>), 55(3):1–38, 2022. 1
569 570 571	Nicholas Carlini and David A. Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. <i>Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security</i> , 2017. 4
572 573 574 575 576	Eoin Delaney, Derek Greene, and Mark T. Keane. Instance-based counterfactual explanations for time series classification. In <i>Case-Based Reasoning Research and Development - 29th Inter-</i> <i>national Conference, ICCBR</i> , volume 12877 of <i>Lecture Notes in Computer Science</i> , pp. 32–47. Springer, 2021. 9
577 578 579 580	 Amit Dhurandhar, Pin-Yu Chen, Ronny Luss, Chun-Chen Tu, Pai-Shun Ting, Karthikeyan Shan-mugam, and Payel Das. Explanations based on the missing: Towards contrastive explanations with pertinent negatives. In <i>Neural Information Processing Systems, NeurIPS</i>, pp. 590–601, 2018. 3
581 582 583	Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. http://archive.ics.uci.edu/ml. 2,5
584 585	Simon Duane, A.D. Kennedy, Brian J. Pendleton, and Duncan Roweth. Hybrid monte carlo. <i>Physics Letters B</i> , 195(2):216–222, 1987. 3
586 587 588	Kelwin Fernandes, Pedro Vinagre, Paulo Cortez, and Pedro Sernadela. Online News Popularity. UCI Machine Learning Repository, 2015. https://doi.org/10.24432/C5NS3V. 2, 5
589	Yarin Gal. Uncertainty in Deep Learning. PhD thesis, University of Cambridge, 2016. 9
590 591 592	Riccardo Guidotti. Counterfactual explanations and how to find them: literature review and bench- marking. <i>Data Min. Knowl. Discov.</i> , 38(5):2770–2824, 2024. 1, 9
593	Mark Hopkins, Erik Reeber, George Forman, and Jaap Suermondt. Spambase. UCI Machine Learn- ing Repository, 1999. https://doi.org/10.24432/C53G6X. 2, 5

594 595 596 597	Junqi Jiang, Francesco Leofante, Antonio Rago, and Francesca Toni. Robust counterfactual ex- planations in machine learning: A survey. In <i>Proceedings of the Thirty-Third International Joint</i> <i>Conference on Artificial Intelligence, IJCAI 2024</i> , pp. 8086–8094. International Joint Conferences on Artificial Intelligence Organization, IJCAI, 2024. 1, 10
599 600 601	Amir-Hossein Karimi, Gilles Barthe, Bernhard Schölkopf, and Isabel Valera. A survey of algorithmic recourse: Contrastive explanations and consequential recommendations. ACM Comput. Surv., 55(5):95:1–95:29, 2023. 1, 3, 9, 10
602 603 604 605	Eoin M. Kenny and Mark T. Keane. On generating plausible counterfactual and semi-factual expla- nations for deep learning. In AAAI Conference on Artificial Intelligence, pp. 11575–11585. AAAI Press, 2021. 1
606 607 608 609	Thibault Laugel, Marie-Jeanne Lesot, Christophe Marsala, Xavier Renard, and Marcin Detyniecki. The dangers of post-hoc interpretability: Unjustified counterfactual explanations. In Sarit Kraus (ed.), <i>Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence,</i> <i>IJCAI</i> , pp. 2801–2807. ijcai.org, 2019. 1
610 611	Y. LeCun, C. Cortes, and C. J. Burges. The mnist database of handwritten digits, 1998. 2, 5
612 613	Francesco Leofante and Nico Potyka. Promoting counterfactual robustness through diversity. In <i>AAAI Conference on Artificial Intelligence, AAAI</i> , pp. 21322–21330. AAAI Press, 2024. 3, 6
615 616 617	Dan Ley, Umang Bhatt, and Adrian Weller. Diverse, global and amortised counterfactual explana- tions for uncertainty estimates. In <i>AAAI Conference on Artificial Intelligence</i> , pp. 7390–7398. AAAI Press, 2022. 9
618 619 620	David J. C. MacKay. A practical bayesian framework for backpropagation networks. <i>Neural Comput.</i> , 4(3):448–472, 1992. 9
621 622 623	Tim Miller. Explanation in artificial intelligence: Insights from the social sciences. <i>Artif. Intell.</i> , 267:1–38, 2019. 1
624 625 626	Kiarash Mohammadi, Amir-Hossein Karimi, Gilles Barthe, and Isabel Valera. Scaling guarantees for nearest counterfactual explanations. In <i>AAAI/ACM Conference on AI, Ethics, and Society</i> , pp. 177–187. ACM, 2021. 2
627 628 629	Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and Anan- thram Swami. The limitations of deep learning in adversarial settings. In 2016 IEEE European Symposium on Security and Privacy (EuroS&P), pp. 372–387, 2016. 8
630 631 632 633 634	Natraj Raman, Daniele Magazzeni, and Sameena Shah. Bayesian hierarchical models for coun- terfactual estimation. In Francisco J. R. Ruiz, Jennifer G. Dy, and Jan-Willem van de Meent (eds.), <i>International Conference on Artificial Intelligence and Statistics</i> , <i>AISTATS</i> , volume 206 of <i>Proceedings of Machine Learning Research</i> , pp. 1115–1128. PMLR, 2023. 10
635 636 637 638 639	Lisa Schut, Oscar Key, Rory McGrath, Luca Costabello, Bogdan Sacaleanu, Medb Corcoran, and Yarin Gal. Generating interpretable counterfactual explanations by implicit minimisation of epis- temic and aleatoric uncertainties. In Arindam Banerjee and Kenji Fukumizu (eds.), <i>International</i> <i>Conference on Artificial Intelligence and Statistics, AISTATS</i> , volume 130 of <i>Proceedings of Ma-</i> <i>chine Learning Research</i> , pp. 1756–1764. PMLR, 2021. 5, 8, 9, 10
640 641 642	Mohammed Yousef Shaheen. Applications of artificial intelligence (ai) in healthcare: A review. <i>ScienceOpen Preprints</i> , 2021. 1
643 644 645	Dylan Slack, Anna Hilgard, Himabindu Lakkaraju, and Sameer Singh. Counterfactual explanations can be manipulated. In <i>Neural Information Processing Systems, NeurIPS</i> , pp. 62–75, 2021. 2, 3
646 647	Jack Smith, J. Everhart, W. Dickson, W. Knowler, and Richard Johannes. Using the adap learning algorithm to forcast the onset of diabetes mellitus. <i>Proceedings - Annual Symposium on Computer Applications in Medical Care</i> , 10, 11 1988. 2, 5

648 649 650	Gabriele Tolomei, Fabrizio Silvestri, Andrew Haines, and Mounia Lalmas. Interpretable predic- tions of tree-based ensembles via actionable feature tweaking. In <i>ACM SIGKDD International</i> <i>Conference on Knowledge Discovery and Data Mining</i> , pp. 465–474. ACM, 2017. 9
651 652	Berk Ustun, Alexander Spangher, and Yang Liu. Actionable recourse in linear classification. In <i>Conference on Fairness, Accountability, and Transparency, FAT*</i> , pp. 10–19. ACM, 2019. 3, 9
653 654 655 656	Sandra Wachter, Brent D. Mittelstadt, and Chris Russell. Counterfactual explanations without open- ing the black box: Automated decisions and the GDPR. <i>Harv. JL & Tech.</i> , 31:841, 2017. 2, 3, 9
657 658	
659 660	
662 663	
664 665	
666 667	
669 670	
671 672	
673 674	
675 676 677	
678 679	
680 681	
682 683 684	
685 686	
687 688	
690 691	
692 693	
694 695	
696 697 698	
699 700	
701	