

Probabilistically Robust Counterfactual Explanations under Model Changes

Luca Marzari^{a,*}, Francesco Leofante^b, Ferdinando Cicalese^a, Alessandro Farinelli^a

^a*Department of Computer Science, University of Verona, Str. le Grazie, 15,
Verona, 37134, Italy*

^b*Department of Computing, Imperial College London, 180 Queen's Gate,
London, SW7 2AZ, United Kingdom*

Abstract

We study the problem of generating robust counterfactual explanations for deep learning models subject to model changes. We focus on *plausible model changes* altering model parameters and propose a novel framework to reason about the robustness property in this setting. To motivate our solution, we begin by showing for the first time that computing the robustness of counterfactuals with respect to model changes is NP-hard. As this (practically) rules out the existence of scalable algorithms for exactly computing robustness, we propose a novel probabilistic approach which is able to provide tight estimates of robustness with strong guarantees while preserving scalability. Remarkably, and differently from existing solutions targeting plausible model changes, our approach does not impose requirements on the network to be analysed, thus enabling robustness analysis on a wider range of architectures, including state-of-the-art tabular transformers. A thorough experimental analysis on four binary classification datasets reveals that our method improves the state of the art in generating robust explanations, outperforming existing methods.

Keywords: Explainable AI, Counterfactual Explanations, Algorithmic Recourse, Robustness of Explanations

*Corresponding author.

Email addresses: luca.marzari@univr.it (Luca Marzari),
f.leofante@imperial.ac.uk (Francesco Leofante), ferdinando.cicalese@univr.it
(Ferdinando Cicalese), alessandro.farinelli@univr.it (Alessandro Farinelli)

1. Introduction

Deep Neural Networks (DNNs) have emerged as a groundbreaking technology revolutionizing several fields ranging from autonomous navigation [1, 2] to image classification [3] and robotics for medical applications [4]. However, despite remarkable successes, their vulnerability to adversarial attacks [5, 6], i.e., imperceptible modifications to input data that can lead to wrong and potentially catastrophic decisions when deployed, has raised crucial safety concerns. Consequently, understanding and explaining the decisions of black-box deep learning models has become a dominant goal in AI research. In this paper, we focus on counterfactual explanations (CFX), a popular class of explanation methods that aim to demystify the decision-making of a DNN by showing how an input needs to be changed to yield a different, typically more desirable, decision (see [7, 8] for recent surveys).

To understand what makes CFXs useful, consider the widely used example of a loan application, where a mortgage applicant represented by an input x with features *unemployed* status, 25 years of age, and *low* credit rating applies for a loan and is rejected by the bank’s AI. A CFX for this decision could be a slightly modified input, where increasing credit rating to *medium* would result in the loan being granted. Ideally, a counterfactual explanation should be as close as possible to the original input to ensure that the changes it suggests are feasible. The approach of [9], showed how this requirement can be mathematically achieved by generating a counterfactual as close as possible to the decision boundary of a DNN. However, producing explanations in this way raises critical concerns about their reliability (see [10] for a survey). For example, as illustrated in Figure 1, fine-tuning the model with additional data can significantly alter its decision boundary, potentially invalidating previously generated counterfactuals. This sensitivity to the model changes for the counterfactuals poses critical questions about the reliability of explanations and long-term usability in dynamic settings.

In particular, recent work has highlighted issues related to the robustness of CFXs against *Plausible Model Changes* (PMC) [11, 12], showing that the validity of CFXs is likely to be compromised when bounded perturbations are applied to the parameters of a DNN, e.g., as a result of fine-tuning [11, 13, 14, 12, 15]. Consider the loan example: if retraining occurs while the applicant is working toward improving their credit rating, without robustness, their

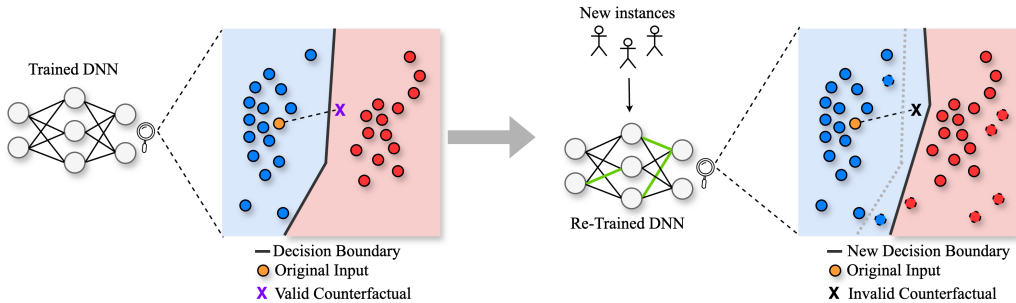


Figure 1: Vignette illustrating the problem of robustness under model changes. A counterfactual explanation is initially generated for a trained model (left). Then, the model is updated to include new data (right). This step might induce slight changes in the decision boundary of the model, ultimately invalidating the counterfactual explanation generated in the first step.

36 modified case may still result in a rejected application, leaving the bank
 37 liable due to their conflicting statements.

38 In this paper, we focus on this troubling phenomenon and advance the
 39 state of the art in CFX robustness research in several directions. More specif-
 40 ically, we start by studying the computational complexity of exactly deter-
 41 mining whether a CFX is robust to PMC in § 3. Our result formally shows for
 42 the first time that this is an NP-hard problem, thus providing new insights
 43 into algorithmic developments in this area.

44 As our hardness results rule out the existence of practical algorithms to
 45 compute the CFX robustness in an exact fashion, we argue that probabilistic
 46 approaches are needed to obtain answers on the CFX robustness under model
 47 changes. Notably, the work by Hamman et al. [15], proposes a probabilis-
 48 tic approach to compute the robustness of CFX under *Naturally-Occurring*
 49 *Model Changes* (NOMC).¹ Even though both PMC and NOMC notions are
 50 commonly used in the literature, very little is known about their potential
 51 interplay, and whether robustness to NOMC subsumes robustness to PMC
 52 is still unresolved. In § 4, we report a complete study of the two notions and
 53 formally prove that these two notions capture profoundly different scenarios.

¹In this work, we primarily use the term “model changes”, following the notation used in recent surveys on the topic [10]. An alternative term “model shifts”, with similar meaning, has also been used in related literature, as in [16]. The two terms will be used interchangeably throughout.

54 As a result, we demonstrate that robustness guarantees given for NOMC do
55 not directly extend to PMC. Having settled this, in § 5, we present an ex-
56 tended overview of our $\text{AP}\Delta\text{S}$, a novel sampling-based certification algorithm
57 that allows us to determine a provable probabilistic bound on the maximum
58 shift a CFX can tolerate under PMC. Unlike existing solutions for robustness
59 under PMC, our approach comes with significantly reduced computational
60 requirements and does not make any assumption on the underlying DNN,
61 thus making it applicable to a wider range of architectures, including state-
62 of-the-art transformer architectures.

63 To confirm this aspect, in § 6, we present a thorough experimental evalua-
64 tion analysing the performance of $\text{AP}\Delta\text{S}$, providing a comprehensive compar-
65 ison of the proposed approach against several state-of-the-art methodologies
66 for CFX robustness and different ablation studies. Crucially, we show that
67 our approach outperforms existing methods on several metrics from the CFX
68 literature, including validity, proximity, and plausibility.

69 The paper is structured as follows. In § 2, we cover the related work,
70 and in § 3, we introduce background notions on computing robust CFXs un-
71 der model changes. § 3 presents our complexity analysis and offers complete
72 proof of NP-hardness for both PMC and NOMC. Motivated by this result, in
73 § 4, we study existing approaches to generate probabilistically robust CFXs
74 and analyze their interplay. Then, in § 5, we introduce our method to gen-
75 erate robust CFXs, $\text{AP}\Delta\text{S}$, and evaluate it extensively in § 6. The core
76 contributions of this work can be summarised as follows:

- 77 • We prove, for the first time, that determining whether a CFX is robust
78 to model changes in a deep neural network is an NP-complete problem,
79 for both existing notions of NOMC and PMC. This finding highlights
80 the need for further research into probabilistic methods to address this
81 problem effectively.
- 82 • We analyse existing approaches to generate probabilistic guarantees
83 for CFXs under NOMC and demonstrate that these guarantees do not
84 extend to PMC.
- 85 • We present $\text{AP}\Delta\text{S}$, a scalable procedure that is able to generate prov-
86 ably robust CFXs. This approach introduces an iterative algorithm
87 to generate probabilistically robust CFXs, which are demonstrated to
88 have superior performance against four robust baselines.

- 89 • To confirm the scalability and effectiveness of our solution, we employ
90 $\text{AP}\Delta\text{S}$ to certify the robustness of CFXs for state-of-the-art transformer
91 architectures [17] employed in tabular data classification. To the best
92 of our knowledge, we are the first to consider models of this size within
93 the robust CFX literature [10].

94 This paper builds upon our previous work [16] with significant extensions.
95 Specifically, § 3 non trivially extends the corresponding section in [16] and
96 offers a full hardness proof for the problem of deciding robustness of coun-
97 terfactual explanations under PMC. As a corollary of this result, we are also
98 able to show the hardness with respect to NOMC, thus providing a rigorous
99 characterization of the complexity of verifying CFX robustness under existing
100 notions of model changes. § 4 is also extended with a thorough experimental
101 evaluation, complementing our theoretical findings of [16] and showing that
102 PMC and NOMC capture very different robustness requirements in practice.
103 Our experimental analysis in § 6 is also extended considerably. In particular,
104 we present a novel analysis of the impact that the main hyper-parameters of
105 $\text{AP}\Delta\text{S}$ can have on the quality of CFXs it generates. Moreover, we demon-
106 strate the scalability of our approach by presenting new results obtained on
107 large-scale tabular transformers. To the best of our knowledge, this is the
108 first time a method for robust CFXs has been shown to scale to state-of-the-
109 art transformer models. These results complement our previous analysis and
110 demonstrate the versatility of $\text{AP}\Delta\text{S}$, as well as its effectiveness in solving
111 robustness issues in state-of-the-art machine learning models.

112 2. Related Work

113 Various methods for generating CFXs for DNNs have been proposed.
114 The seminal work of [9] framed the task of generating CFXs as a gradient-
115 based optimization problem and proposed a loss that promotes CFX *validity*
116 (i.e., the CFX successfully changes the classification outcome of the network)
117 and *proximity* (i.e., the CFX is as close as possible to the original input
118 for some distance metric). In addition to these metrics, other important
119 properties have been highlighted as crucial for the practical applicability of
120 CFXs. Prominent examples include *plausibility* (i.e., the CFX must lie on the
121 data manifold) [18, 19] and *actionability* (i.e., the changes suggested by the
122 CFX must be achievable by the user in practice) [20]. Differently from these
123 works, here we focus on the robustness property of CFXs.

124 Several forms of CFX robustness have been studied in the literature [10].
 125 Robustness to input changes is the focus of, e.g. [21, 22, 23, 24, 25], where
 126 solutions are devised to ensure that explanation algorithms return similar
 127 CFXs for similar inputs. In another line of work, [26, 27, 28, 29] considered
 128 the problem of generating adversarially robust CFXs that preserve validity
 129 under imperfect (or noisy) execution. Robustness to model multiplicity is
 130 instead considered in, e.g. [19, 30, 31], where CFXs that preserve validity
 131 across sets of models are sought. However, the study of these forms of ro-
 132 bustness is outside the scope of this paper as our focus is on model changes.
 133 Robustness to model changes has been studied in, e.g. [32, 13, 14, 12, 15, 33].
 134 Of these, the approaches of [11] and [33] are the most closely related to our
 135 work. The former presents an approach to generate robust CFXs under PMS
 136 using techniques from continuous optimization, which is able to guarantee
 137 robustness in the average-case scenario. The latter instead solves the same
 138 problem using abstraction techniques and discrete optimization tools, obtain-
 139 ing robustness guarantees that hold under worst-case conditions. Given their
 140 relevance, both approaches will be considered for an extensive experimental
 141 comparison in § 6.

142 3. Background

143 **Neural networks and classification tasks.** Let $\mathcal{X} \subseteq \mathbb{R}^d$ denote the in-
 144 put space of a *classifier* $\mathcal{M}_\theta : \mathcal{X} \rightarrow [0, 1]$ mapping an input $x \in \mathcal{X}$ to an
 145 output probability between 0 and 1. We consider classifiers implemented
 146 by feed-forward DNNs parameterized by a (*parameter*) *vector* $\theta \in \Theta \subseteq \mathbb{R}^k$.
 147 Given two parameter vectors $\theta, \theta' \in \Theta$, we refer to the corresponding clas-
 148 sifiers \mathcal{M}_θ and $\mathcal{M}_{\theta'}$ as *instantiations* of the same parametric classifier \mathcal{M}_Θ .
 149 We assume concrete valuations of θ are learned from a set of labeled inputs
 150 as customary in supervised learning settings [34]. Once θ has been learned,
 151 the classifier can be used for inference. Without any loss of generality, we
 152 focus on binary classification tasks, i.e., the classification decision produced
 153 by \mathcal{M}_θ for an unlabeled input x is 1 if $\mathcal{M}_\theta(x) \geq 0.5$, and 0 otherwise.

154 **Counterfactual explanations.** Existing methods in the literature define
 155 CFXs as follows.

156 **Definition 1.** *Consider an input $x \in \mathcal{X}$ and a classifier \mathcal{M}_θ s.t. $\mathcal{M}_\theta(x) <$
 157 0.5 . Given a distance metric $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^+$, a (valid) counterfactual
 158 explanation is any x' such that:*

$$x' = \underset{\hat{x} \in \mathcal{X} : \mathcal{M}_\theta(\hat{x}) \geq 0.5}{\operatorname{arg\,min}} d(x, \hat{x})$$

159 Intuitively, given an input x for which the classifier produces a negative
 160 outcome, a counterfactual explanation is a new input x' which is similar
 161 to x , e.g., in terms of some specified distance between features values, and
 162 for which the classifier predicts a different outcome. Common choices for d
 163 include the ℓ_1 and ℓ_∞ norms [9], which will also be used in this work.

164 **Robustness to model changes.** Among several notions of robustness,
 165 recent work has placed emphasis on generating CFXs that remain valid under
 166 (slight) changes in the classifier they were generated for. While existing
 167 approaches rely on a diverse range of techniques to solve this problem, they
 168 all share a common understanding of what constitutes a model shift, which
 169 we present next.

170 **Definition 2** (Jiang et al. [12]). Let \mathcal{M}_θ and $\mathcal{M}_{\theta'}$ be two instantiations of a
 171 parametric classifier \mathcal{M}_Θ . For $0 \leq p \leq \infty$, the p -distance between \mathcal{M}_θ and
 172 $\mathcal{M}_{\theta'}$ is defined as $d_p(\mathcal{M}_\theta, \mathcal{M}_{\theta'}) = \|\theta - \theta'\|_p$.

173 **Definition 3** (Jiang et al. [12]). A model shift (w.r.t. a fixed p -distance) is
 174 a function S mapping a classifier \mathcal{M}_θ into another classifier $\mathcal{M}_{\theta'} = S(\mathcal{M}_\theta)$
 175 such that:

- 176 • \mathcal{M}_θ and $\mathcal{M}_{\theta'}$ are instantiations of the same \mathcal{M}_Θ ;
- 177 • $d_p(\mathcal{M}_\theta, \mathcal{M}_{\theta'}) > 0$.

178 Informally, a model shift captures changes in the parameters of a DNN,
 179 but does not affect its architecture. Based on this definition, we can formalize
 180 the robustness property for a CFX as follows.

181 **Definition 4.** Consider an input $x \in \mathcal{X}$ and a classifier \mathcal{M}_θ s.t. $\mathcal{M}_\theta(x) <$
 182 0.5 . Let x' be a counterfactual explanation computed for x s.t. $\mathcal{M}_\theta(x') \geq 0.5$.
 183 Given a set of model changes Δ , we say that the counterfactual x' is Δ -robust
 184 if $S(\mathcal{M}_\theta)(x') \geq 0.5$ for all $S \in \Delta$.

185 The definition of a model shift can be specialized to better characterize
 186 how θ is allowed to change under S . In the following, we report two most com-
 187 monly studied notions of model changes: *Naturally-Occurring Model Changes*
 188 and *Plausible Model Changes*.

189 **Definition 5** (Hamman et al. [15] (NOMC)). Consider a classifier \mathcal{M}_θ . A
 190 set of model changes Δ is said to be naturally occurring if for a (randomly)
 191 chosen model change S from Δ and $\mathcal{M}_{\theta'} = S(\mathcal{M}_\theta)$ being the new classifier
 192 obtained after applying S to \mathcal{M}_θ the following hold:

- 193 • $\mathbb{E}[\mathcal{M}_{\theta'}(x)] = \mathcal{M}_\theta(x)$; where the expectation is over the randomness of
 194 $\mathcal{M}_{\theta'}$ given a fixed value of x ;
- 195 • $\text{Var}[\mathcal{M}_{\theta'}(x)] = \nu_x$, where ν_x represents the maximum variance of the
 196 prediction of $\mathcal{M}_{\theta'}(x)$, and whenever x lies on the data manifold \mathcal{X} , ν_x
 197 is upper bounded by a small constant ν ;
- 198 • If \mathcal{M}_θ is Lipschitz continuous for some γ_1 , then $\mathcal{M}_{\theta'}(x)$ is also Lipschitz
 199 continuous for some γ_2 .

200 Broadly speaking, a naturally-occurring model shift allows the application
 201 of arbitrary changes to θ as long as the resulting model remains part of a class
 202 of models that are expected to have the same behaviour. This is in contrast
 203 with the notion of plausible model shift [11, 12], which requires changes to
 204 be bounded.

205 **Definition 6** (Jiang et al. [12] (PMC)). Consider a classifier \mathcal{M}_θ and a
 206 new classifier $\mathcal{M}_{\theta'} = S(\mathcal{M}_\theta)$ obtained after applying a model shift S to \mathcal{M}_θ .
 207 Given some $\delta \in \mathbb{R}_{>0}$ and $0 \leq p \leq \infty$, S is said to be plausible (w.r.t. the
 208 choice of parameters δ and p) if $d_p(\mathcal{M}_\theta, S(\mathcal{M}_\theta)) \leq \delta$.

209 Therefore, for any choice of parameters p, δ , and any instantiation \mathcal{M}_θ of
 210 a parametric classifier \mathcal{M}_Θ we define the set of PMC Δ_p obtained by consid-
 211 ering all changes S that satisfy Definition 6, i.e. $\Delta = \{S \mid d_p(\mathcal{M}_\theta, S(\mathcal{M}_\theta)) \leq$
 212 $\delta\}$.

213 In the following, we will refer to any instantiation $\mathcal{M}_{\theta'}$ of \mathcal{M}_Θ which is
 214 obtainable by applying a model change in Δ to \mathcal{M}_θ as a realisation of Δ .
 215 Moreover, whenever it is not explicitly specified, we will tacitly assume that
 216 the underlying distance $d_p(\cdot, \cdot)$ is the ∞ -norm.

217 Jiang et al. [12] proposed to reason about robustness under PMC using
 218 an Interval Neural Network (INN) [35] as an intermediate representation.

219 **Definition 7.** An interval neural network \mathcal{I} is a neural network where on
 220 each edge e is associated with an interval $I_e = [a_e, b_e]$. A realisation of the
 221 INN \mathcal{I} is a neural network having the same topology of \mathcal{I} and such that the
 222 weight w_e on edge e satisfies $w_e \in I_e$, i.e., it is taken from the interval
 223 associated to the same arc in \mathcal{I} .

224 Jiang et al. [12] exploits the fact that the interval weights of an INN
 225 allow to represent an over-approximation of all the possible models obtain-
 226 able under a set of PMC Δ , thus providing a compact representation of the
 227 problem. Similarly, in our work, we use the INN representation to model
 228 the PMC concept. However, instead of analysing the robustness of counter-
 229 factual explanations through the entire INN, we focus directly on reasoning
 230 regarding the potential realisations within the set Δ . This results in several
 231 computational improvements as we will discuss in section 5.

232 In this section, we study the computational complexity of deciding whether
 233 a given counterfactual explanation is robust in the presence of model shifts.
 234 Our aim here is to better understand the computational challenges arising
 235 from this problem and to use these results to guide the development of novel,
 236 more efficient certification procedures. Without loss of generality, we first fo-
 237 cus on PMC and show the NP-hardness of verifying CFX robustness with
 238 respect to this definition of model changes. Later, we show that the set of
 239 PMC used by our reduction also constitutes a set of NOMC, which implies
 240 that CFX robustness is, in general, also hard to verify with respect to NOMC.

241 Deciding whether for a given \mathcal{M}_θ a CFX x' is robust with respect to a
 242 set of PMC Δ requires to check if, for at least one model shift in Δ , there
 243 exists a realisation $\mathcal{M}_{\theta'}$ which classifies CFX x' differently from \mathcal{M}_θ , i.e.,
 244 $\mathcal{M}_{\theta'}(x) < 0.5 \leq \mathcal{M}_\theta$. This question is encoded in the following problem.

DISTINCT-REALISATIONS PROBLEM (DRP)

Input: an instantiation \mathcal{M}_{θ_1} of a parametric classifier \mathcal{M}_Θ , an input x such that $\mathcal{M}_{\theta_1}(x) \geq 0.5$, and a set Δ of PMC.

Output: yes \iff there exists an instantiation \mathcal{M}_{θ_2} of \mathcal{M}_Θ which is a realisation of Δ and such that $\mathcal{M}_{\theta_2}(x) < 0.5$

245

246 To prove the hardness of the problem, we show a reduction from a simple
 247 variant of 3-SAT, which we refer to as 3-NAF-SAT.

3-NOTALLFALSE-SAT (3-NAF-SAT)

Input: a 3-CNF ϕ such that the assignment of all false values is not satisfying, i.e., $\phi(\text{false}, \text{false}, \dots, \text{false}) = \text{false}$.

Output: yes \iff there exists an assignment \mathbf{a} such that $\phi(\mathbf{a}) = \text{true}$.

248

249 The NP -completeness of 3-NAF-SAT immediately follows from the NP-
 250 completeness of 3-SAT. We provide a proof of this fact for the sake of self-
 251 containment of the paper.

252 **Theorem 1.** 3-NAF-SAT is NP -complete.

253 *Proof.* We show a reduction from 3-SAT. Let ψ be a 3-CNF formula over
 254 n variables x_1, \dots, x_n . Consider the 3-CNF formula ϕ over $n + 1$ variables
 255 defined by $\phi(x_1, \dots, x_n, x_{n+1}) = \psi(x_1, \dots, x_n) \wedge (x_{n+1} \vee \neg x_{n+1} \vee \neg x_{n+1})$. Clearly
 256 for the assignment \mathbf{a} such that $a_i = \text{false}$ for each $i = 1, \dots, n + 1$ we
 257 have $\phi(\mathbf{a}) = \text{false}$, hence ϕ is a proper instance of 3-NAF-SAT, which is
 258 obtainable in polynomial time from the instance ψ of 3-SAT. Moreover, $\mathbf{a} =$
 259 $(a_1, \dots, a_n, a_{n+1})$ is a satisfying assignment for ϕ if and only if $a_{n+1} = \text{true}$
 260 and $\psi(a_1, \dots, a_n) = \text{true}$, i.e., if and only if a_1, \dots, a_n is satisfying for ψ . \square

261 **Theorem 2.** Deciding DRP is NP -complete.

262 *Proof.* The inclusion of DRP in NP is trivial. A certificate is a \mathcal{M}_{θ_2} which is
 263 of the same size as \mathcal{M}_{θ_1} , hence polynomial in the input size. The verification
 264 of such a certificate, consists of a forward propagations of x through \mathcal{M}_{θ_2} in
 265 order to check that $\mathcal{M}_{\theta_2}(x) < 0.5$. This is clearly doable in time polynomial
 266 in the size of the classifier, i.e., polynomial in the input.

267 For the hardness of DRP we show a reduction from 3-NAF-SAT. In
 268 particular, we show that there is a $\delta \in (0, 1]$ such that, given a 3-CNF
 269 formula ϕ , not satisfied by the all-false assignment, we can construct an INN
 270 \mathcal{I} whose edge intervals are all of the width 2δ and an input x such that

- 271 1. for \mathcal{M}_{θ_1} being the DNN with the same topology of \mathcal{I} and such that for
 272 each edge e the weight w_e is taken as the central point of the interval
 273 assigned to e in \mathcal{I} , we have $\mathcal{M}_{\theta_1}(x) \geq 0.5$;
- 274 2. ϕ is satisfiable if and only if there exists another DNN \mathcal{M}_{θ_2} which is
 275 also a realisation of \mathcal{I} and such that $\mathcal{M}_{\theta_2}(x) < 0.5$.

276 Note that we are using the INN \mathcal{I} to represent both the parametric clas-
 277 sifier \mathcal{M}_{Θ} and the set of PMC Δ , consisting of all the possible DNN being
 278 a realisation of \mathcal{I} .

279 We start by analysing several gadgets that will be used as building blocks
 280 of \mathcal{I} . These gadgets are shown in Figures 2, 3, 4, 5.

281 Lemmas 3-8 provide the key properties of such gadgets which will be used
 282 in the reduction. The parameter δ is a number in $(0, 1)$ whose value will be
 283 fixed by the analysis.

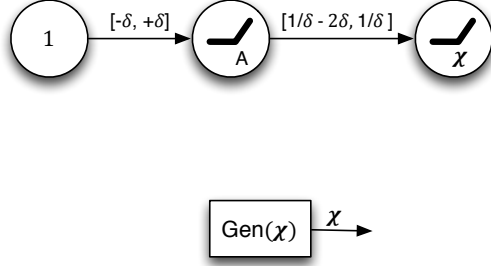


Figure 2: Generating-gadget. The input to this gadget is the constant 1 represented by the leftmost node. The output is the value χ computed in the rightmost node, that depends on the weights chosen in the intervals on the two edges.

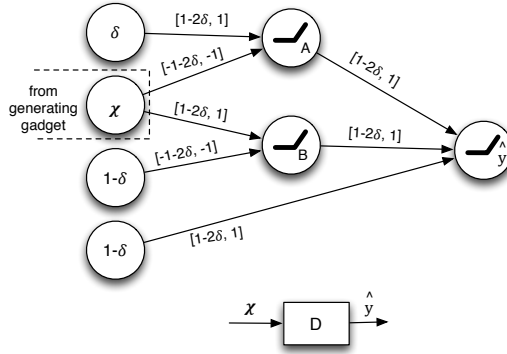


Figure 3: Discretizer-gadget. The only non-constant input is the value computed in the node labelled χ . The output is the value computed in the node labelled \hat{y} .

284 **Lemma 3** (Generating-gadget). *The value χ computed in the leftmost node*
 285 *of the Generating-gadget in Fig. 2, satisfies $\chi \in [0, 1]$.*

286 *Proof.* The value computed by the first node satisfies $A \in [0, \delta]$. Hence, since
 287 $\chi = \max\{0, A \cdot w\}$ with $w \in [\frac{1}{\delta} - 2\delta, \frac{1}{\delta}]$ we have $\chi \in [0, 1]$. \square

288 **Lemma 4** (Discretizer-gadget). *Consider the Discretizer-gadget in figure 3*
 289 *with χ being the the leftmost node of a Generating-gadget, i.e., the corre-*
 290 *sponding value satisfies $\chi \in [0, 1]$. Then, for the value \hat{y} the following holds:*

- 291 1. *if $\hat{y} > 1 - \delta$ then $\chi \in [0, \delta] \cup [1 - \delta, 1]$;*

- 292 2. if $\chi \in \{0, 1\}$ then there are possible choices of the weights yielding
 293 $\hat{y} = 1$;
 294 3. if $\hat{y} \neq 1$ then $\chi \notin \{0, 1\}$.
 295 4. $\hat{y} \in [0, 1]$.

296 *Proof.* The claims are a direct consequence of the following observations (re-
 297 fer to Fig. 3 for the notation):

- 298 (a) if $\chi \in (\delta, 1 - \delta)$ the $A = 0$ and $B = 0$, hence $\hat{y} \in [(1 - \delta)(1 - 2\delta), 1 - \delta]$.
 (b) if $0 < \chi \leq \delta$ then $B = 0$ and $A \in [\max\{0, \delta(1 - 2\delta) - \delta(1 + 2\delta)\}, \delta - \chi] \subseteq [0, \delta]$. Hence

$$\hat{y} \in [(1 - \delta)(1 - 2\delta), (1 - \delta) + (\delta - \chi)] \subseteq [(1 - \delta)(1 - 2\delta), 1).$$

- (c) if $(1 - \delta) \leq \chi < 1$ then $A = 0$ and $B \in [\max\{0, (1 - \delta)(1 - 2\delta) - (1 - \delta)(1 + 2\delta)\}, \chi - (1 - \delta)] \subseteq [0, \delta]$. Hence,

$$\hat{y} \in [(1 - \delta)(1 - 2\delta), (1 - \delta) + \chi - (1 - \delta)] \subseteq [(1 - \delta)(1 - 2\delta), 1).$$

- (d) if $\chi = 0$ then $B = 0$ and $A \in [\delta(1 - 2\delta), \delta]$, hence

$$\hat{y} \in [\delta(1 - 2\delta)^2 + (1 - \delta)(1 - 2\delta), 1].$$

299 In particular, for the realisations of \mathcal{I} where the weights on the topmost
 300 edges and on the bottommost edge are chosen to be 1 we have $\hat{y} = 1$.

- (e) if $\chi = 1$ then $A = 0$ and $B \in [\max\{0, (1 - 2\delta) - (1 + 2\delta)(1 - \delta)\}, \delta] = [0, \delta]$,
 hence

$$\hat{y} \in [(1 - \delta)(1 - 2\delta), (1 - \delta) + \delta] = [(1 - \delta)(1 - 2\delta), 1].$$

301 In particular, for the realisations of \mathcal{I} where all the weights on the edges
 302 are chosen to be the maximum possible value, we have $\hat{y} = 1$.

303 Item 1 in the statement follows directly from (a). Item 2 in the statement
 304 follows from (d) and (e). Item 3 in the statement follows from (b) and (c).
 305 Finally, Item 4 follows from the (a)-(e). \square

306 **Lemma 5** (Negation-gadget). *With reference to the Negation-gadget in Fig.4,*
 307 *for any $0 < \delta < \frac{\sqrt{6}}{2} - 1$, and $\chi \in [0, \delta] \cup [1 - \delta, 1]$, the following holds:*

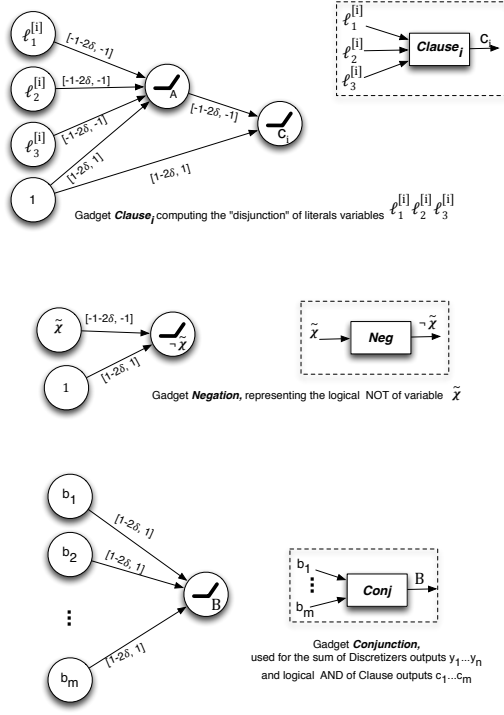


Figure 4: LOGICALPORTS

- 308 1. $\neg\chi \in [0, \delta] \iff \chi \in [1 - \delta, 1]$;
- 309 2. $\neg\chi \in [1 - 3\delta - 2\delta^2, 1] \iff \chi \in [0, \delta]$.
- 310 3. *if $\chi \in \{0, 1\}$ then there is a choice of the weights of the Negation-gadget*
- 311 *such that $\neg\chi = 1 - \chi$. In other words, if the input value is binary, then*
- 312 *there is a choice of the weights such that the Negation-gadget computes*
- 313 *the boolean NOT of the input χ .*

314 *Proof.* We have $\neg\chi = \max\{0, w_2 - |w_1|\chi\}$ where $w_1 \in [-1 - 2\delta]$ and $w_2 \in$

315 $[1 - 2\delta, 1]$. Therefore

(i) if $\chi \in [0, \delta]$ it follows that

$$\neg\chi \in [(1 - 2\delta) - \delta(1 + 2\delta), 1 - 0] = [1 - 3\delta - 2\delta^2, 1];$$

(ii) if $\chi \in [1 - \delta, 1]$ then

$$\neg\chi \in [\max\{(1 - 2\delta) - (1 + 2\delta), 0\}, 1 - (1 - \delta)] = [0, \delta].$$

Moreover, because of the hypothesis $0 < \delta < \frac{\sqrt{6}}{2} - 1$, we have that $1 - 3\delta - 2\delta^2 > \delta$ and

$$[1 - 3\delta - 2\delta^2, 1] \cap [0, \delta] = \emptyset,$$

316 which implies that the implications hold also in the opposite direction.

317 For the third claim of the lemma, it is enough to consider the weight in
318 the associated interval for each edge whose absolute value is equal to 1. \square

319 **Lemma 6** (Clause-gadget). *For the Clause-gadget in Fig.4, the following*
320 *holds: Let $0 < \delta \leq \frac{1}{12}$, and for each $t = 1, 2, 3$, let $\ell_t^{[i]} \in [0, \delta] \cup [1 - 3\delta - 2\delta^2, 1]^2$*
321 *We have that,*

- 322 1. $c_i \in [0, 5\delta + 6\delta^2]$ if and only if for all $t = 1, 2, 3$, $\ell_t^{[i]} \in [0, \delta]$;
- 323 2. $c_i \in [1 - 5\delta - 8\delta^2 - 4\delta^3, 1]$ if and only if there is $t \in \{1, 2, 3\}$ such that
324 $\ell_t^{[i]} \in [1 - 3\delta - 2\delta^2, 1]$.
- 325 3. if for each $t = 1, 2, 3$, $\ell_t^{[i]} \in \{0, 1\}$ then there is a choice of the weights
326 of the Clause-gadget such that $c_i \in \{0, 1\}$ and $c_i = 0$ if and only if
327 $\ell_1^{[i]} = \ell_2^{[i]} = \ell_3^{[i]} = 0$. In other words, if the input values are binary, then
328 there is a choice of the weights such that the Clause-gadget computes
329 the boolean OR of the inputs $\ell_t^{[i]}$.

330 *Proof.* Consider a realisation of the Clause-gadget. Let us denote by w_t^L the
331 weight taken from the interval on the edge connecting $\ell_t^{[i]}$ to A . Moreover,
332 let w_1 denote the weight taken from the interval on the edge connecting the
333 fixed value node 1 to A . Let w_2 be the weight taken from the interval on
334 the edge connecting the fixed value node 1 to the output node of the gadget.
335 Finally, let w_A be the weight taken from the interval associated with the edge
336 connecting the node A to the output node of the gadget.

- (i) If for all $t = 1, 2, 3$, it holds that $\ell_t^{[i]} \in [0, \delta]$ then, using $A = \max\{0, w_1 - \sum_{t=1}^3 |w_t^L| \ell_t^{[i]}\}$, we have that $A \in [\max\{0, 1 - 2\delta - 3\delta(1 + 2\delta)\}, 1] \subseteq [1 - 5\delta - 6\delta^2, 1]$. Since $c_i = \max\{0, w_2 - |w_A| \cdot A\}$, we have

$$c_i \in [\max\{0, (1 - 2\delta - (1 + 2\delta))\}, 1 - (1 - 5\delta - 6\delta^2)] \subseteq [0, 5\delta - 6\delta^2].$$

337 This shows the sufficiency of the condition in the first item of the
338 statement.

²Note that this corresponds to the case when $\ell_t^{[i]}$ is either the output $\neg\chi$ of a Negation-gadget or the output χ of a Generating-gadget such that $\hat{y} > 1 - \delta$

(ii) Assume there exists $\hat{t} \in \{1, 2, 3\}$ such that $\ell_{\hat{t}}^{[i]} \in [1 - 3\delta - 2\delta^2, 1]$. Then

$$A = \max\{0, w_1 - |w_{\hat{t}}^L| \ell_{\hat{t}}^{[i]} + \sum_{t \neq \hat{t}} |w_t^L| \ell_t^{[i]}\}.$$

It follows that

$$A \geq \max\{0, (1 - 2\delta) - (1 + 2\delta) \cdot 1 - 2(1 + 2\delta) \cdot 1\} = 0,$$

and

$$A \leq \max\{0, 1 - 1 \cdot (1 - 3\delta - 2\delta^2) - 2 \cdot (1) \cdot (0)\} = \max\{0, 1 - (1 - 3\delta - 2\delta^2)\} = 3\delta + 2\delta^2.$$

Hence $A \in [0, 3\delta + 2\delta^2]$. Therefore,

$$c_i \in [\max\{0, (1 - 2\delta) - (1 + 2\delta)(3\delta + 2\delta^2)\}, \max\{0, 1 - (1) \cdot 0\}] = [1 - 5\delta - 8\delta^2 - 4\delta^3, 1].$$

339 This shows the sufficiency of the condition in the second item of the
340 statement.

Finally, because of the assumption $\delta \leq \frac{1}{12}$ we have that $5\delta - 6\delta^2 < 1 - 5\delta - 8\delta^2 - 4\delta^3$ hence

$$[0, 5\delta - 6\delta^2] \cap [1 - 5\delta - 8\delta^2 - 4\delta^3, 1] = \emptyset,$$

341 which implies that the conditions in both items of the statement are also
342 necessary.

343 For the third claim of the lemma, it is enough to consider the weight in
344 the associated interval for each edge whose absolute value is equal to 1. \square

345 **Lemma 7** (Conjunction-gadget). *Consider the Conjunction-gadget in Fig.4.*

346 *Let \tilde{n} denote the output of a Conjunction-gadget whose inputs are the*
347 *values $\hat{y}_1, \dots, \hat{y}_n$ output by the n Discretizer-gadgets, with input χ_1, \dots, χ_n ,*
348 *respectively, such that $\chi_j \in [0, 1]$.*

349 *Let \tilde{c} denote the output of a Conjunction-gadget whose inputs are values*
350 *c_1, \dots, c_m . Assume also that for each $i = 1, \dots, m$, c_i is the output of a*
351 *Clause-gadget whose inputs are either the output χ_j of a Generating-gadget*
352 *or the output of a Negation-gadget whose input is the output of a Generating-*
353 *gadget.*

- 354 1. If $\tilde{n} > n - \delta$ then for each $i = 1, \dots, m$ it holds that $\hat{y}_i \in (1 - \delta, 1]$, and
 355 $\chi_i \in [0, \delta] \cup [1 - \delta, 1]$.
 356 2. If $\tilde{c} > m - (5\delta + 8\delta^2 + 4\delta^3)$ then for each $i = 1, \dots, m$ it holds that
 357 $c_i \in (1 - (5\delta + 8\delta^2 + 4\delta^3), 1]$.

358 *Proof.* The first claim follows from Lemma 4. In particular, by Lemma 4,
 359 the condition on the values χ_j implies $\hat{y}_i \in [0, 1]$. Moreover, by $\tilde{n} > n - \delta$, it
 360 follows that for each i we have $\hat{y}_i > 1 - \delta$. Again, by Lemma 4, this implies
 361 that $\chi_i \in [0, \delta] \cup [1 - \delta, 1]$.

362 For the second claim, we first observe that the hypotheses on the Clause-
 363 gadget whose outputs are the values c_1, \dots, c_m , imply that the input to such
 364 gadgets satisfies the hypotheses of Lemma 6. Therefore, for each $i = 1, \dots, m$,
 365 it holds that $c_i \in [0, 5\delta + 6\delta^2] \cup [1 - 5\delta - 8\delta^2 - 4\delta^3, 1]$. It follows that, if $\tilde{c} >$
 366 $m - 5\delta - 8\delta^2 - 4\delta^3$ for each $i = 1, \dots, m$, it holds that $c_i > 1 - 5\delta - 8\delta^2 - 4\delta^3$. \square

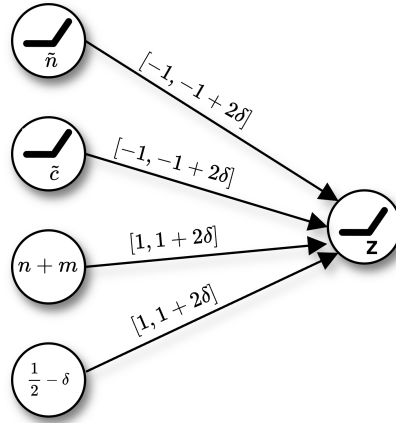


Figure 5: End-gadget

367 **Lemma 8** (End-gadget). *Consider the End-gadget in Fig. 5). For any choice*
 368 *of edge weights, it holds that if $z < 1/2$ then*

- 369 • $\tilde{n} > n - \delta$;
 370 • $\tilde{c} > m - (5\delta + 8\delta^2 + 4\delta^3)$.

Proof. We have that

$$z \geq \max\{0, n + m + \frac{1}{2} - \delta - \tilde{n} - \tilde{c}\}.$$

371 We show that if one of the inequalities in the statement is violated, then
 372 $z \geq 1/2$.

373 Suppose that $\tilde{n} \leq n - \delta$. Then, since $\tilde{c} \leq m$, it follows that $z \geq n + m +$
 374 $1/2 - \delta - n + \delta - m = 1/2$.

375 Suppose now that $\tilde{c} \leq m - (5\delta + 8\delta^2 + 4\delta^3) \leq m - \delta = m - \delta$. Then, since
 376 $\tilde{n} \leq n$, it follows that $z \geq n + m + 1/2 - \delta - n - m + \delta = 1/2$. \square

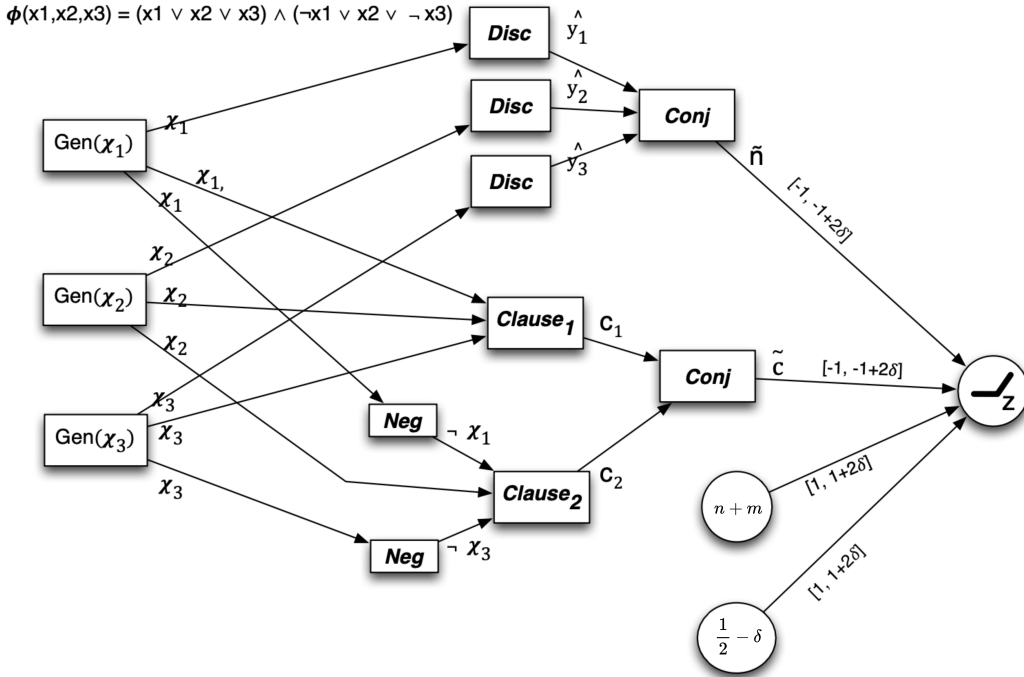


Figure 6: A complete example of the reduction on a simple formula, with $n = 3$ variables and $m = 2$ clauses. All the interval weights not explicitly given are $[1 - 2\delta, 1]$

377 **The reduction** $\mathcal{R} : \phi \mapsto I^\phi = (\mathcal{M}_{\theta_1}^\phi, x^\phi, \Delta^\phi)$. Fix a 3-CNF $\phi(x_1, \dots, x_n)$,
 378 such that $\phi(\mathbf{a}) = false$, for the assignment $\mathbf{a} = (false, \dots, false)$. Fix a

379 positive rational number $\delta \leq \frac{1}{12}$. Consider the INN $\mathcal{I} = \mathcal{I}^\phi$ built as follows
 380 (refer to Fig. 6 for an example of this construction):

- 381 1. For each variable x_i add to the network a copy of the *Generating-gadget*
 382 (and refer to it as $Gen(\chi_i)$) and a copy of the *Discretizer-gadget* (and
 383 refer to it as $Disc_i$).
- 384 2. For each $i = 1, \dots, n$, connect $Gen(\chi_i)$ to $Disc_i$ by identifying the
 385 output node χ_i of $Gen(\chi_i)$ with the non-constant input node χ of $Disc_i$.
 386 Refer to the output node/value of $Disc_i$ as \hat{y}_i (see also Fig. 3).
- 387 3. For each clause $C_j = (\lambda_1^{(j)} \vee \lambda_2^{(j)} \vee \lambda_3^{(j)})$ ($j = 1, \dots, m$) of ϕ add a
 388 *Clause-gadget*, henceforth referred to as $Clause_j$. For each $t = 1, 2, 3$,
 - 389 • if $\lambda_t^{(j)}$ corresponds to the positive variable x_i then create a connec-
 390 tion so that the input of $Clause_j$ that is labelled $\ell_t^{(j)}$ is the output
 391 χ_i of the generating $Gen(\chi_i)$ associated to x_i .
 - 392 • if $\lambda_t^{(j)}$ corresponds to the negated variable $\neg x_i$ then make a connec-
 393 tion so that the input of $Clause_j$ that is labelled $\ell_t^{(j)}$ is the output
 394 of a negation gadget, and the input of such negation gadget is the
 395 output χ_i of the Generating-gadget $Gen(\chi_i)$.
- 396 4. Add a conjunction gadget such that its inputs are the outputs \hat{y}_i ($i =$
 397 $1, \dots, n$) of the *Discretizer-gadgets*. Let \tilde{n} denote the output of such
 398 conjunction gadget.
- 399 5. Add a conjunction gadget such that its inputs are the outputs c_j
 400 ($i = 1, \dots, m$) of the *Clause-gadgets*. Let \tilde{c} denote the output of such
 401 conjunction gadget.
- 402 6. Finally, add an *End-gadget* (Fig. 5) and connect it to the rest of the
 403 network by making the output \tilde{n}, \tilde{c} of the above conjunction gadgets
 404 (defined in items 4, 5) coincide with the *End-gadget* inputs marked
 405 with \tilde{n} and \tilde{c} , respectively.

406 The above construction defines the topology of the DNN, representing the
 407 parametric classifier \mathcal{M}_Θ . The classifier $\mathcal{M}_{\theta_1}^\phi$ is chosen to be the realisation
 408 of \mathcal{I}^ϕ obtained by setting the weight on each edge e to the middle point of
 409 the interval associated to e . Such a classifier takes as input x a vector whose
 410 components are

- 411 • the value in the leftmost node of each Generating-gadget. In the input
 412 x^ϕ defined for our reduction these values are set to 1 as in Fig.7 ;

- 413 • the values in the first (top) and the two last (bottom) nodes in each
414 Discretizer-gadget. In the input x^ϕ defined for our reduction these
415 values are set to δ , $(1 - \delta)$, and $(1 - \delta)$, respectively as in Fig.3;
- 416 • the values in the lowest node of each Clause-gadget. In the input x^ϕ
417 defined for our reduction these values are set to 1 as in Fig.4;
- 418 • the values in the lowest node of each Negation-gadget. In the input x^ϕ
419 defined for our reduction these values are set to 1 as in Fig.4;
- 420 • the values in two bottom nodes of the End-gadget. In the input x^ϕ
421 defined for our reduction these values are set to $n + m$ and $\frac{1}{2} - \delta$,
422 respectively, as in Fig.5.

423 Finally Δ^ϕ is defined as the set of realisations of \mathcal{I} .

424 It is easy to see that by fixing the value δ so that it can be encoded
425 by number of bits polynomial in the size of ϕ , **the instance I^ϕ can be**
426 **constructed from ϕ in polynomial time**, since each gadget has a constant
427 size, the number of gadgets is polynomial in the size of the formula, and the
428 input vector x can be described by a number of bits polynomial in the size
429 of δ and the size of \mathcal{I}^ϕ .

430 We first prove a lemma that characterises realisations of \mathcal{I} such that the
431 output of each χ_i is binary.

432 **Lemma 9.** *The following two claims characterise the realisations of \mathcal{I} such*
433 *that for each $i = 1, \dots, n$, it holds that $\chi_i \in \{0, 1\}$.*

- 434 1. *Fix a truth assignment \mathbf{a} such that $\phi(\mathbf{a}) = \text{false}$. For any realisation*
435 *\mathcal{M}_θ of \mathcal{I} such that for each $i = 1, \dots, n$ it holds that $\chi_i = 0$ if $a_i = \text{false}$*
436 *and $\chi_i = 1$ if $a_i = \text{true}$ it holds that $\mathcal{M}_\theta(x) \geq \frac{1}{2}$.*
- 437 2. *Fix a truth assignment \mathbf{a} such that $\phi(\mathbf{a}) = \text{true}$. Then, there exists a*
438 *realisation \mathcal{M}_θ of \mathcal{I} such that for each $i = 1, \dots, n$ it holds that $\chi_i = 0$*
439 *if $a_i = \text{false}$ and $\chi_i = 1$ if $a_i = \text{true}$, and $\mathcal{M}_\theta(x) < \frac{1}{2}$.*

440 *Proof.* We show the two claims separately.

- 441 1. For the first claim, we observe that
442 (a) $\tilde{n} \leq n$.

443 (b) Since $\phi(\mathbf{a}) = false$, there exists $i \in [m]$ such that the assignment
 444 \mathbf{a} makes all the literals in the i th clause to be false. We also have
 445 that the values $\ell_t^{[i]}$ s, input to the clause gadget encoding the i th
 446 clause, will satisfy $\ell_t^{[i]} \in [0, \delta]$ —in particular, we have $\ell_t^{[i]} = \chi_j = 0$
 447 if the literal corresponds to some variable $x_j = false$; and, if the
 448 the literal correspond to the negation of some variable $x_j = true$,
 449 hence $\ell_t^{[i]} = \neg\chi_j$ with $\chi_j = 1$ and by Lemma 5 $\neg\chi_j \in [0, \delta]$.

Therefore, by Lemma 6, we have $c_i \in [0, 5\delta + 6\delta^2]$. It follows that
 $\tilde{c} < m - \delta$ and

$$z \geq \max\{0, \frac{1}{2} - \delta + n + m - \tilde{n} - \tilde{c}\} \geq \frac{1}{2} - \delta + n + m - n - m + \delta = \frac{1}{2}$$

450 2. For the second claim, consider the realisation obtained by setting the
 451 weights as follows:

- 452 • in the i -th generating gadget (the one associated to x_i) the weights
 453 are chosen in order to have output $\chi = 1$ if $a_i = true$ and $\chi = 0$
 454 if $a_i = false$;
- 455 • in all the other gadgets, the weights are set to the value w such
 456 that $|w| = 1$.

457 Because of the correspondence $a_i = true \rightarrow \chi_i = 1$ and $a_i = false \rightarrow$
 458 $\chi_i = 0$, by Lemma 4, we have $\hat{g}_i = 1$ for each $i = 1, \dots, n$. Hence $\tilde{n} = n$.
 459 Because of the choice of the weights being all of the absolute value one,
 460 it is also easy to see that, interpreting $true$ as 1 and $false$ as 0, for each
 461 $i = 1, \dots, m$ and $t = 1, 2, 3$, we have an exact correspondence between
 462 the truth value assigned by \mathbf{a} to the t th literal of the i th clause and
 463 the value $\ell_t^{[i]}$. Hence, by the assumption that $\phi(\mathbf{a}) = true$, we also have
 464 that $c_i = 1$ for each $i = 1, \dots, m$. It follows that $\tilde{c} = m$.

Therefore,

$$z = \max\{0, \frac{1}{2} - \delta + n + m - \tilde{n} - \tilde{c}\} = \frac{1}{2} - \delta < \frac{1}{2}$$

465

□

466 Let \mathcal{M}_{θ_1} be the realisation of \mathcal{I} obtained by setting the weight on each
 467 edge e to the middle point of the interval associated to e . Let \mathbf{a} be the
 468 assignment for ϕ such that $a_i = false$ for each $i = 1, \dots, n$. Therefore, \mathcal{M}_{θ_1}

469 coincides with the realisation of \mathcal{I} such that for each $i = 1, \dots, n$ it holds
 470 that $\chi_i = 0$ and the hypotheses of Lemma 9 are satisfied Hence, by Lemma
 471 9, it holds that $\mathcal{M}_\theta(x) \geq \frac{1}{2}$. We have shown the following.

472 **Lemma 10.** *For each instance ϕ of 3-NAF-SAT, the reduction \mathcal{R} produces*
 473 *in polynomial time a proper instance $(\mathcal{M}_{\theta_1}, \mathcal{I}, x)$ of DRP.*

474 In order to complete the proof of the Theorem, the following remains to
 475 be shown.

476 **Lemma 11.** *The formula ϕ is satisfiable if and only if there is a realisation*
 477 *\mathcal{M}_{θ_2} of \mathcal{I} , such that $\mathcal{M}_{\theta_2}(x) < \frac{1}{2}$.*

478 *Proof.* The sufficiency of the condition directly follows from the second claim
 479 of Lemma 9, which shows that: *If there exists an assignment \mathbf{a} such that*
 480 *$\phi(\mathbf{a}) = true$ then there is a realisation \mathcal{M}_{θ_2} of \mathcal{I} , such that $\mathcal{M}_{\theta_2}(x) < \frac{1}{2}$.*

481 Let us now focus on the other direction. Assume that there is a realisation
 482 \mathcal{M}_{θ_2} such that $\mathcal{M}_{\theta_2}(x) < \frac{1}{2}$. By Lemma 8, it follows that for the realisation
 483 $\mathcal{M}_{\theta_2}(x)$ it holds that $\tilde{n} > n - \delta$ and $\tilde{c} > m - (5\delta + 8\delta^2 + 4\delta^3)$.

484 Then, by Lemma 7 it follows that

- 485 1. for each $i = 1, \dots, n$ it holds that $\hat{y}_i \in (1 - \delta, 1]$, and $\chi_i \in [0, \delta] \cup [1 - \delta, 1]$;
- 486 2. for each $i = 1, \dots, m$ it holds that $c_i \in (1 - (5\delta + 8\delta^2 + 4\delta^3), 1]$.

487 These two conditions together with $\delta < 1/12$ imply, by Lemmas 4, 5 and 6,
 488 that for each $i = 1, \dots, m$, there is $t \in \{1, 2, 3\}$ such that one of the following
 489 holds

- 490 1. $\ell_t^{[i]} \in [1 - 3\delta - 2\delta^2, 1]$ and $\ell_t^{[i]}$ coincides with some output $\neg\chi_j$ of a
 491 negation-gadget, and the input value satisfies $\chi_j \in [0, \delta]$; moreover by
 492 construction, then literal $\lambda_t^{[i]}$ is $\neg x_j$;
- 493 2. $\ell_t^{[i]} \in [1 - \delta, 1]$ and $\ell_t^{[i]}$ coincides with some output χ_j of a generating-
 494 gadget, whence by construction, then literal $\lambda_t^{[i]}$ is x_j

Let $\mathbf{a} = (a_1, \dots, a_n)$ be the truth assignment defined by

$$a_i = \begin{cases} true & \text{if } \chi_j \geq 1 - \delta \\ false & \text{if } \chi_j \leq \delta. \end{cases}$$

495 Then, for each clause $i = 1, \dots, m$ at least one literal is set to *true*, and
 496 $\phi(\mathbf{a}) = true$. □

497 The proof is complete. □

498 From Theorem 2, it follows that deciding whether a CFX x' is not robust
 499 to a set of PMC Δ is NP-complete. We now show that the above reduction
 500 can be used to prove the NP-completeness of deciding robustness with respect
 501 to a set Δ of NOMC models. In particular, we have the following:

502 **Theorem 12** (Hardness of DRP for NOMC). *Given classifier \mathcal{M}_{θ_1} , an input*
 503 *x and a set Δ of NOMC, deciding whether $\exists \mathcal{M}_{\theta_2} \in \Delta$ s.t $\mathbb{E}[\mathcal{M}_{\theta_2}(x)] < \frac{1}{2} \leq$*
 504 *$\mathcal{M}_{\theta_1}(x)$ is NP-complete.*

505 *Proof.* The proof of the inclusion in NP is analogous to the one of DRP for
 506 PMC models.

507 For the hardness, we use again a reduction from 3-NAF-SAT: given a
 508 3-CNF ϕ that is not satisfied by the all-false assignment, build an interval
 509 neural network exactly like in Theorem 2 but for one difference consisting in
 510 the interval of weights on the first edge of the Generating-gadget, which are
 511 now set to $[0, 2\delta]$ as in Fig. 7.

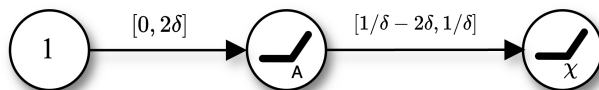


Figure 7: *Generating-gadget* used in this proof.

512 We denote by \mathcal{I}_{NOMC} this interval neural network. obtained from this
 513 reduction starting from a 3-CNF ϕ . Note that the new generating-gadget
 514 can also produce any value in $[0, 1]$. More generally, we have the following
 515 important remark.

516 **Remark 1.** *Lemmas 3-8 also hold for the interval neural network \mathcal{I}_{NOMC} .*

517 We define \mathcal{M}_{θ_1} to be the realisation of \mathcal{I}_{NOMC} obtained by setting the
 518 weight on each edge e to the middle point of the interval associated with e .
 519 The input value x is defined as in the proof of Theorem 2. We also let Δ to
 520 be the set of realisation of \mathcal{I}_{NOMC} .

521 From the above remark and the proof of Theorem 2, it follows that the
 522 3-CNF ϕ has a satisfying assignment if and only if there exists a realisation
 523 \mathcal{M}_{θ_2} in Δ such that $\mathcal{M}_{\theta_2}(x) < 0.5$.

524 Then, in order to complete the proof, we only need to show that $\mathcal{M}_{\theta_1}(x) \geq$
525 $\frac{1}{2}$ and Δ properly defines a set of NOMC for \mathcal{M}_{θ_1} (Def. 5), which we recall
526 here for readability purposes:

- 527 1. $\mathbb{E}[\mathcal{M}_\theta(x)] = \mathcal{M}_{\theta_1}(x)$; where the expectation is over the randomness³ of
528 \mathcal{M}_θ ;
- 529 2. $\text{Var}[\mathcal{M}_\theta(x)] = \nu_x$, where ν_x represents the maximum variance of the
530 prediction of $\mathcal{M}_\theta(x)$, and whenever x lies on the data manifold \mathcal{X} , ν_x
531 is upper bounded by a small constant ν ;
- 532 3. If \mathcal{M}_{θ_1} is Lipschitz continuous for some γ_1 , then \mathcal{M}_θ is also Lipschitz
533 continuous for some γ_2 .

534 For a (random or fixed) realisation \mathcal{M}_θ of \mathcal{I}_{NOMC} and a node ν , let
535 us denote by ν_θ the value computed in the node ν by \mathcal{M}_θ on input x . In
536 accordance with the analysis in Theorem 2, we assume $\delta = 1/12$.

537 To show that $\mathcal{M}_{\theta_1}(x) \geq \frac{1}{2}$ and that Δ satisfies property 1 for being a set
538 of NOMC, we prepare the following.

539 **Lemma 13.** *Let \mathcal{M}_θ be a random realisation of \mathcal{I}_{NOMC} . It holds that*

- 540 1. *for the output node χ of each Generating-gadget, we have $\mathbb{E}[\chi_\theta] = \chi_{\theta_1} =$
541 $\frac{1}{2} - \delta^2$.*
- 542 2. *for the the output node $\neg\chi$ of each Negation-gadget we have $\mathbb{E}[\neg\chi] =$
543 $\neg\chi_{\theta_1} = \frac{1}{2} - \frac{3}{2}\delta + \delta^2 + \delta^3$.*
- 544 3. *for the the output node \hat{y} of each Discretizing-gadget we have $\mathbb{E}[\hat{y}_\theta] =$
545 $\hat{y}_{\theta_1} = 1 - \delta$*
- 546 4. *for the output node c of each Clause-gadget we have $\mathbb{E}[c_\theta] = c_{\theta_1} = 1 - \delta$*
- 547 5. *for the output node \tilde{c} of the Conjunction-gadget collecting the outputs
548 of the Clause-gadgets we have $\mathbb{E}[\tilde{c}_\theta] = \tilde{c}_{\theta_1} = m(1 - \delta)^2$*
- 549 6. *for the output node \tilde{n} of the Conjunction-gadget collecting the outputs
550 of the Discretizing-gadgets we have $\mathbb{E}[\tilde{n}_\theta] = \tilde{n}_{\theta_1} = n(1 - \delta)^2$*
- 551 7. *for the output node z of the End-gadget, we have $\mathbb{E}[z_\theta] = z_{\theta_1} \geq \frac{1}{2}$.*

552 *Proof.*

³In our case this is a *random realisation of \mathcal{I}_{NOMC}* , i.e., a realisation of \mathcal{I}_{NOMC} obtained by independently choosing the weight on each edge sampling uniformly at random from the interval associated to that edge.

1. Let w_1, w_2 denote the weights on the edges of the Generating-gadget, respectively, in order from left to right. Because of the independence of the choices of the weights of the random realisation \mathcal{M}_θ , we have

$$\mathbb{E}[\chi_\theta] = \mathbb{E}[w_2]\mathbb{E}[A_\theta] = \mathbb{E}[w_2]\mathbb{E}[w_1] = \left(\frac{1}{\delta} - \delta\right) \delta = \frac{1}{2} - \delta^2.$$

553 It is immediate to verify that this value is equal to χ_{θ_1} .

2. Let w_1, w_2 denote the weights on the top edge and the bottom edge, respectively, of the Negation-gadget. Using 1. the independence in the choice of the weights, and the fact that with $\delta = \frac{1}{12}$ the argument of the ReLU is always non-negative, we have that

$$\mathbb{E}[-\chi_\theta] = \mathbb{E}[\chi_\theta] \cdot \mathbb{E}[w_1] + \mathbb{E}[w_2].$$

554 The claim then follows by the fact that the expected values of the
555 weights are given by the middle point of the intervals from which they
556 are respectively taken.

3. Item 1. and the first claim of Lemma 4, together with $\delta = 1/12$ imply
557 that both for a random realisation and for the realisation \mathcal{M}_{θ_1} , the
558 (expected) values computed in nodes A and B of the discretizing-gadget
559 are both 0. Hence, we have $\mathbb{E}[\hat{y}] = \mathbb{E}[w_\theta]$, where w denotes the weight
560 on the lowest edge of the gadget. By noticing that this expected value
561 is equal to the middle point of the interval, we have the desired result.
562
4. Because of 1. and 2. we have that the expected value (as well as the
563 value computed by \mathcal{M}_{θ_1} on x) of the input nodes $\ell_t^{[i]}$ of each clause-
564 gadgets are from the set $\{\frac{1}{2} - \delta^2, \frac{1}{2} - \frac{3}{2}\delta + \delta^2 + \delta^3\}$. It follows that the
565 argument of the ReLU function computed in the node A is negative.
566 Hence, we have $\mathbb{E}[c] = \mathbb{E}[w_\theta]$, where w denotes the weight on the lowest
567 edge of the gadget. Again, noticing that this expected value is equal
568 to the middle point of the interval gives the desired result.
569
5. For a realisation \mathcal{M}_θ let $w_{\theta,i}$ denote the weight on the i th edge (counting
from top to bottom) of the conjunction-gadget collecting the outputs of
the Clause-gadgets, and $c_{\theta,i}$ the output value of the i th clause gadget,
as computed by \mathcal{M}_θ on input x . Then, we have

$$\mathbb{E}[\tilde{c}_\theta] = \sum_{i=1}^m \mathbb{E}[c_{\theta,i}] \cdot \mathbb{E}[w_{\theta,i}].$$

570 The results follow from 4. and the fact that the expected value of a
 571 uniformly sampled weight is equal to the middle point of the interval
 572 from which it is taken.

- 573 6. The proof of this point is analogous to the proof of 5.
 7. For a realisation \mathcal{M}_θ let $w_{\theta,i}$ denote the weight on the i th edge (counting from top to bottom) of the End-gadget. We start by observing that from the results of the previous points, it follows that the argument of the RELU function in node z is always non-negative. Hence, we have

$$\mathbb{E}[z_\theta] = \mathbb{E}[\tilde{n}_\theta] \cdot \mathbb{E}[w_{\theta,1}] + \mathbb{E}[\tilde{c}_\theta] \cdot \mathbb{E}[w_{\theta,2}] + (n+m) \cdot \mathbb{E}[w_{\theta,3}] + \left(\frac{1}{2} - \delta\right) \cdot \mathbb{E}[w_{\theta,4}]$$

574 Then, the equality $\mathbb{E}[z_\theta] = z_{\theta_1}$ in the claim follows again from the fact
 575 that the expected values of the weights of a random realisation are
 576 equal to the middle point of the interval, i.e., the value of the weight
 577 on the edge in the realisation \mathcal{M}_{θ_1} . The inequality in the claim follows
 578 from Lemma 8 since, with $\delta = 1/12$, it holds that $n(1 - \delta)^2 < n - \delta$.

579 □

580 Claim 7 of the lemma directly implies that the first property of an NOMC
 581 is satisfied by Δ , i.e., $\mathbb{E}[\mathcal{M}_\theta(x)] = \mathcal{M}_{\theta_1}(x)$. The same claim also proves that
 582 $\mathcal{M}_{\theta_1}(x) \geq \frac{1}{2}$.

583 For the second property, namely $\text{Var}[\mathcal{I}(x)] = \nu_x$, we use the uniform
 584 continuity of the function computed by the realisations of \mathcal{I}_{NOMC} , which is a
 585 direct consequence of being linear combinations of RELU functions which are
 586 Lipschitz continuous functions, hence uniform continuous. By the Extreme
 587 Value Theorem (see, e.g., [36, Thm. 4.16]) any realisation of \mathcal{I}_{NOMC} will be
 588 bounded and achieve its minimum and maximum on the compact domain,
 589 and thus the variance will be indeed bounded.

590 Finally, the last property follows from the fact that the linear composition
 591 of Lipschitz continuous operations (ReLU) is also Lipschitz continuous, which
 592 is indeed the case of any realisation of \mathcal{I}_{NOMC} . □

593 Summarizing, we have shown the following.

594 **Corollary 14.** *Given a model \mathcal{M}_θ , a CFX x' and a set of either NOMC or*
 595 *PMC model changes Δ , the problem of verifying the Δ -robustness of x' is*
 596 *NP-complete.*

597 These hardness results motivate the introduction of novel approximate
 598 solutions to estimate the robustness of a counterfactual under a set of PMC
 599 Δ .

600 4. Probabilistic Guarantees for Existing Notions of Model Changes

601 As we have established in the previous section, exact methods for com-
 602 puting robustness under model changes are bound to lack scalability. This
 603 motivates the design of approximate and/or probabilistic approaches to solve
 604 the problem. Previous work by Hamman et al. [15] presented an approach
 605 to obtain counterfactual explanations that are probabilistically robust under
 606 NOMC. A natural question that arises then is whether guarantees obtained
 607 for NOMC also transfer to the PMC setting. As we show for the first time
 608 below, this is not the case in general.

609 **Lemma 15.** *Naturally-Occurring model changes may not be Plausible, and*
 610 *vice-versa.*

611 *Proof.* Consider the DNN \mathcal{M}_θ depicted in Fig. 8 (a) with two input nodes,
 612 one hidden layer with two ReLU nodes⁴ and one single output. The param-
 613 eters $\theta = [w_1, \dots, w_6]$ are the weights on the edges listed top-bottom and
 614 left-right.



Figure 8: (a) The model \mathcal{M}_θ used as an example to prove the lemma. (b) An interval neural network representing the realisations that can be obtained from \mathcal{M}_θ considering a set of PMC Δ_δ with $\delta = 0.3$.

Propagating an input vector $x = [x_1, x_2]^T$ through \mathcal{M}_θ , we obtain $\mathcal{M}_\theta(x) = y = w_5 \cdot \max\{0, w_1 \cdot x_1 + x_2 \cdot w_3\} + w_6 \cdot \max\{0, w_2 \cdot x_1 + x_2 \cdot w_4\}$. Now assume

⁴In this proof, we consider a DNN with only ReLU activation functions. However, we notice that it is possible to have a similar counterexample even with other activations, e.g., Tanh, Sigmoid.

an input vector $x = [0.9, 0.9]^T$ and weights $w_1 = 1, w_2 = 0, w_3 = 0, w_4 = 0.6, w_5 = 1, w_6 = -1$. The corresponding output generated by the DNN is $\mathcal{M}_\theta(x) = 0.46$. A counterfactual for x could be given as a new input vector $x' = [1, 0.8]^T$, for which we obtain $\mathcal{M}_\theta(x') = 0.52 > 0.5$. Now, following Definition 6, we consider a set of plausible model changes obtained for $\delta = 0.3$. This can be captured by defining on each weight w_i the corresponding interval in $[w_i - \delta, w_i + \delta]$ depicted in Fig. 8 (b) that represents the set of all the possible models obtained from \mathcal{M}_θ , replacing each w_i with a weight in the interval $[w_i - \delta, w_i + \delta]$. We then have that the expected result of a model $\mathcal{M}_{\theta'}$ sampled uniformly from such a set satisfies:

$$\begin{aligned}
\mathbb{E}[\mathcal{M}_{\theta'}(x')] &= \mathbb{E}[w_5] \cdot \mathbb{E}[\text{ReLU}(x_1 \cdot w_1 + x_2 \cdot w_3)] + \\
&\quad \mathbb{E}[w_6] \cdot \mathbb{E}[\text{ReLU}(x_1 \cdot w_2 + x_2 \cdot w_4)] \\
&= \mathbb{E}[[0.7, 1.3]] \cdot \mathbb{E}[\max\{0, x_1 \cdot [0.7, 1.3] + \\
&\quad x_2 \cdot [-0.3, 0.3]\}] + \mathbb{E}[-1.7, -0.3] \cdot \\
&\quad \mathbb{E}[\max\{0, x_1 \cdot [-0.3, 0.3] + x_2 \cdot [0.3, 0.9]\}] \\
&> 0.52 \neq \mathcal{M}_\theta(x')
\end{aligned}$$

615 Definition 5 states that a model change is naturally occurring if $\mathbb{E}[\mathcal{M}_{\theta'}(x)] =$
616 $\mathcal{M}_\theta(x)$. This implies that Δ contains models that cannot be characterized
617 as naturally occurring model changes. Vice versa, the existence of Naturally-
618 Occurring model changes not being plausible is implicit in the definition, and
619 for the sake of completeness, we provide an example network in Fig. 9.

620 Consider a DNN having a single input value x and a single parameter θ
621 and computing the function $\mathcal{M}_\theta(x) = \text{ReLU}(0.5 - \text{ReLU}(x - \theta))$

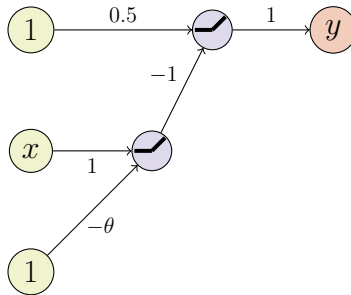


Figure 9: The DNN considered in this proof

Table 1: Empirical evaluation across model perturbations of increasing magnitude δ and different sample sizes n .

	<i>Credit</i>				<i>Spam</i>				<i>News</i>			
	$n = 1000$		$n = 10000$		$n = 1000$		$n = 10000$		$n = 1000$		$n = 10000$	
	Avg diff.	Rej. (%)	Avg diff.	Rej. (%)	Avg diff.	Rej. (%)	Avg diff.	Rej. (%)	Avg diff.	Rej. (%)	Avg diff.	Rej. (%)
$\delta = 0.05$	0.008	90	0.022	90	0.018	50	0.017	70	0.034	70	0.033	80
$\delta = 0.1$	0.017	100	0.047	100	0.034	100	0.035	100	0.064	80	0.063	100
$\delta = 0.2$	0.046	100	0.086	100	0.0748	90	0.064	100	0.127	90	0.141	100
$\delta = 0.3$	0.110	100	0.140	90	0.121	100	0.087	100	0.207	90	0.173	100

Fix a data set \mathcal{X} and let $\theta = \max_{x \in \mathcal{X}} x$. Let us consider the set of model changes $\Sigma = \{S_\tau \mid \tau \in \mathbb{R}_+\}$ defined by $S_\tau(\mathcal{M}_\theta) = \mathcal{M}_{\theta+\tau}$. Clearly for any $\tau \geq 0$, we have

$$\mathcal{M}_{\tau+\theta}(x) = \mathcal{M}_\theta(x) = 0.5,$$

622 for any $x \in \mathcal{X}$. This trivially implies that Σ is a set of naturally occurring
623 model changes (all changes considered have exactly the same value in all
624 points in \mathcal{X}).

625 The claim now follows by observing that there is no finite δ such that the
626 corresponding set of plausible model changes $\Delta = \{S \mid d_p(\mathcal{M}_\theta, S(\mathcal{M}_\theta)) \leq \delta\}$
627 contains Σ .

628 □

629 Lemma 15 shows the existence of witnesses proving that Definition 5
630 (NOMC) and Definition 6 (PMC) may capture very different model changes
631 in general. To complement this observation, we also ran experiments to
632 determine how often these definitions disagree empirically. In particular, we
633 considered three binary classification datasets commonly used in Explainable
634 AI:

- 635 • the *credit* dataset [37], which is used to predict the credit risk of a
636 person (good or bad) based on a set of attributes describing their credit
637 history;
- 638 • the *spambase* dataset [38] is used to predict whether an email is to be
639 considered spam or not based on selected attributes of the email;
- 640 • the *online news popularity* dataset [39], referred to as *news* in the fol-
641 lowing, is used to predict the popularity of online articles.

642 We trained a neural network classifier with two hidden layers (20 and 10
 643 neurons, respectively) for each dataset and used a Nearest-Neighbor Counter-
 644 factual Explainer [40] to generate counterfactual explanations for 10 different
 645 inputs. After generating a counterfactual, we produce n different perturba-
 646 tions $\mathcal{M}_{\theta'}$ of the original neural network \mathcal{M}_{θ} for $n \in \{1000, 10000\}$ under
 647 *plausible* model change with $\Delta \in \{0.05, 0.1, 0.2, 0.3\}$. We then considered
 648 two measures:

- 649 • average difference in output between \mathcal{M}_{θ} and $\mathcal{M}_{\theta'}$, for each of the n
 650 model $\mathcal{M}_{\theta'}$ and across all CFXs;
- 651 • for each counterfactual, we perform a one-sided t-test to check whether
 652 the average prediction generated by n models $\mathcal{M}_{\theta'}$ equals the original
 653 prediction of \mathcal{M}_{θ} . We report the percentage of CFXs for which the null
 654 hypothesis was rejected (p-value used 0.05).

655 Table 1 reports our results. We observe that the requirement that the
 656 expected output of perturbed models remains equal to the original predic-
 657 tion is often violated. These results complement the result of Lemma 15,
 658 confirming that the two notions indeed capture two different settings in gen-
 659 eral. In particular, our results show that (probabilistic) methods devised for
 660 NOMC may fail to guarantee robustness under PMC, thus motivating the de-
 661 velopment of dedicated approaches for probabilistic guarantees under PMC.
 662 Indeed, having clarified the relationship between the two notions of model
 663 changes, in the following, we focus on certification approaches for robust-
 664 ness under PMC, presenting a novel approximate solution with probabilistic
 665 guarantees.

666 5. Robustness under PMC with Probabilistic Guarantees

667 Jiang et al. [12, 33] proposed to use INNs to enable a compact represen-
 668 tation of a superset of the models that can be obtained by a perturbation
 669 of the starting model under a set Δ . By exploiting an exact reachable set
 670 computation method, e.g., based on MILP [41], the authors could determine
 671 whether or not a CFX is robust under the chosen Δ via a single forward
 672 propagation of the CFX. However, in view of the NP-hardness of the prob-
 673 lem discussed in the § 3 and the typical non-linear nature of the classifiers,
 674 it presents some computational limitations.

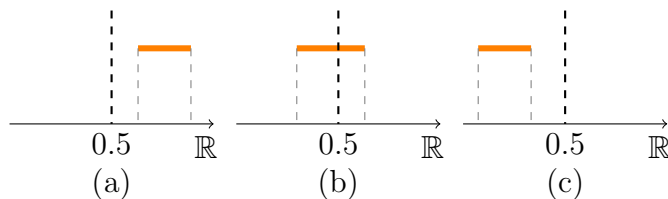


Figure 10: Visual representation of the possible output reachable set for an interval abstraction for a binary classification model. (a) For a given Δ , we classify an input as 1 (robust) if the output range for that input is always greater 0.5. Otherwise, the input is classified as 0, i.e., not robust (b),(c).

675 In general, interval neural networks map inputs to intervals representing
676 an over-approximation of all possible outcomes that can be produced by any
677 shifted model $\mathcal{M}_{\theta'}$ obtained under Δ . Given this property, if the output
678 reachable set is completely disjoint from the decision threshold 0.5, then one
679 can assert – in a sound and complete fashion – whether or not a given CFX
680 is robust (Fig. 10 (a,c)). On the other hand, if we run into a situation such as
681 the one depicted in Fig. 10 (b), one cannot assert robustness with certainty.
682 In this scenario, Jiang et al. [12] propose to classify the CFX as not robust,
683 which preserves the soundness of their result. Nonetheless, this might lead to
684 discarding a CFX even when the actual probability that after retraining, we
685 incur in plausible model changes for which the CFX is not robust is extremely
686 low. As we will show in § 6, this worst-case notion of robustness affects the
687 CFXs generated by [12], which may end up being unnecessarily expensive
688 (in terms of proximity) and having low plausibility. Additionally, computing
689 the exact output reachable set of an interval abstraction may be costly (e.g.,
690 MILP is known to be NP-hard). This is expected: Theorems 2 and 3 show
691 that there is no polynomial time algorithm able to return an exact estimate
692 of the fraction of plausible changes for which the CFX is robust (hence a
693 fortiori deciding whether it is Δ -robust), unless P=NP. In the following, we
694 propose a novel certification approach that aims to alleviate this problem.

695 5.1. A Provable Probabilistic Approach

696 One possible idea to avoid exact reachable set computation to determine
697 the robustness of a CFX under PMC is to use naive interval propagation.
698 Given an input CFX, we propagate this input through the network, keeping
699 track of all the possible activation values that can be obtained under Δ until

700 the output layer is reached. However, the non-linear and non-convex nature
 701 of DNNs may result in a significant overestimation of the actual reachable set,
 702 thus resulting in a spurious decision of non-robustness. In such cases, a CFX
 703 may end up being labeled as non-robust even though the CFX is actually
 704 robust. Additionally, even with exact methods, a CFX may be discarded
 705 even though the fraction of plausible model changes in Δ for which the CFX
 706 is not robust is negligible.

707 To avoid these problems, we propose an approximate certification ap-
 708 proach based on Monte-Carlo sampling that draws sample realisations di-
 709 rectly from Δ to obtain an underestimation of the space of possible classifi-
 710 cations under PMC. The idea of using a sample-based approach stems from
 711 the fact that the Δ set, representing all the plausible model changes, ab-
 712 stracts an infinite number of models to test. As testing this infinite number
 713 of models may be impossible in practice, efficient sampling-based solutions
 714 hold great promise. In detail, given a CFX x' we can compute an underesti-
 715 mation of the output reachable set under Δ by sampling n random realisa-
 716 tions $\mathcal{M}_{\theta_1}, \dots, \mathcal{M}_{\theta_n}$ from Δ , and compute the output reachable set by taking,
 717 respectively, the $\min_i \mathcal{M}_{\theta_i}(x')$ and the $\max_i \mathcal{M}_{\theta_i}(x')$ for $i \in \{1, \dots, n\}$.

718 This approach is very effective and allows us to obtain an estimate of the
 719 output reachable set without using an exact solver. Nonetheless, the number
 720 n of realisation to sample in order to achieve a good reachable set estimation
 721 remains unclear, as well as what kind of guarantees one could obtain from
 722 this approach. To answer these questions, we leverage previous results on
 723 the *statistical prediction of tolerance limits* [42]. Indeed, we observe that
 724 for each realisation \mathcal{M}_{θ_i} sampled from Δ , the resulting output of the DNN
 725 $\mathcal{M}_{\theta_i}(x')$ can be interpreted as an instantiation of a random variable X whose
 726 tolerance interval we are trying to estimate. Following this observation, we
 727 can derive a probabilistic bound on the correctness of the solution returned
 728 from n samples, using the following lemma based on [42]:

729 **Lemma 16.** *Fix an integer $n > 0$ and an approximation parameter $R \in$
 730 $(0, 1)$. Given a sample of n models $\mathcal{M}_{\theta_1}, \dots, \mathcal{M}_{\theta_n}$ from the (continuous) set
 731 of possible realisations Δ , the probability that for at least a fraction R of the
 732 models in a further possibly infinite sequence of samples $\mathcal{M}_{\theta_1}^{(2)}, \dots, \mathcal{M}_{\theta_m}^{(2)}$ from
 733 Δ we have*

$$\min_i \mathcal{M}_{\theta_i}^{(2)}(x) \geq \min_i \mathcal{M}_{\theta_i}(x) \tag{1}$$

$$(respectively \max_i \mathcal{M}_{\theta_i}^{(2)}(x) \leq \max_i \mathcal{M}_{\theta_i}(x))$$

734 is given by $\alpha = n \cdot \int_R^1 x^{n-1} dx = 1 - R^n$.

735 Informally, Lemma 16 allows us to derive the minimum number n of
736 realisations that it is enough to sample and check in order to guarantee
737 that with probability α at least a fraction R of the models in Δ satisfy
738 the robustness property. More precisely, from these n realisations, we can
739 obtain an underestimation of the reachable set of any realisation in Δ that
740 is guarantee to be correct with confidence α for at least a fraction R of a
741 possibly infinite further sample of realisations from Δ . In practice, if we set,
742 e.g. $\alpha = 0.999$ and $R = 0.995$, we can derive n as $n = \log_R(1 - \alpha) = 1378$.
743 After having selected 1378 random realisations from Δ , if the lower bound of
744 the underestimated reachable set computed as $\min_i \mathcal{M}_{\theta_i}(x')$ is greater than
745 0.5, then with probability $\alpha = 0.999$, R is a lower bound on the fraction
746 of plausible model changes in Δ for which x' is robust. In other words,
747 Lemma 16 allows us to assert with a confidence α that x' is not Δ -robust for
748 at most a fraction $(1 - R) = 0.05$ of models from Δ .

749 5.2. The *AP Δ S* Algorithm

750 Using the result of Lemma 16, we now present our approximation method
751 *AP Δ S* to generate probabilistic robustness guarantees. The procedure, shown
752 in Algorithm 1, receives as input a model \mathcal{M}_θ , a CFX x' for which robust-
753 ness guarantees are sought, and the two confidence parameters α, R . The
754 algorithm then searches for the largest δ_{max} such that, with probability α ,
755 the CFX x' is robust for at least a fraction R of the set of plausible model
756 changes $\Delta = \{S \mid d_p(\mathcal{M}_\theta, S(\mathcal{M}_\theta)) \leq \delta_{max}\}$.

757 The algorithm starts by computing the size n of a sample of realisations
758 that is sufficient to guarantee the condition in Lemma 16 (line 3). *AP Δ S* then
759 initializes a small δ_{init} and checks if x' is at least robust to a small model
760 shift. To this end, it employs `realisations`($\mathcal{M}_\theta, x', \delta, n$) which samples n
761 realisations, perturbing each model parameter by at most a factor δ and
762 checks if for each of these realisation $\mathcal{M}_{\theta_i}(x') \geq 0.5$, thus computing a ro-
763 bustness rate. If not all these realisations result in a robust outcome, thus
764 achieving a final rate not equal to 1, the algorithm discards the CFX x' as
765 non-robust (lines 6-8). Otherwise, it combines an exponential search (lines
766 9-12) and a binary search (lines 13-24) to find δ_{max} . At each step of this
767 search, the procedure checks whether for each of the n realisations from
768 $\Delta = \{S \mid d_p(\mathcal{M}_\theta, S(\mathcal{M}_\theta)) \leq \delta_{max}\}$ the condition $\mathcal{M}_{\theta_i}(x') \geq 0.5$ is verified.

Algorithm 1 Approximate Plausible Δ -Shift (**AP Δ S**)

```
1: Input: Model  $\mathcal{M}_\theta$ , set of PMC  $\Delta$ , CFX  $x'$ ,  $\alpha$ ,  $R$ ,  $\delta_{init}$ 
2: Output:  $\delta_{max}$ 
3:  $n \leftarrow \log_R(1 - \alpha)$  ▷ number of samples
4:  $rate \leftarrow \mathbf{realisations}(\mathcal{M}_\theta, x', \delta_{init}, n)$ 
5: if  $rate \neq 1$  then
6:   return 0 ▷ not robust for  $\delta_{init}$ 
7: end if
8:  $\delta \leftarrow \delta_{init}$ 
9: while  $rate = 1$  do
10:   $\delta \leftarrow 2\delta$ 
11:   $rate \leftarrow \mathbf{realisations}(\mathcal{M}_\theta, x', \delta, n)$ 
12: end while
    ▷ we exit from the while because we have found at least one model in the
    realisations with an output  $< 0.5$ , and we have  $[\delta/2, \delta)$  to search for a  $\delta_{max}$ .
13:  $\delta_{max} \leftarrow \delta/2$ 
14: while True do
15:   if  $|\delta - \delta_{max}| \leq \delta_{init}$  then
16:     return  $\delta_{max}$ 
17:   end if
18:    $\delta_{new} \leftarrow (\delta_{max} + \delta)/2$ 
19:    $rate \leftarrow \mathbf{realisations}(\mathcal{M}_\theta, x', \delta_{new}, n)$ 
20:   if  $rate = 1$  then
21:      $\delta_{max} \leftarrow \delta_{new}$ 
22:   else
23:      $\delta \leftarrow \delta_{new}$ 
24:   end if
25: end while
```

769 **Proposition 17.** Fix $\delta_{init} > 0$. Given a model \mathcal{M}_θ and a CFX x' , let δ^* be
770 the (exact) maximum magnitude of model changes such that x' is robust with
771 respect to the set of PMC $\Delta_{\delta^*} = \{S \mid d_p(\mathcal{M}_\theta, S(\mathcal{M}_\theta)) \leq \delta^*\}$. Then, with
772 probability α , **AP Δ S** returns a $\delta_{max} \geq \delta^* - \delta_{init}$ such that the CFX x' is robust
773 for at least a fraction R of the set of PMC $\Delta_{\delta_{max}}$. Moreover, the computation
774 of δ_{max} is polynomial.

775 *Proof sketch.* The δ_{max} returned by the algorithm is obtained by iteratively
776 increasing δ , sampling n models from the corresponding Δ_δ and verifying that

777 $\mathcal{M}_{\theta_i}(x) \geq 0.5$ for each model \mathcal{M}_{θ_i} sampled. By definition, δ^* is the actual
778 value we are trying to estimate. When the algorithm stops, with values δ_{\max}
779 and δ such that $\delta - \delta_{\max} \leq \delta_{init}$, we have that for n models in $\Delta_{\delta_{\max}}$ the
780 CFX x' is robust and for at least one model in Δ_{δ} the x' is not robust. Since
781 for each model in Δ_{δ^*} the CFX is robust, it must hold that $\delta > \delta^*$, hence
782 $\delta_{\max} \geq \delta^* - \delta_{init}$.

783 Once the exponential search ends, by exploiting Lemma 16, we can state
784 that with probability α , the CFX x' is robust for at least R of any infinite
785 further realisations from $\Delta_{\delta_{\max}}$. The time complexity of the algorithm cor-
786 responds to $n \cdot m$ forward propagations, with n being the sample size and
787 $m = \log \frac{\delta_{\max}}{\delta_{init}}$ being the number of iterations of the exponential search, which
788 is polynomial in the input size of the problem. \square

789 6. Experimental Analysis

790 Section 5 laid the theoretical foundations of a novel sampling-based method
791 that allows the obtaining of provable probabilistic guarantees on the robust-
792 ness of CFXs. In this section, we evaluate our approach by considering five
793 experiments:

- 794 • In § 6.1 we show how to instantiate AP Δ S in practice using a synthetic
795 example. Specifically, we first demonstrate the interplay of parameters
796 n , α , and R used to obtain a probabilistic guarantee. Then, using
797 the maximum δ_{\max} discovered by AP Δ S, we precisely characterize the
798 subsets $\hat{\Delta}$ of the set of PMC $\Delta_{\delta_{\max}}$ for which the given CFX x' cannot
799 be proved to be robust. In our experiments at most a fraction $(1 - R)$
800 of $\Delta_{\delta_{\max}}$ is in $\hat{\Delta}$, so complementing empirically our theoretical results.
- 801 • In § 6.2 we compare our certification approach with the one proposed
802 in [12]. In particular, we focus on the difference between the worst-
803 case guarantees offered by their approach and compare them with the
804 average-case guarantees of AP Δ S in terms of maximum changes that can
805 be certified. These experiments confirm our intuition that worst-case
806 guarantees might be too conservative in practice, leading to a larger
807 number of CFXs being discarded.
- 808 • In § 6.3, we consider the problem of generating robust CFXs and com-
809 pare with two state-of-the-art approaches for robustness under PMC, [12]

810 and [11]. We show that our approach produces CFXs that are less ex-
 811 pensive (in terms of ℓ_1 distance from the original input) and more
 812 plausible, without sacrificing robustness.

813 • In § 6.4, we perform an in-depth analysis of the impact that the two
 814 main hyper-parameters of $\text{AP}\Delta\text{S}$, α and R , have on the quality of gener-
 815 ated CFXs. We show that higher values of these parameters typically
 816 lead to tighter estimates that result in improved robustness. These
 817 results also align with existing literature on CFXs in revealing that
 818 improved robustness appears to be correlated with higher plausibility
 819 and cost.

820 • Finally, in § 6.5 we analyse the scalability of $\text{AP}\Delta\text{S}$. We consider
 821 tabular transformer architectures, such as *TabNet* [17] and show that
 822 $\text{AP}\Delta\text{S}$ scales well even when employed in recent architectures employed
 823 at the state of the art and containing hundreds of thousands of param-
 824 eters, thus confirming the wide applicability of our method.

825 An implementation of $\text{AP}\Delta\text{S}$ is integrated in the *RobustX* library [43]
 826 available at <https://github.com/RobustCounterfactualX/RobustX.git>.
 827 Additional material is available at <https://github.com/lmarza/APAS>.

828 6.1. $\text{AP}\Delta\text{S}$ in Action

829 This experiment is designed to demonstrate how the three main param-
 830 eters of $\text{AP}\Delta\text{S}$, i.e., n , α , and R , can be used to obtain probabilistic robust-
 831 ness guarantees. To this end, we focus on the synthetic example depicted in
 832 Fig. 11. Weights for the original network \mathcal{M}_θ , as well as the input used for
 833 testing robustness, are generated randomly.

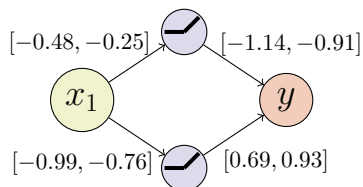


Figure 11: The interval neural network used for exact enumeration.

834 Considering a random input $x = -2.57$, we use $\text{AP}\Delta\text{S}$ to estimate a
 835 δ_{max} for which we seek the guarantee that for at least $R = 90\%$ of the

Algorithm 2 *Exact CFX Δ -Robustness*

```
1: Input: An INN  $\mathcal{N}$  and a CFX  $x'$  and an maximum  $\epsilon$ -precision for the splitting
   phase
2: Output: set of INNs for which  $x'$  is robust.
3: robust_INNs  $\leftarrow \emptyset$ 
4: non-robust_INNs  $\leftarrow \emptyset$ 
5: unknown  $\leftarrow \text{Push}(\mathcal{N})$ 
6: while (unknown  $\neq \emptyset$ ) or ( $\epsilon$ -precision not reached) do
7:    $\mathcal{I} \leftarrow \text{GetINNToVerify}(\text{unknown})$ 
8:    $\mathcal{R}_{\mathcal{I}} \leftarrow \text{ComputeReachableSet}(\mathcal{I}, x')$ 
9:   if  $\text{lower}(\mathcal{R}_{\mathcal{I}}) \geq 0.5$  then
10:    robust_INNs  $\leftarrow \text{Push}(\mathcal{I})$ 
11:    unknown  $\leftarrow \text{Pop}(\mathcal{I})$ 
12:   else if  $\text{upper}(\mathcal{R}_{\mathcal{I}}) < 0.5$  then
13:    non-robust_INNs  $\leftarrow \text{Push}(\mathcal{I})$ 
14:    unknown  $\leftarrow \text{Pop}(\mathcal{I})$ 
15:   else
16:     $\mathcal{I}', \mathcal{I}'' \leftarrow \text{ChooseIntervalToSplit}(\mathcal{I})$ 
17:    unknown  $\leftarrow \text{Push}(\mathcal{I}', \mathcal{I}'')$ 
18:   end if
19: end while
20: return robust_INNs
```

836 plausible model changes induced by such δ_{max} the CFX x' is robust. Following
837 Proposition 17, we set a confidence level $\alpha > 1 - 10^{-40}$ (i.e., with certainty, in
838 practice), which yields $n = 100k$ realisations. For this setting, AP Δ S identifies
839 a $\delta_{max} = 0.115$.

840 To validate this result, we define a procedure to exactly characterize,
841 following the intuitions of [44, 45], the models within $\Delta_{\delta_{max}}$ for which the
842 robustness property does not hold. The interval abstraction proposed by
843 [12] can be used to exactly compute the portion of the model changes from
844 Δ for which a CFX x' is not robust. In fact, it is possible to build an interval
845 neural network using the δ_{max} value identified by AP Δ S, setting each weight
846 w_i in θ to $[w_i - \delta_{max}, w_i + \delta_{max}]$. Then, recursively splitting each interval
847 weight of the network in half allows to identify portions of Δ that are not
848 robust. Employing the following strategy (reported in Algorithm 2), after
849 $s = 7$ splits, we obtain that for $\sim 92\%$ of sub-interval networks, the CFX is
850 robust. The remaining 8% produced an *unknown* answer (i.e., the situation

Table 2: Comparison on the robustness of CFXs using five state-of-the-art methods and AP Δ S proposed in this work.

	<i>Diabetes</i>				<i>no2</i>				<i>SBA</i>				<i>Credit</i>			
	VM1 $\delta = 0.11$	VM2 $\delta_e = 0.27$	ℓ_1	lof	VM1 $\delta = 0.02$	VM2 $\delta_e = 0.07$	ℓ_1	lof	VM1 $\delta = 0.11$	VM2 $\delta = 0.25$	ℓ_1	lof	VM1 $\delta = 0.05$	VM2 $\delta_e = 1.28$	ℓ_1	lof
Wacht-R	100%	100%	0.122	1.00	100%	100%	0.084	1.00	92%	92%	0.023	-0.78	-	-	-	-
Proto-R	100%	96%	0.104	1.00	100%	100%	0.069	1.00	90%	88%	0.011	-0.02	32%	30%	0.300	-1.00
MILP-R	100%	100%	0.212	-0.48	100%	100%	0.059	1.00	100%	100%	0.018	-0.88	100%	100%	0.031	1.00
ROAR	82%	14%	0.078	0.95	88%	34%	0.074	1.00	82%	78%	0.031	-0.80	62%	60%	0.047	1.00
AP Δ S	100%	100%	0.072	1.00	100%	100%	0.042	1.00	100%	100%	0.009	0.44	100%	94%	0.028	1.00

851 depicted in Fig. 10(b)) that would require further splits, corresponding to
 852 only ten nodes to explore in the next iteration. In the worst case, even
 853 considering all the remaining ten nodes left to explore as non-robust, we
 854 would have a maximum percentage of non-robustness still lower than the
 855 desired upper bound $(1 - R) = 10\%$, confirming that the guarantees produced
 856 by AP Δ S indeed hold in practice.

857 6.2. Worst-case vs Average-case Guarantees

858 This set of experiments aims to compare the probabilistic guarantees
 859 offered by AP Δ S with the worst-case guarantees offered by [12]. What we
 860 aim to show here is that adopting an average-case certification perspective
 861 may be more practical in some circumstances, as worst-case guarantees may
 862 be unnecessarily conservative. Our approach aims to obtain a δ_{max} for which
 863 the CFX is robust with confidence α for at least a fraction R of model changes
 864 in Δ . This is in stark contrast with the worst-case reasoning of [12], where
 865 even a single realisation of Δ for which the CFX is not robust results in the
 866 corresponding δ being discarded.

867 To show why such strict guarantees may not be needed, we use an anal-
 868 ogous experimental setup and the training process of [12], which considers
 869 four datasets: *Diabetes* (continuous) [46], *Credit* (heterogeneous) [37], *no2*
 870 (continuous) [47] and *Small Business Administration* (SBA) (continuous fea-
 871 tures) [48]. In detail, for the training procedure of the classifier, we randomly
 872 shuffle each dataset and split it into two halves, denoted \mathcal{D}_1 and \mathcal{D}_2 . First, we
 873 use \mathcal{D}_1 to train a base neural network; then we use both \mathcal{D}_1 and \mathcal{D}_2 to train
 874 a shifted model. We then generate 50 robust CFXs for the base network
 875 using the MILP-R and the same δ values as in [12] for a fair comparison.

Algorithm 3 Provable Plausible Δ -Shift

```
1: Input: Model  $\mathcal{M}_\theta$ , CFX  $x'$ ,  $\alpha$ ,  $R$ 
2: Output:  $\delta_{max}$ 
3:  $\delta_{init} \leftarrow 0.0001$ 
4: rate  $\leftarrow$  MILP( $\mathcal{M}_\theta, x', \delta_{init}$ )
5: if rate  $\neq 1$  then
6:   return 0 ▷ no robustness
7: end if
8:  $\delta \leftarrow \delta_{init}$ 
9: while rate = 1 do
10:   $\delta \leftarrow 2\delta$ 
11:  rate  $\leftarrow$  MILP( $\mathcal{M}_\theta, x', \delta$ )
12: end while
13:  $\delta_{max} \leftarrow \delta/2$ 
14: while True do
15:  if  $|\delta - \delta_{max}| \leq \delta_{init}$  then
16:    return  $\delta_{max}$ 
17:  end if
18:   $\delta_{new} \leftarrow (\delta_{max} + \delta)/2$ 
19:  rate  $\leftarrow$  MILP( $\mathcal{M}_\theta, x', \delta_{new}$ )
20:  if rate = 1 then
21:     $\delta_{max} \leftarrow \delta_{new}$ 
22:  else
23:     $\delta \leftarrow \delta_{new}$ 
24:  end if
25: end while
```

876 Specifically, we use $\delta = 0.11$ for *Diabetes*, $\delta = 0.02$ for *no2*, $\delta = 0.11$ for *SBA*
877 and $\delta = 0.05$ for *Credit*. Subsequently, we evaluate the resulting CFXs by
878 looking at two metrics: (i) **VM1**, the percentage of CFXs that are valid on
879 the base neural network and (ii) **VM2**, the percentage of CFXs that remain
880 valid for the shifted neural network trained using both \mathcal{D}_1 and \mathcal{D}_2 . Table 2
881 reports the results we obtained for this experiment.

882 As previously observed by Jiang et al. [12], the training procedure used to
883 generate shifted models may result in changes that exceed the δ used to gener-
884 ate provably robust CFXs. Indeed, after inspecting the networks obtained,
885 we noted that the maximum empirical difference observed after retraining
886 (denoted as δ_e) is well above the δ values used during CFX generation. In

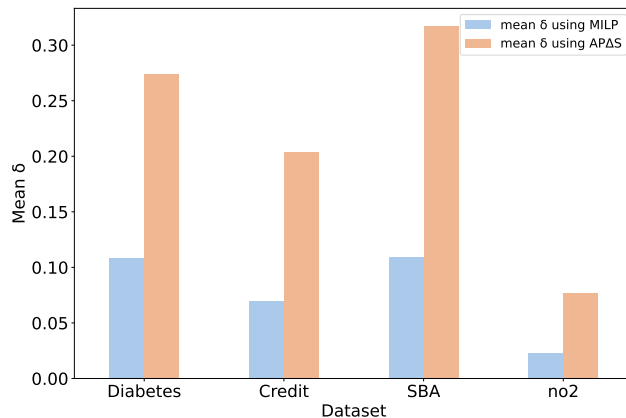


Figure 12: Average robust δ obtained using MILP-based certification and APDS .

887 particular, we recorded $\delta_e = 0.27$ for *Diabetes*, $\delta_e = 0.07$ for *no2*, $\delta = 0.25$
 888 for *SBA* and $\delta_e = 1.28$ for *Credit*. Given the magnitude of these changes,
 889 the robustness of the CFXs generated by MILP-R cannot be guaranteed in
 890 practice. However, the results show a rather intriguing picture: the **VM2**
 891 metric appears to be unaffected by retraining, and all CFXs remain valid on
 892 the respective final models.

893 These results suggest that certification approaches based on worst-case
 894 reasoning may be too strict in practical scenarios. To further understand the
 895 implications of worst-case vs average-case reasoning, we adapted Algorithm 1
 896 to use the certification procedure of Jiang et al., i.e., a MILP solver instead
 897 of a sampled-based approach, and compute the maximum provable δ^* for
 898 which the previously generated CFXs are robust (Algorithm 3).

899 Fig. 12 shows a comparison between the average maximum provable δ
 900 obtained by this procedure and APDS . As we can observe, our average-case
 901 guarantees allow to obtain δ values that are much higher, exceeding the
 902 MILP-certified in all instances. This is expected, given the results discussed
 903 in Proposition 17. However, what remains unclear is how these differences
 904 may affect the cost and plausibility of CFXs when certification procedures
 905 are embedded in procedures to generate CFXs.

906 6.3. Generating Robust CFXs using APDS

907 The results discussed in the previous section have important implications
 908 on algorithms for the generation of robust CFXs. Recent works, e.g. [32, 12,

Algorithm 4 Generation of Robust CFXs

```
1: Input: Model  $\mathcal{M}$ , input  $x$  such that  $\mathcal{M}(x) = c$ , set of plausible model changes  
    $\Delta$ , maximum iteration number  $\tau$   
2: Output:  $\Delta$ -robust CFX  $x'$   
3:  $t \leftarrow 0$  ▷ iteration number  
4: while  $t < \tau$  do  
5:    $x' \leftarrow \text{ComputeCFX}(x, \mathcal{M})$   
6:    $\text{rate} \leftarrow \text{AP}\Delta\text{S}(\mathcal{M}, x', \Delta)$   
7:   if  $\text{rate} = 1$  then  
8:     return  $x'$  ▷  $x'$  is approx.  $\Delta$ -robust  
9:   else  
10:    increase allowed distance of next CFX  
11:    increase iteration number  $t$   
12:   end if  
13: end while  
14: return no robust CFX can be found
```

909 15], have proposed iterative procedures that generate provably robust CFXs
910 by alternating two phases. First, a CFX is generated solving (variations
911 of) Definition 1; then, a robustness certification procedure is invoked on
912 the CFX. If the CFX is robust, then it is returned to the user; otherwise,
913 the search continues, allowing for CFXs of increasing distance to be found.
914 Clearly, the certification step has the potential to affect the CFXs computed
915 in several ways. A robustness test that is too conservative may discard
916 potentially good explanations and keep relaxing the distance constraint until
917 the CFX is deemed robust. Ultimately, this may result in CFXs that exhibit
918 poor proximity and plausibility.

919 To test this hypothesis, we adapt the CFX generation algorithm of [12]
920 and replace their Δ -robustness test with the one performed by **AP** $\Delta**S**. The
921 complete procedure is shown in Algorithm 4. In detail, after some initial-
922 ization steps, we compute the first CFX using **ComputeCFX**(x, \mathcal{M}) (line 5),
923 which employs the solution proposed in [12] and presented above. Given a
924 CFX x' and a plausible model shift Δ , at line 6, we employ **AP** Δ **S** setting
925 $\alpha = 0.999$ and $R = 0.995$, thus obtaining 1378 realisations to perform in
926 the robustness test. If the CFX x' returned by our approximation results
927 robust for all these realisations, then we return it to the user. Otherwise, we
928 increased the allowed distance for the next CFX generation and the iteration$

929 number t (lines 10-11).

930 We then compare the resulting procedure with the four generation algo-
931 rithm studied in [12]: Wacht-R, Proto-R, MILP-R, and finally, ROAR [11].
932 Notably, ROAR is specifically designed to generate robust CFXs under plau-
933 sible model changes using average-case certification. Using the same datasets
934 and training procedures of § 6.2, we generate 50 CFXs for each dataset. We
935 evaluate CFXs based on their proximity, measured by the ℓ_1 distance, and
936 plausibility, measured by the local outlier factor (**lof**) which determines if
937 an instance is within the data manifold by quantifying the local data den-
938 sity [49] (+1 for inliers, -1 otherwise). We average ℓ_1 and **lof** over the
939 generated CFXs. We also report **VM1** and **VM2** for completeness. The
940 results obtained, which we report in Table 2, confirm our hypothesis. In-
941 deed, **AP Δ S** produces the best results across all datasets, always generating
942 CFXs with high plausibility and better proximity. Notably, **AP Δ S** outper-
943 forms ROAR as well, producing CFXs that retain a higher degree of validity
944 after retraining.

945 6.4. Impact of hyper-parameters on validity, plausibility and cost

946 The previous set of experiments demonstrated that **AP Δ S** is able to out-
947 perform existing approaches and generate CFXs that are robust, but also
948 plausible and less expensive than other robust approaches. What remains
949 unclear is the role that the main hyper-parameters of our algorithm, α and
950 R , might play in obtaining these results. We therefore conducted additional
951 experiments to evaluate the interplay between the tightness of the probabilis-
952 tic guarantees offered by **AP Δ S** and the quality of resulting explanations. In
953 particular, focusing on the same datasets used in previous experiments, we
954 started by checking the influence that α and R have on the validity of CFXs
955 after retraining. We generated 50 CFXs for each dataset using an instantia-
956 tion of Algorithm 4 that uses MILP encodings to generate candidate CFXs
957 as done in [12]. Figures 13-16 report the results obtained for the *Diabetes*
958 and *SBA* datasets. Additional results for the two remaining datasets are
959 reported in the appendix for this first set of experiments.

960 Our intuition is that lower values for α and R should result in coarser
961 robustness guarantees (i.e., larger δ values) and, thus, lower validity rates.
962 As we can observe, our intuition is confirmed across all datasets, further
963 clarifying the nature of the probabilistic guarantees that **AP Δ S** can offer.
964 Next, we investigate the impact that α and R have on the plausibility and
965 cost of CFXs generated by **AP Δ S**. As per our previous experiments, we

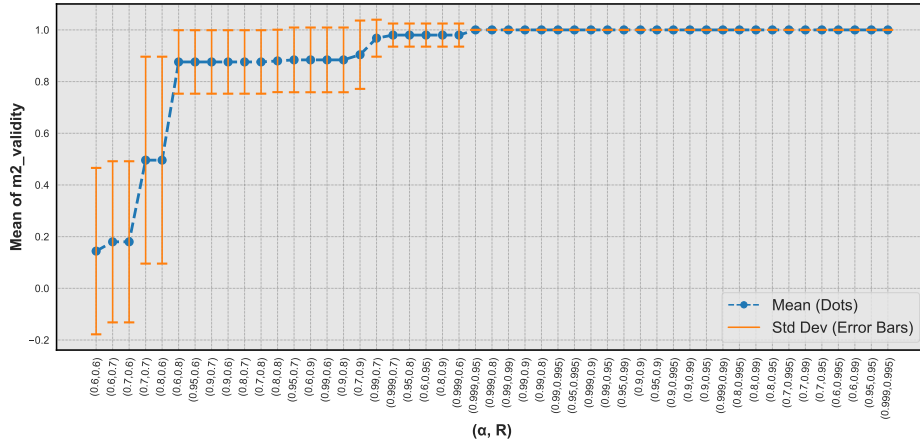


Figure 13: Mean validity after retraining visualised for increasing α, R values using the *Diabetes* dataset.

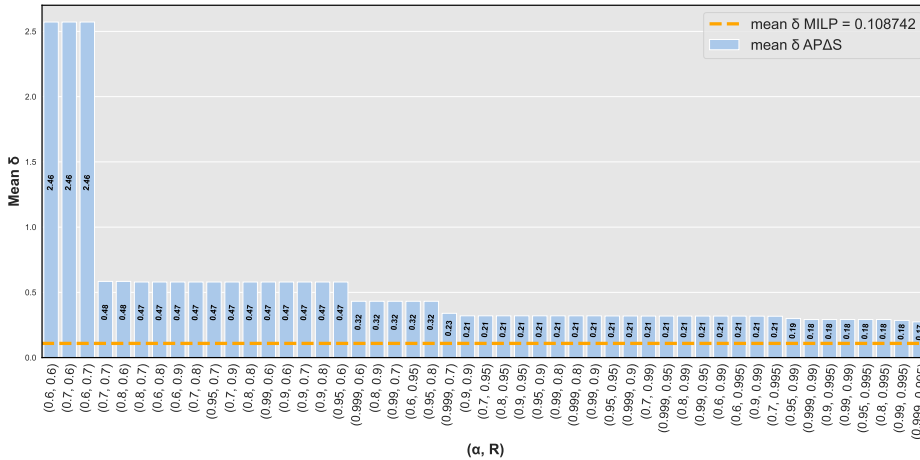


Figure 14: Mean certifiable δ obtained for increasing α, R values using the *Diabetes* dataset.

966 measure plausibility using the *LOF* score, and we use ℓ_0, ℓ_1 and ℓ_∞ norms to
 967 measure the proximity of CFXs. For conciseness, we only report results for
 968 the *Diabetes* dataset in Figure 17 below and delegate additional results to the
 969 appendix. To improve the readability of our results, we decided to separate
 970 CFXs that achieved 100% validity after retraining for the rest. Overall we
 971 can observe a clear trend, whereby increasing α, R results in CFXs that are
 972 further away from the decision boundary and thus more plausible. These

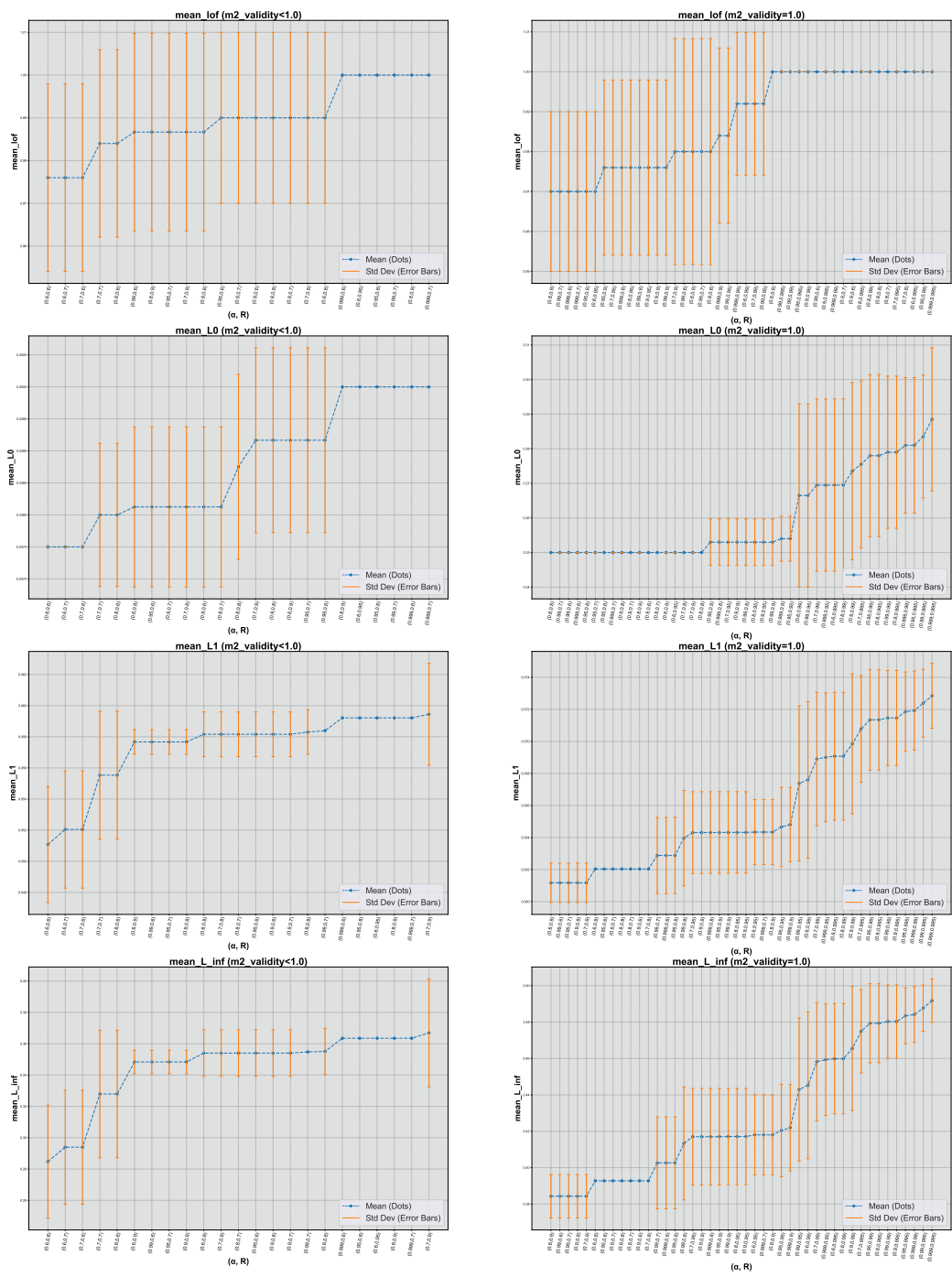


Figure 17: Mean LOF , l_0, l_1, l_{∞} metrics for increasing α, R values using the *Diabetes* dataset. CFXs with 100% validity after retraining are shown on the right, while the remaining CFXs are shown on the left.

980 focus on *TabNet* [17], a tabular transformer recently introduced that lever-
981 ages sparse attention and sequential feature selection to learn interpretable
982 feature representations. At its core, *TabNet* processes data in a series of de-
983 cision steps, with each step using a learned attention mask to select a subset
984 of features, which allows the model to focus on the most relevant attributes
985 at each stage. This sparse attention mechanism makes *TabNet* computationally
986 efficient and helps to improve interpretability in tabular datasets. From
987 our perspective, this architecture is interesting as it comprises an attention
988 mechanism, with encoder-decoder components typical of other recent trans-
989 former architectures and a consequent significant number of parameters to
990 test the scalability of AP Δ S . To the best of our knowledge, this is the first
991 time that CFXs with robustness guarantees are generated for such a complex
992 architecture with tens of thousands of parameters.

993 Before considering the robustness property, we analyse the accuracy in
994 the training and testing phases of *TabNet*. To this end, we employ a su-
995 pervised training approach, splitting the datasets employed in the previous
996 evaluation, namely *Diabetes*, *No2*, *SBA* and *Credit*, into training and testing
997 datasets, and we first compare the accuracy obtained using this architecture
998 with standard MLPs employed in § 6.2. To ensure statistical significance of
999 our results, we consider, for each dataset tested, the mean of the accuracies
1000 obtained using ten random initializations of the transformer architecture. As
1001 highlighted in the first two columns of Tab. 3, with *TabNet*, we have an
1002 increased number of parameters in the model but similar or even higher ac-
1003 curacy with respect to the classical MLP, confirming the potential of this
1004 novel architecture in selecting important features to get more precise final
1005 accuracy in the prediction.

1006 As our results show a similar level of accuracy between MLP and *TabNet*,
1007 we move on to how to generate robust CFXs for this transformer architecture.
1008 Given the significantly higher number of parameters in *TabNet*, we replace
1009 the MILP-based procedure used in Section 6.3, Algorithm 4 with a Nearest
1010 Neighbors Counterfactual Explainer (NNCFX) [40] to ensure scalability of
1011 our generation procedure. More specifically, line 5 in Algorithm 4 (reported
1012 in the appendix) now implements the following strategy. Given an input x for
1013 which a robust CFX is sought, we identify the nearest data point belonging
1014 to the dataset for which *TabNet* produces a different classification outcome.
1015 Our implementation uses k-d trees to improve the efficiency of this nearest-

Table 3: Scalability experiments of $\text{AP}\Delta\text{S}$.

Diabetes				
	# Parameters	Mean Accuracy	Mean δ_{max}	Mean Comp. Time
<i>MLP</i>	81	79%	0.32	0.01s
<i>TabNet</i>	30992	82%	0.48	5.21s
no2				
	# Parameters	Mean Accuracy	Mean δ_{max}	Mean Comp. Time
<i>MLP</i>	145	64%	0.11	0.008s
<i>TabNet</i>	30676	68%	0.35	9.2s
SBA				
	# Parameters	Mean Accuracy	Mean δ_{max}	Mean Comp. Time
<i>MLP</i>	199	99%	0.53	0.02s
<i>TabNet</i>	30992	100%	0.16	10.3s
Credit				
	# Parameters	Mean Accuracy	Mean δ_{max}	Mean Comp. Time
<i>MLP</i>	371	74%	0.34	0.01s
<i>TabNet</i>	40946	74%	1.8	11.3s

1016 neighbor search.⁵

1017 For the following experiments, we consider the same datasets used in Sec-
 1018 tion 6.3 and the same 50 original inputs employed in those experiments. The
 1019 last two columns of Table 3 report the results we obtained when generating
 1020 CFX with robustness guarantees for *TabNet*. As we can observe, $\text{AP}\Delta\text{S}$ is still
 1021 able to compute robust CFXs within tens of seconds, even when employed
 1022 in transformer-based architecture with $\sim 400x$ times parameters, showing
 1023 a linear growth time computation. Similar runtimes are observed across all
 1024 four datasets, thus confirming the high scalability of our approach. To fur-
 1025 ther confirm this aspect, we ran an in-depth scalability study by training a

⁵We perform a further experiment reported in the Appendix B to understand the difference between the two CFX-generation approaches, namely MILP and NNCFX.

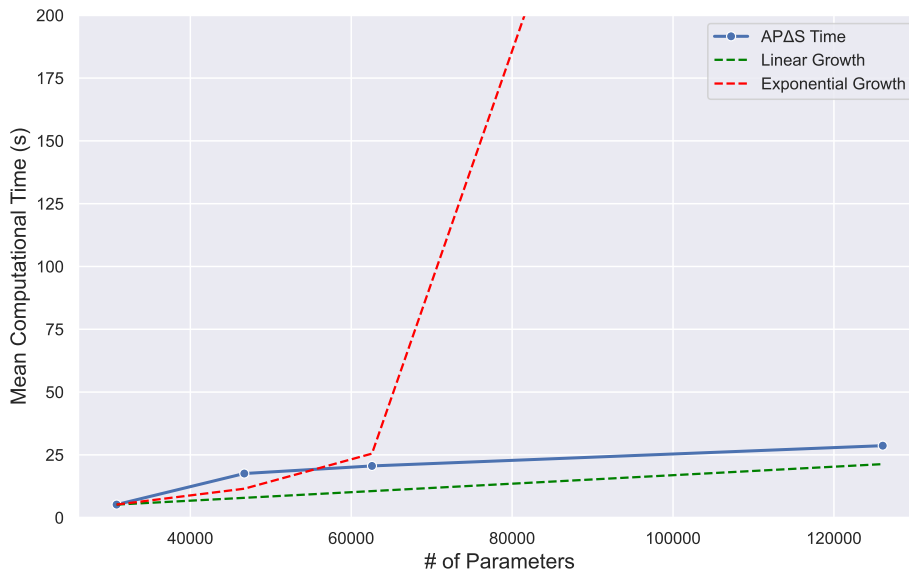


Figure 18: Mean computation time of $\text{AP}\Delta\text{S}$ applied to *TabNet* architectures with increasing number of parameters.

1026 set of 4 *TabNet* models with an increasing number of parameters. Using the
 1027 diabetes dataset, we trained models containing [30834, 40946, 62578, 126066]
 1028 respectively and generated 50 robust CFXs for each. We stored the runtimes
 1029 for each robust CFX and report the mean computation time for each model
 1030 in Figure 18. As we can observe, the runtime increase follows a linear trend,
 1031 thus highlighting the effectiveness and applicability of our proposed solution
 1032 even when targetting complex architectures.

1033 7. Conclusions

1034 We studied the problem of generating robust CFXs with respect to plausi-
 1035 ble model changes. We proved for the first time that certifying the robustness
 1036 of CFX with respect to this notion of robustness is an NP-hard problem, and
 1037 also extended this result to show that the same complexity results apply to
 1038 naturally-occurring model changes. These results motivate the quest for new
 1039 scalable algorithms to certify robustness under plausible model changes. To
 1040 this end, we investigated existing methods to generate robust CFXs with
 1041 probabilistic guarantees and showed that these approaches may not be di-
 1042 rectly applicable to our setting. We then introduced $\text{AP}\Delta\text{S}$, a novel scalable

1043 approach for probabilistic robustness certification, and used it to generate
1044 robust CFXs under plausible model changes. We carried out an extensive
1045 experimental analysis, demonstrating the advantages of $\text{AP}\Delta\text{S}$ and outper-
1046 forming SOTA methods on a range of metrics, including validity, plausibility,
1047 and cost. Crucially, we also applied our method to certify CFXs’ robustness
1048 for tabular transformers containing thousands of parameters. To the best
1049 of our knowledge, we are the first to consider models of this size within the
1050 robust CFX literature [10], further demonstrating the scalability and wide
1051 applicability of our approach. We see these outcomes as important contribu-
1052 tions towards complementing existing formal approaches for Explainable AI
1053 and making them applicable in practice.

1054 Acknowledgements

1055 Leofante was supported by an Imperial College Research Fellowship. Any
1056 views or opinions expressed herein are solely those of the authors listed.

1057 References

- 1058 [1] L. Tai, G. Paolo, M. Liu, Virtual-to-real drl: Continuous control of
1059 mobile robots for mapless navigation, in: IROS, 2017.
- 1060 [2] L. Marzari, D. Corsi, E. Marchesini, A. Farinelli, Curriculum learning
1061 for safe mapless navigation, in: Proceedings of the 37th ACM/SIGAPP
1062 Symposium on Applied Computing, 2022, pp. 766–769.
- 1063 [3] K. O’Shea, R. Nash, An introduction to convolutional neural networks,
1064 arXiv preprint arXiv:1511.08458 (2015).
- 1065 [4] D. Corsi, L. Marzari, A. Pore, A. Farinelli, A. Casals, P. Fiorini,
1066 D. Dall’Alba, Constrained reinforcement learning and formal verifica-
1067 tion for safe colonoscopy navigation, in: 2023 IEEE/RSJ International
1068 Conference on Intelligent Robots and Systems (IROS), IEEE, 2023.
- 1069 [5] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfel-
1070 low, R. Fergus, Intriguing properties of neural networks, arXiv preprint
1071 arXiv:1312.6199 (2013).
- 1072 [6] G. Amir, D. Corsi, R. Yerushalmi, L. Marzari, D. Harel, A. Farinelli,
1073 G. Katz, Verifying learning-based robotic navigation systems, in: 29th
1074 International Conference, TACAS 2023, Springer, 2023, pp. 607–627.

- 1075 [7] I. Stepin, J. M. Alonso, A. Catalá, M. Pereira-Fariña, A survey of
1076 contrastive and counterfactual explanation generation methods for ex-
1077 plainable artificial intelligence, *IEEE Access* 9 (2021) 11974–12001.
- 1078 [8] A. Karimi, G. Barthe, B. Schölkopf, I. Valera, A survey of algorithmic
1079 recourse: Contrastive explanations and consequential recommendations,
1080 *ACM Comput. Surv.* 55 (2023) 95:1–95:29.
- 1081 [9] S. Wachter, B. D. Mittelstadt, C. Russell, Counterfactual explana-
1082 tions without opening the black box: Automated decisions and the
1083 GDPR, *CoRR* abs/1711.00399 (2017). URL: [http://arxiv.org/abs/](http://arxiv.org/abs/1711.00399)
1084 [1711.00399](http://arxiv.org/abs/1711.00399). arXiv:1711.00399.
- 1085 [10] J. Jiang, F. Leofante, A. Rago, F. Toni, Robust counterfactual explana-
1086 tions in machine learning: A survey, in: *Proceedings of the Thirty-Third*
1087 *International Joint Conference on Artificial Intelligence, IJCAI 2024,*
1088 *Jeju, South Korea, August 3-9, 2024*, ijcai.org, 2024, pp. 8086–8094.
1089 URL: <https://www.ijcai.org/proceedings/2024/894>.
- 1090 [11] S. Upadhyay, S. Joshi, H. Lakkaraju, Towards robust and reliable algo-
1091 rithmic recourse, in: *Advances in Neural Information Processing Sys-*
1092 *tems 34 (NeurIPS21)*, 2021, pp. 16926–16937.
- 1093 [12] J. Jiang, F. Leofante, A. Rago, F. Toni, Formalising the robustness of
1094 counterfactual explanations for neural networks, in: *Proceedings of the*
1095 *37th AAAI Conference on Artificial Intelligence (AAAI23)*, 2023, pp.
1096 14901–14909.
- 1097 [13] E. Black, Z. Wang, M. Fredrikson, Consistent counterfactuals for deep
1098 models, in: *Proceedings of the 10th International Conference on Learn-*
1099 *ing Representations (ICLR22)*, OpenReview.net, 2022.
- 1100 [14] T.-D. H. Nguyen, N. Bui, D. Nguyen, M.-C. Yue, V. A. Nguyen, Ro-
1101 bust Bayesian recourse, in: *Proceedings of the 38th Conference on*
1102 *Uncertainty in AI (UAI22)*, 2022, pp. 1498–1508.
- 1103 [15] F. Hamman, E. Noorani, S. Mishra, D. Magazzeni, S. Dutta, Robust
1104 counterfactual explanations for neural networks with probabilistic guar-
1105 antees, in: *Proceedings of the International Conference on Machine*
1106 *Learning (ICML23)*, 2023, pp. 12351–12367.

- 1107 [16] L. Marzari, F. Leofante, F. Cicalese, A. Farinelli, Rigorous probabilistic
1108 guarantees for robust counterfactual explanations, in: ECAI 2024, IOS
1109 Press, 2024, pp. 1059–1066.
- 1110 [17] S. Ö. Arik, T. Pfister, Tabnet: Attentive interpretable tabular learn-
1111 ing, in: Proceedings of the AAAI conference on artificial intelligence,
1112 volume 35, 2021, pp. 6679–6687.
- 1113 [18] R. Poyiadzi, K. Sokol, R. Santos-Rodríguez, T. D. Bie, P. A. Flach,
1114 FACE: feasible and actionable counterfactual explanations, in: Proceed-
1115 ings of the AAAI/ACM Conference on AI, Ethics, and Society (AIES20),
1116 2020, pp. 344–350.
- 1117 [19] M. Pawelczyk, K. Broelemann, G. Kasneci, Learning model-agnostic
1118 counterfactual explanations for tabular data, in: Proceedings of the
1119 Web Conference (WWW20), ACM / IW3C2, 2020, pp. 3126–3132.
- 1120 [20] B. Ustun, A. Spangher, Y. Liu, Actionable recourse in linear classifica-
1121 tion, in: Proceedings of the Conference on Fairness, Accountability, and
1122 Transparency (FAT*19), 2019, pp. 10–19. URL: [https://doi.org/10.](https://doi.org/10.1145/3287560.3287566)
1123 [1145/3287560.3287566](https://doi.org/10.1145/3287560.3287566). doi:10.1145/3287560.3287566.
- 1124 [21] A. Artelt, V. Vaquet, R. Velioglu, F. Hinder, J. Brinkrolf, M. Schilling,
1125 B. Hammer, Evaluating robustness of counterfactual explanations, in:
1126 2021 IEEE Symposium Series on Computational Intelligence (SSCI),
1127 IEEE, 2021, pp. 01–09.
- 1128 [22] D. Slack, A. Hilgard, H. Lakkaraju, S. Singh, Counterfactual
1129 explanations can be manipulated, in: Advances in Neural
1130 Information Processing Systems 34 (NeurIPS21), 2021, pp.
1131 62–75. URL: [https://proceedings.neurips.cc/paper/2021/hash/](https://proceedings.neurips.cc/paper/2021/hash/009c434cab57de48a31f6b669e7ba266-Abstract.html)
1132 [009c434cab57de48a31f6b669e7ba266-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/009c434cab57de48a31f6b669e7ba266-Abstract.html).
- 1133 [23] R. Dominguez-Olmedo, A. Karimi, B. Schölkopf, On the adversarial ro-
1134 bustness of causal algorithmic recourse, in: Proceedings of the Interna-
1135 tional Conference on Machine Learning (ICML22), 2022, pp. 5324–5342.
- 1136 [24] S. Zhang, X. Chen, S. Wen, Z. Li, Density-based reliable and robust
1137 explainer for counterfactual explanation, *Expert Syst. Appl.* 226 (2023)
1138 120214.

- 1139 [25] F. Leofante, N. Potyka, Promoting counterfactual robustness through
1140 diversity, in: Thirty-Eighth AAAI Conference on Artificial Intelli-
1141 gence, AAAI 2024, February 20-27, 2024, Vancouver, Canada, AAAI
1142 Press, 2024, pp. 21322–21330. URL: <https://doi.org/10.1609/aaai.v38i19.30127>. doi:10.1609/AAAI.V38I19.30127.
- 1144 [26] V. Guyomard, F. Fessant, T. Guyet, T. Bouadi, A. Termier, Gener-
1145 ating robust counterfactual explanations, in: Machine Learning and
1146 Knowledge Discovery in Databases: Research Track - European Confer-
1147 ence, ECML PKDD 2023, volume 14171 of *Lecture Notes in Computer
1148 Science*, Springer, 2023, pp. 394–409.
- 1149 [27] M. Virgolin, S. Fracaros, On the robustness of sparse counterfactual
1150 expl. to adverse perturbations, *Artif. Intell.* 316 (2023). URL: <https://doi.org/10.1016/j.artint.2022.103840>. doi:10.1016/j.artint.2022.103840.
- 1153 [28] F. Leofante, A. Lomuscio, Towards robust contrastive explanations for
1154 human-neural multi-agent systems, in: N. Agmon, B. An, A. Ricci,
1155 W. Yeoh (Eds.), Proceedings of the 2023 International Conference on
1156 Autonomous Agents and Multiagent Systems (AAMAS23), ACM, 2023,
1157 pp. 2343–2345.
- 1158 [29] M. Pawelczyk, T. Datta, J. van den Heuvel, G. Kasneci, H. Lakkaraju,
1159 Probabilistically robust recourse: Navigating the trade-offs between
1160 costs and robustness in algorithmic recourse, in: Proceedings of the
1161 11th International Conference on Learning Representations (ICLR23),
1162 OpenReview.net, 2023.
- 1163 [30] F. Leofante, E. Botoeva, V. Rajani, Counterfactual explanations and
1164 model multiplicity: a relational verification view, in: Proceedings of the
1165 20th International Conference on Principles of Knowledge Representa-
1166 tion and Reasoning (KR23), 2023, pp. 763–768.
- 1167 [31] J. Jiang, A. Rago, F. Leofante, F. Toni, Recourse under model multi-
1168 plicity via argumentative ensembling, in: Proceedings of the 2024 In-
1169 ternational Conference on Autonomous Agents and Multiagent Systems
1170 (AAMAS24), 2024.

- 1171 [32] S. Dutta, J. Long, S. Mishra, C. Tilli, D. Magazzeni, Robust coun-
1172 terfactual explanations for tree-based ensembles, in: Proceedings of
1173 the International Conference on Machine Learning (ICML22), 2022, pp.
1174 5742–5756.
- 1175 [33] J. Jiang, F. Leofante, A. Rago, F. Toni, Interval abstractions for ro-
1176 bust counterfactual explanations, *Artif. Intell.* 336 (2024) 104218. URL:
1177 <https://doi.org/10.1016/j.artint.2024.104218>. doi:10.1016/J.
1178 ARTINT.2024.104218.
- 1179 [34] I. J. Goodfellow, Y. Bengio, A. C. Courville, *Deep Learning, Adaptive*
1180 *computation and machine learning*, MIT Press, 2016. URL: [http://](http://www.deeplearningbook.org/)
1181 www.deeplearningbook.org/.
- 1182 [35] P. Prabhakar, Z. R. Afzal, Abstraction based output range analysis
1183 for neural networks, in: *Advances in Neural Information Processing*
1184 *Systems 32 (NeurIPS19)*, 2019, pp. 15762–15772.
- 1185 [36] W. Rudin, et al., *Principles of mathematical analysis, volume 3*,
1186 McGraw-hill New York, 1964.
- 1187 [37] D. Dua, C. Graff, UCI machine learning repository, [http://archive.](http://archive.ics.uci.edu/ml)
1188 [ics.uci.edu/ml](http://archive.ics.uci.edu/ml), 2017. Accessed: 2022-08-30.
- 1189 [38] M. Hopkins, E. Reeber, G. Forman, J. Suermondt, Spam-
1190 base, UCI Machine Learning Repository, 1999. DOI:
1191 <https://doi.org/10.24432/C53G6X>.
- 1192 [39] K. Fernandes, P. Vinagre, P. Cortez, P. Sernadela, Online
1193 News Popularity, UCI Machine Learning Repository, 2015. DOI:
1194 <https://doi.org/10.24432/C5NS3V>.
- 1195 [40] R. Guidotti, Counterfactual explanations and how to find them: litera-
1196 ture review and benchmarking, *Data Mining and Knowledge Discovery*
1197 (2022) 1–55.
- 1198 [41] V. Tjeng, K. Xiao, R. Tedrake, Evaluating robustness of neural networks
1199 with mixed integer programming, arXiv preprint arXiv:1711.07356
1200 (2017).

- 1201 [42] S. S. Wilks, Statistical prediction with special reference to the problem
1202 of tolerance limits, *The annals of mathematical statistics* 13 (1942)
1203 400–409.
- 1204 [43] J. Jiang, L. Marzari, A. Purohit, F. Leofante, Robustx: Robust counter-
1205 factual explanations made easy, arXiv preprint arXiv:2502.13751 (2025).
- 1206 [44] L. Marzari, D. Corsi, F. Cicalese, A. Farinelli, The #DNN-Verification
1207 Problem: Counting Unsafe Inputs for Deep Neural Networks, in: Inter-
1208 national Joint Conference on Artificial Intelligence (IJCAI), 2023, pp.
1209 217–224.
- 1210 [45] L. Marzari, D. Corsi, E. Marchesini, A. Farinelli, F. Cicalese, Enumer-
1211 ating safe regions in deep neural networks with provable probabilistic
1212 guarantees, in: Proceedings of the AAAI Conference on Artificial Intel-
1213 ligence, volume 38, 2024, pp. 21387–21394.
- 1214 [46] J. W. Smith, J. E. Everhart, W. Dickson, W. C. Knowler, R. S. Jo-
1215 hannes, Using the adap learning algorithm to forecast the onset of dia-
1216 betes mellitus, in: Proceedings of the annual symposium on computer
1217 application in medical care, American Medical Informatics Association,
1218 1988, p. 261.
- 1219 [47] J. Vanschoren, J. N. van Rijn, B. Bischl, L. Torgo, Openml: networked
1220 science in machine learning, *SIGKDD Explor.* 15 (2013) 49–60. URL:
1221 <https://doi.org/10.1145/2641190.2641198>. doi:10.1145/2641190.
1222 2641198.
- 1223 [48] M. Li, A. Mickel, S. Taylor, “should this loan be approved or denied?”:
1224 A large dataset with class assignment guidelines, *Journal of Statistics*
1225 *Education* 26 (2018) 55–66.
- 1226 [49] M. M. Breunig, H. Kriegel, R. T. Ng, J. Sander, LOF: identifying
1227 density-based local outliers, in: Proceedings of the 2000 ACM SIGMOD
1228 International Conference on Management of Data, ACM, 2000, pp. 93–
1229 104.

1230 Appendix A. Additional results on *Credit* and *No2* datasets for
 1231 § 6.4 experiments.

1232 *Credit*

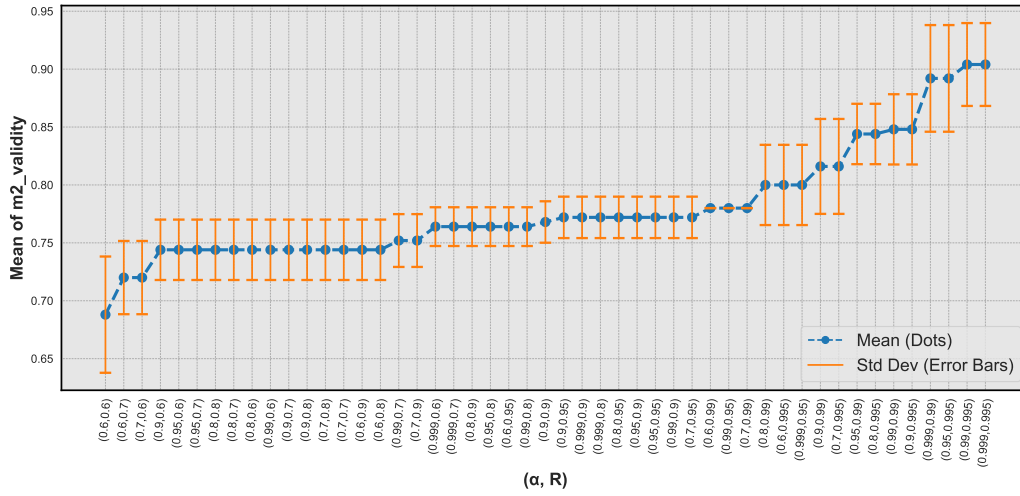


Figure A.19: Mean validity after retraining visualised for increasing α , R values using the *Credit* dataset.

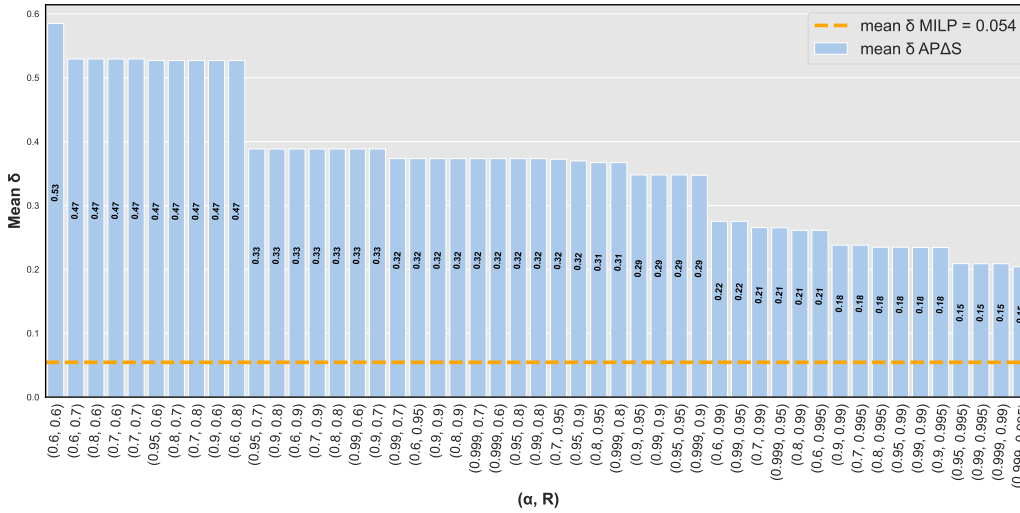


Figure A.20: Mean certifiable δ obtained for increasing α , R values using the *Credit* dataset.

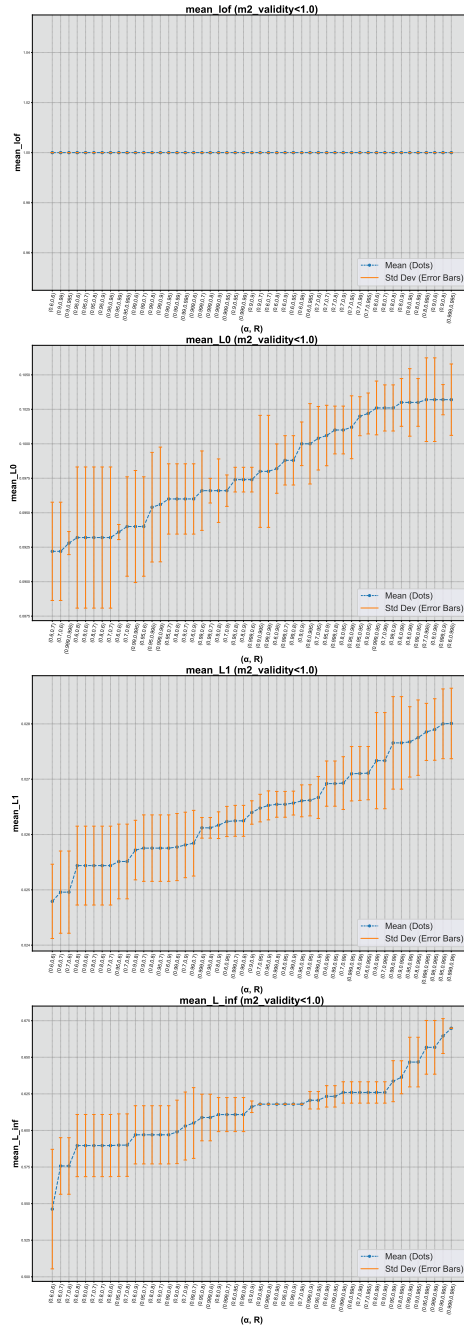


Figure A.21: Mean LOF , ℓ_0 , ℓ_1 , ℓ_∞ metrics for increasing α , R values using the *Credit* dataset. CFXs never reach 100% validity after retraining on this dataset.

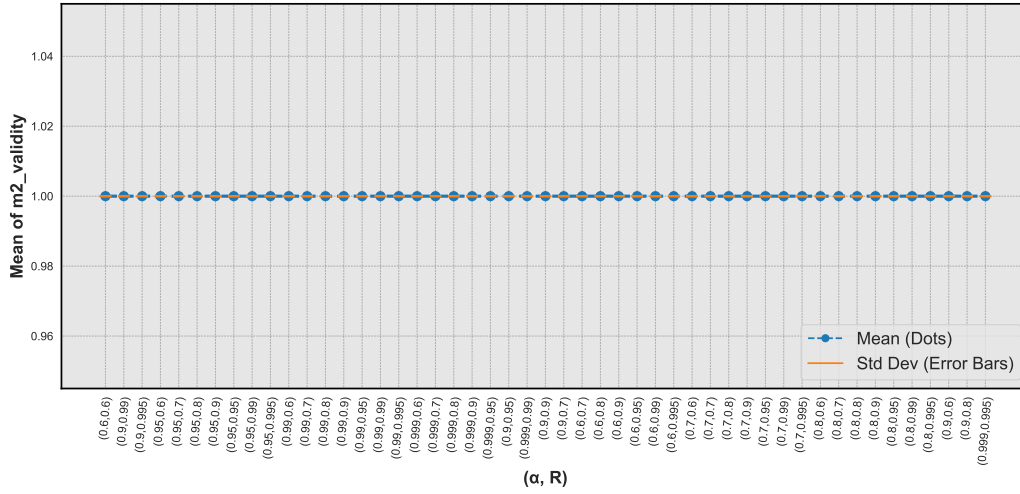


Figure A.22: Mean validity after retraining visualised for increasing α, R values using the *No2* dataset.

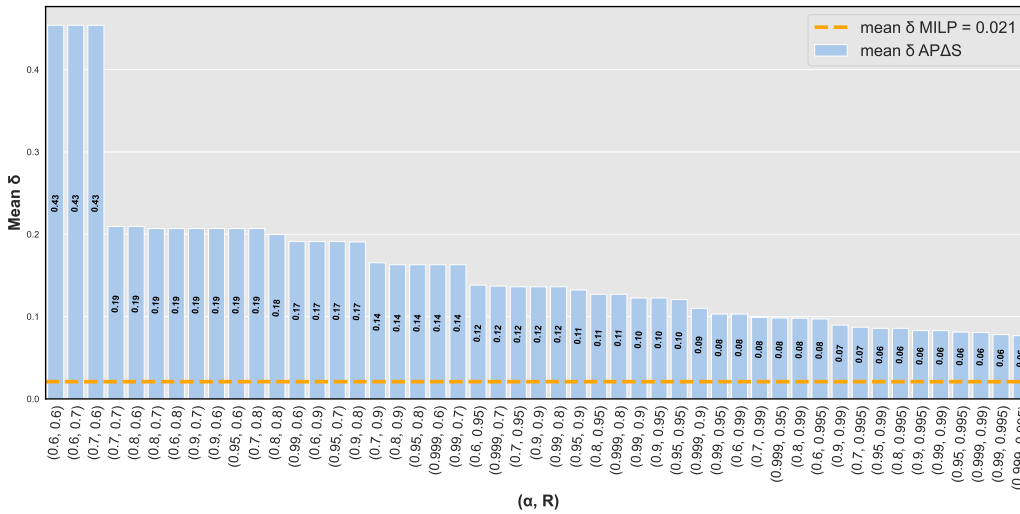


Figure A.23: Mean certifiable δ obtained for increasing α, R values using the *No2* dataset.

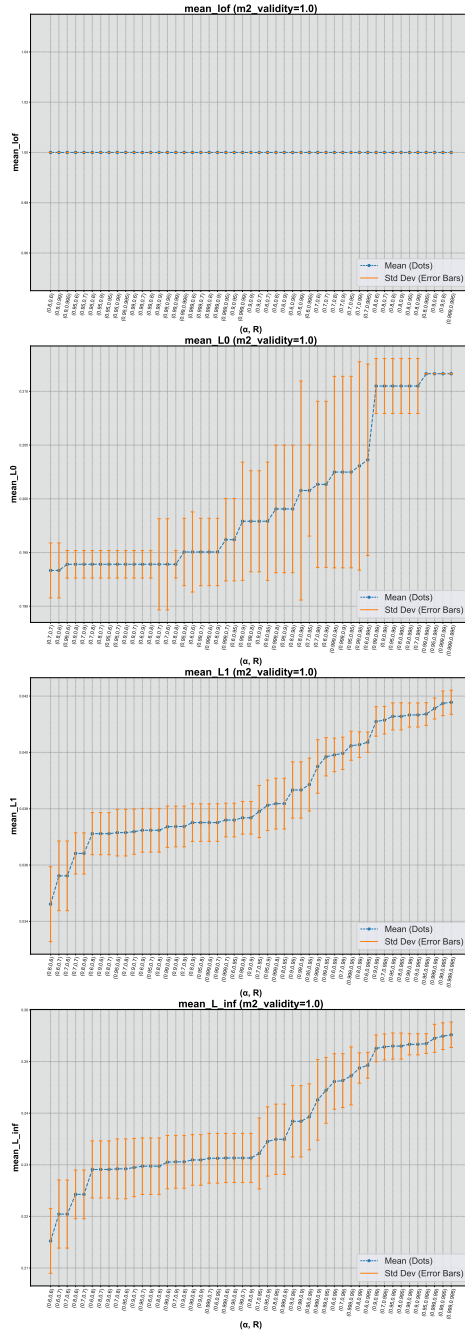


Figure A.24: Mean LOF , l_0 , l_1 , l_∞ metrics for increasing α , R values using the *Credit* dataset. CFXs always obtain 100% validity after retraining on this dataset.

1234 **Appendix B. Comparison of two CFXs generation approaches of**
 1235 **§ 6.5 experiments.**

1236 As stated in § 6.5 in the main paper, given the significantly higher number
 1237 of parameters in *TabNet*, we replace the MILP-based procedure used in
 1238 Section 6.3, Algorithm 4 with a Nearest Neighbors Counterfactual Explainer
 1239 (NNCFX) [40] to ensure scalability of our generation procedure. More specifically,
 1240 line 5 in Algorithm 4 now implements the following strategy.

Algorithm 5 *Nearest Neighbors Counterfactual Explanation*

```

1: Input: Dataset  $\mathbf{d}$ , a k-d tree built from dataset features,  $\mathbf{x}$  set of original
   inputs,  $\mathbf{y}$  set of original outcomes
2: Output:  $\mathbf{x}'$  set of nearest counterfactual explanation.
3:  $\mathbf{x}' \leftarrow \emptyset$ 
4: for  $i$  in  $len(\mathbf{x})$  do
5:    $x \leftarrow \mathbf{x}[i]$  ▷ original input
6:    $y \leftarrow \mathbf{y}[i]$  ▷ original output
7:    $y' \leftarrow 1 - y$  ▷ desired outcome
8:    $idx, distance \leftarrow \text{k-d tree.query}(x, len(features))$  ▷ already sorted per
   distance
9:   if  $\mathbf{d}[idx][\textit{outcome}'] == y'$  then
10:     $\mathbf{x}' \leftarrow \mathbf{d}[idx]$ 
11:   else
12:     $\mathbf{x}' \leftarrow None$ 
13:   end if
14: end for
15: return  $\mathbf{x}'$ 

```

1241 This function iteratively searches for a neighboring data point with the
 1242 opposite outcome by evaluating distances between features and selecting the
 1243 nearest as possible. Clearly, this approach and the one of [12] can produce
 1244 different explanations. We perform a further experiment to understand the
 1245 difference between the two CFX-generation approaches. Hence, we consider
 1246 the same datasets used in Section 6.3 and the same 50 original inputs employed
 1247 in those experiments. In the NNCFX approach, once a valid counterfactual
 1248 is found, the mean feature-wise distance between the identified
 1249 counterfactual and the MILP-generated counterfactual is calculated. This
 1250 distance serves as a measure of similarity between the counterfactuals identified
 1251 by the TabNet-based approach and those obtained via MILP in an MLP

1252 setting. Our results are reported in Fig. B.25. On the x-axis, we report the
1253 index of CFX, while on the y-axis, the mean and standard deviation distance
1254 between our CFX and the one generated with MILP in each dataset.

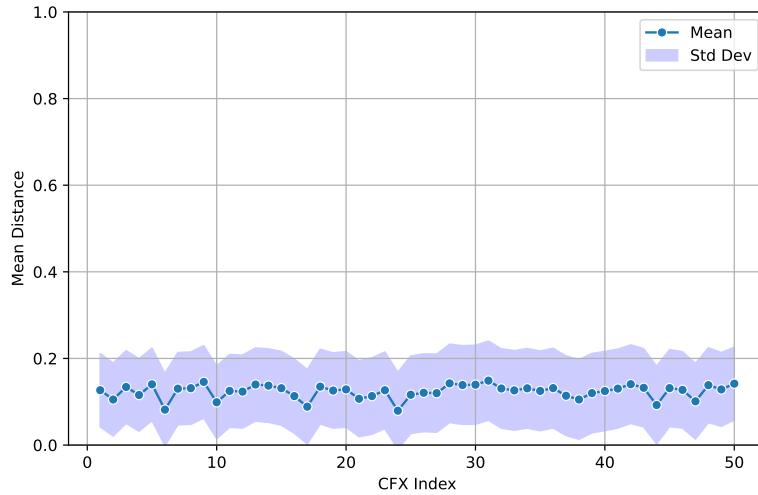


Figure B.25: Mean and standard deviation distance between CFXs generated with NNCFX and MILP in *Diabetes*, *Credit*, *SBA*, *no2* datasets.

1255 As we can notice, since the values in the datasets are typically normalized
1256 in a range $[0, 1]$, the CFXs generated with the two approaches are consistently
1257 close. In fact, there is a mean feature distance between the CFX generated
1258 with NNCFX and MILP of ~ 0.12 for the 50 inputs selected. This result
1259 shows the correctness and efficiency of the NNCFX generation approach in
1260 the transformers-based setting.