

---

# Rethinking Counterfactual Explanations as Local and Regional Counterfactual Policies

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        Among the challenges not yet resolved for Counterfactual Explanations (CE),  
2        there are stability, synthesis of the various CE and the lack of plausibility/sparsity  
3        guarantees. From a more practical point of view, recent studies show that the  
4        prescribed counterfactual recourses are often not implemented exactly by the  
5        individuals and demonstrate that most state-of-the-art CE algorithms are very likely  
6        to fail in this noisy environment. To address these issues, we propose a probabilistic  
7        framework that gives a sparse local counterfactual rule for each observation: we  
8        provide rules that give a range of values that can change the decision with a given  
9        high probability instead of giving diverse CE. In addition, the recourses derived  
10       from these rules are robust by construction. These local rules are aggregated  
11       into a regional counterfactual rule to ensure the stability of the counterfactual  
12       explanations across observations. Our local and regional rules guarantee that the  
13       recourses are faithful to the data distribution because our rules use a consistent  
14       estimator of the probabilities of changing the decision based on a Random Forest.  
15       In addition, these probabilities give interpretable and sparse rules as we select  
16       the smallest set of variables having a given probability of changing the decision.  
17       Codes for computing our counterfactual rules are available, and we compare their  
18       relevancy with standard CE and recent similar attempts.

## 19    1 Introduction

20    In recent years, many explanations methods have been developed for explaining machine learning  
21    models, with a strong focus on local analysis, i.e., generating explanations for individual prediction  
22    (see [Molnar, 2022] for a survey). Among this plethora of methods, one of the most prominent and  
23    active techniques are Counterfactual Explanations [Wachter et al., 2017b]. Unlike popular local  
24    attribution methods, e.g., SHAP [Lundberg et al., 2020] and LIME [Ribeiro et al., 2016], which  
25    highlight the importance score of each feature, Counterfactuals Explanations (CE) describe the  
26    smallest modification to the feature values that changes the prediction to a desired target. Although  
27    CE are intuitive and user-friendly by giving recourse in some scenarios (e.g., loan application), they  
28    have many shortcomings in practice. Indeed, most counterfactual methods rely on a gradient-based  
29    algorithm or heuristics approaches [Karimi et al., 2020b], thus can fail to identify the most natural  
30    explanations and lack guarantees. Most algorithms either do not guarantee sparse counterfactuals  
31    (changes in the smallest number of features) or do not generate in-distribution samples (see [Verma  
32    et al., 2020, Chou et al., 2022] for a survey on counterfactuals methods). Although some works  
33    [Parmentier and Vidal, 2021, Poyiadzi et al., 2019, Looveren and Klaise, 2019] try to solve the  
34    plausibility/sparsity problem, the suggested solutions are not entirely satisfactory.  
35    In another direction, many papers [Mohtilal et al., 2020, Karimi et al., 2020a, Russell, 2019] encour-  
36    ages the generation of diverse counterfactuals in order to find actionable recourse [Ustun et al., 2019].  
37    Actionability is a vital desideratum, as some features may be non-actionable, and generating many

38 counterfactuals increases the chance of getting actionable recourse. However, the diversity of CE  
 39 makes the explanations less intelligible, and the synthesis of various CE or local explanations, in  
 40 general, is yet to be comprehensively solved [Lakkaraju et al., 2022]. In addition, recently Pawelczyk  
 41 et al. [2022] highlights a new problem of local CE called: *noisy responses to prescribed recourses*.  
 42 Indeed, in real-world scenarios, some individuals may not be able to implement exactly the prescribed  
 43 recourses, and they show that most CE methods fail in this noisy environment. Therefore, we propose  
 44 to reverse the usual way of explaining with counterfactual by computing *Counterfactual rules*. We  
 45 introduce a new line of counterfactuals: we build interpretable policies for changing a decision with  
 46 a given probability that ensure the stability of the deduced recourse. These policies are optimal (in  
 47 sparsity) and faithful to the data distribution. Their computation comes with statistical guarantees  
 48 as they use a consistent estimator of the conditional distribution. Our proposal is to find a general  
 49 policy or rule that permits changing the decision while fixing some features instead of generating  
 50 many counterfactual samples. One of the main challenges is to identify the (minimal) set of features  
 51 that provide the best promising directions for changing the decision to the desired output. We also  
 52 show this approach can be extended for finding a collection of regional counterfactuals, such that we  
 53 have a global counterfactual policy for analyzing a model. An example of the counterfactual rules  
 that we introduce is given in figure 1.

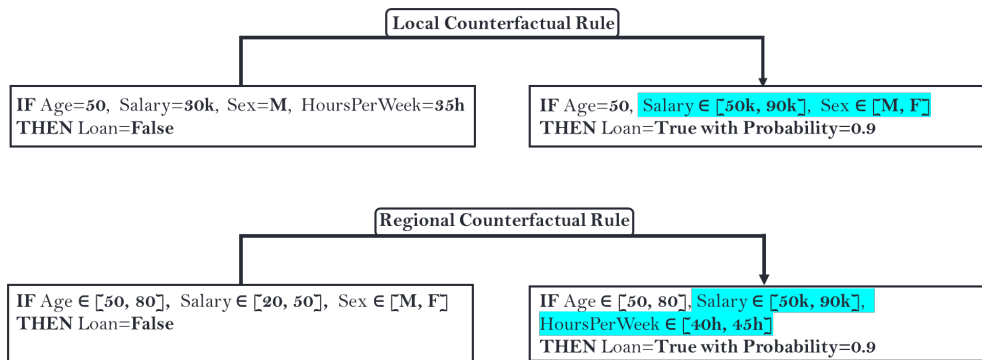


Figure 1: Illustration of the local and regional Counterfactuals Rules that we introduced on a dataset with 4 variables: Age, Salary, Sex, and HoursPerWeek. The Counterfactual Rules define intervals on the minimal subset of features to change the decision of a model prediction in the local counterfactual rule or the decision of a rule that applies on a sub-population in the regional counterfactual rule. In Blue, we have the proposed rules to change the decision.

54

## 55 2 Motivation and Related works

56 Most of the methods that propose Counterfactuals Explanations are based on the approach of the  
 57 seminal work of Wachter et al. [2017a]: the counterfactuals are generated by optimizing a cost, but  
 58 this procedure does not account directly the plausibility of the counterfactual examples (see [Verma  
 59 et al., 2020] for classification of CE methods). Indeed, a major shortcoming is that the adverse  
 60 decision needed for obtaining the counterfactual is not designed to be feasible or representative of the  
 61 underlying data distribution. However, some recent studies proposed ad-hoc plausibility constraint  
 62 in the optimization, using for instance an outlier score [Kanamori et al., 2020], an Isolation Forest  
 63 [Parmentier and Vidal, 2021] or a density-weighted metrics [Poyiadzi et al., 2019] to generate in-  
 64 distribution samples. In another direction, Looveren and Klaise [2019] proposes to use an autoencoder  
 65 that penalizes out-of-distribution candidates. Instead of relying on ad-hoc constraints, we propose CE  
 66 that gives plausible explanations by design. Indeed, for each observation, we identify the variables  
 67 and associated ranges of values that have the highest probability of changing the prediction. We can  
 68 compute this probability with a consistent estimator of the conditional distribution  $P(Y|X_S)$ . As a  
 69 consequence, the sparsity of the counterfactuals is not encouraged indirectly by adding a penalty term  
 70 ( $\ell_0$  or  $\ell_1$ ) as existing works [Mothilal et al., 2020]. Our approach is inspired by the concept of *Same*  
 71 *Decision Probability (SDP)* (introduced in [Chen et al., 2012]) that can be used for identifying the  
 72 smallest subset of features to guarantee (with a given probability) the stability of a prediction. This  
 73 minimal subset is called *Sufficient Explanations*. In [Amoukou and Brunel, 2021], it has been shown

74 that the *SDP* and the *Sufficient Explanations* can be estimated and computed efficiently for identifying  
75 important local variables in any classification and regression models. For counterfactuals, we are  
76 interested in the dual set: we want the minimal subset of features that have a high probability of  
77 changing the decision (when the other features are fixed). Another limitation of the current CE is their  
78 local nature and the multiplicity of the explanations produced. While some papers [Mothilal et al.,  
79 2020, Karimi et al., 2020a, Russell, 2019] promote the generation of diverse counterfactual samples  
80 to ensure actionable recourse, such diverse explanations should be summarized to be intelligible  
81 [Lakkaraju et al., 2022], but the compilation of local explanations is often a very difficult problem. To  
82 address this problem, we do not generate counterfactual samples, but we build a rule *Counterfactual*  
83 *Rules* (CR) from which we can derive counterfactuals. Contrary to classic CE which gives the nearest  
84 instances with a desired output, we find the most effective rule for each observation (or group of similar  
85 observations) that changes the prediction to the desired target. This local rule easily aggregates similar  
86 counterfactuals. For example, if  $x = \{\text{Age}=20, \text{Salary}=35\text{k}, \text{HoursWeek}=25\text{h}, \text{Sex}=\text{M}, \dots\}$   
87 with  $\text{Loan}=\text{False}$ , fixing the variables Age and Sex and changing the Salary and HoursWeek  
88 change the decision. Therefore, instead of given multiples combination of Salary and HoursWeek  
89 (e.g. 35k and 40h or 40k and 55h, ...) that result in many instances, the counterfactual  
90 rule gives the range of values: IF  $\text{HoursWeek} \in [35\text{h}, 50\text{h}]$ ,  $\text{Salary} \in [40\text{k}, 50\text{k}]$ , and  
91 the **remaining features are fixed** THEN  $\text{Loan}=\text{True}$ . It can be extended at a regional scale,  
92 e.g., given a rule  $\mathbf{R} = \{\text{IF } \text{Salary} \in [35\text{k}, 20\text{k}], \text{Age} \in [20, 80] \text{ THEN } \text{Loan}=\text{False}\}$ ,  
93 the regional Counterfactual Rule (CR) could be  $\{\text{IF } \text{Salary} \in [40\text{k}, 50\text{k}], \text{HoursWeek} \in$   
94  $[35\text{h}, 50\text{h}]$  and the **remaining rules are fixed** THEN  $\text{Loan}=\text{True}\}$ . The main difference be-  
95 tween a local and a global CR is that the Local-CR explain a single instance by fixing the remaining  
96 feature values (not used in the CR) ; while a regional-CR is defined by keeping the remaining variables  
97 in a given interval (not used in the regional-CR). Moreover, by giving ranges of values that guarantee  
98 a high probability of changing the decision, we partly answer the problem of *noisy responses to*  
99 *prescribed recourses* [Pawelczyk et al., 2022] so long as the perturbations are within our ranges.

100 Although the *Local Counterfactual Rule* is new, the *Regional Counterfactual Rule* can be related to  
101 some recent works. Indeed, Rawal and Lakkaraju [2020] proposed Actionable Recourse Summaries  
102 (AREs), a framework that constructs global counterfactual recourses in order to have a global insight  
103 of the model and detect unfair behavior. While AREs is similar to the Regional Counterfactual  
104 Rule, we emphasize some significant differences. Our methods can address regression problems and  
105 deal with continuous features. Indeed, AREs needs to discretize the continuous features, inducing a  
106 trade-off between speed and performance as noticed by [Ley et al., 2022]. Thus, too few bins result  
107 in unrealistic recourse, while too many bins result in excessive computation time. In addition, AREs  
108 uses a greedy heuristic search approach to find global recourse, which might produce sub-optimal  
109 recourse. As we have already mentioned, the changes we provide overcome these two limitations  
110 because the consistency of our counterfactual is controlled by an estimation of the probability of  
111 changing the decision, and because we favor changes of a minimum number of features. Another  
112 global CE framework has been introduced in [Kanamori et al., 2022] to ensure transparency: the  
113 Counterfactual Explanation Tree (CET) partitions the input space with a decision tree and assigns  
114 an appropriate action for changing the decision of each subspace. Therefore, it gives a unique  
115 action for changing the decision of multiple instances. In our case, we offer more flexibility in the  
116 counterfactual explanations because we provide a range of possible values that guarantee a change  
117 with a given probability. In our approach, we do not make any assumption about the cost of changing  
118 the feature nor the causal structure. If we have such information, then we can add it as additional  
119 post-processing such that it can be made more explicit and more transparent for the final user as  
120 required for trustworthy AI.

### 121 3 Minimal Counterfactual Rules

122 We assume that we have an i.i.d sample  $\mathcal{D}_n = \{(\mathbf{X}_i, Y_i)_{i=1, \dots, n}\}$  such that  $(\mathbf{X}, Y) \sim P_{(\mathbf{X}, Y)}$  where  
123  $\mathbf{X} \in \mathcal{X}$  (typically  $\mathcal{X} = \mathbb{R}^p$ ) and  $Y \in \mathcal{Y}$ . The output  $\mathcal{Y}$  can be discrete or continuous. We want to  
124 explain the predictor  $f : \mathbb{R}^p \mapsto \mathcal{Y}$ , that has been learned with the dataset  $\mathcal{D}_n$ . We use uppercase  
125 letters for random variables and lowercase letters for their value assignments. For a given subset  
126  $S \subset [p]$ ,  $\mathbf{X}_S = (X_i)_{i \in S}$  denotes a subgroup of features, and we write  $\mathbf{x} = (x_S, \mathbf{x}_{\bar{S}})$  (with some  
127 abuse of notation).

128 For an observation  $(x, y = f(x))$ , we have a target set  $\mathcal{Y}^* \subset \mathcal{Y}$ , such that  $y \notin \mathcal{Y}^*$ . For the simple  
129 case of classification problem,  $\mathcal{Y}^* = \{y^*\}$  is the standard singleton such that  $y^* \in \mathcal{Y}$  is different of  
130  $y$ . Contrary to standard approaches, our definition of the counterfactual deals also with the regression  
131 case by considering  $\mathcal{Y}^* = [a, b] \subset \mathbb{R}$ ; our definitions and computations of counterfactuals are the  
132 same for both classification and regression. We remind that the classic CE problem (defined only for  
133 classification) is to find a function  $\mathbf{a} : \mathcal{X} \mapsto \mathcal{X}$ , such that for all observations  $x \in \mathcal{X}$ ,  $f(x) \neq y^*$ ,  
134 and we have  $f(\mathbf{a}(x)) = y^*$ . With standard CE, the function is defined point-wise by solving an  
135 optimisation program. Most often  $\mathbf{a}(\cdot)$  is not a real function, as  $\mathbf{a}(x)$  may be in fact a collection of  
136 (random) values  $\{\mathbf{x}_1^*, \dots, \mathbf{x}_p^*\}$ . A more recent point of view was proposed by Kanamori et al. [2022],  
137 and it defines  $\mathbf{a}$  as a decision tree, where in each leaf  $L$ , the best perturbation  $a_L$  is predicted and add  
138 it to all the instances  $x \in L$ .

139 Our approach is hybrid, because we do not propose a single action for each subspace of  $\mathcal{X}$  or sub-group  
140 of population, but we give sets of possible perturbations. Indeed, a *Local Counterfactual Rule* (Local-  
141 CR) for  $\mathcal{Y}^*$  and observation  $x$  (with  $f(x) \notin \mathcal{Y}^*$ ) is a rectangle  $C_S(x; \mathcal{Y}^*) = \prod_{i \in S} [a_i, b_i]$ ,  $a_i, b_i \in$   
142  $\overline{\mathbb{R}}$  such that for all perturbations of  $x = (x_S, x_{\bar{S}})$  obtained as  $x^* = (z_S, x_{\bar{S}})$  with  $z_S \in C_S(x; \mathcal{Y}^*)$   
143 and  $x^*$  an in-distribution sample, then  $f(x^*)$  is in  $\mathcal{Y}^*$  with a high probability.

144 Similarly, a *Regional Counterfactual Rule* (Regional-CR)  $C_S(\mathbf{R}; \mathcal{Y}^*)$  is defined for  $\mathcal{Y}^*$  and a  
145 rectangle  $\mathbf{R} = \prod_{i=1}^d [a_i, b_i]$ ,  $a_i, b_i \in \overline{\mathbb{R}}$ , if for all observations  $x = (x_S, x_{\bar{S}}) \in \mathbf{R}$ , the perturbations  
146 obtained as  $x^* = (z_S, x_{\bar{S}})$  with  $z_S \in C_S(\mathbf{R}, \mathcal{Y}^*)$  and  $x^*$  an in-distribution sample are such that  
147  $f(x^*)$  is in  $\mathcal{Y}^*$  with high probability.

148 We build such rectangles sequentially, first, we propose to find the best directions  $S \subset [p]$  that offers  
149 the best probability of change. Then, we find the best intervals  $[a_i, b_i]$ ,  $i \in S$  that change the decision  
150 to the desired target. A central tool in this approach is the Counterfactual Decision Probability.

151 **Definition 3.1. Counterfactual Decision Probability (CDP).** The Counterfactual Decision Prob-  
152 ability of the subset  $S \subset [1, p]$ , w.r.t  $x = (x_S, x_{\bar{S}})$  and the desired target  $\mathcal{Y}^*$  (s.t.  $f(x) \notin \mathcal{Y}^*$ )  
153 is

$$CDP_S(\mathcal{Y}^*; x) = P(f(\mathbf{X}) \in \mathcal{Y}^* | \mathbf{X}_{\bar{S}} = x_{\bar{S}}).$$

154 The CDP of the subset S is the probability that the decision changes to the desired target  $\mathcal{Y}^*$   
155 by sampling the features  $\mathbf{X}_S$  given  $\mathbf{X}_{\bar{S}} = x_{\bar{S}}$ . It is related to the Same Decision Probability  
156  $SDP_S(\mathcal{Y}; x) = P(f(\mathbf{X}) \in \mathcal{Y} | \mathbf{X}_S = x_S)$  used in [Amoukou and Brunel, 2021] for solving the  
157 dual problem of selecting the most local important variables for obtaining and maintaining the decision  
158  $f(x) \in \mathcal{Y}$  (where  $f(x) \in \mathcal{Y} \subset \mathcal{Y}$ ). The set  $S$  is called the Minimal Sufficient Explanation. Indeed,  
159 we have  $CDP_S(\mathcal{Y}^*; x) = SDP_{\bar{S}}(\mathcal{Y}^*; x)$ . The computation of these probabilities is challenging  
160 and discussed in Section 4. We now focus on the minimal subset of features  $S$  such that the model  
161 makes the desired decision with a given probability  $\pi$ .

162 **Definition 3.2. (Minimal Divergent Explanations).** Given an instance  $x$  and a desired target  $\mathcal{Y}^*$ ,  
163  $S$  is a Divergent Explanation for probability  $\pi > 0$ , if  $CDP_S(\mathcal{Y}^*; x) \geq \pi$ , and no subset  $Z$  of  $S$   
164 satisfies  $CDP_Z(\mathcal{Y}^*; x) \geq \pi$ . Hence, a Minimal Divergent Explanation is a Divergent Explanation  
165 with minimal size.

166 The set minimizing this probability is not unique, and we can have several Minimal Divergent  
167 Explanations. Note that the probability  $\pi$  represents the minimum level required for a set to be chosen  
168 for generating counterfactuals, and its value should be as high as possible and depends on the use  
169 case. We have now enough material to define our main criterion for building a Local Counterfactual  
170 Rule (Local-CR):

171 **Definition 3.3. (Local Counterfactual Rule).** Given an instance  $x$ , a desired target  $\mathcal{Y}^* \not\ni f(x)$ , a  
172 Minimal Divergent Explanation  $S$ , the rectangle  $C_S(x; \mathcal{Y}^*) = \prod_{i \in S} [a_i, b_i]$ ,  $a_i, b_i \in \overline{\mathbb{R}}$  is a Local  
173 Counterfactual Rule with probability  $\pi_C$  if

$$CRP_S(\mathcal{Y}^*, x, C_S(x; \mathcal{Y}^*)) \triangleq P(f(\mathbf{X}) \in \mathcal{Y}^* | \mathbf{X}_S \in C_S(x; \mathcal{Y}^*), \mathbf{X}_{\bar{S}} = x_{\bar{S}}) \geq \pi_C. \quad (3.1)$$

174 The  $CRP_S$  is the Counterfactual Rule Probability.

175 The higher the probability  $\pi_C$  is, the better the relevance of the rule  $C_S(x; \mathcal{Y}^*)$  is, for this instance.  
176 Given a set  $S$ , we seek for the maximal rectangle in the direction  $S$  satisfying Definition 3.1.

177 In practice, we can observe that the Local-CR  $C_S(\cdot; \mathcal{Y}^*)$  for neighbors  $x, x'$  are often quite close, be-  
178 cause the Minimal Divergent Explanations are similar and the corresponding rectangles often overlaps.

179 Hence, this motivates a generalisation of these Local-CR to hyperrectangle  $\mathbf{R} = \prod_{i=1}^d [a_i, b_i]$ ,  $a_i, b_i \in$   
180  $\mathbb{R}$  regrouping similar observations. We denote  $\text{supp}(\mathbf{R}) = \{i : [a_i, b_i] \neq \mathbb{R}\}$  the support of the  
181 rectangle, and we extend the Local-CR to Regional Counterfactual Rules (Regional-CR). In order  
182 to do it, we denote  $\mathbf{R}_{\bar{S}} = \prod_{i \in \bar{S}} [a_i, b_i]$  as the rectangle with intervals of  $\mathbf{R}$  in  $\text{supp}(\mathbf{R}) \cap \bar{S}$  and we  
183 also defines the corresponding Counterfactual Decision Probability CDP (Definition 3.1) for rule  $\mathbf{R}$   
184 and subset  $S$  as  $CDP_S(\mathcal{Y}^*; \mathbf{R}) = P(f(\mathbf{X}) \in \mathcal{Y}^* | \mathbf{X}_{\bar{S}} \in \mathbf{R}_{\bar{S}})$ . Therefore, we can also compute  
185 the Minimal Divergent Explanation for rule  $\mathbf{R}$  using Definition 3.2 with the CDP for rules.

186 **Definition 3.4. (Regional Counterfactual Rule).** Given any rectangle  $\mathbf{R}$ , a desired target  $\mathcal{Y}^*$ ,  
187 a Minimal Divergent Explanation  $S$  of  $\mathbf{R}$ , the rectangle  $C_S(\mathbf{R}; \mathcal{Y}^*) = \prod_{i \in S} [a_i, b_i]$  is a Regional  
188 Counterfactual Rule with probability  $\pi_C$  if

$$CRP_S(\mathcal{Y}^*; \mathbf{R}, C_S(\mathbf{R}, \mathcal{Y}^*)) \triangleq P(f(\mathbf{X}) \in \mathcal{Y}^* | \mathbf{X}_S \in C_S(\mathbf{R}, \mathcal{Y}^*), \mathbf{X}_{\bar{S}} \in \mathbf{R}_{\bar{S}}) \geq \pi_C. \quad (3.2)$$

189  $CRP_S(\mathcal{Y}^*; \mathbf{R}, C_S(\mathbf{R}))$  is the corresponding Counterfactual Rule Probability for rule  $\mathbf{R}$ .

190 **Remarks:** Local-CR and regional-CR differ slightly: for local, we condition by  $\mathbf{X}_{\bar{S}} = \mathbf{x}_{\bar{S}}$  in Eq.  
191 3.1, while for regional, we condition by  $\mathbf{X}_{\bar{S}} \in \mathbf{R}_{\bar{S}}$ . For computing regional-CR, we can start for a  
192 rectangle generated by any method, such as [Wang et al., 2017, Lin et al., 2020]. The only condition  
193 is that it contains a homogeneous group, i.e. with almost the same output. However, by default  
194 we use as initial rules the Sufficient Rules derived in [Amoukou and Brunel, 2021] as it handles  
195 regression problem. The Sufficient Rules are minimal support rectangles define for a given output  $\mathcal{Y}$   
196 as  $C_S(\mathcal{Y}) = \prod_{i \in S} [a_i, b_i]$  such that  $\forall \mathbf{x} \in \mathcal{X}, \mathbf{x}_S \in C_S(\mathcal{Y}), P(f(\mathbf{X}) \in \mathcal{Y} | \mathbf{X}_S = \mathbf{x}_S) \geq \pi$ .

## 197 4 Estimation of the CDP and CRP

198 In order to compute the probabilities  $CDP_S$  and  $CRP_S$  for any  $S$ , we use a dedicated Random  
199 Forest (RF)  $m_{k,n}$  that learns the model  $f$  to explain. Indeed, the conditional probabilities  $CDP_S$   
200 and  $CRP_S$  can be easily computed from a RF by combining the Projected Forest algorithm [Bénard  
201 et al., 2021a] and the Quantile Regression Forest [Meinshausen and Ridgeway, 2006]: hence we can  
202 estimate consistently the probabilities  $CDP_S(\mathcal{Y}^*; \mathbf{x})$ . We adapt the approach used in [Amoukou and  
203 Brunel, 2021] and remind for the sake of completeness, the computation of the estimate of  $SDP_S$ .

### 204 4.1 Projected Forest and $CDP_S$

205 The estimator of the  $SDP_S$  is built upon a learned Random Forest [Breiman et al., 1984]. A Random  
206 Forest (RF) is a predictor consisting of a collection of  $k$  randomized trees (see [Loh, 2011] for a  
207 detailed description of decision tree). For each instance  $\mathbf{x}$ , the predicted value of the  $j$ -th tree is  
208 denoted  $m_n(\mathbf{x}, \Theta_j)$  where  $\Theta_j$  represents the resampling data mechanism in the  $j$ -th tree and the  
209 successive random splitting directions. The trees are then averaged to give the prediction of the forest  
210 as:

$$m_{k,n}(\mathbf{x}, \Theta_{1:k}, \mathcal{D}_n) = \frac{1}{k} \sum_{l=1}^k m_n(\mathbf{x}; \Theta_l, \mathcal{D}_n) \quad (4.1)$$

211 However, the RF can also be view as an adaptive nearest neighbor predictor. For every instance  $\mathbf{x}$ ,  
212 the observations in  $\mathcal{D}_n$  are weighted by  $w_{n,i}(\mathbf{x}; \Theta_{1:k}, \mathcal{D}_n)$ ,  $i = 1, \dots, n$ . Therefore, the prediction  
213 of RF can be rewritten as

$$m_{k,n}(\mathbf{x}, \Theta_{1:k}, \mathcal{D}_n) = \sum_{i=1}^n w_{n,i}(\mathbf{x}; \Theta_{1:k}, \mathcal{D}_n) Y_i.$$

214 This emphasizes the central role played by the weights in the RF's algorithm, see [Meinshausen and  
215 Ridgeway, 2006, Amoukou and Brunel, 2021] for detailed description of the weights. Therefore,  
216 it naturally gives estimators of other quantities e.g., Cumulative hazard function [Ishwaran et al.,  
217 2008], Treatment effect [Wager and Athey, 2017], conditional density [Du et al., 2021]. For instance,  
218 Meinshausen and Ridgeway [2006] showed that we can used the same weights to estimate the  
219 Conditional Distribution Function with the following estimator:

$$\hat{F}(y | \mathbf{X} = \mathbf{x}, \Theta_{1:k}, \mathcal{D}_n) = \sum_{i=1}^n w_{n,i}(\mathbf{x}; \Theta_{1:k}, \mathcal{D}_n) \mathbb{1}_{Y_i \leq y} \quad (4.2)$$

220 In another direction, Bénard et al. [2021a] introduced the Projected Forest algorithm [Bénard et al.,  
221 2021c,a] that aims to estimate  $E[Y | \mathbf{X}_S]$  by modifying the RF's prediction algorithm.

222 **Projected Forest:** To estimate  $E[Y|\mathbf{X}_S = \mathbf{x}_S]$  instead of  $E[Y|\mathbf{X} = \mathbf{x}]$  using a RF, [Bénard et al.](#)  
223 [\[2021b\]](#) suggests to simply ignore the splits based on the variables not contained in  $S$  from the  
224 tree predictions. More formally, it consists of projecting the partition of each tree of the forest on  
225 the subspace spanned by the variables in  $S$ . The authors also introduced an algorithmic trick that  
226 computes the projected partition efficiently without modifying the initial tree structures. We drop  
227 observations down in the initial trees, ignoring the splits which use a variable not in  $S$ : when a  
228 split involving a variable outside of  $S$  is met, the observations are sent both to the left and right  
229 children nodes. Therefore, each instance falls in multiple terminal leaves of the tree. We drop the  
230 new query point  $\mathbf{x}_S$  down the tree, following the same procedure, and gather the set of terminal  
231 leaves where  $\mathbf{x}_S$  falls. Next, we collect the training observations which belong to every terminal leaf  
232 of this collection, in other words, we keep only the observations that fall in the intersection of the  
233 leaves where  $\mathbf{x}_S$  falls. Finally, we average the outputs  $Y_i$  of the selected training points to generate  
234 the estimation of  $E[Y|\mathbf{X}_S = \mathbf{x}_S]$ . Notice that this algorithm converges asymptotically to the true  
235 projected conditional expectation  $E[Y|\mathbf{X}_S = \mathbf{x}_S]$ .

236 As the RF, the PRF gives also a weight to each observation. The associated PRF is denoted  
237  $m_{k,n}^{(\mathbf{x}_S)}(\mathbf{x}_S) = \sum_{i=1}^n w_{n,i}(\mathbf{x}_S)Y_i$ . Therefore, as the weights of the original forest was used to  
238 estimate the CDF in equation 4.2, [Amoukou and Brunel \[2021\]](#) used the weights of the Projected  
239 Forest Algorithm to estimate the  $SDP$  as  $\widehat{SDP}_S(\mathcal{Y}; \mathbf{x}) = \sum_{i=1}^n w_{n,i}(\mathbf{x}_S)\mathbb{1}_{Y_i \in \mathcal{Y}}$ . The idea is  
240 essentially to replace  $Y_i$  by  $\mathbb{1}_{Y_i \in \mathcal{Y}}$  in the Projected Forest equation defined above. The authors also  
241 show that this estimator converges asymptotically to the true  $SDP_S$ . Therefore, we can estimate the  
242  $CDP$  with the following estimator

$$\widehat{CDP}_S(\mathcal{Y}^*; \mathbf{x}) = \sum_{i=1}^n w_{n,i}(\mathbf{x}_S)\mathbb{1}_{Y_i \in \mathcal{Y}^*}. \quad (4.3)$$

243 **Remarks:** Note that we only give the estimator of the  $CDP_S$  of an instance  $\mathbf{x}$ . The estimator of the  
244  $CDP_S$  of a rule  $R$  will be discussed in the next section as it is related to the estimator of the  $CRP_S$ .

## 245 4.2 Regional RF and $CRP_S$

246 In this section, we focus on the estimation of the  $CRP_S(\mathcal{Y}^*, \mathbf{x}, C_S(\mathbf{x}; \mathcal{Y}^*)) = P(f(\mathbf{X}) \in$   
247  $\mathcal{Y}^* | \mathbf{X}_S \in C_S(\mathbf{x}; \mathcal{Y}^*), \mathbf{X}_{\bar{S}} = \mathbf{x}_{\bar{S}})$  and  $CRP_S(\mathcal{Y}^*, \mathbf{R}, C_S(\mathbf{R}; \mathcal{Y}^*)) = P(f(\mathbf{X}) \in \mathcal{Y}^* | \mathbf{X}_S \in$   
248  $C_S(\mathbf{R}; \mathcal{Y}^*), \mathbf{X}_{\bar{S}} \in \mathbf{R}_{\bar{S}})$ . For simplicity, we remove the dependency of the rectangles in  $\mathcal{Y}^*$ . Based  
249 on the previous Section, we already know that the estimators using the RF will be in the form of  
250  $\widehat{CRP}_S(\mathcal{Y}^*, \mathbf{x}, C_S(\mathbf{x})) = \sum_{i=1}^n w_{n,i}(\mathbf{x})\mathbb{1}_{Y_i \in \mathcal{Y}^*}$ , thus we only need to find the right weighting.  
251 The main challenge is that we have a condition based on a region, e.g.,  $\mathbf{X}_S \in C_S(\mathbf{x})$  or  $\mathbf{X}_{\bar{S}} \in \mathbf{R}_{\bar{S}}$   
252 (regional-based) instead of condition of type  $\mathbf{X}_S = \mathbf{x}_S$  (fixed value-based) as usually. However, we  
253 introduced a natural generalization of the RF algorithm to make predictions when the conditions  
254 are both regional-based and fixed value-based. Thus, the case where there are only regional-based  
255 conditions are naturally derived.

256 **Regional RF to estimate  $CRP_S(\mathcal{Y}^*, \mathbf{x}, C_S(\mathbf{x})) = P(f(\mathbf{X}) \in \mathcal{Y}^* | \mathbf{X}_S \in C_S(\mathbf{x}), \mathbf{X}_{\bar{S}} = \mathbf{x}_{\bar{S}})$ :**  
257 The algorithm is based on a slight modification of RF. Its works as follow: we drop the observations  
258 in the initial trees, if a split used variable  $i \in \bar{S}$ , i.e., fixed value-based condition, we use the  
259 classic rules of RF, if  $x_i \leq t$ , the observations go to the left children, otherwise the right children.  
260 However, if a split used variable  $i \in S$ , i.e., regional-based condition, we use the rectangles  $C_S(\mathbf{x}) =$   
261  $\prod_{i=1}^{|\bar{S}|} [a_i, b_i]$ . The observations are sent to the left children if  $b_i \leq t$ , right children if  $a_i > t$  and  
262 if  $t \in [a_i, b_i]$  the observations are sent both to the left and right children. Therefore, we use the  
263 weights of the Regional RF algorithm to estimate the  $CRP_S$  as in equation 4.3, the estimator is  
264  $\widehat{CRP}_S(\mathcal{Y}^*; \mathbf{x}, C_S(\mathbf{x})) = \sum_{i=1}^n w_{n,i}(\mathbf{x})\mathbb{1}_{Y_i \in \mathcal{Y}^*}$ . A more detailed version of the algorithm is provided  
265 and discussed in Appendix.

266 To estimate the  $CDP$  of a rule  $CDP_S(\mathcal{Y}^*; \mathbf{R}) = P(f(\mathbf{X}) \in \mathcal{Y}^* | \mathbf{X}_{\bar{S}} \in \mathbf{R}_{\bar{S}})$ , we just have to  
267 apply the projected Forest algorithm to the Regional RF, i.e., when a split involving a variable outside  
268 of  $\bar{S}$  is met, the observations are sent both to the left and right children nodes, otherwise we use the  
269 Regional RF split rule, i.e., if an interval of  $\mathbf{R}_{\bar{S}}$  is below  $t$ , the observations go to the left children,  
270 otherwise the right children and if  $t$  is in the interval, the observations go to the left and right children.

271 The estimator of the  $CRP_S(\mathcal{Y}^*; \mathbf{R}, C_S(\mathbf{R}))$  for rule is also derived from the Regional RF. Indeed, it  
 272 is a special case of the Regional RF algorithm where there are only regional-based conditions.

## 273 5 Learning the Counterfactual Rules

274 We compute the Local and Regional CR using the estimators of the previous section. First, we find  
 275 the Minimal Divergent Explanation in the same way as Minimal Sufficient Explanation can be found  
 276 [Amoukou and Brunel, 2021]. As the exploration of all possible subsets is exponential, we search  
 277 the Minimal Divergent Subset among the  $K = 10$  most frequently selected variables in the RF  $m_{k,n}$   
 278 used to estimate the probabilities  $CDP_S, CRP_S$  ( $K$  is an hyper-parameter to select according to the  
 279 use case and computational power). We can also use any importance measure.

280 Given an instance  $\mathbf{x}$  or rectangle  $\mathbf{R}$  (and set  $\mathcal{Y}^*$ ) and their corresponding Minimal Divergent  
 281 Explanation  $S$ , we want to find a rule  $C_S(\mathbf{x}) = \prod_{i \in S} [a_i, b_i]$  s.t. given  $\mathbf{X}_{\bar{S}} = \mathbf{x}_{\bar{S}}$  or  $\mathbf{X}_{\bar{S}} \in \mathbf{R}_{\bar{S}}$  and  
 282  $\mathbf{X}_S \in C_S(\mathbf{x})$ , the probability that  $Y \in \mathcal{Y}^*$  is high. More formally, we want:  $P(f(\mathbf{X}) \in \mathcal{Y}^* | \mathbf{X}_S \in$   
 283  $C_S(\mathbf{x}), \mathbf{X}_{\bar{S}} = \mathbf{x}_{\bar{S}})$  or  $P(f(\mathbf{X}) \in \mathcal{Y}^* | \mathbf{X}_S \in C_S(\mathbf{x}), \mathbf{X}_{\bar{S}} \in \mathbf{R}_{\bar{S}})$  above  $\pi_C$ .

284 The computation of the rectangles  $C_S(\mathbf{x}) = \prod_{i \in S} [a_i, b_i]$  relies heavily on our use of RF and on the  
 285 algorithmic trick of the projected RF. Indeed, the rectangles defining the rules arise naturally from RF,  
 286 while AREs [Rawal and Lakkaraju, 2020] relies on binned variables to generate candidate rules and  
 287 tests all these possible rules for choosing an optimal one. We overcome the computational burden  
 288 and the challenge of choosing the number of bins.

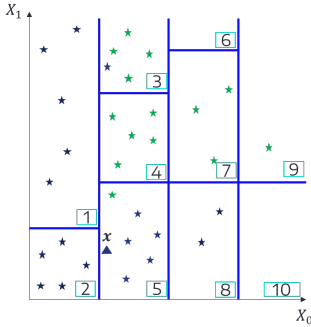


Figure 2: The partition of the RF learned to classify the toy data (Green/Blue stars). It has 10 leaves. The explainee  $\mathbf{x}$  is the Blue triangle in leaf 5.

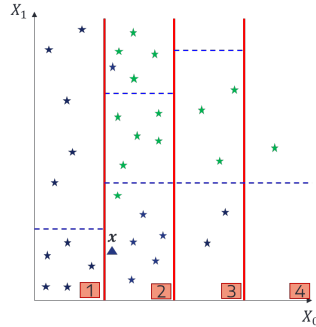


Figure 3: The partition of the projected Forest when we condition on  $X_0$ , i.e., ignoring the splits based on  $X_1$  (the dashed lines).

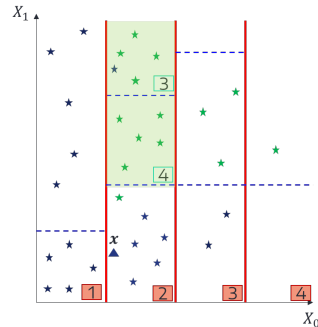


Figure 4: The optimal CR for  $\mathbf{x}$  when we condition given  $X_0 = x_0$  is the Green region, its corresponds to the union of leaf 3 and 4 of the forest

289 To illustrate the idea, we use a two-dimensional data  $(X_0, X_1)$  with label  $Y$  represented as Green/Blue  
 290 stars in figure 2. We fit a Random Forest to classify this dataset and show its partition in figure 2. The  
 291 explainee  $\mathbf{x}$  is the Blue triangle observation.

292 By looking at the different cells/leaves of the RF, we can guess that the Minimal Divergent Explanation  
 293 of  $\mathbf{x}$  is  $S = X_1$ . Indeed, in figure 3, we observe the leaves of the Projected Forest when we do not  
 294 condition on  $S = X_1$ , thus projected the RF's partition only on the subspace  $X_0$ . It consists of  
 295 ignoring all the splits in the other directions (here the  $X_1$ -axis), thus  $\mathbf{x}$  falls in the projected leaf 2  
 296 (see figure 3) and its  $CDP$  is  $CDP_{X_1}(\text{Green}; \mathbf{x}) = \frac{10 \text{ Green}}{10 \text{ Green} + 17 \text{ Blue}} = 0.58$ .

297 Finally, the problem of finding the optimal rectangle  $C_S(\mathbf{x}) = [a_i, b_i]$  in the direction of  $X_1$  s.t. the  
 298 decision changes can be easily solved by using the leaves of the RF. In fact, by looking at the leaves  
 299 of the RF (figure 2) of the observations that belong in the projected RF leaf 2 (figure 3) where  $\mathbf{x}$  falls,  
 300 we see in figure 4 that the optimal rectangle to change the decision given  $X_0 = x_0$  or being in the  
 301 projected RF leaf 2 is the union of the intervals on  $X_1$  of the leaf 3 and 4 of the RF (see the Green  
 302 region of figure 4).

303 Given an instance  $\mathbf{x}$  and its Minimal Divergent Explanation  $S$ , the first step is the collect of the  
 304 observations which belong to the leaf of the Projected Forest given  $\bar{S}$  where  $\mathbf{x}$  falls. It corre-  
 305 sponds to the observations that has positive weights in the computation of the  $CDP_S(\mathcal{Y}^*; \mathbf{x}) =$

306  $\sum_{i=1}^n w_{n,i}(\mathbf{x}_{\bar{S}}) \mathbb{1}_{Y_i \in \mathcal{Y}^*}$ , i.e.,  $\{\mathbf{x}_i : w_{n,i}(\mathbf{x}_{\bar{S}}) > 0\}$ . Then, we used the partition of the original forest  
 307 to find the possible leaves  $C_S(\mathbf{x})$  in the direction  $S$ . The possible leaves is among the RF’s leaves  
 308 of the collected observations  $\{\mathbf{x}_i : w_{n,i}(\mathbf{x}_{\bar{S}}) > 0\}$ . Let denote  $L(\mathbf{x}_i)$  the leaves of the observations  
 309  $\mathbf{x}_i$  with  $w_{n,i}(\mathbf{x}_{\bar{S}}) > 0$ . A possible leaf is a leaf  $L(\mathbf{x}_i)$  s.t.  $CRP_S(\mathcal{Y}^*, \mathbf{x}, L(\mathbf{x}_i)_S) = P(f(\mathbf{X}) \in$   
 310  $\mathcal{Y}^* | \mathbf{X}_S \in L(\mathbf{x}_i)_S, \mathbf{X}_{\bar{S}} = \mathbf{x}_{\bar{S}}) \geq \pi_C$ . Finally, we merge all the neighboring possible leaves to get  
 311 the largest rectangle, and this maximal rectangle is the counterfactual rule. Note that the union of the  
 312 possible leaves is not necessary a connected space, thus we can have multiple counterfactual rules.

313 We apply the same idea to find the regional CR. Given a rule  $\mathbf{R}$  and its Minimal Divergent Explanation  
 314  $S$ , we used the Projection given  $\mathbf{X}_{\bar{S}} \in \mathbf{R}_{\bar{S}}$  to find the compatible observations and their leaves  
 315 and combine the possible ones to obtain the regional CR that has  $CRP_S(\mathcal{Y}^*, \mathbf{R}, C_S(\mathbf{R})) \geq \pi_C$ .  
 316 For example, if we consider the leaf 5 of the original forest as a rule: If  $\mathbf{X} \in \text{Leaf 5}$ , then  
 317 predict Blue. Its Minimal Divergent Explanation is also  $S = X_1$ . The R-CR would also be the  
 318 Green region in figure 4. Indeed, if we satisfy the  $X_0$  condition of the leaf 5 and  $X_1$  condition of the  
 319 leaf 3 and 4, then the decision change to Green.

## 320 6 Experiments

321 To demonstrate the performance of our framework, we conduct two experiments on real-world  
 322 datasets. The first consists of showing how we can use the *Local Counterfactual Rules* for explaining  
 323 a regression model. In the second experiment, we compare our approaches with the 2 baselines  
 324 methods in classification problem: (1) **CET** [Kanamori et al., 2022], which partition the input  
 325 space using a decision tree and associate a vector perturbation for each leaf, (2) **ARes** [Rawal and  
 326 Lakkaraju, 2020] performs an exhaustive search for finding global counterfactual rules, but we used  
 327 the implementation of Kanamori et al. [2022] that adapts the algorithm for returning counterfactuals  
 328 samples instead of rules. We compare the methods only in classification problem as most prior works  
 329 do not deal regression problem. In all experiments, we split our dataset into train (75%) - test (25%),  
 330 and we learn a model  $f$ , a LightGBM (*estimators=50, nb leaves=8*), on the train set that is the  
 331 explaine. We learn  $f$ ’s predictions on the train set with an approximating RF  $m_{nb,n}$  (*estimators=20,*  
 332 *max depth=10*): **that** will be used to generate the CR with  $\pi = 0.9$ . The used parameters for **ARes**,  
 333 **CET** are *max rules=8, bins=10* and *max iterations=1000, max leaf=8, bins=10* respectively. Due to  
 334 page limitation, the detailed parameters of each method are provided in Appendix.

335 **Sampling CE using the Counterfactual Rules:** Notice that our approaches cannot be directly  
 336 compare with the baseline methods since they all return counterfactual samples while we give rules  
 337 (range of vector values) that permit to change the decision with high probability. However, we adapt  
 338 the CR to generate also counterfactual samples using a generative model. For example, given an  
 339 instance  $\mathbf{x} = (\mathbf{x}_S, \mathbf{x}_{\bar{S}})$ , target  $\mathcal{Y}^*$  and its counterfactual rule  $C_S(\mathbf{x}; \mathcal{Y}^*)$ , we want to find a sample  
 340  $\mathbf{x}^* = (\mathbf{z}_S, \mathbf{x}_{\bar{S}})$  with  $\mathbf{z}_S \in C_S(\mathbf{x}, \mathcal{Y}^*)$  s.t  $\mathbf{x}^*$  is an in-distribution sample and  $f(\mathbf{x}^*) \in \mathcal{Y}^*$ . Instead  
 341 of using a complex conditional generative model as [Xu et al., 2019, Patki et al., 2016] that can be  
 342 difficult to calibrate, we use an energy-based generative approach [Grathwohl et al., 2020, Lecun et al.,  
 343 2006]. The core idea is to find  $\mathbf{z}_S \in C_S(\mathbf{x}, \mathcal{Y}^*)$  s.t  $\mathbf{x}^*$  maximize a given energy score to ensure that  
 344 it is an in-distribution sample. As an example of an energy function, we use the negative outlier score  
 345 of an Isolation Forest [Liu et al., 2008]. We use Simulated Annealing (see [Guilmeau et al., 2021]  
 346 for a review) to maximize the negative outlier score using the information of the counterfactual rules  
 347  $C_S(\mathbf{x}; \mathcal{Y}^*)$ . In fact, the range values given by the CR  $C_S(\mathbf{x}; \mathcal{Y}^*)$  reduce the search space for  $\mathbf{z}_S$   
 348 drastically. We used the training set  $\mathcal{D}_n$  to find the possible values i.e., we defined  $P_i, P_S$  as the list of  
 349 values of the variable  $i \in S$  found in  $\mathcal{D}_n$  and  $P_S = \{\mathbf{z}_S = (z_1, \dots, z_S) : \mathbf{z}_S \in C_S(\mathbf{x}, \mathcal{Y}^*), z_i \in P_i\}$   
 350 the possible values of  $\mathbf{z}_S$  respectively. Then, we sample  $\mathbf{z}_S$  in the set  $P_S$  and use Simulated Annealing  
 351 to find a  $\mathbf{x}^*$  that maximizes the negative outlier score. Note that the algorithm is the same for sampling  
 352 CE with the Regional-CR. A more detailed version of the algorithm is provided in Appendix.

353 Finally, we compare the methods on unseen observations using three criteria. *Correctness* is the aver-  
 354 age number of instances for which acting as prescribed change to the desired prediction. *Plausibility*  
 355 is the average number of inlier (predict by an Isolation Forest) in the counterfactual samples. *Sparsity*  
 356 is the average number of features that have been changed, and especially for the global counterfactual  
 357 methods (ARes, Regional-CR) that do not ensure to cover all the instances, we compute *Coverage*  
 358 that corresponds to the average number of unseen observations we cover.



359 **Local counterfactual rules for regression:** We give recourse for the **California House Price**  
 360 dataset [Kelley Pace and Barry, 1997] derived from the 1990 U.S. census. We have information about  
 361 each district (demography, ...), and the goal is to predict the median house value of each district.

362 To illustrate the efficiency of the Local-CR, we select all the observations in the test set having a price  
 363 lower than  $100k$  (1566 houses), and we aim to find the recourse that permit to increase their price  
 364 : we want the price  $y$  to be in the interval  $\mathcal{Y}^* = [200k, 250k]$ . For each instance  $x$ , we compute  
 365 the Minimal Divergent Explanation  $S$ , the Local-CR  $C_S(x; [200k, 250k])$  and a CE using the  
 366 Simulated Annealing as described above. We succeed in changing the decision of all the observations  
 367 ( $Correctness = 1$ ) and most of them passed the outlier test with  $Plausibility = 0.92$ . On top of that,  
 368 our Local-CR have sparse support ( $Sparsity = 4.45$ ). For example, the Local-CR of the instance  $x =$   
 369 (Longitude=-118.2, latitude=33.8, housing median age=26, total rooms=703,  
 370 total bedrooms=202, population=757, households=212, median income=2.52) is  
 371  $C_S(x, [200k, 250k]) = (\text{total room} \in [2132, 3546], \text{total bedrooms} \in [214, 491])$ . It  
 372 means if total room and total bedrooms satisfy the conditions in  $C_S(x, [200k, 250k])$  and  
 373 the remaining features of  $x$  is fixed, then the probability that the price is in  $[200k, 250k]$  is 0.97.

374 **Comparisons of Local-CR and Regional-CR with baselines (ARes, CET):** We use 3 real-world  
 375 datasets: **Diabetes** [Kaggle, 2016] contains diagnostic measurements and aims to predict whether  
 376 or not a patient has diabetes, **Breast Cancer Wisconsin (BCW)** [Dua and Graff, 2017] consists of  
 377 predicting if a tumor is benign or not using the characteristic of the cell nuclei, and **Compas** [Larson  
 378 et al., 2016] was used to predict recidivism, and it contains information about the criminal history,  
 379 demographic attributes. During the evaluation, we observe that **ARes, CET** are very sensitive to the  
 380 number of bins and the maximal number of rules or actions as noticed by [Ley et al., 2022]. A bad  
 381 parameterization gives completely useless explanations. Moreover, a different model needs to be  
 382 trained for each class to be accurate, while we only need to have a RF that has good precision.

383 In table 1, we notice that the Local and Regional-CR succeed in changing decisions with a high  
 384 accuracy in all datasets, outperforming **ARes** and **CET** with a large margin on **BCW**, and **Diabetes**.  
 385 Moreover, we notice that the baselines struggle to change at the same time the positive and negative  
 386 class, (e.g. CET has  $Acc=1$  in the positive class, and 0.21 for the negative class on **BCW**) or when  
 387 they have a good  $Acc$ , the CE are not plausible. For instance, CET has  $Acc=0.98$  and  $Psb=0$  on  
 388 **Compas**, meaning that all the CE are outlier. Regarding the coverage of the global CE, CET covers  
 389 all the instances as it partitions the space, but we observe that **ARes** has a smaller  $Coverage =$   
 390  $\{0.43, 0.44, 0.81\}$  than the Regional-CR which has  $\{1, 0.7, 1\}$  for **BCW, Diabetes, and Compas**  
 391 respectively. To sum up, the CR is easier to train and provides more accurate and plausible rules than  
 392 the baselines methods.

Table 1: Results of the  $Correctness$  ( $Acc$ ),  $Plausibility$ , and  $Sparsity$  ( $Sprs$ ) of the different methods. We compute each metric according to the positive (Pos) and negative (Neg) class.

	COMPAS						BCW						Diabetes					
	Acc		Psb		Sps		Acc		Psb		Sps		Acc		Psb		Sps	
	Pos	Neg	Pos	Neg	Pos	Neg	Pos	Neg	Pos	Neg	Pos	Neg	Pos	Neg	Pos	Neg	Pos	Neg
<b>L-CR</b>	1	0.9	0.87	0.73	2	4	1	1	0.96	1	9	7	0.97	1	0.99	0.8	3	4
<b>R-CR</b>	0.9	0.98	0.74	0.93	2	3	0.89	0.9	0.94	0.93	9	9	0.99	0.99	0.9	0.87	3	4
<b>ARes</b>	0.98	1	0.8	0.61	1	1	0.63	0.34	0.83	0.80	4	3	0.73	0.60	0.77	0.86	1	1
<b>CET</b>	0.85	0.98	0.7	0	2	2	1	0.21	0.6	0.80	8	2	0.84	1	0.60	0.20	6	6

## 393 7 Conclusion

394 Most current works that generate CE are implicit through an optimization process or a brunch of  
 395 random samples, thus lacking guarantees. For this reason, we rethink CE as *Counterfactual Rules*.  
 396 For any individual or sub-population, it gives the simplest policies that change the decision with  
 397 high probability. Our approach learns robust, plausible, and sparse adversarial regions where the  
 398 observations should be moved. We make central use of Random Forests, which give consistent  
 399 estimates of the interest probabilities and naturally give the counterfactual rules we want to extract.  
 400 In addition, it permits us to deal with regression problems and continuous features. Consequently,  
 401 our methods are suitable for all datasets where tree-based model performs well (e.g., tabular data). A  
 402 prospective work is to evaluate the robustness of our methods to noisy human responses, i.e., when  
 403 the prescribed recourse is not implemented exactly, and to refine the methodology for selecting the  
 404 threshold probabilities  $\pi$  and  $\pi_C$ .

## 405 References

- 406 Salim I Amoukou and Nicolas JB Brunel. Consistent sufficient explanations and minimal local rules  
407 for explaining regression and classification models. *arXiv preprint arXiv:2111.04658*, 2021.
- 408 Clément Bénard, Gérard Biau, Sébastien Da Veiga, and Erwan Scornet. Shaff: Fast and consistent  
409 shapley effect estimates via random forests. *arXiv preprint arXiv:2105.11724*, 2021a.
- 410 Clément Bénard, Gérard Biau, Sébastien Veiga, and Erwan Scornet. Interpretable random forests  
411 via rule extraction. In *International Conference on Artificial Intelligence and Statistics*, pages  
412 937–945. PMLR, 2021b.
- 413 Clément Bénard, Sébastien Da Veiga, and Erwan Scornet. Mda for random forests: inconsistency,  
414 and a practical solution via the sobol-mda. *arXiv preprint arXiv:2102.13347*, 2021c.
- 415 Leo Breiman, Jerome Friedman, Richard Olshen, and Charles Stone. Classification and regression  
416 trees. *wadsworth int. Group*, 37(15):237–251, 1984.
- 417 S. Chen, Arthur Choi, and Adnan Darwiche. The same-decision probability: A new tool for decision  
418 making. 2012.
- 419 Yu-Liang Chou, Catarina Moreira, Peter Bruza, Chun Ouyang, and Joaquim Jorge. Counterfactuals  
420 and causability in explainable artificial intelligence: Theory, algorithms, and applications. *Informa-*  
421 *tion Fusion*, 81:59–83, 2022. ISSN 1566-2535. doi: <https://doi.org/10.1016/j.inffus.2021.11.003>.  
422 URL <https://www.sciencedirect.com/science/article/pii/S1566253521002281>.
- 423 Qiming Du, Gérard Biau, François Petit, and Raphaël Porcher. Wasserstein random forests and appli-  
424 cations in heterogeneous treatment effects. In *International Conference on Artificial Intelligence*  
425 *and Statistics*, pages 1729–1737. PMLR, 2021.
- 426 Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL [http://archive.ics.](http://archive.ics.uci.edu/ml)  
427 [uci.edu/ml](http://archive.ics.uci.edu/ml).
- 428 Will Grathwohl, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi,  
429 and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like  
430 one. In *International Conference on Learning Representations*, 2020.
- 431 Thomas Guilmeau, Emilie Chouzenoux, and Víctor Elvira. Simulated annealing: a review and a new  
432 scheme. pages 101–105, 07 2021. doi: 10.1109/SSP49050.2021.9513782.
- 433 Hemant Ishwaran, Udaya B Kogalur, Eugene H Blackstone, and Michael S Lauer. Random survival  
434 forests. *The annals of applied statistics*, 2(3):841–860, 2008.
- 435 Kaggle. Pima indians diabetes database, 2016. URL [https://www.kaggle.com/datasets/](https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database)  
436 [uciml/pima-indians-diabetes-database](https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database).
- 437 Kentaro Kanamori, Takuya Takagi, Ken Kobayashi, and Hiroki Arimura. Dace: Distribution-aware  
438 counterfactual explanation by mixed-integer linear optimization. In *IJCAI*, 2020.
- 439 Kentaro Kanamori, Takuya Takagi, Ken Kobayashi, and Yuichi Ike. Counterfactual explanation trees:  
440 Transparent and consistent actionable recourse with decision trees. In *Proceedings of The 25th*  
441 *International Conference on Artificial Intelligence and Statistics*, PMLR 151:1846-1870, 2022.
- 442 Amir-Hossein Karimi, Gilles Barthe, Borja Balle, and Isabel Valera. Model-agnostic counterfactual  
443 explanations for consequential decisions. *ArXiv*, abs/1905.11190, 2020a.
- 444 Amir-Hossein Karimi, Gilles Barthe, Bernhard Schölkopf, and Isabel Valera. A survey of algorithmic  
445 recourse: definitions, formulations, solutions, and prospects. *CoRR*, abs/2010.04050, 2020b. URL  
446 <https://arxiv.org/abs/2010.04050>.
- 447 R. Kelley Pace and Ronald Barry. Sparse spatial autoregressions. *Statistics, Probability Letters*, 33  
448 (3):291–297, 1997. ISSN 0167-7152. doi: [https://doi.org/10.1016/S0167-7152\(96\)00140-X](https://doi.org/10.1016/S0167-7152(96)00140-X). URL  
449 <https://www.sciencedirect.com/science/article/pii/S016771529600140X>.

- 450 Himabindu Lakkaraju, Dylan Slack, Yuxin Chen, Chenhao Tan, and Sameer Singh. Rethinking  
 451 explainability as a dialogue: A practitioner’s perspective. *CoRR*, abs/2202.01875, 2022. URL  
 452 <https://arxiv.org/abs/2202.01875>.
- 453 Jeff Larson, Surya Mattu, Lauren Kirchner, , and Julia Angwin. How we analyzed  
 454 the compas recidivism algorithm, 2016. URL [https://www.propublica.org/article/  
 455 how-we-analyzed-the-compas-recidivism-algorithm](https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm).
- 456 Yann Lecun, Sumit Chopra, and Raia Hadsell. *A tutorial on energy-based learning*. 01 2006.
- 457 Dan Ley, Saumitra Mishra, and Daniele Magazzeni. Global counterfactual explanations: Investiga-  
 458 tions, implementations and improvements, 2022. URL <https://arxiv.org/abs/2204.06917>.
- 459 Jimmy J. Lin, Chudi Zhong, Diane Hu, Cynthia Rudin, and Margo I. Seltzer. Generalized optimal  
 460 sparse decision trees. *ArXiv*, abs/2006.08690, 2020.
- 461 Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In *2008 eighth ieee international  
 462 conference on data mining*, pages 413–422. IEEE, 2008.
- 463 Wei-Yin Loh. Classification and regression trees. *Wiley Interdisciplinary Reviews: Data Mining and  
 464 Knowledge Discovery*, 1, 2011.
- 465 Arnaud Van Looveren and Janis Klaise. Interpretable counterfactual explanations guided by proto-  
 466 types. *CoRR*, abs/1907.02584, 2019. URL <http://arxiv.org/abs/1907.02584>.
- 467 Scott M. Lundberg, Gabriel Erion, Hugh Chen, Alex DeGrave, Jordan M. Prutkin, Bala Nair, Ronit  
 468 Katz, Jonathan Himmelfarb, Nisha Bansal, and Su-In Lee. From local explanations to global  
 469 understanding with explainable ai for trees. *Nature Machine Intelligence*, 2(1):2522–5839, 2020.
- 470 Nicolai Meinshausen and Greg Ridgeway. Quantile regression forests. *Journal of Machine Learning  
 471 Research*, 7(6), 2006.
- 472 Christoph Molnar. *Interpretable Machine Learning*. 2 edition, 2022. URL [https://christophm.  
 473 github.io/interpretable-ml-book](https://christophm.github.io/interpretable-ml-book).
- 474 Ramaravind K. Mothilal, Amit Sharma, and Chenhao Tan. Explaining machine learning classifiers  
 475 through diverse counterfactual explanations. In *Proceedings of the 2020 Conference on Fairness,  
 476 Accountability, and Transparency, FAT\* ’20*, page 607–617, New York, NY, USA, 2020. Associa-  
 477 tion for Computing Machinery. ISBN 9781450369367. doi: 10.1145/3351095.3372850. URL  
 478 <https://doi.org/10.1145/3351095.3372850>.
- 479 Axel Parmentier and Thibaut Vidal. Optimal counterfactual explanations in tree ensembles. *CoRR*,  
 480 abs/2106.06631, 2021. URL <https://arxiv.org/abs/2106.06631>.
- 481 N. Patki, R. Wedge, and K. Veeramachaneni. The synthetic data vault. In *2016 IEEE International  
 482 Conference on Data Science and Advanced Analytics (DSAA)*, pages 399–410, Oct 2016. doi:  
 483 10.1109/DSAA.2016.49.
- 484 Martin Pawelczyk, Teresa Datta, Johannes van-den Heuvel, Gjergji Kasneci, and Himabindu  
 485 Lakkaraju. Algorithmic recourse in the face of noisy human responses, 2022. URL [https:  
 486 //arxiv.org/abs/2203.06768](https://arxiv.org/abs/2203.06768).
- 487 Rafael Poyiadzi, Kacper Sokol, Raúl Santos-Rodriguez, Tijl De Bie, and Peter A. Flach. FACE:  
 488 feasible and actionable counterfactual explanations. *CoRR*, abs/1909.09369, 2019. URL [http:  
 489 //arxiv.org/abs/1909.09369](http://arxiv.org/abs/1909.09369).
- 490 Kaivalya Rawal and Himabindu Lakkaraju. Beyond individualized recourse: Interpretable and  
 491 interactive summaries of actionable recourses. *Advances in Neural Information Processing Systems*,  
 492 33:12187–12198, 2020.
- 493 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the  
 494 predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference  
 495 on knowledge discovery and data mining*, pages 1135–1144, 2016.

- 496 Chris Russell. Efficient search for diverse coherent explanations. In *Proceedings of the Conference*  
 497 *on Fairness, Accountability, and Transparency*, FAT\* '19, page 20–28, New York, NY, USA, 2019.  
 498 Association for Computing Machinery. ISBN 9781450361255. doi: 10.1145/3287560.3287569.  
 499 URL <https://doi.org/10.1145/3287560.3287569>.
- 500 Berk Ustun, Alexander Spangher, and Yang Liu. Actionable recourse in linear classification. *Pro-*  
 501 *ceedings of the Conference on Fairness, Accountability, and Transparency*, 2019.
- 502 Sahil Verma, John P. Dickerson, and Keegan Hines. Counterfactual explanations for machine learning:  
 503 A review. *CoRR*, abs/2010.10596, 2020. URL <https://arxiv.org/abs/2010.10596>.
- 504 Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening  
 505 the black box: Automated decisions and the gdpr. *Harv. JL & Tech.*, 31:841, 2017a.
- 506 Sandra Wachter, Brent Daniel Mittelstadt, and Chris Russell. Counterfactual explanations without  
 507 opening the black box: Automated decisions and the gdpr. *Cybersecurity*, 2017b.
- 508 Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using  
 509 random forests, 2017.
- 510 Tong Wang, Cynthia Rudin, Finale Doshi-Velez, Yimin Liu, Erica Klampfl, and Perry MacNeille. A  
 511 bayesian framework for learning rule sets for interpretable classification. *J. Mach. Learn. Res.*, 18:  
 512 70:1–70:37, 2017.
- 513 Lei Xu, Maria Skoularidou, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. Modeling tabular  
 514 data using conditional gan. In *NeurIPS*, 2019.

## 515 Checklist

- 516 1. For all authors...
- 517 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
 518 contributions and scope? **[Yes] In Section 3 the two new explanation methods:**  
 519 **The Local and Regional Counterfactual rules. Section 4 shows that our methods**  
 520 **are consistent, contrary to prior works. Then, in Section 6 we demonstrate the**  
 521 **performance of our new explanations w.r.t SOTA on real-world datasets. Finally,**  
 522 **we provide a Python Package that computes our methods. Additional experiments**  
 523 **can be found in Appendix.**
- 524 (b) Did you describe the limitations of your work? **[Yes] In conclusion, we emphasize**  
 525 **that as our estimators are based on a Random Forest, our methods are suitable**  
 526 **for all datasets on which tree-based models perform well. Therefore, it works well**  
 527 **on tabular data, but it is not adapted for big computer vision models for example.**
- 528 (c) Did you discuss any potential negative societal impacts of your work? **[No] Our con-**  
 529 **tributions are fully dedicated to the positive societal impacts. Indeed, we propose**  
 530 **new and better explanation methods. For instance, our regional counterfactual**  
 531 **rules permit us to detect unfair behavior of model as AReS but more accurately.**  
 532 **On the other hand, the local counterfactual rules permit to give more robust**  
 533 **recourse in real-world scenarios.**
- 534 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
 535 them? **[Yes] We conform to the ethics review as we are not concerned by potential**  
 536 **negative impacts (methodologies, data,...), as the general ethical conduct.**
- 537 2. If you are including theoretical results...
- 538 (a) Did you state the full set of assumptions of all theoretical results? **[Yes] See Section 4.**
- 539 (b) Did you include complete proofs of all theoretical results? **[Yes] Our methods are**  
 540 **based on the theoretical results of previous work. Thus, the complete proofs can**  
 541 **be found in the given references.**
- 542 3. If you ran experiments...

- 543 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
544 perimental results (either in the supplemental material or as a URL)? **[Yes]** **All**  
545 **the codes to reproduce the results are given in** [https://github.com/anoxai/](https://github.com/anoxai/counterfactual_rules)  
546 [counterfactual\\_rules](https://github.com/anoxai/counterfactual_rules)
- 547 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
548 were chosen)? **[Yes]** **See Appendix.**
- 549 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
550 ments multiple times)? **[Yes]** **We run the experiments multiple times.**
- 551 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
552 of GPUs, internal cluster, or cloud provider)? **[Yes]** **See Appendix.**
- 553 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 554 (a) If your work uses existing assets, did you cite the creators? **[Yes]** **We use the**  
555 **implementation of AReS, and CET provided by Kanamori et al. at** <https://github.com/kelicht/cet>. **We have cited the data we used in the experiment**  
556 **Section.**
- 557 (b) Did you mention the license of the assets? **[Yes]**
- 558 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**  
559 **At** [https://github.com/anoxai/counterfactual\\_rules](https://github.com/anoxai/counterfactual_rules)
- 560 (d) Did you discuss whether and how consent was obtained from people whose data you're  
561 using/curating? **[No]** **Not relevant**
- 562 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
563 information or offensive content? **[No]** **Not relevant**
- 564
- 565 5. If you used crowdsourcing or conducted research with human subjects...
- 566 (a) Did you include the full text of instructions given to participants and screenshots, if  
567 applicable? **[No]** **Not relevant**
- 568 (b) Did you describe any potential participant risks, with links to Institutional Review  
569 Board (IRB) approvals, if applicable? **[No]** **Not relevant**
- 570 (c) Did you include the estimated hourly wage paid to participants and the total amount  
571 spent on participant compensation? **[No]** **Not relevant**