

---

# Supplementary materials: Rethinking Counterfactual Explanations as Local and Regional Counterfactual Policies

---

Anonymous Author(s)

Affiliation

Address

email

1	<b>Contents</b>	
2	<b>A Regional RF detailed</b>	<b>2</b>
3	<b>B Additional experiments</b>	<b>3</b>
4	<b>C Simulated annealing to generate counterfactual samples using the Counterfactual Rules</b>	<b>3</b>
5	<b>D Parameters detailed</b>	<b>5</b>

## 6 A Regional RF detailed

7 In this section, we give a simple application of the Regional RF algorithm to better understand how  
 8 it works. Recall that the regional RF is a generalization of the RF's algorithm to give prediction  
 9 even when we condition given a region, e.g., to estimate  $E(f(\mathbf{X}) | \mathbf{X}_S \in C_S(\mathbf{x}), \mathbf{X}_{\bar{S}} = \mathbf{x}_{\bar{S}})$  with  
 10  $C_S(\mathbf{x}) = \prod_{i=1}^{|\bar{S}|} [a_i, b_i]$ ,  $a_i, b_i \in \mathbb{R}$  a hyperrectangle. The algorithm works as follows: we drop the  
 11 observations in the initial trees, if a split used variable  $i \in \bar{S}$ , a fixed value-based condition, we used  
 12 the classic rules i.e., if  $x_i \leq t$ , the observations go to the left children, otherwise the right children.  
 13 However, if a split used variable  $i \in S$ , regional-based condition, we used the hyperrectangle  
 14  $C_S(\mathbf{x}) = \prod_{i=1}^{|\bar{S}|} [a_i, b_i]$ . The observations are sent to the left children if  $b_i \leq t$ , right children if  $a_i > t$   
 15 and if  $t \in [a_i, b_i]$  the observations are sent both to the left and right children.

16 To illustrate how it works, we use a two dimensional variables  $\mathbf{X} \in \mathbb{R}^2$ , a simple decision tree  $f$   
 17 represented in figure 1, and want to compute for  $\mathbf{x} = [1.5, 1.9]$ ,  $E(f(\mathbf{X}) | \mathbf{X}_1 \in [2, 3.5], \mathbf{X}_0 = 1.5)$ .  
 18 We assume that  $P(X_1 \in [2, 3.5] | X_0 = 1.5) > 0$  and denoted  $T_1$  as the set of the values of the splits  
 19 based on variables  $X_1$  of the decision tree. One way of estimating this conditional mean is by using  
 20 Monte Carlo sampling. Therefore, there are two cases :

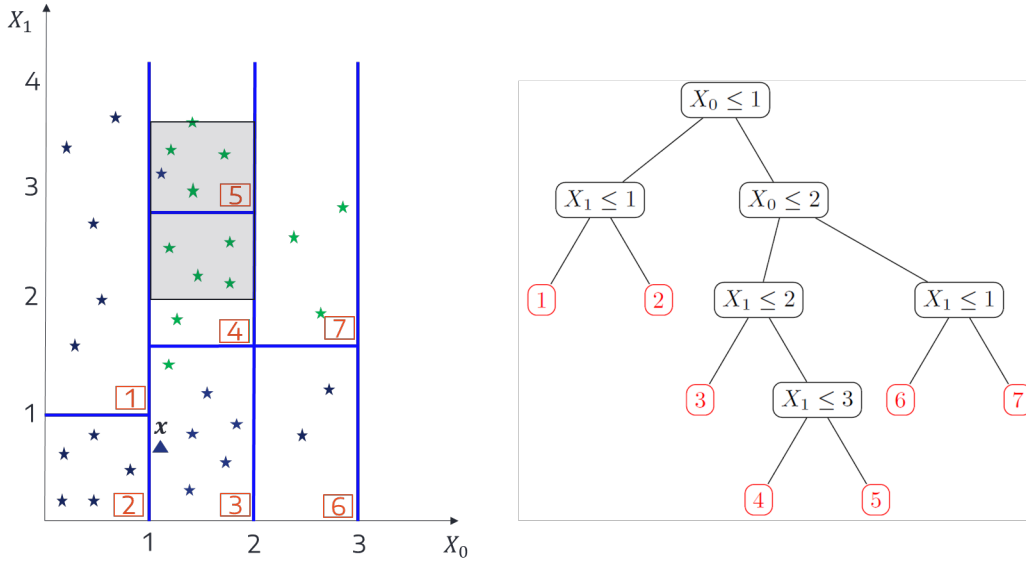


Figure 1: Representation of a simple decision tree (right figure) and its associated partition (left figure). The gray part in the partition corresponds to the region  $[2, 3.5] \times [1, 2]$

21 • If  $\forall t \in T_1, t \leq 2$  or  $t > 3$ , then all the observations sampled s.t.  $\tilde{X}_i \sim \mathcal{L}(\mathbf{X} | \mathbf{X}_1 \in$   
 22  $[2, 3.5], \mathbf{X}_0 = 1.5)$  follow the same path and fall in the same leaf. The Monte Carlo  
 23 estimator of the decision tree  $E(f(\mathbf{X}) | \mathbf{X}_1 \in [2, 3.5], \mathbf{X}_0 = 1.5)$  is equal to the output of  
 24 the Regional RF algorithm.

25 – For instance, a special case of the case above is: if  $\forall t \in T_1, t \leq 2$ , and we sample using  
 26  $\mathcal{L}(\mathbf{X} | \mathbf{X}_1 \in [2, 3.5], \mathbf{X}_0 = 1.5)$ , then all the observations go to the right children  
 27 when they encounters a node using  $X_1$  and fall in the same leaf.

28 • If  $\exists t \in T_1$  and  $t \in [2, 3.5]$ , then the observations sampled s.t.  $\tilde{X}_i \sim \mathcal{L}(\mathbf{X} | \mathbf{X}_1 \in$   
 29  $[2, 3.5], \mathbf{X}_0 = 1.5)$  can fall in multiple terminal leaf depending on if their coordinates  
 30  $x_1$  is lower than  $t$ . Following our example, if we generate samples using  $\mathcal{L}(\mathbf{X} | \mathbf{X}_1 \in$   
 31  $[2, 3.5], \mathbf{X}_0 = 1.5)$ , the observations will fall in the gray region of figure 1, and thus can  
 32 fall in node 4 or 5. Therefore, the true estimate is:

$$\begin{aligned} E(f(\mathbf{X}) | \mathbf{X}_1 \in [2, 3.5], \mathbf{X}_0 = 1.5) \\ = p(X_1 \leq 2.9 | X_0 = 1.5) * E[f(\mathbf{X}) | \mathbf{X} \in L_4] + p(X_1 > 2.9 | X_0 = 1.5) * E[f(\mathbf{X}) | \mathbf{X} \in L_5] \end{aligned} \quad (\text{A.1})$$

33 Concerning the last case ( $t \in [2, 3.5]$ ), we need to estimate the different probabilities  $p(X_1 \leq$   
34  $2.9 | X_0 = 1.5), p(X_1 > 2.9 | X_0 = 1.5)$  to compute  $E(f(\mathbf{X}) | \mathbf{X}_1 \in [2, 3.5], \mathbf{X}_0 = 1.5)$ , but  
35 these probabilities are difficult to estimate in practice. However, we argue that we can ignore these  
36 splits, and thus do not need to fragment the query region using the leaves of the tree. Indeed, as we  
37 are no longer interested in a point estimate but regional (population mean) we do not need to go to  
38 the level of the leaves. We propose to ignore the splits of the leaves that divide the query region.  
39 For instance, the leaves 4 and 5 split the region  $[2, 3.5]$  in two cells, by ignoring these splits we  
40 estimate the mean of the gray region by taking the average output of the leaves 4 and 5 instead of  
41 computing the mean weighted by the probabilities as in Eq. A.1. Roughly, it consists to follow  
42 the classic rules of a decision tree (if the region is above or below a split) and ignore the splits  
43 that are in the query region, i.e., we average the output of all the leaves that are compatible with  
44 the condition  $\mathbf{X}_1 \in [2, 3.5], \mathbf{X}_0 = 1.5$ . We think that it leads to a better approximation for two  
45 reasons. First, we observe that the case where  $t$  is in the region and thus divides the query region  
46 does not happen often. Moreover, the leaves of the trees are very small in practice, and taking the  
47 mean of the observations that fall in the union of leaves that belong to the query region is more  
48 reasonable than computing the weighted mean and thus trying to estimate the different probabilities  
49  $p(X_1 \leq 2.9 | X_0 = 1.5), p(X_1 > 2.9 | X_0 = 1.5)$ .

## 50 B Additional experiments

51 In table 1, we compare the *Correctness* (Acc), *Plausibility* (Psb), and *Sparsity* (Sprs) of the different  
52 methods on additional real-world datasets: FICO [FICO, 2018], NHANESI [CDC, 1999-2022].

53 We observe that the L-CR, and R-CR outperform the baseline methods by a large margin on *Correct-*  
54 *ness* and *Plausibility*. The baseline methods still struggle to change at the same time the positive and  
55 negative class. In addition, AReS and CET give better sparsity, but their counterfactual samples are  
56 less plausible than the ones generated by the CR.

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.

	FICO						NHANESI					
	Acc		Psb		Sps		Acc		Psb		Sps	
	Pos	Neg	Pos	Neg	Pos	Neg	Pos	Neg	Pos	Neg	Pos	Neg
<b>L-CR</b>	0.98	0.94	0.98	0.99	5	5	0.99	0.98	0.98	0.97	5	6
<b>R-CR</b>	0.90	0.94	0.98	0.99	9	8.43	0.86	0.95	0.96	0.99	7	7
<b>AReS</b>	0.34	0.01	0.85	0.86	2	1	0.06	1	0.87	0.92	1	1
<b>CET</b>	0.76	0	0.76	0.60	2	2	0	0.40	0.82	0.56	0	5

## 57 C Simulated annealing to generate counterfactual samples using the 58 Counterfactual Rules

```

59 1 import numpy as np
60 2
61 3 def generate_candidate(x, S, x_train, C_S, n_samples):
62 4     """
63 5     Generate sample by sampling marginally between the features value
64 6     of the training observations.
65 7     Args:
66 8     x (numpy.ndarray): 1-D array, an observation
67 9     S (list): contains the indices of the variables on which to
68 10    condition
69 11    x_train (numpy.ndarray): 2-D array represent the training
70 12    samples
71 13    C_S (numpy.ndarray): 3-D (#variables x 2 x 1) representing
72 14    the hyper-rectangle on which to condition
73 15    n_samples (int): number of samples
74 16    Returns:
75 17    The generated samples

```

```

7614     """
7715     x_poss = [x_train[(C_S[i, 0] <= x_train[:, i]) * (x_train[:, i] <=
78         C_S[i, 1]), i] for i in S]
7916     x_cand = np.repeat(x.reshape(1, -1), repeats=n_samples, axis=0)
8017
8118     for i in range(len(S)):
8219         rdm_id = np.random.randint(low=0, high=x_poss[i].shape[0],
83             size=n_samples)
8420         x_cand[:, S[i]] = x_poss[i][rdm_id]
8521
8622     return x_cand
8723
8824
8925 def simulated_annealing(outlier_score, x, S, x_train, C_S, batch,
90     max_iter, temp, max_iter_convergence):
9126     """
9227     Generate sample X s.t. X_S \in C_S using simulated annealing and
93     outlier score.
9428     Args:
9529         outlier_score (lambda functon): outlier_score(X) return a
96     outlier score. If the value are negative, then the observation is
97     an outlier.
9830         x (numpy.ndarray): 1-D array, an observation
9931         S (list): contains the indices of the variables on which to
100     condition
10132         x_train (numpy.ndarray): 2-D array represent the training
102     samples
10333         C_S (numpy.ndarray): 3-D (#variables x 2 x 1) representing
104     the hyper-rectangle on which to condition
10534         batch (int): number of sample by iteration
10635         max_iter (int): number of iteration of the algorithm
10736         temp (double): the temperature of the simulated annealing
108     algorithm
10937         max_iter_convergence (double): minimum number of iteration to
110     stop the algorithm if it find an in-distribution observation
11138
11239     Returns:
11340         The generated sample, and its outlier score
11441     """
11542
11643     best = generate_candidate(x, S, x_train, C_S, n_samples=1)
11744     best_eval = outlier_score(best)[0]
11845     curr, curr_eval = best, best_eval
11946
12047     it = 0
12148     for i in range(max_iter):
12249
12350         x_cand = generate_candidate(curr, S, x_train, C_S, batch)
12451         score_candidates = outlier_score(x_cand)
12552
12653         candidate_eval = np.max(score_candidates)
12754         candidate = x_cand[np.argmax(score_candidates)]
12855
12956         if candidate_eval > best_eval:
13057             best, best_eval = candidate, candidate_eval
13158             it = 0
13259         else:
13360             it += 1
13461
13562         # check convergence
13663         if best_eval > 0 and it > max_iter_convergence:
13764             break
13865
13966         diff = candidate_eval - curr_eval
14067         t = temp / np.log(float(i + 1))

```

```

14168     metropolis = np.exp(-diff / t)
14269
14370     if diff > 0 or rand() < metropolis:
14471         curr, curr_eval = candidate, candidate_eval
14572
14673     return best, best_eval

```

Listing 1: The simulated annealing algorithm to generate samples that satisfy the condition CR

## 147 D Parameters detailed

148 In this section, we give the different parameters of each method. For all methods and datasets, we first  
 149 used a greedy search given a set of parameters. For AReS, we use the following set of parameters:

- 150 • max rule = {4, 6, 8}, max rule length = {4, 8}, max change num = {2, 4, 6},
- 151 • minimal support = 0.05, discretization bins = {10, 20},
- 152 •  $\lambda_{acc} = \lambda_{cov} = \lambda_{cst} = 1$ .

153 For CET, we search in the following set of parameters:

- 154 • max iterations = {500, 1000},
- 155 • max leaf size = {4, 6, 8, -1},
- 156 •  $\lambda = 0.01, \gamma = 1$ .

157 Finally, for the Counterfactual Rules, we used the following parameters:

- 158 • nb estimators = {20, 50}, max depth= {8, 10, 12},
- 159 •  $\pi = 0.9, \pi_C = 0.9$ .

160 We obtained the same optimal parameters for all datasets:

- 161 • AReS: max rule = 4, max rule length= 4, max change num = 4, minimal support = 0.05,  
 162 discretization bins = 10,  $\lambda_{acc} = \lambda_{cov} = \lambda_{cst} = 1$
- 163 • CET: max iterations = 1000, max leaf size = -1,  $\lambda = 0.01, \gamma = 1$
- 164 • CR: nb estimators= 20, max depth= 10,  $\pi = 0.9, \pi_C = 0.9$

165 The code and the results can be found at [https://github.com/anoxai/counterfactual\\_](https://github.com/anoxai/counterfactual_rules)  
 166 [rules](https://github.com/anoxai/counterfactual_rules).

167 **References**

- 168 CDC. National health and nutrition examination survey, 1999-2022. URL [https://www.cdc.](https://www.cdc.gov/Nchs/Nhanes/Default.aspx)  
169 [gov/Nchs/Nhanes/Default.aspx](https://www.cdc.gov/Nchs/Nhanes/Default.aspx).
- 170 FICO. Fico. explainable machine learning challenge, 2018. URL [https://community.fico.com/](https://community.fico.com/s/explainable-machine-learning-challenge)  
171 [s/explainable-machine-learning-challenge](https://community.fico.com/s/explainable-machine-learning-challenge).