

Mining Pareto-Optimal Counterfactual Antecedents With A Branch-And-Bound Model-Agnostic Algorithm

Marcos M. Raimundo (✉ marcosmrai@gmail.com)

Fundação Getúlio Vargas <https://orcid.org/0000-0003-0499-2564>

Luis Gustavo Nonato

Brazil and New York University

Jorge Poco

Fundação Getúlio Vargas

Research Article

Keywords: Counterfactual, Interpretability, Machine learning, Knowledge Mining, Multi-objective optimization, Crime Data.

Posted Date: September 20th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-551661/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Mining Pareto-Optimal Counterfactual Antecedents with a Branch-and-Bound Model-Agnostic Algorithm

Marcos M. Raimundo · Luis Gustavo Nonato · Jorge Poco

Received: date / Accepted: date

Abstract Mining *counterfactual antecedents* became a valuable tool to discover knowledge and explain machine learning models. It consists of generating synthetic samples from an original sample to achieve the desired outcome in a machine learning model thus helping to understand the prediction. An insightful methodology would explore a broader set of counterfactual antecedents to reveal multiple possibilities while operating on any classifier. Thus, we create a tree-based search that requires monotonicity from the objective functions (a.k.a. cost functions); it allows pruning branches that will not improve the objective functions. Since monotonicity is only required for the objective function, this method can be used for any family of classifiers (e.g., linear models, neural networks, decision trees). However, additional classifier properties speed up the tree-search when it foresees branches that will not result in feasible actions. Moreover, the proposed optimization generates a diverse set of *Pareto-optimal* counterfactual antecedents by relying on multi-objective concepts. The results show an algorithm with working guarantees that enumerates a wide range of counterfactual antecedents. It helps the decision-maker understand the machine learning decision and finds alternatives to achieve the desired outcome. The user can inspect these multiple counterfactual antecedents to find the most suitable one and have a broader understanding of the prediction.

Keywords Counterfactual · Interpretability · Machine learning · Knowledge Mining · Multi-objective optimization · Crime Data.

Marcos M. Raimundo
Fundação Getúlio Vargas, Brazil, E-mail: marcos.raimundo@fgv.br,

Luis Gustavo Nonato
ICMC-USP, São Carlos, Brazil and New York University, USA, E-mail: gnonato@icmc.usp.br

Jorge Poco
Fundação Getúlio Vargas, Brazil and Universidad Católica San Pablo, E-mail: jorge.poco@fgv.br,

1 Introduction

Given a sample with an undesired outcome from a machine learning model, a **counterfactual antecedent** is a synthetic sample that achieves the desired outcome with minimal changes compared to the original sample. It helps on explaining the prediction of such sample by observing what could be changed to achieve another outcome. For instance, let us suppose that a person is diagnosed with a high risk of diabetes by a machine learning technique. The attributes involved in the decision include an **Insulin** level of 285, **BMI** (Body Mass Index) of 44, **BloodPressure** of 80, **Skin Thickness** of 34, and **Glucose** level of 114. A **counterfactual antecedent mining** system would suggest some counterfactual antecedents to move from high to low risk. — e.g., this person should decrease his/her **BMI** from 44 to 36 or decrease **Glucose** level from 114 to 99. Counterfactual mining is also known by other names such as Actionable knowledge (Yang et al, 2003), Inverse Classification (Yang et al, 2012), Counterfactual Explanations (Wachter et al, 2018), and other variations. People use this technique to recommend medical treatments (Yang et al, 2012; Krause et al, 2016), marketing strategies to retain customers (Yang et al, 2003), and personal profile changes to have credit approved (Ustun et al, 2019). Recently, counterfactual antecedent mining has become popular to explain the behavior of learning models (Wachter et al, 2018).

Counterfactual mining methods explore counterfactual antecedent with a single (Krause et al, 2016) or multiple (Wachter et al, 2018) feature changes to achieve the desired outcome. However, the methods capable of finding a diverse set of counterfactual antecedents (Karimi et al, 2020; Ustun et al, 2019; Mothilal et al, 2020) rely on a single objective function. Although this function should capture the cost of making changes to create a counterfactual antecedent, it is still unyielding to design and rely only on a single objective function to fulfill the preferences of a user (Rudin, 2019). To clarify the importance of having multiple counterfactual antecedents, Fig. 1 shows the original sample (see the **Orig** column) and a set of counterfactual antecedents (C1, . . . , C12 columns) for a diabetes example. For instance, to change a patient from high to low risk of diabetes:

- C1 suggests increasing **Insulin** from 285 to 468, decrease **BMI** from 44 to 42, and decrease **Glucose** from 114 to 113.
- C8 suggests increasing **Insulin** from 285 to 468, decrease **BMI** from 44 to 42, and decrease **SkinThickness** from 34 to 24.
- C10 gives another alternative, namely, to increase **Insulin** from 285 to 416, decrease **BMI** from 44 to 42, and decrease **SkinThickness** from 34 to 12.

This example shows that it is possible to explore changes in distinct features (C1 vs. C8) and changes in the values of the same features (C8 vs. C10). This diversity is achievable by using multiple objectives, especially using multi-objective optimization. Moreover, most existing methods usually operate on a limited family of classification methods (Ustun et al, 2019; Mothilal et al, 2020); thus, making it impossible to use them with some types of classifiers.

Therefore, there is a need for methods to operate in any family of classifiers (model-agnostic method) and capable of mining a diverse set of counterfactual antecedents with multiple objectives.

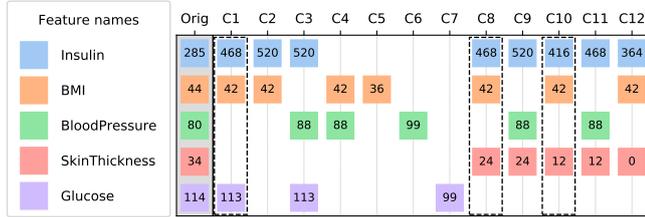


Fig. 1 A multi-objective enumeration of counterfactual antecedents that reduces the risk of diabetes. Orig column is the original sample and the columns C1 to C12 are counterfactual antecedents.

This paper’s main contribution is a methodology called MAPOCAM (Model-Agnostic Pareto-Optimal Counterfactual Antecedent Mining), which relies on a tree-based search to enumerate all Pareto-optimal counterfactual antecedents using a given set of objective functions. The proposed methodology only requires the monotonicity of the objective functions to prune branches that do not improve the solution’s quality. Moreover, MAPOCAM can also explore optional conditions to improve performance. For example, when the decision function (that evaluates if we achieve the result) is monotonic (an increase in the value of a feature increases the probability of achieving the target outcome), the search can be speed up. It happens because it is possible to foresee branches that take to unfeasible actions thus it can be pruned.

In summary, the contributions of this research are:

- A model agnostic methodology that enumerates all counterfactual antecedents considering the trade-off of multiple objective functions;
- A mechanism to create an overview of how much the variables need to change in order to give rise to counterfactual antecedents;
- A comprehensive set of experiments that explore multiple machine learning models — logistic regression, decision trees, and neural networks — to validate our methodology on real applications;
- An in-depth experiment using urban and social-economic variables to explain how those variables can alleviate to some extent the ever-decreasing crime incidents.

2 Related work

Counterfactual antecedent mining is approached in the literature under various names depending on the application: Actionable knowledge (Cui et al, 2015; Yang et al, 2007; Lv et al, 2018; Lu et al, 2017; Yang et al, 2003), Actionable Plans (Lu et al, 2017), Actionable recourse (Ustun et al, 2019), Actionable

feature tweaking (Tolomei et al, 2017), Inverse Classification (Yang et al, 2012; Aggarwal et al, 2010), Counterfactual Explanations (Wachter et al, 2018), Why-not questions (Gao et al, 2015; He and Lo, 2012; Chen et al, 2015). The coined name *counterfactual antecedent* tries to homogenize those terms. The term counterfactual already explains the nature of the method, in which we create an antecedent (distinct from the original sample) that achieves the consequent (the target outcome, which is fixed).

Actionable knowledge methods (plans/recourse/feature tweaking), usually, aim to mine counterfactual antecedent for specific machine learning models, trying to explore their structure to find counterfactual patterns. They usually explore a single counterfactual antecedent with a single objective function, differing mainly on the optimization procedure. Heuristic approaches typically require an already fitted decision tree or ensemble of trees and use greedy algorithms to find groups of samples on which actions can change the outcome (Yang et al, 2003, 2007; Subramani et al, 2016). Heuristics can also search on every positive path of a decision tree (i.e., changes that result in a swap of outcome), searching positive paths that change the ensemble’s outcome with the lowest cost (Tolomei et al, 2017). Another heuristic consists of making greedy moves (i.e., changes that increase the probability) on a KNN classifier (Yang et al, 2012). Other approaches try to explore counterfactual antecedents with more theoretical support but using speed-up heuristics. In that vein, some approaches explore A^* -like search to find an optimal change to additive tree models using a heuristic that considers the probability achieved by the action and the action’s cost (Lu et al, 2017). Other works use the desired state distance as a heuristic to refine an A^* -like search (Lv et al, 2018). Some formulations use mixed linear-integer modeling (enforcing optimality guarantees) to swap between leaves of the trees, searching for a feature change with lower cost (Cui et al, 2015). Those formulations can also search through a set of actions to find a counterfactual antecedent on linear classifiers (Ustun et al, 2019).

Inverse classification defines a set of features (also interpreted as actions) that achieves the desired class. An example of inverse classification consists of searching for a set of feature changes whose correspondent subset of samples have a high Gini index (Aggarwal et al, 2010). Therefore, the inverse classification does not rely on a machine learning model; it mines sets of features that achieve the target outcome.

Similarly, the why-not questions approach focuses on determining the absence of an object in a database query. Knowing a user preference (usually expressed using a weight vector), a query returns a set of objects. If the user was expecting another answer, then it can use a quadratic optimization to find a weight vector with the desired object (Gao et al, 2015); creating a sampling procedure to find new weight vectors with the desired object (He and Lo, 2012); or using a tree-based bound and prune algorithm (Chen et al, 2015). This family of methods is quite limited since it only uses a weight vector to classify the samples, similar to linear classifiers.

Counterfactual explanations are usually capable of exploring broader families of machine learning models. Convex optimization can be applied to linear

classifiers and neural networks to find a new outcome that is as close as possible to a new target and still close enough to the original sample (Wachter et al, 2018), and add constraints to ensure plausibility (Artelt and Hammer, 2020). The method that uses the Nelder-Mead optimization finds counterfactual antecedent on black-box models with an iterative approach to determine how close it is to the original sample (Grath et al, 2018). Then it uses differentiable approximations of tree ensembles to keep the convexity of the problem (Lucic et al, 2019). Other approaches include using the density of samples to create a graph and find a path to a counterfactual that both are feasible and actionable (Poyiadzi et al, 2020); and using binary search to verify the counterfactual antecedent’s satisfiability on a model-agnostic oracle of logic formula’s representation of machine learning models (Karimi et al, 2020).

There are some strategies to explore multiple counterfactual antecedents: (i) a mixed-integer model can exclude the combination of features that already find a counterfactual antecedent (Ustun et al, 2019); (ii) a convex optimization problem searches for multiple solutions at once and uses distance metrics between the counterfactual antecedents to induce diversity (Mothilal et al, 2020); or (iii) an iterative algorithm that finds new solutions with a minimum distance from the already found solution, i.e., using l_0 larger than one to have at least one new change (Karimi et al, 2020). The main limitation of those methods consists of using distances like l_0 , which does not allow exploring the same set of features but increasing one and decreasing others. Furthermore, even using l_2 distance could induce a new counterfactual antecedent with all features slightly worse. In contrast, multi-objective optimization allows creating diversity without relying on a diversity-inducing function. A proper representation of the Pareto-front promoted by multi-objective optimization will result in solutions with different trade-offs between objectives (as presented in the example of Fig. 1 in the Introduction), thus generating counterfactual antecedents sufficiently different to promote diversity.

Our method has a set of characters not concomitantly present in the methods described above: (i) it is model-agnostic; (ii) is capable of exploring multiple objective functions; and (iii) it is capable of exploring multiple counterfactual antecedents, property that is provided by exploring multiple trade-offs by using multi-objective optimization. Moreover, our approach has solid theoretical guarantees, and it is capable of dealing with binary features—quality not present in methods that rely on convex optimization.

3 Proposed method

This section defines the goals and premises of our methodology, a branch-and-bound algorithm, and it details the essential properties of the proposed methodology.

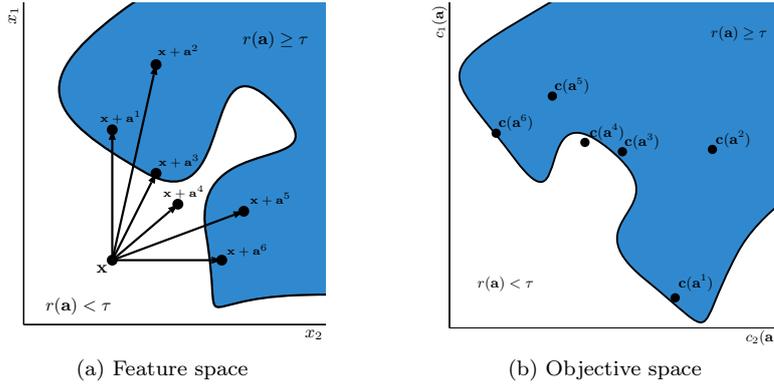


Fig. 2 Representation of feature space (a) and objective space (b), taking two features and two objectives. The blue area represents the space where the samples are classified with the target outcome. In the feature space (a), the goal is to find a plan \mathbf{a}^i that makes \mathbf{x} achieve a desired outcome $r(\mathbf{x} + \mathbf{a}^i) \geq \tau$, thus \mathbf{a}^i is a feasible action (e.g., $\mathbf{a}^1, \mathbf{a}^2, \mathbf{a}^3, \mathbf{a}^5, \mathbf{a}^6$). In the objective space (b), the goal is to exclude all feasible actions \mathbf{a}^i that have other actions \mathbf{a}^j with lower (thus better) values of objectives $\mathbf{c}(\mathbf{a}^i) \geq \mathbf{c}(\mathbf{a}^j)$ (e.g., \mathbf{a}^2 has higher values than \mathbf{a}^1 and \mathbf{a}^3). The feasible actions with no other feasible action with higher values are called Pareto-optimal (e.g., $\mathbf{a}^1, \mathbf{a}^3, \mathbf{a}^6$).

3.1 Background

In the following, we consider a binary classification problem where each sample comprises a vector of features $\mathbf{x} \in \mathcal{R}^d$, where d is the number of features, and a binary outcome $y \in \{0, 1\}$. We represent the classifier by a **decision function** $r(\bullet)$ that is 1 when $r(\mathbf{x}) \geq \tau$ and 0 when $r(\mathbf{x}) < \tau$, where τ is a given threshold. For the sake of simplicity, 1 will be the target outcome.

An **action** \mathbf{a} on a sample is a set of changes that can be positive (increase the value of the feature) or negative (decrease the value of the feature), but not both for the same feature. Without losing generality, we will focus only on non-negative changes $a_i \geq 0, \forall i \in \{1, \dots, d\}$. In this context, a feasible action $\mathbf{x} + \mathbf{a}$ that achieves the target outcome $r(\mathbf{x} + \mathbf{a}) \geq \tau$ is called a counterfactual antecedent. Fig. 2(a) illustrates these concepts.

Definition 1 - Feasible action.

An action $\mathbf{a} \in \mathbb{R}^d$ belongs to the feasible set of solutions \mathcal{A} if and only if it achieves the desired outcome $r(\mathbf{x} + \mathbf{a}) \geq \tau$.

We propose using **multiple objective** functions $c_1(\bullet), \dots, c_m(\bullet)$ where each objective function $c_i(\bullet), i \in \{1, \dots, m\}$ is designed to search for a specific desire from the user. When the objectives are conflicting, it might be impossible to find a single feasible action with the best performance since one objective's optimization might negatively impact others. However, it is possible to define a partial order in multiple objectives and define a dominant action — an action with all objectives better than other actions.

Definition 2 - Partial order.

Ordering relation on partially ordered sets occurs when all components are ordered in the same sense. We use the symbols \preceq and \succeq to describe the ordering relations on partially ordered sets; for example, $\mathbf{x} \preceq \mathbf{y}$ is equivalent to $x_i \leq y_i, \forall i \in \{1, \dots, m\}$.

Definition 3 - Dominant action. A feasible action $\mathbf{a} : r(\mathbf{x} + \mathbf{a}) \geq \tau$ dominates \mathbf{a}' if and only if $\mathbf{c}(\mathbf{a}) \preceq \mathbf{c}(\mathbf{a}')$.

The concept of a dominant action used here is based on the dominance concept in multi-objective optimization where given two solutions \mathbf{a} and \mathbf{a}' , \mathbf{a} dominates \mathbf{a}' iff $c_i(\mathbf{a}) \leq c_i(\mathbf{a}'), \forall i \in \{1, \dots, m\}$ and exists j such that $c_j(\mathbf{a}) < c_j(\mathbf{a}')$ (Miettinen, 1999). We call Pareto-optimal action a feasible action that no other feasible action dominates. Fig. 2(b) illustrated this concept that is mathematically defined below.

Definition 4 - Pareto-optimal action (Miettinen, 1999).

Consider an objective function vector $\mathbf{c}(\bullet) : \mathbf{R}^d \Rightarrow \mathbf{R}^m$ that we want to minimize, and a feasible set of solutions \mathcal{A} . An action \mathbf{a}^* is Pareto-optimal iff there is no action $\mathbf{a} \in \mathcal{A}$ that dominates \mathbf{a}^* .

Finally, another important concept for the proposed algorithm is monotonicity.

Definition 5 - Monotonicity w.r.t. a partial order.

Given any two actions $\mathbf{a} \in \mathbf{R}^d$ and $\mathbf{a}' \in \mathbf{R}^d$ such that $a_i \geq a'_i, \forall i \in \{1, \dots, d\}$, a function vector $\mathbf{f}(\bullet) : \mathbf{R}^d \Rightarrow \mathbf{R}^m$ is monotone if only if $f_j(\mathbf{a}) \geq f_j(\mathbf{a}'), \forall j \in \{1, \dots, m\}$.

Monotonicity helps our branch-and-bound (B&B) algorithm to foresee if a search space is fruitful or not, avoiding unnecessary computation. Monotonicity is a requirement for all objective functions; however, it is optional for the decision functions (classifiers), although a monotone decision function improves the algorithm's performance.

3.2 Model-Agnostic Pareto-optimal Counterfactual Antecedents Mining

The proposed algorithm is a tree-based branch-and-bound (Lawler and Wood, 1966) algorithm called MAPOCAM (Model-Agnostic Pareto-Optimal Counterfactual Antecedents Mining); we detail it in Algorithm 1. In this recursive algorithm, each parameter's call is a node that can create other nodes with a branching procedure or avoid exploring subsequent nodes on a pruning procedure. In essence, MAPOCAM can examine every combination of changes. However, it cleverly bounds nodes if it preemptive knows that no subsequent branching will improve the optimal set of solutions. This branching procedure is the cornerstone of the algorithm's efficiency.

Algorithm 1 Model-Agnostic Pareto-optimal Counterfactual Antecedents Mining

Require: A sample \mathbf{x} , an objective function $c(\bullet)$, a decision rule $r(\bullet)$, a threshold τ , and a number of allowed changes k .

```

1: procedure ENUMERATE( $\mathbf{a}, \mathcal{D}$ )
2:   if  $|i : a_i \neq 0 \forall i \in \mathcal{D}| > k$  or  $\exists \mathbf{a}' \in \mathcal{A} : c(\mathbf{a}) \succeq c(\mathbf{a}')$  then
3:     return
4:   end if
5:   if  $r(\bullet)$  is monotone and  $r(\mathbf{x} + \bar{\mathbf{a}}^*) < \tau$  then
6:     return
7:   end if
8:   if  $r(\mathbf{x} + \mathbf{a}) \geq \tau$  then
9:      $\mathcal{A} = \mathcal{A} \cup \{\mathbf{a}\}$ 
10:    return
11:  end if
12:   $i = \text{SELECT\_FEATURE}(\forall i : i \notin \mathcal{D})$ 
13:  for  $\forall \mathbf{a}' : a'_i \geq a_i$  do
14:    ENUMERATE( $\mathbf{a}', \mathcal{D} \cup \{i\}, \mathcal{A}$ )
15:  end for
16: end procedure
17:  $\mathcal{A} = \{\}, \mathcal{D}^0 = \{\}$ 
18:  $\mathbf{a}_i^0 = 0 \forall i \in \{1, \dots, d\}$ .
19: ENUMERATE( $\mathbf{a}^0, \mathcal{D}^0$ )
20: return  $\mathcal{A}$ 

```

A **node** consists of an **action** $\mathbf{a} : a_i \geq 0, \forall i \in \{1, \dots, d\}$ and a set of fixed features \mathcal{D} (line 1). The **branching** procedure consists of selecting a non-fixed feature $i \notin \mathcal{D}$ (line 12), calling the recursive function (line 14) with different actions' values $a'_i \geq 0$ for that feature i that now is fixed $\mathcal{D}' \equiv \mathcal{D} \cup i$ for the subsequent calls. All features with no made decision are considered null $a_i = 0, \forall i \notin \mathcal{D}$ and might have their value altered in the subsequent branching.

The **bounding step** consists of finding any branch that will not improve the optimal set of solutions; thus, the node should not go deeper. This algorithm works with three possible conditions: (i) bounding any node that has more than k changes ($|i : a_i \neq 0 \forall i \in \mathcal{D}| > k$ in line 2); (ii) bounding any node that will not improve the objective functions ($\mathbf{a}', \mathcal{D}'$) : $\mathbf{c}(\mathbf{a}) \preceq \mathbf{c}(\mathbf{a}'), \forall \mathbf{a} \in \mathcal{A}$ when we already have a set of feasible solutions \mathcal{A} (line 2); and (iii) bounding nodes that any subsequent bounding will not generate a feasible solution (line 5).

In Appendix A.1, we show two essential properties of the branching procedure. First, by enforcing monotonicity in the objective functions, it is possible to see that any subsequent branching of a node with dominated action will be dominated (Corollary 1). Second, by enforcing monotonicity on the classifier, it is possible to prove when all subsequent node branches are unfeasible because the maximal achievable action $\bar{\mathbf{a}}^*$ is not feasible (Theorem 1). In Appendix A.1, we also prove two properties of the algorithm: First, Algorithm A.1 finds all Pareto-optimal solutions (Theorem 2). Second, the time complexity of Algorithm 1 is $T(d, k) = \mathcal{O}((bd)^k)$ (Theorem 3), where d is the number of variables, b is the maximum number of possible states of each variable, and k is the maximal number of allowed changes.

Worth mentioning that the procedure `SELECT_FEATURE` in Algorithm 1 is responsible for indicating the next feature to be explored. Ideally, the procedure `SELECT_FEATURE` will return the feature that most likely helps to find a feasible solution. We give further details in Appendix B.

3.3 Objective functions

The monotonicity of the objective functions is the main property required by the algorithm — in Appendix A.2, we define some basic properties of these monotonic functions. To define the monotonic objective functions used in this paper, we need to know that we have a set of N samples $\mathbf{x}^i \in \mathcal{R}^d$, $i \in \{1, \dots, N\}$, where d is the number of features. We also define the working sample $\mathbf{x} \in \mathcal{R}^d$ that we want to obtain a different outcome. We present three objective functions: percentile change, feature change, and number of changes.

Definition 6 Maximal percentile change (MPC).

First, let's define the percentile for each feature j :

$$l_j(\mathbf{x}) = \frac{|\{i | \hat{x}_j^i \geq x_j, \forall i \in \{1, \dots, N\}\}|}{N} \times 100 \quad (1)$$

where $|\bullet|$ is the cardinality of the set.

Given that, a percentile change of an action \mathbf{a} for the feature j is the absolute difference in percentiles $|l_j(\mathbf{x}) - l_j(\mathbf{x} + \mathbf{a})|$. Thus, the maximal percentile change $c(\mathbf{a})$ of an action \mathbf{a} is $\max(|l_j(\mathbf{x} + \mathbf{a}) - l_j(\mathbf{x})|, j \in \{1, \dots, d\})$.

Definition 7 j -th feature change.

The feature change for the feature j consists on the magnitude of an action \mathbf{a} for the feature j : $c_j(\mathbf{a}) = a_j$.

Definition 8 Number of changes.

Consists on counting the number of changes (non-zero values) of an action \mathbf{a} : $c(\mathbf{a}) = |\{\mathbf{a}_j | \mathbf{a}_j \neq 0, \forall j \in \{1, \dots, d\}\}|$.

Theorem 5 in Appendix A.2 recognizes that any Pareto-optimal solution of monotonic objective function from any nature will be a Pareto-optimal solution using the feature changes as objectives. Thus, resorting to feature changes as objectives will create a pool of counterfactual antecedents that will satisfy any preference, only needing to filter other sets of objectives to find the desired antecedents.

3.4 Monotonicity on classifiers

Another characteristic that a classifier can have is its monotonicity. It helps to avoid branches with no feasible actions, but it is an optional property to MAPOCAM. In Appendix A.3, we show that some linear classifiers and ensembles of monotone classifiers are monotone.

3.4.1 Enforcing monotonicity

Some methods enforce the machine learning model to have a monotone behavior w.r.t. to any desired feature – the increase of that feature will always increase (or decrease) the outcome. In this category, we have LightGBM¹ (Ke et al, 2017) and Gupta’s research (Gupta et al, 2016).

3.4.2 Dealing with non-monotone classifiers

Despite being possible to enforce/guarantee the monotonicity of the classifiers/decision functions, it is also possible to anticipate that some high-quality classifiers cannot hold monotonicity. However, it is possible to use Algorithm 1 to find Pareto-optimal actions for non-monotone classifiers. We can see that conditions on Line 2 and Line 8 hold despite the lack of monotonicity in the classifier as long as the objective function is monotone: any actions with higher costs are not interesting. However, the second condition in Line 5 cannot hold since it can exist a feasible action $\mathbf{a}' : \mathbf{a} \succeq \mathbf{a}' \succeq \bar{\mathbf{a}}^*$ in non-monotone classifiers. Given that, the first condition ($r(\bullet)$ is monotone) in Line 5 ensures the second condition ($r(\mathbf{x} + \bar{\mathbf{a}}^*) < \tau$) only being applied if the decision rule is monotonic.

3.4.3 Other similar properties

The main goal of monotonicity for classifiers is to assess if the desired probability is achievable in the subsequent branching in the B&B optimization. It is possible to obtain that without enforcing monotonicity in decision trees. Each node in B&B constrains the value of actions a_i for all features i in the set of taken decisions \mathcal{D} . With that, it is possible to filter which leaves are reachable in the decision trees by the counterfactuals on that node. Thus, finding that no leaf in decision trees reaches the desired probability can rule out the correspondent node in B&B optimization. We also use this methodology to speedup the algorithm in the experiments.

3.5 Counterfactual antecedent index

Let us suppose the method generated a set of solutions $\mathcal{A} \equiv \{\mathbf{a}^1, \dots, \mathbf{a}^{|\mathcal{A}|}\}$. An interesting property would be to observe how much each feature $i \in \{1, \dots, d\}$ helps in finding a counterfactual antecedent. To do that, we propose a counterfactual antecedent index (CA_i) for each feature i :

$$CA_i = \frac{1}{\nu_i} \sum_{\mathbf{a}^j \in \mathcal{A}} \frac{a_i^j - \max a_i}{\max a_i} \times \nu_j \quad (2)$$

¹ It can be set with the `monotone_constraints` parameters at lightgbm.readthedocs.io/en/latest/Parameters.html

where $\nu_i = |\mathbf{a}_i^j : \mathbf{a}_i^j \neq 0, \mathbf{a}_i^j \in \mathcal{A}|$ is the number of counterfactual antecedents with non-zero value in the feature i , and $\nu_j = |\mathbf{a}_i^j : \mathbf{a}_i^j \neq 0, j \in \{1, \dots, d\}|$ is the number of non-zero values in the action \mathbf{a}_i .

With this definition, it is possible to estimate how much, on average, a variable would need to change to create a counterfactual antecedent. The index penalizes features that need help from other features by multiplying by the number ν_j . Given that, if this number is bigger than 1, this variable would have to surpass its limit to create a counterfactual antecedent, indicating the need from other features to find a counterfactual antecedent.

This index helps summarize many counterfactual antecedents and aids the user in understanding the importance of each feature.

4 Experiments

The following experiments aim at evaluating MAPOCAM on two aspects: (Experiment 1) the importance of multiple antecedents to give broader reasoning of counterfactuals, as well as showing the importance of exploring counterfactuals on different machine learning models; and (Experiment 2) the correctness and computational performance of the method on monotonic and non-monotonic classifiers with single and multiple objectives.

To evaluate the capability of finding counterfactual antecedents, we consider the three credit datasets used in (Ustun et al, 2019): German Credit Dataset (labeled as *german*), Give Me Some Credit (labeled as *giveme*), and Taiwan Credit (labeled as *taiwan*). These datasets have 1000, 120269, and 30000 samples and 27, 10, and 17 features, respectively. All datasets have socio-economic and credit/bank history as independent variables. The dependent variable indicates whether a person deserves credit or not. Given that, a counterfactual goal is to provide a series of changes suggestion to allow a person classified as a high risk to have their credit granted.

4.1 Experiment 1 - Pareto-optimal counterfactual antecedents

In this experiment, we create Pareto-optimal counterfactual antecedents for a single sample of *credit* to show the capability of MAPOCAM to work on a wide range of scenarios. In Section 4.1.1, we evaluate how MAPOCAM behaves on extracting counterfactual antecedents on monotone classifiers., Section 4.1.2 discusses how some important classifiers (such as decision trees and neural networks) cannot hold the monotonicity property and how different the mined counterfactual antecedents are.

We limit the counterfactual antecedents to change at most three features for any case in this experiment. We choose value three because it can attain a good variety of antecedents without demanding excessive computational effort. Fig. 3, Fig. 4, Fig. 7, Fig. 8, and Fig. 9 depict the enumeration of counterfactual antecedents. Following the same pattern of Fig. 1, the first column (Orig column)

shows the original values inside the squares, and the following columns (C1, C2, and others) illustrate a set of counterfactual antecedents with the new value inside the squares. We can observe the changes' sparsity (only a few squares per column); we also use the colors to differentiate the features easily. We sort the counterfactual antecedents by MPC cost: lower cost at left and higher cost at right.

4.1.1 Monotone classifiers

In this experiment, we show the algorithm's capability of enumerating a representation of all Pareto-optimal solutions. This experiment uses the *credit* dataset to train an l_2 penalized Logistic Regression using Scikit-learn (Pedregosa et al, 2011) and selects a sample to extract counterfactual antecedents using MAPOCAM. We use two multi-objective formulations to show the behavior in more than one scenario: (i) minimize the action for every feature; (ii) minimize the maximal percentile change (MPC) cost and the number of changes — Fig. 3 and Fig. 4 show the results for formulation (i) and (ii), respectively.

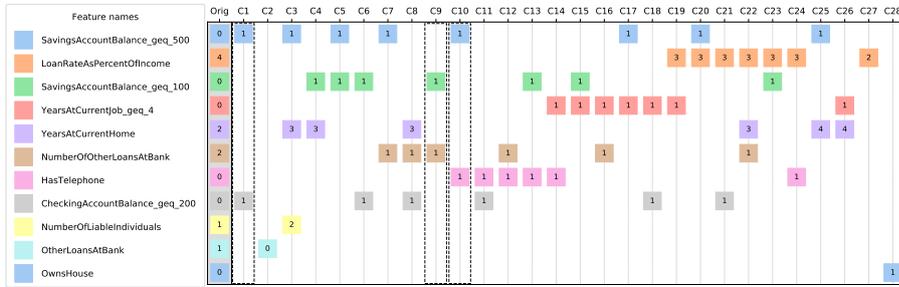


Fig. 3 Enumeration for Logistic Regression of counterfactual antecedents with Pareto-optimal values when each feature is considered as an objective function. Orig column is the original sample and the columns C1 to C28 are counterfactual antecedents.

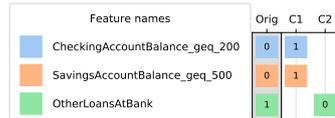


Fig. 4 Enumeration for Logistic Regression of counterfactual antecedents with Pareto-optimal values when the MPC cost and the number of changes are the objectives. Orig column is the original sample and the columns C1 and C2 are counterfactual antecedents.

Fig. 3 shows a vast set of counterfactual antecedents allowing multiple combinations. Despite having the lowest MPC cost, increasing the balance above \$500 in the savings account suggested in (C1) might be hard for the

user. However, other antecedents might be more accessible to the user: (C9) increasing savings account balance above \$100 and close a loan in other banks; or (C10) increasing checking account above \$200 and signing in to have a telephone. Our technique, using this formulation, works as an enumeration of every possible combination of features' changes without depicting solutions with an even more profound feature change. This modeling gives a wide range of possibilities for users to find the most suitable counterfactual antecedent.

Using the second formulation, MAPOCAM gets two counterfactual antecedents (See Fig. 4). (C1) shows that increasing checking balance above \$200 simultaneously on increase savings above \$500 is a counterfactual antecedent with the lowest MPC cost. Meanwhile, (C2) shows that ending loans at other banks is an alternative with higher MPC costs but with a single change. Both antecedents represent the designer's preferences also enumerated in Fig. 3. It shows that no matter user preferences, the resulting trade-offs can be filtered from the enumeration using feature changes as objectives, as said in Section 3.3.

4.1.2 Non-monotone classifiers

In this experiment, we explore the capability of MAPOCAM to find counterfactual antecedents on non-monotonic classifiers. Fig. 5 shows the impact of increasing the variables `LoanAmount`, `LoanDuration`, and `Age` on the probability of having the credit granted for an adjusted multilayer perceptron (MLP). Here we can see that MLP does not have monotonic behavior since the probability might increase or decrease for the same variable.

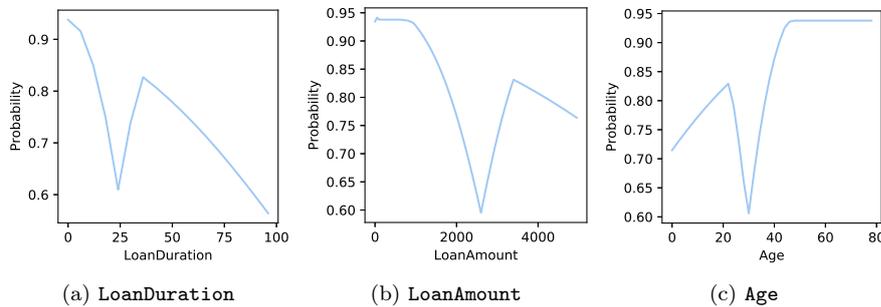


Fig. 5 Representation of impact of increasing each variable in the probability of having credit granted in a multilayer perceptron.

Fig. 6 shows an adjusted tree with a depth of three. To understand its non-monotonicity we rely on a sample with the following attributes: `CriticalAccountOrLoansElsewhere=1`, `OtherLoansAtBank=1`, `Age=41`, `LoanDuration=16`, and `LoanAmount=7000`. If we change `CriticalAccountOrLoansElsewhere` from 1 to 0, its probability of getting credit will decrease from 45.5% to 14.3%. However, if we change `LoanAmount` from 7000 to 6000, changing

`CriticalAccountOrLoansElsewhere` from 1 to 0, would cause the probability to increase from 45.5% to 76.1%. This change of behavior, depending on other attributes, characterizes `CriticalAccountOrLoansElsewhere` as a non-monotone attribute.

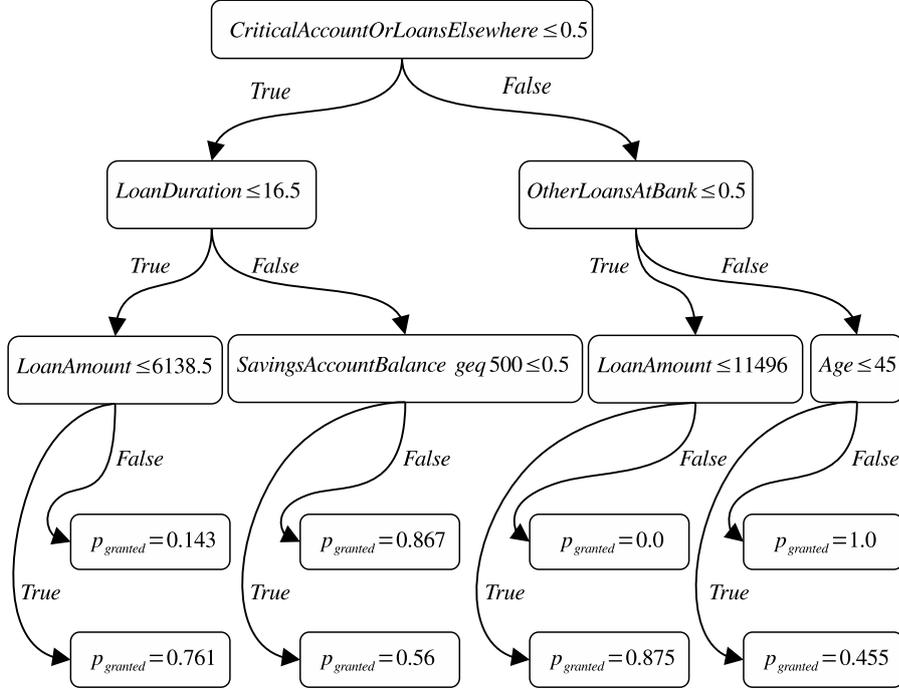


Fig. 6 Representation of a decision tree for the German dataset with depth of 3.

Both examples show relevant machine learning models that do not preserve the monotonicity property, but they are still in demand of a robust methodology capable of finding counterfactual antecedents for those complex models. Selecting the same sample investigated in Fig. 3 and Fig. 4, we used MAPOCAM to enumerate counterfactual antecedents to show how different they are and how different machine learning models uncover distinct counterfactual antecedents. We chose three machine learning models: (a) the multilayer perceptron, which has the probability behavior depicted in Fig. 5; (b) the decision tree illustrated in Fig. 6; and (c) a gradient boosting trees (LightGBM) to complement the decision trees with a more robust classifier.

Fig. 7, Fig. 8, and Fig. 9 show valuable suggestions to change the prediction given by non-monotone classifiers. It worth observing that new features are essential to the decision depending on the classifier: (1) changing `CriticalAccountOrLoansElsewhere` is an antecedent on Decision Trees (C1 in Fig. 8) and LightGBM (C2, C3, and others, in Fig. 9), but it is not in Logistic



Fig. 7 Enumerations for a multilayer perceptron (whose monotonicity is depicted in Fig. 5) with Pareto-optimal values when each feature is considered as an objective function. Orig column is the original sample and the columns C1 to C12 are counterfactual antecedents.

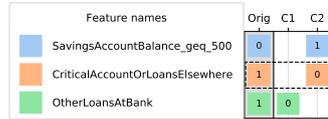


Fig. 8 Enumerations for a decision tree (depicted in Fig. 6) with Pareto-optimal values when each feature is considered as an objective function. Orig column is the original sample and the columns C1 and C2 are counterfactual antecedents.

Regression (Fig. 3) and Neural Networks (Fig. 7). (2) changing **RentsHouse** is an antecedent on LightGBM (C3 and C7 in Fig. 9) and Neural Networks (C4 and C5 in Fig. 7), but it is not in Logistic Regression (Fig. 3) and Decision Trees (Fig. 8). The variations found by enumeration in distinct learning machines show the importance of using a model-agnostic algorithm. MAPOCAM compares the classifiers’ counterfactual antecedents and promotes richer analysis to understand the problem in real life.

4.2 Experiment 2 - Comparing correctness and performance

In this experiment, we rely on the three datasets (german, giveme, and taiwan) to assess the capability of finding the optimal solution (on single-objective) or set of Pareto-optimal solutions (on multi-objective). Moreover, we analyze the computational time of our approach. No matter the experiment, we first train a machine learning model (which can be Logistic Regression or lightGBM classifier) in a training set and then perform a test on a separate data set of 100 samples. Finally, we use MAPOCAM (and the baselines) to extract counterfactual antecedents from samples with an undesired outcome on the test set, allowing us to access the algorithms’ correctness and performance.

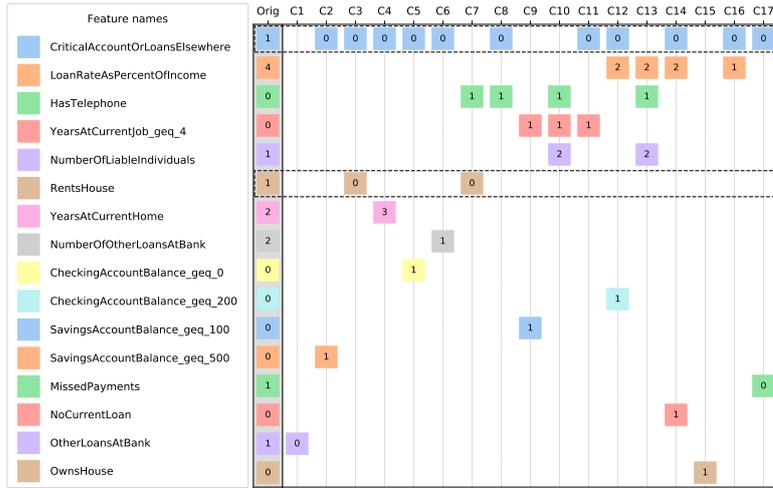


Fig. 9 Enumerations for a LightGBM with Pareto-optimal values when each feature is considered as an objective function. Orig column is the original sample and the columns C1 to C17 are counterfactual antecedents.

4.2.1 Baselines

We have different sets of baselines for single and multi-objective experiments. For single-objective experiments (using the objective of MPC cost described in Section 3.3), we compare MAPOCAM with a mixed-integer formulation (to compare with an exact approach) and a greedy approach (to have an algorithmic reference). The mixed-integer formulations consist of optimization models that extract properties from specific classifiers to find counterfactual antecedents, thus demanding individualized approaches: *mi-logistic* for logistic regression (based on Ustun et al. (Ustun et al, 2019)) and *mi-trees* for the ensemble of decision trees (based on Cui et al. (Cui et al, 2015)). The *greedy* approach is not individualized. Given an intermediate action, the greedy algorithm finds a new action with the best performance to participate in the next iteration. This action should have a change in a single feature, and the performance is the ratio between the MPC cost of making such action and the improvement in the decision rule. This iterative procedure continues until it finds counterfactual antecedents with the desired outcome, or no other action is capable of making improvements in the decision rule.

For multi-objective experiments (using the objective of MPC cost vs. the number of changes and j -th feature change described in Section 3.3), we also compare MAPOCAM with a mixed-integer formulation (to compare with an exact approach) and a brute-force approach (to have an algorithmic reference). The mixed-integer formulations in multi-objective context consist of approaches that execute the mixed-integer optimization several times: *po-mi-logistic* for logistic regression and *po-mi-trees* for the ensemble of decision trees. At each time, we add a constraint that forbids the optimization to find counterfactual

antecedents dominated by the previously found solutions—resorting to a procedure used in general mixed-integer problems (Sylva and Crema, 2004) that is further detailed in Algorithm 2 in Appendix B. We proposed both approaches (*po-mi-logistic* and *po-mi-trees*) as simple extensions to Ustun et al. (Ustun et al, 2019) and Cui et al. (Cui et al, 2015) because there is no Pareto-optimal approach for counterfactual antecedents in the literature. The *brute-force* approach is not individualized and is designed similarly to Algorithm 1 without the code from Line 2 to 11 to avoid any early bounding, storing every counterfactual antecedent, even the dominates. Thus, it enumerates all possible actions without trying to stop the explorations early. Both multi-objective formulations and methods have the number of changes limited to three to attain a good variety of antecedents without demanding excessive computational resources.

Worth mentioning that Ustun et al. (Ustun et al, 2019) indicate a method to find multiple counterfactuals; however, we decided to promote a modification to find all Pareto-optimal counterfactual antecedents instead of the suggested method to be comparable with the proposed method. We did not include other methods that find multiple counterfactual antecedents because they can not find Pareto-optimal counterfactual antecedents. Also, they use different norms (to enforce diversity (Karimi et al, 2020; Mothilal et al, 2020)), generating antecedents of different nature and even generating dominated antecedents. Thus, the execution time would not be comparable to the proposed method.

4.2.2 Experiment 1 - Comparing performance on logistic regression

In this experiment, we consider the three credit datasets and train a l_2 penalized logistic regression using Scikit-learn (Pedregosa et al, 2011) library. After cross-validation training, we define a threshold τ to guarantee the probability of credit granted to be at least 80%. We use three optimization approaches to assess the correctness and performance: MAPOCAM, *mi-logistic*, and *greedy*. Table 1 shows the time performance of such algorithms. Worth mentions that MAPOCAM and *mixed-integer* always find the best solution while *greedy* algorithms find 38%, 24%, and 1% optimal solutions for the datasets.

Table 1 Average execution time, in seconds, for finding an action optimizing the MPC cost.

	german	giveme	taiwan
mi-logistic	0.04	0.01	0.02
greedy	0.05	0.03	0.03
MAPOCAM	0.04	0.01	0.11

This experiment verifies that the proposed method can find the best result in all situations, showing MAPOCAM’s correctness. The computational cost is comparable to other approaches in the two datasets, being costly only in the taiwan dataset but still quite acceptable. The higher computational cost, mainly when compared against mixed-integer, can be justified by two factors.

(i) *Data structure and programming language*: our approach is implemented in Python with a NumPy array as a data structure, while the mixed-integer uses a CPLEX² package implemented in C. (ii) *The mixed-integer approach relies on more information to tackle the problem*: the mixed-integer formulation makes use of additional information in every node to find a tighter projection of the “future” of the node. In contrast, our approach does not rely on any of this information. This fact renders our approach more general and capable of dealing with any decision function such as decision trees or multilayer perceptron.

4.2.3 Experiment 2 - Enumeration of Pareto-optimal counterfactual antecedents

In this experiment, we show the algorithm’s capability of enumerating a representation of all Pareto-optimal solutions. We use the same training procedure and threshold of the previous experiment to explore two multi-objective formulations: (i) minimize the action for every feature; (ii) minimize the MPC cost and the number of changes.

Table 2 shows the time performance of the MAPOCAM, *mixed-integer*, and *brute-force* approach; we also confirm that both searches find the same number of solutions for all of the samples. The proposed method tends to perform similarly to the *mixed-integer* approach, being faster than a *brute-force* approach in both scenarios: with two objectives (MPC cost vs. # changes) and more than two objectives (features). Notice that MAPOCAM is costly than *mixed-integer* with two objectives (MPC cost vs. # changes), but it performs better with more than two objectives (features). This difference is mainly caused by evaluating the MPC cost — while this computation is costly in our approach, it is not relevant in the mixed-integer since the model itself embeds this information — since the cost of comparing the features is low.

Table 2 Average execution time, in seconds, for finding a representation of a Pareto-front.

conflict	method	german	giveme	taiwan
MPC vs. #changes	<i>po-mi-logistic</i>	0.01	0.01	0.01
	<i>brute-force</i>	0.26	59.30	21.90
	MAPOCAM	0.03	0.01	0.32
features	<i>po-mi-logistic</i>	0.02	0.05	0.54
	<i>brute-force</i>	0.23	12.30	2.90
	MAPOCAM	0.03	0.01	0.13

4.2.4 Experiment 3 - Comparing performance on non-monotonic classifiers

We know that the proposed method works properly on non-monotonic classifiers. This experiment aims to examine the performance of the proposed algorithm in a forest of decision trees. We consider the three credit datasets and train a

² Available at ibm.com/analytics/cplex-optimizer

Random Forest from Scikit-learn (Pedregosa et al, 2011), with five estimators, depth of five, and 31 leaves. We define a threshold of τ to guarantee the probability of credit granted to be at least 80% and use three optimization approaches to assess the correctness and performance.

Table 3 shows the time performance of such algorithms with one (MPC cost), two (MPC cost vs. # changes), and multiple objectives (feature). Except for the *greedy* algorithm, all methods were capable of finding the Pareto-optimal set of solutions. In this experiment, we show two versions of the proposed algorithm: (i) MAPOCAM-N is a version with no unfeasibility branching, and (ii) MAPOCAM-P explores the possibility to evaluate all leaves from the random forest and find the maximal probability of the optimization tree, considering that it already took some actions. This second approach can behave similarly from a monotone classifier, because it avoids unpromising branching, thus improving time performance.

Table 3 Average execution time, in seconds, for finding solutions on a random forest (*greedy not always find the best solution).

conflict	method	german	giveme	taiwan
MPC cost	<i>mi-trees</i>	0.02	0.02	0.03
	<i>greedy*</i>	0.01	0.02	0.02
	MAPOCAM-N	11.90	0.21	41.40
	MAPOCAM-P	0.10	0.02	0.10
MPC vs. #changes	<i>po-mi-trees</i>	0.04	0.04	0.27
	<i>brute-force</i>	4.20	1.10	17.10
	MAPOCAM-N	0.38	0.13	0.80
	MAPOCAM-P	0.08	0.01	0.10
feature	<i>po-mi-trees</i>	0.05	0.06	0.27
	<i>brute-force</i>	0.44	0.57	5.00
	MAPOCAM-N	0.21	0.36	0.70
	MAPOCAM-P	0.07	0.03	0.12

In this experiment, the MAPOCAM can find optimal solutions in all instances. When the predictor can foresee unfeasibility in the nod (MAPOCAM-P version), the performance is quite competitive. However, without resorting to methods that foresee unfeasible nodes (MAPOCAM-N version), the technique tends to be time-consuming. Worth mentioning that the results closer to the *brute-force* approach are usually related to the fact that these datasets have many unfeasible instances. It forces the proposed algorithm to cover all the search trees, making it similar to *brute-force*.

Since it is possible to observe gain on MAPOCAM-P performance, we investigate the monotonicity’s impact on the optimization procedure. In that regard, we explore the monotonicity constraints on the LGBMClassifier class from the lightGBM³ library. Given that, we train LGBMClassifier with ten estimators, depth of ten, and 63 leaves. We define a threshold of τ to guarantee the probability of credit granted to be at least 80% and use three optimization

³ `lightgbm.readthedocs.io`

approaches to assess the correctness and performance. We compare the same contenders from the previous comparison, but now we have three variations from the proposed algorithm. (i) MAPOCAM-N is a version with no unfeasibility branching. (ii) MAPOCAM-P explores the possibility of evaluating all leaves from the random forest and finding the maximal probability of the optimization tree, considering that it already took some actions. And (iii) MAPOCAM-M explores the monotonicity to avoid unfeasible branches. Table 4 shows the time performance of such algorithms with one (MPC cost), two (MPC cost vs. # changes), and multiple objectives (feature).

Table 4 Average execution time, in seconds, for finding solutions on a monotone lightgbm (*greedy not always find the best solution).

conflict	method	german	giveme	taiwan
MPC cost	<i>mi-trees</i>	0.02	0.02	0.10
	<i>greedy*</i>	0.03	0.08	0.09
	MAPOCAM-N	93.40	9.30	23.60
	MAPOCAM-P	0.02	7.10	8.10
	MAPOCAM-M	0.02	0.01	0.45
MPC vs. #changes	<i>po-mi-trees</i>	0.02	0.02	1.18
	<i>brute-force</i>	0.30	0.43	180.50
	MAPOCAM-N	0.30	0.42	12.50
	MAPOCAM-P	0.01	0.07	1.70
	MAPOCAM-M	0.01	0.01	0.72
feature	<i>po-mi-trees</i>	0.02	0.02	0.61
	<i>brute-force</i>	0.29	0.43	68.90
	MAPOCAM-N	0.29	0.43	13.10
	MAPOCAM-P	0.01	0.07	2.50
	MAPOCAM-M	0.01	0.01	0.89

This experiment shows that MAPOCAM can find optimal solutions in all of the instances. MAPOCAM has a competitive performance when it explores the monotonicity property (MAPOCAM-M version) of the classifier. Still, it has performance gain when it explores other tree classifiers’ properties to estimate the maximum probability on a counterfactual antecedent (MAPOCAM-P version). The difference in performance is due to the difference in computational effort requested by the foreseeing step on MAPOCAM-P and MAPOCAM-M. However, without resorting to methods that foresee unfeasible nodes (MAPOCAM-N version), the method tends to be time-consuming. Worth mentioning that the results closer to the *brute-force* approach are usually related to the fact that these datasets have many unfeasible instances. It forces the proposed algorithm to cover all the search trees, making the search similar to *brute-force*.

4.2.5 Discussion

In general, MAPOCAM achieves good performance when: (i) the classifiers are monotonic, or (ii) it explores the constraints of achievable leaves in decision trees.

The execution time is also manageable without any information about the model (MAPOCAM-N)—generally taking less than one minute. This application of the method leads to a flexible plug-and-play model-agnostic method capable of working with any classifier without providing any model-related procedure. Any other method in the literature needs model-related procedures to work, such as (i) creating a new mixed-integer model (Ustun et al, 2019; Cui et al, 2015), (ii) requiring differentiable functions (Mothilal et al, 2020), or (iii) transforming the model to another representation (Karimi et al, 2020).

Worth mentions that MAPOCAM can create a diverse set of counterfactual antecedents using multi-objective optimization. Theorem 5 in Appendix A.2 shows that MAPOCAM enumerates all proper counterfactual antecedents without resorting to norms and adjusting their parameters (Mothilal et al, 2020; Karimi et al, 2020). Also, *po-mi-logistic* and *po-mi-trees* are modifications from mixed-integer models (Ustun et al, 2019; Cui et al, 2015) that we adapted to interactively find Pareto-optimal counterfactual antecedents. Moreover, this paper is a pioneer in considering multi-objective optimization into counterfactual antecedent mining.

5 Case study - Enumeration of counterfactual antecedents for crimes in São Paulo

In this experiment, we consider the impact of counterfactual antecedents on Police Decision making on crime. As a case study, we investigate crime patterns in São Paulo city in Brazil. We gathered socioeconomic, urban, and crime information from census regions of São Paulo. The Center for the Study of Violence from the University of São Paulo (NEV-USP)⁴ provided us the criminal records. The Center for Metropolitan Studies (CEM)⁵ provided us with data about schools, bus stops, and bars. With this data, we ranked census regions by the rate of passerby crimes normalized by each region’s total population. Then, we labeled the 10% regions with higher rates as dangerous and fitted an l_2 regularized Logistic Regression to classify if an area is dangerous or not.

Given a census region classified as dangerous, MAPOCAM can create counterfactuals that give insights into what a decision-maker might do to convert that region into a safer area. Fig. 10 describes counterfactuals in Region 1. This region is known for having a high incidence of violent crimes in the 2000s. The original attributes of this region are depicted in the first column (Orig) of the figure, showing an impoverished region with a relatively low `HighIncomeHolder` (percentage of householders), `PermanentHousing` (percentage), `WaterSupply` and `SewageCollection` (service coverage). The other columns (C1, C2, and others) show possible counterfactuals to reduce criminality. For instance, C2 shows that increasing `PermanentHousing` from 0.843 to 0.876 is a counterfactual. On the other hand, C5 shows that a minor increase in `PermanentHousing`

⁴ nevsp.org

⁵ centrodametropole.fflch.usp.br/pt-br

mentioned disabling. This new scenario shows more realistic planning, such as modifying the bus routes (to reduce the number of passengers passing through this region), reducing the traveling time, improving the water supply, and incentivizing population growth. Note that we labeled the regions considering the criminality rate, so the suggestion to increase the population might reduce the rate, explaining this proposed change.

Fig. 10 and Fig.11 show a set of highly insightful counterfactual antecedents. These options help understand the prediction, but when the number of counterfactual antecedents is too high, it is hard to select the best suitable one and understand each feature’s impact on the learning. To better understand the overview of counterfactual antecedents before selecting one, Fig. 12 (a) presents the feature importance in a logistic regression (higher the value, more important is the variable), and Fig. 12(b)-Fig. 12(d) present the CA_i index described in Section 3.5 for three different regions in São Paulo: the already mentioned Region 1, the central and commercial Region 2, and the wealthy Region 3 (with headquarters of many corporations). The CA_i index indicates how much a feature has to change to create a counterfactual antecedent on relative average (near 0 – small change, near 1 – change to its maximum, between 1 and 3 – need another variable to create a counterfactual antecedent, and gray bar – unachievable).

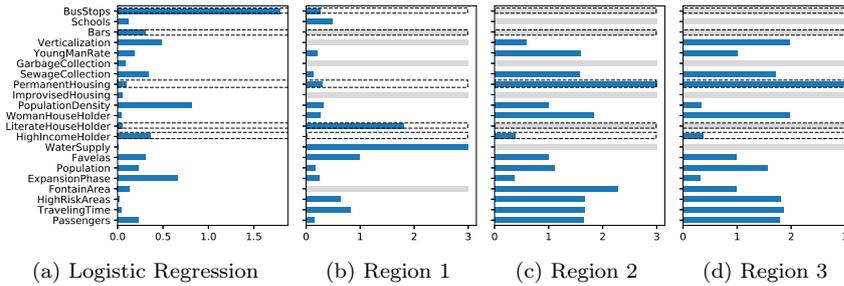


Fig. 12 This figure shows the representation of the importance of each feature. (a) represents the feature importance for the logistic regression learning process (higher the better); (b-d) represent how much a feature has to change to participate in a counterfactual antecedent (lower the better). The value is normalized by how much the variable can change; greater than one means that the variable usually needs other features to participate in a counterfactual antecedent. Grey bars represent features that never participate in a counterfactual antecedent, thus we represented them as having the largest cost – 3.

Fig. 12 shows exciting properties: (i) the CA_i index for each feature is different from the importance calculated with counterfactuals antecedents; (ii) the indexes on each region are different from each other. From these results, we can say that reducing crime requires different approaches for each region. Also, we observe that reducing the number of **BusStops** and **Bars** does not help achieve counterfactuals on Region 2 and 3 because these regions do not have those amenities. Socio-economic variables (**PermanentHousing**,

HighIncomeHolder, and LiterateHouseHolding) need more change to help Region 2 and Region 3 become safer, mainly because those regions already have better socio-economic conditions.

This case study depicts the importance of relying on counterfactual antecedents in Police Decision making. We could be able to identify particularities of each region pointing to more effort on solving socio-economic issues in more vulnerable regions, which would not necessarily be pointed by a single explanation for the whole model (e.g., feature importance in Logistic Regression). We envision a real impact on criminality control using our methodology with a richer dataset (e.g., information about road signaling, sidewalk preservation, wall painting, and urban lighting), which would allow suggestions for departments on urban administration help with criminality control.

6 Discussion and limitations

This research’s final product achieves a milestone in constructing a model-agnostic algorithm that finds diverse counterfactual antecedents using multi-objective optimization concepts. However, there are some limitations and future work that can be addressed in follow-up work.

Performance. Compared to other approaches that explore mixed-integer solvers (Cui et al, 2015; Ustun et al, 2019), we notice that our method is rarely better but has competitive performance. The mixed-integer approach has more available information (e.g., linear relaxations of the problem), more computational resources (can explore parallelism), and more cost-efficient programming language and data structure. Thus, it indicates that having a competitive performance shows the real value of our proposal, that is also more general and can deal with any classifier.

Therefore, the following steps can improve our approach:

- Use model information before the branching procedure on each node, such as: (1) explore the use of a linear relaxation for logistic regression; (2) use convex optimization to find heuristic solutions on neural networks; and (3) use satisfiability modulo theories (SMT) to verify the feasibility of some actions (Karimi et al, 2020).
- As a tree search, a node does not depend on information from other nodes. Thus, we could implement a parallel version that only needs to share the found counterfactual antecedents.
- As a first approach, we used the NumPy array as a data structure for counterfactual antecedents, but it can be improved using more sophisticated solutions.
- Another option to improve the performance is to use other search algorithms such as A^* . This strategy can take advantage of heuristics designed for every objective function.

Discretization. Given a set of possible actions, we proved that our algorithm could find every Pareto-optimal counterfactual antecedent. The optimality is

globally valid for decision trees we create a step for every interval between decisions on the features promoted by decision trees (or ensemble trees), thus covering all possibilities of changes for that decision tree. We use a simple stepwise discretization with absolute or relative steps for numerical features for linear classifiers and neural networks. However, this decision does not guarantee the optimal result for all discretizations. One solution could be to use a binary search for further refining when a counterfactual antecedent is found (Karimi et al, 2020). Improving the discretization can also help to find better-suited solutions and optimal value with lower costs. On the other hand, for non-monotone classifiers, it helps to find more counterfactual antecedents. Nevertheless, it can drastically change the search structure or require a specific design for every family of classifier methods.

Opposite direction actions. Another alternative to deal with non-monotone classifiers is to create two variables for each feature, i.e., one for the positive and one for the negative values. This strategy removes any conflict between the negative and positive directions. Also, it allows the design of a monotone objective function to deal with negative actions. We plan to work on it as an immediate future work; however, it needs additional theoretical and experimental developments to ensure a stable algorithm behavior.

Categorical variables. Another limitation in our proposal consists of dealing with categorical features. We did not consider this feature in the algorithm because it depends on the encoding strategy. However, it is possible to encode constraints in the algorithm that would allow using categorical variables. It might be an immediate future work, but we need more theoretical and experimental developments to ensure a stable algorithm’s stable behavior.

Visualization. All the experiments take care of choosing examples that do not have many possible actions; we did it to simplify the analysis. However, we notice cases with more than a hundred actions, which could be overwhelming to anyone analyzing the possible counterfactual antecedents. We are now working on creating a visual analytic system to explore the counterfactuals interactively. This system should also allow selecting and filtering counterfactuals with some properties, e.g., having (or not having) a feature of a specific change.

7 Conclusion

The main finding of this research is a tree-based algorithm, called MAPOCAM, that enumerates all Pareto-optimal counterfactual antecedents. Our technique helps to interpret model prediction and to mine counterfactual scenarios. Moreover, our algorithm possesses properties validated by empirical evidence. First, it is optimal; we found the same outcomes as a mixed-integer algorithm in single-objective problems and with a brute force approach in multi-objective situations. Second, it is model-agnostic; we tested our algorithms with linear models, multilayer perceptron, and decision trees. Third, it is diverse; we explore multiple counterfactual antecedents to analyze real-world applications.

Another exciting feature of this research relies on the counterfactual antecedent index. This index helps in exploring the real-world applications when the set of antecedents is numerous. For instance, in the case study of São Paulo, our tools gave a good overview of the usefulness of diverse counterfactual antecedents on crime analysis.

The use of multi-objective concepts by MAPOCAM diversifies the counterfactual antecedents, helping to provide information on the model's behavior for a given sample. Modeling each feature as an objective helps to find counterfactual antecedents with distinct sets of features and find trade-offs among them. Both aspects and the Pareto-optimality (no other counterfactual antecedent is better in all objectives) are keys to mine counterfactual antecedents. This research makes an essential step in this matter.

Acknowledgments

This work was supported by CNPq-Brazil (grants #312483/20180) and Getúlio Vargas Foundation.

References

- Aggarwal CC, Chen C, Han J (2010) The inverse classification problem. *Journal of Computer Science and Technology* 25(3):458–468
- Artelt A, Hammer B (2020) Convex density constraints for computing plausible counterfactual explanations. In: *International Conference on Artificial Neural Networks*, Springer, pp 353–365
- Breiman L (2001) Random forests. *Machine learning* 45(1):5–32
- Chen L, Lin X, Hu H, Jensen CS, Xu J (2015) Answering why-not questions on spatial keyword top-k queries. *Proceedings - International Conference on Data Engineering 2015-May*:279–290, DOI 10.1109/ICDE.2015.7113291
- Cui Z, Chen W, He Y, Chen Y (2015) Optimal action extraction for random forests and boosted trees. In: *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Association for Computing Machinery, New York, NY, USA, pp 179–188
- Gao Y, Liu Q, Chen G, Zheng B, Zhou L (2015) Answering why-not questions on reverse Top-k queries. *Proceedings of the VLDB Endowment* 8(7):738–749, DOI 10.14778/2752939.2752943
- Grath RM, Costabello L, Van CL, Sweeney P, Kamiab F, Shen Z, Lecue F (2018) Interpretable Credit Application Predictions With Counterfactual Explanations. *arXiv preprint arXiv:181105245* 1:1–9
- Gupta M, Cotter A, Pfeifer J, Voevodski K, Canini K, Mangylov A, Moczydlowski W, Van Esbroeck A (2016) Monotonic calibrated interpolated look-up tables. *The Journal of Machine Learning Research* 17(1):3790–3836
- He Z, Lo E (2012) Answering why-not questions on top-k queries. *Proceedings - International Conference on Data Engineering* pp 750–761, DOI 10.1109/ICDE.2012.8

- Karimi AH, Barthe G, Balle B, Valera I (2020) Model-agnostic counterfactual explanations for consequential decisions. In: International Conference on Artificial Intelligence and Statistics, PMLR, pp 895–905
- Ke G, Meng Q, Finley T, Wang T, Chen W, Ma W, Ye Q, Liu TY (2017) LightGBM: A highly efficient gradient boosting decision tree. *Advances in Neural Information Processing Systems 2017-Decem(Nips)*:3147–3155
- Krause J, Perer A, Ng K (2016) Interacting with predictions: Visual inspection of black-box machine learning models. In: *Conference on Human Factors in Computing Systems - Proceedings*, Association for Computing Machinery, New York, NY, USA, pp 5686–5697
- Lawler EL, Wood DE (1966) Branch-And-Bound Methods : A Survey. *Operations Research* 14(4):699–719
- Lu Q, Cui Z, Chen Y, Chen X (2017) Extracting optimal actionable plans from additive tree models. *Frontiers of Computer Science* 11(1):160–173
- Lucic A, Oosterhuis H, Haned H, de Rijke M (2019) Focus: Flexible optimizable counterfactual explanations for tree ensembles. *arXiv preprint arXiv:191112199*
- Lv Q, Chen Y, Li Z, Cui Z, Chen L, Zhang X, Shen H (2018) Achieving data-driven actionability by combining learning and planning. *Frontiers of Computer Science* 12(5):939–949
- Miettinen K (1999) *Nonlinear Multiobjective Optimization*. Springer, New York
- Mothilal RK, Sharma A, Tan C (2020) Explaining machine learning classifiers through diverse counterfactual explanations. In: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pp 607–617
- Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E (2011) Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* 12:2825–2830
- Poyiadzi R, Sokol K, Santos-Rodriguez R, De Bie T, Flach P (2020) Face: feasible and actionable counterfactual explanations. In: *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pp 344–350
- Rudin C (2019) Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence* 1(5):206–215
- Subramani S, Wang H, Balasubramaniam S, Zhou R, Ma J, Zhang Y, Whittaker F, Zhao Y, Rangarajan S (2016) Mining Actionable Knowledge Using Reordering Based Diversified Actionable Decision Trees. In: *Web Information Systems Engineering – WISE 2016*, Springer International Publishing, Cham, pp 553–560
- Sylva J, Crema A (2004) A method for finding the set of non-dominated vectors for multiple objective integer linear programs. *European Journal of Operational Research* 158(1):46–55, DOI 10.1016/S0377-2217(03)00255-8
- Tolomei G, Silvestri F, Haines A, Lalmas M (2017) Interpretable predictions of tree-based ensembles via actionable feature tweaking. *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data*

Mining Part F1296:465–474

- Ustun B, Spangher A, Liu Y (2019) Actionable recourse in linear classification. In: FAT* 2019 - Proceedings of the 2019 Conference on Fairness, Accountability, and Transparency, Association for Computing Machinery, New York, NY, USA, pp 10–19
- Wachter S, Mittelstadt B, Russell C (2018) Counterfactual Explanations Without Opening the Black Box: Automated Decisions and the GDPR. *Harvard Journal of Law & Technology* 31(2):841–887
- Yang C, Street WN, Robinson JG (2012) 10-year CVD risk prediction and minimization via inverse classification. In: IHI'12 - Proceedings of the 2nd ACM SIGHIT International Health Informatics Symposium, Association for Computing Machinery, New York, NY, USA, pp 603–609
- Yang Q, Yin J, Ling CX, Chen T (2003) Postprocessing decision trees to extract actionable knowledge. *Proceedings - IEEE International Conference on Data Mining, ICDM 1*:685–688, DOI 10.1109/icdm.2003.1251008
- Yang Q, Yin J, Ling C, Pan R (2007) Extracting actionable knowledge from decision trees. *IEEE Transactions on Knowledge and Data Engineering* 19(1):43–55

Appedices

A Proofs of the Theorems

In this appendix we present some theoretical proofs of properties concerning the proposed methodology.

A.1 Algorithm properties

To help on branching nodes in Algorithm 1 it is simple to see that from Definition 3 subsequent nodes of a node with dominated action will be dominated.

Corollary 1 *Let \mathbf{a} and \mathbf{a}' as in Definition 3 and consider that \mathbf{a}' and a set of taken decisions \mathcal{D}' defines a node. It is possible to see that no subsequent action in this node $\bar{\mathbf{a}} \succeq \mathbf{a}'$ will improve any monotone objective functions.*

To further improve the search, if the chosen decision rule $r(\mathbf{x})$ is monotone, we can define an **unfeasible node**.

Definition 9 - Unfeasible node. Given a node with its related actions \mathbf{a} and taken decisions \mathcal{D} , and given a monotone decision function $r(\bullet)$; the maximal achievable action $\bar{\mathbf{a}}^*$ for this node is given by:

$$\bar{a}_i^* = \begin{cases} a_i, & \text{if } i \in \mathcal{D} \\ \max a_i, & \text{otherwise} \end{cases} \quad (3)$$

Given that, a **node is unfeasible** iff $r(\mathbf{x} + \bar{\mathbf{a}}^*) < \tau$.

Theorem 1 *Let \mathbf{a} and \mathcal{D} as in Definition 9. We want to prof that any action $\bar{\mathbf{a}} : \bar{a}_i \geq a_i \forall i \notin \mathcal{D}$, is infeasible $r(\bar{\mathbf{a}}) < \tau$.*

Proof Let's suppose by contradiction that $r(\bar{\mathbf{a}}) \geq \tau$, since $\bar{\mathbf{a}}^* \geq \bar{\mathbf{a}}$, we have, by monotonicity of $r(\bullet)$ that $r(\bar{\mathbf{a}}^*) \geq r(\bar{\mathbf{a}}) \geq \tau$, what contradicts $r(\bar{\mathbf{a}}^*) < \tau$.

Notice that is only possible to define a node as unfeasible if the decision rule respects the monotone property. But when we have this property, we can preemptively discard this node.

There are some guarantees that Algorithm 1 is capable of detecting all optimal solutions and evaluate its time complexity.

Theorem 2 *All Pareto-optimal solutions are found.*

Proof We know that Algorithm 1 would search through all possible actions if no pruning condition (lines 2, 5 and 8 in Algorithm 1) is activated. Any Pareto-optimal solution $\mathbf{a}^* : \mathbf{c}(\mathbf{a}^*) \leq \mathbf{c}(\mathbf{a}), \forall \mathbf{a} \in \mathcal{A}$ would not be found iff the conditions in lines 2, 5, and 8 are activated by \mathbf{a}^* or a predecessor $\mathbf{a} \preceq \mathbf{a}^*$.

We can split the proof by contradiction into 3 cases:

- Line 2 would be activated if a predecessor $\mathbf{a} \preceq \mathbf{a}^*$ is dominated by a feasible solution $\mathbf{a}' \in \mathcal{A} : \mathbf{c}(\mathbf{a}') \preceq \mathbf{c}(\mathbf{a})$. This cannot be true since $\mathbf{c}(\mathbf{a}') \preceq \mathbf{c}(\mathbf{a}) \preceq \mathbf{c}(\mathbf{a}^*)$ would contradict the Pareto-optimality of \mathbf{a}^* .
- Line 5 would be activated if $r(\mathbf{x} + \bar{\mathbf{a}}^*) < \tau$. This cannot be true since $\mathbf{a}^* \preceq \bar{\mathbf{a}}^*$ and it is a feasible solution.
- Line 8 would be activated if $r(\mathbf{x} + \mathbf{a}) \geq \tau$. This cannot be true because this condition would indicate that \mathbf{a} is feasible, and $\mathbf{c}(\mathbf{a}) \preceq \mathbf{c}(\mathbf{a}^*)$ would contradict the Pareto-optimality of \mathbf{a}^* .

Since all conditions led to a contradiction the proof is done.

Theorem 3 *The worst-case complexity of the algorithm is $T(d, k) = \mathcal{O}((bd)^k)$, where d is the number of variables, b is the maximum number of possible states of each variable, and k is the maximal number of allowed changes.*

Proof We can see the cost of each recurrence is given by:

$$T(d, k) = \begin{cases} dbT(d-1, k-1) + db + d^2, & k > 0 \\ 1, & k = 0 \end{cases} \quad (4)$$

By substitution we can see, for $k > 0$ that:

$$T(d, k) \leq dbT(d-1, k-1) + db + d^2, \quad (5)$$

$$\leq db(b(d-1))^{(k-1)} + db + d^2, \quad (6)$$

$$\leq b^k(d^k - d) + db + d^2 \quad (7)$$

Given that, we can see that $T(d, k) = \mathcal{O}((bd)^k)$.

A.2 Objective function properties

Since monotonicity is a requirement for the algorithm, in this sections we discuss some properties about monotone functions as well as prove the monotone properties of the functions used in the experiment.

Definition 10 Opposite direction actions do not decrease the cost.

Defining a pair of actions \mathbf{b} , and \mathbf{a} . In which \mathbf{b} have opposite actions and \mathbf{a} is equivalent to \mathbf{b} but with *null* action in such way: $\mathbf{a}, \mathbf{b} : a_j = 0, b_j < 0, \forall j \in \mathcal{J}$ and $a_j = b_j, \forall i \notin \mathcal{J}$.

The objective functions must be defined like that: $c_k(\mathbf{a}) = c_k(\mathbf{b}), \forall k \in \{1, \dots, M\}$.

Theorem 4 *Any feasible solution with opposite action has a correspondent feasible action on monotone decision functions.*

Proof Let's suppose an action $\mathbf{b} : b_j < 0 \in \mathcal{J}$, that is feasible. By Definition 10, it is possible to see that exists an action $\mathbf{a} : a_j = 0, b_j < 0, \forall j \in \mathcal{J}$ and $a_j = b_j, \forall i \notin \mathcal{J}$ with the same cost.

By the monotonicity of $r(\bullet)$ it is possible to see that \mathbf{a} is also feasible.

Also, we want to show that using the j -th feature change for all features leads to an enumeration of Pareto-optimal solutions that will contain any set of Pareto-optimal solutions from other objective functions.

Theorem 5 *Any Pareto-optimal solution for a vector of monotone objective functions is a Pareto-optimal solution if we optimize the feature changes.*

Proof Knowing that $\mathbf{c}(\bullet)$ is the vector of functions for the j -th feature changes, and $\mathbf{c}'(\bullet)$ as a vector of function that is monotone.

Lets assume, by contradiction, that \mathbf{a}^* is a Pareto-optimal for $\mathbf{c}'(\bullet)$, and \mathbf{a} dominates \mathbf{a}^* for $\mathbf{c}'(\bullet)$. We also assume that exists k such that $c'_k(\mathbf{a}^*) < c'_k(\mathbf{a})$ to enforce the solution to have different objective function values.

Thus $c_j(\mathbf{a}^*) \geq c_j(\mathbf{a}), \forall j \in \{1, \dots, d\}$, meaning that $\mathbf{a}^* \succeq \mathbf{a}$. Since $c'_k(\bullet)$ is monotone $c'_k(\mathbf{a}^*) \geq c'_k(\mathbf{a})$. Thus being a contradiction.

A.3 Monotone classifier's properties

Another important trait that a classifier can have is its monotonicity. It helps on avoiding branches with no feasible actions.

A.3.1 Logistic regression and SVM

The decision rules from Logistic regression and SVM can be formulated as $p(\mathbf{x}) = \frac{1}{1+e^{-\theta^\top \mathbf{x}+b}}$ (SVM probability being formulated as in Platt Scaling). Assuming that $\theta_i \geq 0$ for all i mutable (consequently $\mathbf{a}_j = 0$ for all j non-mutable).

Theorem 6 *The decision function $p(\mathbf{x}) = \frac{1}{1+e^{-\theta^\top \mathbf{x}+b}}$ is monotone.*

Proof Since $\mathbf{a} \succeq \mathbf{0}$:

$$\frac{1}{1+e^{-\theta^\top (\mathbf{x}+\mathbf{a})+b}} = \frac{1}{1+e^{-\theta^\top \mathbf{a}} e^{-\theta^\top (\mathbf{x})+b}} \quad (8)$$

Since $e^{-\theta^\top \mathbf{a}} \leq 1$ then $e^{-\theta^\top \mathbf{a}} e^{-\theta^\top (\mathbf{x})+b} \leq e^{-\theta^\top \mathbf{x}+b}$. Consequently $p(\mathbf{x} + \mathbf{a}) \geq p(\mathbf{x})$.

A.3.2 Ensemble of monotone classifiers

Supposing an ensemble of classifiers $p^{ens}(\mathbf{x}) = \sum_{i=1}^E \frac{p^i(\mathbf{x})}{E}$ where $p^i(\mathbf{x})$ is monotone for all $i \in \{1, \dots, E\}$.

Theorem 7 *The decision function $p^{ens}(\mathbf{x})$ is monotone.*

Proof Let's suppose $\mathbf{x}' \succeq \mathbf{x}$. Given that:

$$p^i(\mathbf{x}') \geq p^i(\mathbf{x}), \quad \forall i \in \{1, \dots, E\} \quad (9)$$

$$\sum_{i=1}^E \frac{1}{E} p^i(\mathbf{x}') \geq \sum_{i=1}^E \frac{1}{E} p^i(\mathbf{x}) \quad (10)$$

$$p^{ens}(\mathbf{x}') \geq p^{ens}(\mathbf{x}) \quad (11)$$

Proving that $p^{ens}(\bullet)$ is monotone.

B Implementation details

It is crucial to notice that in all formalization of the method, all actions are additive $\mathbf{a} \succeq \mathbf{0}$ to simplify the notation. However, it is possible to explore actions that can subtract a value of sample \mathbf{x} , and the objective functions and the decision functions will handle only negative actions to sustain the monotone property.

Also, worth mentioning the actions must be discretized to work in the proposed method (as well as other approaches (Ustun et al, 2019; Cui et al, 2015)). To do so, we used the strategy in Ustun et al (2019) to discretize the actions for linear classifiers by taking an absolute step that increases the action by the defined step as well as a relative step that will walk in steps from the value of the sample until it achieves a limit (the maximal of minimal found value in the dataset). Moreover, we used the strategy in Cui et al., 2015 (Cui et al, 2015) to find the actions in sets of decision trees by dividing the space search for each feature using the splits in the decision trees, in the space between the splits is allocated an action for that feature, enumerating all actions to the decision trees.

As previously mentioned, the procedure `SELECT_FEATURE` in Algorithm 1 is responsible for indicating the next feature to be explored. To do that, we resorted to variable importance from Random Forests (Breiman, 2001) that randomly permute the values of a feature and estimate the impact in the prediction to evaluate the feature importance. We select an unexplored feature on a given node with higher importance to find a counterfactual antecedent as fast as possible.

Another implementation detail involves the modification presented in this work to make mixed-integer models (Cui et al, 2015; Ustun et al, 2019) to generate Pareto-optimal

counterfactual antecedents. Given a *model* that finds a counterfactual explanation (Cui et al, 2015; Ustun et al, 2019), the procedure in Algorithm 2 consists of iteratively finding counterfactual antecedents and adding constraints (line 7 and 8) that counts (v) the number of objective components of action with a value larger than a previously found antecedent \mathbf{a}^* and enforces it to be lower than the number of objectives ($v \leq m - 1$). This procedure induces any subsequent optimization of the *model* not to have any new action dominated by previously found antecedents.

Algorithm 2 Pareto-optimal procedure for mixed-integer formulations.

Require: A mixed-integer model builder *build*, a sample \mathbf{x} , a objective function $\mathbf{c}(\bullet)$, a decision rule $r(\bullet)$, a threshold τ , and a number of allowed changes k .

```

1: procedure ENUMERATE_MIP(model)
2:    $\mathcal{A} = \{\}$ 
3:   while model.feasible do
4:      $\mathbf{a}^* = \text{model.optimize}()$ 
5:     if model.feasible then
6:        $\mathcal{A} = \mathcal{A} \cup \{\mathbf{a}^*\}$ 
7:        $v = |i : c_i(\mathbf{a}) \geq c_i(\mathbf{a}^i), \forall i \in \{1, \dots, m\}|$ 
8:       model.addConstr( $v \leq m - 1$ )
9:     end if
10:  end while
11:  return  $\mathcal{A}$ 
12: end procedure
13: model = build( $\mathbf{x}, \mathbf{c}, \tau, \tau$ )
14:  $\mathcal{A} = \text{ENUMERATE}(\text{model})$ 
15: return  $\mathcal{A}$ 

```
