# Calculating and Visualizing Counterfactual Feature Importance Values

Bjorge Meulemeester*a,b, Raphael Mazzine Barbosa De Oliveira*a, David Martensa

*aDepartment of Engineering Management, University of Antwerp, Prinsstraat 13, 2000 Antwerpen, BEL*
*bMax-Planck-Institut für Neurobiologie des Verhaltens – Caesar, Ludwig-Erhard-Allee 2, 53175 Bonn, DE*

## Abstract

Despite the success of complex machine learning algorithms, mostly justified by an outstanding performance in prediction tasks, their inherent opaque nature still represents a challenge to their responsible application. Counterfactual explanations surged as one potential solution to explain individual decision results. However, two major drawbacks directly impact their usability: (1) the isonomic view of feature changes, in which it is not possible to observe *how much* each modified feature influences the prediction, and (2) the lack of graphical resources to visualize the counterfactual explanation. We introduce Counterfactual Feature (change) Importance (CFI) values as a solution: a way of assigning an importance value to each feature change in a given counterfactual explanation. To calculate these values, we propose two potential CFI methods. One is simple, fast, and has a greedy nature. The other, coined Counter-Shapley, provides a way to calculate Shapley values between the factual-counterfactual pair. Using these importance values, we additionally introduce three chart types to visualize the counterfactual explanations: (a) the Greedy chart, which shows a greedy sequential path for prediction score increase up to predicted class change, (b) the CounterShapley chart, depicting its respective score in a simple and one-dimensional chart, and finally (c) the Constellation chart, which shows all possible combinations of feature changes, and their impact on the model's prediction score. For each of our proposed CFI methods and visualization schemes, we show how they can provide more information on counterfactual explanations. Finally, an open-source implementation is offered, compatible with any counterfactual explanation generator algorithm.

*Keywords:* Counterfactual explanations, Counterfactual feature importance method, explainable AI, XAI, visualization

## 1. Introduction

Artificial intelligence (AI) has become capable of capturing, learning, and predicting complex data structures. The performance of GPT-3 [1] models, for example, are so performant that a piece of software can mimic language to a level that is often indistinguishable from the human-written text. Deep learning models can learn a voice's intonation, tone, and pronunciation habits using just a couple of seconds of audio as training data [2, 3, 4]. Learning highly complex data structures, such as speech, generally requires an equally complex model to capture the nuances of the data accurately. While this increase in complexity opens up the possibility of improved performance, it comes with a price: its explainability is greatly reduced [5]. This lack of justification for decisions can decisively impact their application [6]. Current legislation, for example, requires explainability of models that impact high-stake decisions for people's lives [7]. Moreover, understanding why a model works can provide insights to guide improvements in itself [5] and possibly detect erratic patterns learned from data that can cause discrimination or unfair decisions [8, 9, 5]. Models that become inexplicable due to their inherent complex characteristics, such as a high number of parameters or complex feature associations, are referred to as "black box" models, since their inner workings are not clear enough for humans to comprehend.

---

*Authors contributed equally.

*Email addresses:* bjorgemeulemeester@mpinb.mpg.de (Bjorge Meulemeester),
Raphael.MazzineBarbosaDeOliveira@uantwerpen.be (Raphael Mazzine Barbosa De Oliveira*),
david.martens@uantwerpen.be (David Martens)

The field of eXplainable Artificial Intelligence (XAI) dedicates great efforts to shed light on how complex models reach decisions by adopting several methods [10, 9]. The wide variety of XAI methods can generally be categorized by the scope of their explanation, model compatibility, and output format [5, 10]. In terms of scope, local explanation methods aim at uncovering why a model predicts some outcome for a single data instance, while global explanations seek to provide an overall explanation of how the model works on the entire dataset. For model compatibility, we have methods that only work for a specific kind of model - since they often use particular mechanisms to generate explanations - and approaches that are independent of model type i.e. "model-agnostic". Finally, the output format of the explanations can be simple feature importance values or more sophisticated representations such as heatmaps for image explanations.

Popular global explanation approaches include, but are not limited to: using intrinsic model characteristics, such as impurity decrease analysis for tree-based models [11]; rule extraction strategies, such as Trepan [12] and ALBA [13], where a less complex but comprehensible surrogate model is built to mimic the black-box model, and its intrinsic explainable parts are used as a proxy to explain the black-box model; and depicting the marginal effect of features on the prediction score (partial dependence plots, PDP) [14]. As for local explanations, LIME [15] and SHAP [16] are popular feature importance (FI) methods, indicating the most important (subset of) features for a given instance's prediction score. LIME considers data in the vicinity of the instance to be explained and fits a locally linear relationship. The coefficients of this linear fit can be interpreted as the FI values. SHAP calculates approximate or exact Shapley values, compared to the average prediction and average feature value of the entire dataset. Finally, counterfactual explanations provides a set of changes in feature values that lead to a different model decision [17].

In terms of explanation output format, all listed explanation methods (except for counterfactual explanations) include ways to visualize the feature's importance graphically. LIME allows bar charts to visualize the weights of the locally linear decision boundary[1]. SHAP even allows various visual methods to compare the feature Shapley values, such as waterfall plots, force plots, the popular beeswarm and violin plots, and even image heatmaps for image recognition models[2]. And the random forest's mean decrease in impurity (MDI) allows a simple representation of the feature's importance by bar charts [18].

One naive visualization approach that is commonly used to introduce counterfactual explanations is to show the factual and counterfactual points graphically in their input space. However, this is only possible for 2 to 3 dimensions, a rare occurrence for real-life datasets. Moreover, just showing the value changes in Cartesian space does not add considerable value to the explanation, since features usually have vastly different and interdependent distributions. They cannot be directly compared. Therefore, the most practical and adopted solution to represent the output of a counterfactual explanation method is to provide a set of feature names and associated values to be changed. We can illustrate such an explanation with the example below, which considers a machine learning model that approves or rejects credit card applications:

---

**Explanation 1: Counterfactual explanation why a credit card application was rejected.**

---

If your age were 28 years old instead of 20, and
your salary were \$4,000 instead of \$1,200
then you would be approved instead of rejected.

---

The representation does not include any visualization, nor any way to assess which of the feature changes was more important than the other. This is in stark contrast to LIME or SHAP, who have easily interpretable charts that show the importance of each feature. Additionally, the unique objective of counterfactuals explanations is to look at feature *changes*, rather than just features, that lead to a *class change*. This does not allow for a simple translation of LIME or SHAP visualization. Hence, this leads to the following two questions: can we obtain importance values for each feature change in a counterfactual explanation, and can these be used to visualize the explanation?

---

[1]The visualization resources for LIME can be seen at https://github.com/marcotcr/lime

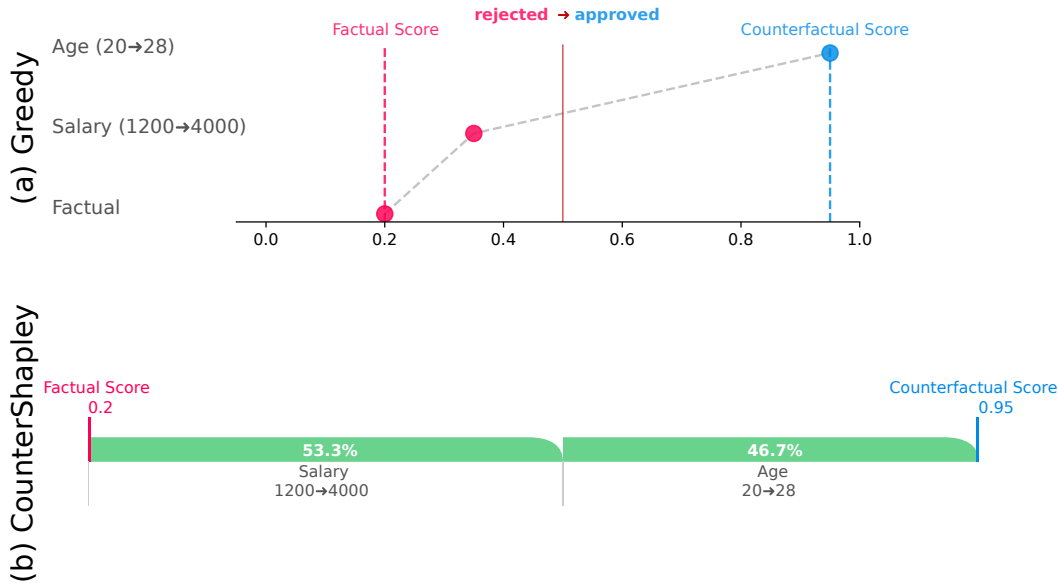[2]SHAP's charts can be found in their original repository: https://github.com/slundberg/shap

Figure 1: Counterfactual explanation visualizations for *explanation 1* using the (a) Greedy and (b) CounterShapley charts.

In this paper, we solve this problem as follows: we generate Counterfactual Feature change Importance (CFI) values, and propose two methods to do so[3]. Figure 1 already illustrates the goal and result of such CFI methods for the same counterfactual explanation we previously mentioned. The *Greedy CFI method and chart (a)* visualize the modification sequence that occurs when changing the feature with the greatest effect on the model prediction score, one at a time. Since each change adapts the feature, each feature change will have a different effect after each modification, thus making these values dependent on the order of the feature changes. The benefit of this approach is its efficiency and simplicity. In this example, we see that the prediction score for the data instance (the factual) is 0.2, which is below the 0.5 threshold, leading to a rejected decision from the model. Changing salary from 1,200 to 1,400 apparently leads to the highest increase in the predicted score, yet still leading to a rejection score. Additionally changing the age from 20 to 28 then leads to the counterfactual instance with a prediction score of 0.95 and approval decision. Note that changing the age only has an impact this high *after* the salary has been changed as well.

The *CounterShapley CFI method and chart* generalizes this concept, by looking at all possible subsets of feature changes, and assessing the average marginal value of each feature change. This naturally leads to Shapley values [19]. Section 5 will describe this in more detail, but the example already indicates that salary change has an importance value of 53.3%, while age is slightly less important, having an importance value of 46.7%. Finally, we will also introduce *Constellation charts*, that aim to visualize all these subsets of feature-value changes. Further demonstration of this concept will be provided in our empirical section.

## 2. Motivation

The research interest in counterfactual explanations is evident when we look at the increasing number of publications in this field [20]: ever since it has been introduced in 2011 [21], over 300 counterfactual generation methods have been proposed [17]. This interest is justified by several desirable properties that counterfactual explanations have: they are simple, model-agnostic, sparse, and decision-oriented [17, 22, 23]. The simplicity is not only related to the output explanation itself, but also refers to the familiarity that humans have to use counterfactual arguments to provide explanations [24]. This process is well-established both in philosophy and psychology [25, 24], showing that

---

[3]Note that we will interchangeably talk about feature changes or feature value changes, as well as their associated CFI values and importance scores. Each feature change has a value in feature space, and an associated CFI value as well.

such explanations come natural to end users. In theory, counterfactual explanations do not rely on any specific model since they consider the prediction model as a black box [20], making them model-agnostic. The explanations are typically also sparse, in the sense that only a small subset of all features occurs in the explanation. Not only does this make the explanation smaller and hence more comprehensible [9] (compared to an explanation using all features), but it also considerably reduces the complexity of the problem of assigning importance values to each feature change in the explanation. This sparsity property thereby removes the need to limit the search or visualization to some smaller subset of features, as is done in LIME and SHAP where the number of features is user-defined.

Finally, the decision-oriented nature is a key point of counterfactuals; they focus on modifications that lead to a prediction change [23]. This characteristic is fundamental for highly complex machine learning methods, since their nonlinearity can lead to nontrivial feature-value relationships that can potentially bring misleading assumptions. We can verify this last point in literature works [26, 23] that compare counterfactual explanations with feature importance-based methods (LIME and SHAP) and show that the most important features do not always lead to a change in classification decision. Hence, counterfactual explanations can be advantageous whenever one is more interested in explaining a predicted decision, rather than a prediction. This decision-oriented nature of counterfactuals aligns well with the transparency requirements in legislations, such as the European GDPR [7].

Despite these numerous benefits of counterfactual explanations, one major drawback, as introduced in the previous section, is the way these explanation results are presented. This issue is not only limited to the lack of visual approaches to plot charts, but also to the inability to assign meaningful values to feature changes. One might wonder: why not simply apply an existing feature importance method, such as SHAP, to the counterfactual explanation? The reason is multifaceted, as we will describe in detail in the next sections: (1) we focus on feature-value *changes* (e.g. age changing from 20 to 28), not just on feature values (e.g. age equals 20); (2) the dimensionality that we have to deal with is generally much smaller, yielding unique visualisation opportunities; (3) the sum of an instance's SHAP values sum up to the difference between the average prediction score and the instance prediction, not to the difference between counterfactual and factual scores; and (4) the reference point and the features considered in SHAP are different from those considered in counterfactual explanations. So, we aim to define Counterfactual Feature change Importance (CFI) methods that enhance the informative value of counterfactual explanations by creating feature change importance values and correspondent visualization resources that are tailored to the cause.

The significance of visualization is well known in the field [17], and it is further reinforced in a recent study in which consumers preferred feature importance methods over counterfactual explanations [27]. The authors of the study conjecture that this preference is likely due to the lack of visualization for counterfactual explanations. Likewise, merely providing what feature values should be changed in order to generate a counterfactual explanation does not show any distinction in their relative importance. This additional information can provide useful information on how the model prediction score works, ultimately leading to a better comprehension of the model's decision. This would allow users and practitioners to critically evaluate whether the model behaves as expected.

The graphical representation of counterfactuals traces back to its origins [28], where the greedy nature of the proposed counterfactual generation algorithm (SEDC) allowed the plotting of sequential changes until the counterfactual class is achieved. In more recent studies, the visualization problem was tackled by multiple studies. For example, GAMUT [29] implements a user interface that shows multiple statistics and methods to assess the relationship between the instance's features and the model's prediction. Here, counterfactual visualisation is limited to a simple model output change given certain modifications, without any assignment of how important each of those changes is. Google's What-If Tool [30] focuses on allowing the user to explore the model's behavior by probing feature changes that are dynamically presented in tables and charts. But in terms of reporting, the counterfactual feature changes do not offer much information about their importance: they are depicted by simple tables that highlight the value changes, and two-dimensional charts with features values or different model's prediction scores. DECE [31] also presents a user interface that includes various statistical analyses of features, and shows multiple counterfactual explanations according to sequential modification steps. But it still does not show the importance of each feature change to prediction scoring. Lastly, VICE [32] represents counterfactual changes in charts together with the data distribution, and it allows customization of which features to include in the explanation. Yet it again does not disclose the influence of each feature change. Although the previous approaches have advantages by allowing users to investigate interactively how counterfactuals work, and they include multiple statistical analyses related to the target explanation and its features, they introduce an extra layer of complexity that users must learn to operate. Furthermore, they all represent counterfactual changes by only showing what features were modified in tables or spatial shifts with two-dimensional

charts. This all emphasizes the current lack of a method to represent how each counterfactual explanation feature change impacts the model's prediction scoring.

## 3. Counterfactual Explanations

The methods presented here do not rely on any specific type of counterfactual generator. A CFI method simply assumes some counterfactual explanation has been provided with corresponding feature changes, for which importance values need to be assigned, irrespective of the algorithm that was used to generate the explanation. Therefore, this section introduces a general concept of counterfactual explanation, which covers any generation algorithm. The formal definition is then used to introduce the CFI methods' theoretical reasoning.

Counterfactual explanations [28, 33] assume four basic components: a dataset ($\mathbf{D}$); an instance to explain ($\mathbf{x}$); a prediction model ($\mathcal{M}$); and, in the case of a classification task, a threshold ($t$).

Consider an $N \times M$ dimension dataset, with a set of $N$ instances

$$\mathbf{D} = [\mathbf{d}^1, \mathbf{d}^2, ..., \mathbf{d}^N] \tag{1}$$

and $M$ features

$$\mathbf{d}^n = [d_1^n, d_2^n, ..., d_M^n] \, , \, n \in [1...N] \tag{2}$$

where $\mathbf{d}^n$ is called a datapoint or an instance: a vector of size $M$.

These instances are used by a prediction model $\mathcal{M}$ to assign a score $s$. Scoring an instance with a model will be denoted by

$$\mathcal{M}(\mathbf{d}^n) = s \tag{3}$$

If we consider a binary classification task where we set a decision threshold $t$, assigning instances to *class 0* if $s < t$ and to *class 1* if $s \geq t$, we can effectively categorize all instances into either class. We will only assume binary classification in this work for simplicity, but this can easily be generalized to multiclass classification with, e.g., a one-vs-one or one-vs-rest paradigm. Since CFI only assumes the existence of some counterfactual and its associated feature values, no matter the prediction task, this can also be generalized to regression tasks.

Given two different datapoints $\boldsymbol{a}$ and $\boldsymbol{b}$, we can define the difference between $\boldsymbol{b}$ and $\boldsymbol{a}$ as the set of indices where their feature values differ:

$$\delta_{a,b} = \{i \mid a_i \neq b_i\} \tag{4}$$

For notational ease, we define the scoring of the instance $\boldsymbol{a}$ with replaced feature values $\delta_{a,b}$ as being:

$$\mathcal{M}(\delta_{a,b}) = s \tag{5}$$

Note that the definition of $\delta$ is symmetric with respect to the two associated instances. By convention, this notation will denote the scoring of $\mathbf{a}$, adapted to have the values of $\mathbf{b}$ on the indices in $\delta$, where $\mathbf{a}$ and $\mathbf{b}$ are the left and right instance of the $\neq$ sign in Equation 4. This becomes important as soon as we consider subsets of $\delta_{a,b}$. In the case of factual and counterfactual instances (see below), this notation will always denote the factual instance, adapted to have feature values of the counterfactual on some indices. Not the other way around.

Using this notation, we can formalize a counterfactual instance $\mathbf{c}$, derived from a factual instance $\mathbf{x} \in \mathbf{D}$, having the following properties:

1. The factual and counterfactual's predicted classes are not the same:

$$(\mathcal{M}(\mathbf{x}) > t) \neq (\mathcal{M}(\mathbf{c}) > t) \tag{6}$$

2. The difference between the factual and counterfactual can be defined by a set of indices where their features differ, i.e., the counterfactual explanation itself:

$$\delta_{x,c} = \{i \mid x_i \neq c_i\} \tag{7}$$

3. The counterfactual evidence is irreducible: there's no subset of valid changes that changes the original class.

$$\nexists \, \delta'_{x,c} \subset \delta_{x,c} : (\mathcal{M}(\delta'_{x,c}) > t) = (\mathcal{M}(\delta_{x,c}) > t) \tag{8}$$

## 4. Assessing Feature Importance Values in Counterfactual Explanations

### 4.1. Greedy CFI method

The greedy CFI method is a simple yet useful approach. It will iteratively select the feature change that makes the largest contribution towards the counterfactual class prediction, as formalized in Algorithm 1. The resulting scores have some useful properties: the sum of the feature change scores is equal to the difference between the counterfactual and factual prediction ($\mathcal{M}(\mathbf{c}) - \mathcal{M}(\mathbf{x})$); it has a simple interpretation that can easily be presented in charts (further discussed in Section 5.1); and it is computationally efficient. Note that Algorithm 1 is not the most efficient implementation of a greedy algorithm, in favor of readability.

This method shares a large similarity with the SEDC [21] counterfactual generation approach since both have a similar objective. However, while SEDC is an algorithm designed to find a counterfactual explanation given a factual instance and model, the Greedy CFI method already assumes a counterfactual instance and assigns a score for each feature change. So Greedy CFI is a generic version of the initially proposed 'Score Evolution' chart of Martens and Provost.

---

**Algorithm 1:** Retrieving the vector with Greedy scores ($\varrho$) for a factual ($\mathbf{x}$) and counterfactual ($\mathbf{c}$) instance and prediction model ($\mathcal{M}$). This code is not optimized as vectorization and other approaches can accelerate the calculation and this is done in our package implementation.

**Input:** Prediction model ($\mathcal{M}$), factual ($\mathbf{x}$) and counterfactual ($\mathbf{c}$) instance
**Output:** Map ($\varrho$) of greedy feature importances for all counterfactual changes in $\delta$

```
 1  ϱ ← {};
 2  current_pred ← M(x);
 3  repeat |δ| times
 4  │   best_pred ← 0;
 5  │   best_feature ← −1;
 6  │   for i in δ do
 7  │   │   x_r ← x;
 8  │   │   x_r[i] ← c[i];
 9  │   │   pred_r ← M(x_r);
10  │   │   if pred_r > best_pred then
11  │   │   │   best_pred ← pred_r;
12  │   │   │   best_feature ← i;
13  │   │   end
14  │   end
15  │   ϱ[best_feature] ← best_pred − current_pred;
    │   // update for next iteration
16  │   x[best_feature] ← c[best_feature];
17  │   current_pred ← best_pred;
18  │   δ ← δ \ best_feature;
19  end
20  return ϱ
```

---

### 4.2. CounterShapley CFI Method

Given the non-linearity of the complex models, different sequences may give different importance values for the same feature change. Therefore, we introduce the CounterShapley CFI method next, which presents a way to calculate importance values independent of a specific modification path. As this relies on Shapley values, we will first provide a brief introduction to this game theoretical concept, followed by a short discussion of the previous use of Shapley in machine learning.

### 4.2.1. Shapley Values

Shapley values quantify the contribution of different variables to a final numeric result ($\varphi$). This method is derived from game theory, where each variable is considered a player $p$ whose combinations can form coalitions $V$ that contribute to a certain outcome according to a gain function $\mathcal{H}(V)$ [19]. Therefore, we can calculate the contribution ($\varphi_i$) of a certain player $p_i$ in a coalition formed by a set of $|Z|$ players by considering its contribution in all possible coalitions:

$$\varphi_i(\mathcal{H}) = \sum_{V \subseteq Z \setminus \{i\}} \frac{|V|!(|Z| - |V| - 1)!}{|Z|!} \left( \mathcal{H}(V \cup \{i\}) - \mathcal{H}(V) \right) \tag{9}$$

The total number of terms Equation 12 will have to add together to calculate the Shapley value for a single player equals:

$$N_{coalitions,i} = \sum_{i=0}^{|Z|-1} \binom{|Z| - 1}{i} = 2^{|Z|-1} - 1 \tag{10}$$

However, to have a complete calculation of all players' contributions, one must consider every possible subset of players, and not just the ones excluding $i$. This then yields the following number of iterations:

$$N_{coalitions} = \sum_{i=0}^{|Z|} \binom{|Z|}{i} = 2^{|Z|} - 1 \tag{11}$$

These Shapley values have several desirable properties [34]:

1. Symmtery: two players' contributions are the same, only if they contribute the same to every coalition that does not contain either player.
2. Efficiency: the sum of all players' contributions is equal to the total worth of the game or the grand coalition.
3. Dummy: a player with zero contribution does not affect the outcome contribution of a coalition, independent of the coalition in which it occurs.
4. Linearity: if one takes the sum of two gain functions, then the Shapley value of any player as scored by this sum of gain functions is the same as when one calculates the Shapley value for each gain function, and sums the results. Similarly, when scaling up the gain function by a factor of $a$, this will simply yield $a$ times the Shapley value as calculated by the original gain function.

### 4.2.2. Shapley Values in Machine Learning - SHAP

The assignment of feature importance values for individual instances traces back to 2010 with Strumbelj and Kononenko work [35] and has gained popularity with LIME in 2016 [15, 36]. Since then, multiple algorithms have iterated upon this explanation approach to improve its stability [37, 38], where the same instance could yield different importance weights. Among the strategies to circumvent this problem, Shapley value properties demonstrate their potential for solving the stability issue demonstrated by LIME. Given the model's non-linear behavior, features are allowed to contribute differently, depending on which coalition of features they are added to. Moreover, the dummy property assures that features that do not contribute to any coalition have no importance either. Also, the efficiency property allows assigning proportional importance values relative to a certain baseline. What exactly this baseline is, depends on what we consider as the model prediction of an "empty set of players".

SHAP [39] adapts the Shapley values calculation to the machine learning context. As a main novelty, they developed a strategy of unifying an adapted Shapley value calculation and feature importance methods. The features are considered to be the players, and the scores obtained by the model are meant to describe, similarly to Shapley's gain function, how each feature contributes to the final prediction score. It has as a baseline the average prediction score. The calculation is performed through a model agnostic approaches (such as KernelSHAP), or model-specific approaches (such as TreeSHAP) to generate more reliable feature importance values. For KernelSHAP, Shapley values are estimated by performing a least squares regression using Lasso normalization. On the other hand, TreeSHAP [40] is data independent: it does not directly need external data since it uses the model's tree nodes information to calculate the features' marginal contributions.

Despite the success of these approximations, they have limitations that substantially impact their application in explaining machine learning models. For TreeSHAP, the obvious implication is its narrow compatibility to tree-based models only. Although KernelSHAP does not have this constraint, its calculation can be computationally expensive for a large number of features and training points: it can take up to $L \cdot 2^M$ calculations, where $L$ is the sample number of (training) instances, and $M$ is the number of total features. This creates a trade-off between stability (large sample) and computational efficiency (small sample). Moreover, even though SHAP yields more stable results over different runs when compared to LIME, its feature importance values can still change depending on the sampling since both the baseline score and the coalitions rely on the randomly selected data. Finally, the coalitions considered by SHAP can be unlikely or even include impossible data points. This is due to the fact that SHAP adapts feature values in order to compare instances to a baseline case. For instance, if we consider one-hot encoded (OHE) features, some SHAP coalitions can consider multiple active features from the same OHE (e.g., concurrently being married and single). This issue has broader consequences with KernelSHAP, since it uses multiple data instances that are not necessarily in the same local region as the instance of interest. For this reason, it can end up giving deceptive importance weights to features [41] for a local context.

Finally, SHAP is unable to assign importance scores to counterfactual explanations. It is not compatible with comparing pairs of instances, nor feature *changes*, because it considers the average prediction as a baseline. Despite that, SHAP could still be useful given its variety in plotting resources, filling the gap that counterfactual explanations generally fail to deliver engaging visual resources.

### 4.2.3. CounterShapley definition

As discussed earlier, counterfactual explanations consist of a factual instance to be explained, and a counterfactual point with a different class. Although this type of explanation has benefits, discussed in Section 2, such as simplicity and sparsity, this representation alone does not precisely show each feature change's influence on this classification change. Consequently, we cannot tell which feature changes contribute more to the score change than others. This directly impacts the quality of the model's explanation since we can only describe the class change and feature modifications. The former is trivially inherent to a counterfactual explanation, while the latter does not directly reflect the impact of each individual feature modification on the model score.

Augmenting the informative value of counterfactual explanations would leverage not just one but multiple explanation methods since counterfactual generation algorithms can have different objectives: SEDC [28], for example, tries to improve sparsity and CADEX [42] minimizes euclidean distance in feature space. Other algorithms target both metrics at the same time, such as DiCE [43] and ALIBI [44].

As described in Section 4.1, we can assign importance values to the features changes in counterfactual explanations by using their impact over the prediction score. However, the Greedy CFI method considers a specific modification pattern that may not be useful when overall relative importance is desired, especially in highly non-linear models. To solve that issue, we can adopt a strategy similar to SHAP and use the Shapley value approach (introduced in Section 4.2.1) to calculate counterfactual explanation feature importance values.

Let us consider the counterfactual evidence set $\delta_{x,c}$ and associated feature values $x_i$ and $c_i$, such that the factual feature values have to be changed from $x_i$ to $c_i$ in order to go from a factual $\mathbf{x}$ to a counterfactual $\mathbf{c}$. If we define the model $\mathcal{M}$ as the scoring function, $\delta_{x,c} = \delta$ and $|\delta| = K$ for notational ease, we can use Equation 9 to define a Shapley value $\varphi_i$ for the features-to-be-changed $\delta$ as:

$$\varphi_i(\mathcal{M}) = \sum_{V \subseteq \delta \setminus \{i\}} \frac{|V|! \, (K - |V| - 1)!}{K!} \left( \mathcal{M} \left( V \cup \{i\} \right) - \mathcal{M} \left( V \right) \right) \tag{12}$$

where $K$ is the number of features, and the sum extends over all sets $V$ that are a subset of $\delta \setminus \{i\}$. This weighted sum is simply the average marginal contribution of $x_i$ to the class switch, as predicted by the model $\mathcal{M}$. In more detail, from right to left, this equation can be understood by:

1. Evaluating the model outcomes of some factual instance $\mathbf{x}$, but adapted to have all features changes in $V$, as well as the modification of the feature of interest ($x_i = c_i$).
2. Evaluating the same model outcome for the same adapted instance, but without changing the feature of interest.
3. Considering the difference in these two model outcomes

4. Weighing this difference, depending on how often this particular coalition of feature changes may occur.
5. Considering all coalitions of all sizes that do not include a particular feature change, defined by index $i$.

If we want to know how much some feature change contributed to the change in model prediction score, we can change these features from the factual to the counterfactual value one by one, each time evaluating the difference in the model outcome. Given the non-linearity of complex models, the order in which these features are changed may have an influence on the difference in the model outcome. For example, changing some feature value of an instance if two previous values have already been changed can yield a different outcome than changing the exact same feature if no previous changes have been made yet. For this reason, we must consider every possible order in which some feature can be changed, justifying the application of Shapley values. This is a major difference of CounterShapley compared to the Greedy approach, as the latter considers a fixed sequence of modifications.

### 4.2.4. Difference between Shapley, SHAP and CounterShapley

Similarly to SHAP, the CounterShapley approach inherits all the desirable properties of Shapley values. It is fundamental to highlight that, despite the similarity with SHAP, CounterShapley cannot be directly compared with it. They have divergent characteristics and objectives, as shown in Table 1, and explained in detail below.

*Baseline.* The baseline in Shapley is originally just the value of the empty coalition where no players are involved. In contrast, for KernelSHAP, the baseline consists of the average model's prediction over all available data points. For CounterShapley, the baseline is the model prediction of the factual instance that is to be explained. This is a completely different baseline, which is simpler to understand and trivial to calculate.

*Coalitions.* The information needed to calculate the importance values is also different for the three methods. Despite their similar goal of calculating the marginal contribution of different coalitions, they require different coalitions. For Shapley values, the coalitions are formed by all possible combinations of players. For KernelSHAP, the coalitions are formed by all possible feature value combinations of the sampled data points; the contribution values are then obtained by calculating the average prediction score of the instances with the same replaced features for the different points being considered. CounterShapley's coalitions are all the possible combinations of feature value differences between the factual and counterfactual instance, i.e. $\delta$.

*Complexity.* With these considerations on what each method calculates, we can estimate their complexity. All have an exponential complexity, but we observe that CounterShapley does not necessarily consider all features (or players). When a counterfactual algorithm is optimized to minimize the number of features that are to be changed in order to create a counterfactual, this can greatly reduce the complexity of assessing the importance of these features. In Appendix A, we illustrate how the calculation of CounterShapley works for one feature change, which also explicitly illustrates its complexity. Apart from the fact that KernelSHAP requires more features than CounterShapley, it also has a higher complexity because it requires more data instances. CounterShapley only needs 2, while KernelSHAP needs a decent number before it can assess the feature importances with some degree of stability.

*Efficiency.* Additionally, their efficiency property is also notably different. Shapley values have their sum equal to the value of the scoring function when all players are being considered. When we sum the KernelSHAP values we have the difference between the average prediction of the sample and the explained instance prediction score. For CounterShapley, the sum of their values is equal to the difference between the counterfactual and factual prediction scores.

*Goal and applicability.* These different baselines, coalitions, calculations, and sums highlight their distinct interpretations and applicability. CounterShapley considerably differs from SHAP since it aims to explain the feature's effect on the factual to the counterfactual change, while SHAP compares the feature of interest with the average prediction. SHAP is more suitable if one wants to understand how every single feature influences the prediction scoring, and CounterShapley is advisable for assigning importance weights to counterfactual explanations. Additionally, CounterShapley does not create the explanation, but rather augments the informative value of pre-existing counterfactual explanations. Therefore, if the speed of generating an explanation is a relevant matter, one must also include the counterfactual generation time, which can vary greatly depending on the counterfactual generator being applied [22].

| **Shapley** | |
|---|---|
| **Objective** | Given a gain function $\mathcal{H}$ and $Z$ players that can form $2^Z$ coalitions (V), calculate each player's contribution to all possible coalitions. |
| **Base value** | Output value of empty coalition $\mathcal{H}(\emptyset)$. |
| **Information needed** | For any coalition $V$, its output value with $(\mathcal{H}(V \cup \{i\}))$ and without $(\mathcal{H}(V))$ player's contribution. |
| **Calculations** | $\mathcal{O}(2^Z)$ (all coalitions). |
| **Values' sum** | Gain function's value when all players are considered $(\mathcal{H}(all\ players))$. |
| **Interpretation** | Marginal contribution of any player over the final result considering all players. |
| **KernelSHAP** | |
| **Objective** | For a prediction model $\mathcal{M}$ and a (sampled) dataset $\mathbf{D}$ with $N$ rows and $M$ columns, calculate each feature importance of a data point $\mathbf{x}_e$. |
| **Base value** | Average model's prediction over all data points $\left(\sum_{i=0}^{N} \frac{\mathcal{M}(\mathbf{x}_i)}{N}\right)$. |
| **Information needed** | Model's prediction value for all combinations between target $(\mathbf{x}_e)$ and data points $(\forall \mathbf{x}_i \in \mathbf{D}_{N,M})$. |
| **Calculations** | $\mathcal{O}(N2^M)$. |
| **Values' sum** | Difference between the average prediction over the sampled data points and the instance point to be explained $\left(\mathcal{M}(\mathbf{x}_e) - \sum_{i=0}^{N} \frac{\mathcal{M}(\mathbf{x}_i)}{N}\right)$. |
| **Interpretation** | Average contribution to the model's prediction of each feature in relation to sampled data points. |
| **CounterShapley** | |
| **Objective** | For a factual $(\mathbf{x})$ and counterfactual $(\mathbf{c})$ instance with $K$ modified features out of $M$ features in total, calculate the contribution of each feature change according to the model $\mathcal{M}$. |
| **Base value** | Factual instance's prediction score $(\mathcal{M}(\mathbf{x}))$. |
| **Information needed** | Model's prediction value for all possible combinations of feature changes $(\delta)$ in the factual instance $(\mathbf{x})$. |
| **Calculations** | $\mathcal{O}(2^K)\ (K \leq M)$. |
| **Values' sum** | Difference between counterfactual and factual predictions $(\mathcal{M}(\mathbf{c}) - \mathcal{M}(\mathbf{x}))$. |
| **Interpretation** | Contribution of all counterfactual features and their respective changes compared to the factual prediction score. |

Table 1: Comparison table of Shapley values and two model agnostic approaches for explaining machine learning predictions, highlighting the differences between Shapley, KernelSHAP and CounterShapley.

*4.2.5. Implementation of CounterShapley*

We continue by describing how CounterShapley values can be calculated, highlighting the methods' theoretical differences and advantages, and show how it lends itself very nicely to visualisation purposes. Algorithm 2 presents a method of generating a map between any-sized coalitions of feature changes and their corresponding model prediction outcomes. This map can be used in conjunction with Algorithm 3 to calculate the Shapley value of the change in some feature $i$. Since Algorithm 2 has already concerned itself with mapping all possible coalitions to their corresponding model prediction outcome, Algorithm 3 does not need to perform any further calls to the model's prediction method and can simply fetch them from the result of Algorithm 2. This is an important aspect of not just this specific implementation, but any implementation, as the algorithmic complexity is defined by the number of calls to this model prediction method. The extension to calculate all Shapley values for all features $\delta$ is now trivial, and can be done by using Algorithm 3 for every change in feature $i$ in $\delta$.

---

**Algorithm 2:** Generating marginal feature contributions for varying coalition sizes

    **Input** : factual ($\mathbf{x}$), counterfactual ($\mathbf{c}$), model ($\mathcal{M}$)
    **Output:** $G = \{coalition : score\}$
1  $G \leftarrow \{\}$;
2  **for** $L \leftarrow 0$ **to** $K$ **do**
3      $subsets \leftarrow$ `getCombinations`($V, L$);
4      **for** $V$ **in** $subsets$ **do**
          `// create adapted instance to score`
          `// start with the factual instance`
5          $\mathbf{c}^* \leftarrow \mathbf{x}$;
          `// change out feature values`
6          $\mathbf{c}^*[V] \leftarrow \mathbf{c}[V]$;
          `// save score of adapted instance`
7          $G[V] \leftarrow \mathcal{M}(\mathbf{c}^*)$;
8      **end**
9  **end**
10  **return** $G$

---

## 5. Visualizing Counterfactual Explanations

In this section, we present three different chart types that increase the informative value of counterfactual explanations. Two of them, Greedy and CounterShapley charts, use the CFI methods described in the previous section. The third chart type, the Constellation chart, uses the same calculations as CounterShapley since it analyzes all possible coalitions with the counterfactual feature modifications.

*5.1. Greedy Chart*

The greedy chart aims to illustrate the sequential nature of the changes made by the greedy CFI method. We argue that this method brings value to graphical representations for two main reasons. First, as said previously, the calculation of scores has lower computational costs if compared to CounterShapley, from an exponential $O(2^K)$ to a lower complexity $O(K^2)$ behavior. Second, this approach gives us a fixed path that can be depicted in charts, governed by a greedy selection of feature changes. The latter can be useful if the receiver of the counterfactual explanation wants to increase, as much as possible in each modification step, the prediction score. Moreover, this greedy CFI method gives the same importance scores as CounterShapley if the model has a linear association between features. Thus, for lower-complexity models, it can give a sufficiently good approximation of the CFI values obtained with the CounterShapley method. To exemplify its use, let us consider the following counterfactual explanation:

---

**Algorithm 3:** Calculating the Shapley value of some feature change with index $i$

---

**Input:** Feature $i$, map of all coalitions $V$ and corresponding scores $G$
**Output:** CounterShapley value $\varphi_i$ of feature's $i$ respective change

1  $\varphi_i \leftarrow 0$;
2  **for** $L \leftarrow 1$ **to** $K$ **do**
3     $C \leftarrow \text{getCoalitionsOfLength}(L, G)$;
4     **for** $V$ **in** $C$ **do**
         `// construct adapted instances`
         `// one that has feature i`
5        $with\_i \leftarrow \mathbf{x}$;
6        $with\_i[V] \leftarrow \mathbf{c}[V]$;
         `// and one that does not`
7        $no\_i \leftarrow \mathbf{x}$;
8        $no\_i[V \setminus \{i\}] \leftarrow \mathbf{c}[V \setminus \{i\}]$;
         `// marginal contribution`
9        $diff \leftarrow G[with\_i] - G[no\_i]$;
         `// add normalised marginal contribution`
10       $\varphi_i \leftarrow \varphi_i + diff \cdot L! \cdot (K - L - 1)!/K!$;
11    **end**
12 **end**
13 **return** $\varphi_i$

---

> **Explanation 2: Explanation for a loan application decision**
>
> ---
>
> If your age were 30 instead of 20, salary were \$2,200 instead of \$1,500, sex were F instead of M, and state were CA instead of NY, the application would be classified as approved instead of rejected.

Figure 2 depicts how *Explanation 2* would be displayed using the Greedy chart. The first point, at the bottom, represents the prediction score of the factual point, which has the features *Sex*, *City*, *Salary*, and *Age* equal to "M", "NY", $1,500$, and $20$ respectively. The next point (with a prediction score of about $0.35$) shows the change in prediction when the feature *Sex* is changed from "M" to "F". This modification (if compared to the three others) is the one that leads the highest increase towards the counterfactual class. The next points follow the same pattern, greedily, showing the best feature change to modify the classification result (from rejected to approved). The chart then has a sequence of changes starting from the factual point (bottom left) and gradually selects (from the bottom to the top) the changes that most positively contribute to the counterfactual class, which is depicted on the right side of the red threshold line. This simple chart allows for the visualization of sequential changes and their impact on the prediction score.

It's fundamental to emphasize that this chart represents scores influences in a specific setting: a greedy search. This may not reflect the CounterShapley values in non-linear models. Note for example how changing the City from "NY" to "CA" has a larger impact on the model score than changing Sex from "M" to "F", but is not the first change in line. This is because this change would not be as influential if the "Sex" wasn't changed first. Likewise for changing "Salary": without changing "Sex" and "City" first, the impact of this feature change would not be as big as it is. Such is the nonlinear nature of the model. Therefore, practitioners must be aware of whether this chart meets their expectations in explaining a decision, and possibly avoid it if they want to provide a precise estimation of feature change impact over the prediction score of complex models.
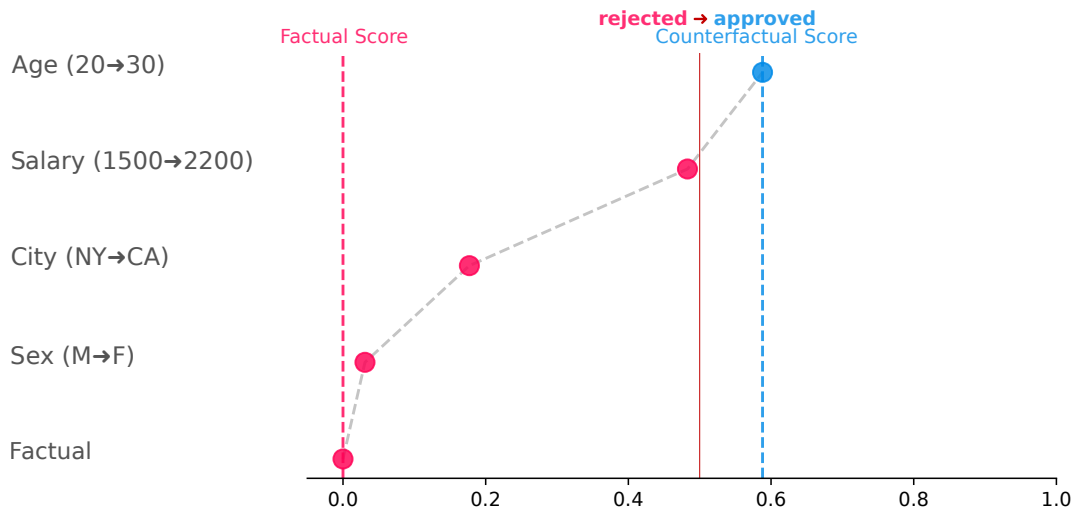
Figure 2: Counterfactual explanation visualization for *Explanation 2* using the Greedy chart.

## 5.2. CounterShapley Chart

The importance scores of counterfactual feature changes allow us to create a new visualization, here referred to as a CounterShapley chart. Our main objective with this chart is to represent counterfactual explanations in a simple but informative way which allows easy user understanding of what features lead to a class change and how important each of those features is.

Figure 3 shows the CounterShapley chart for the counterfactual *Explanation 2*, summarizing diverse information in a single image. The bars represent the CounterShapley values for each feature, and inside them we indicate their contribution in percentage to the total CounterShapley value. Below the bars, we have the feature names and their respective value changes. In addition, the chart shows the original instance, counterfactual instance, and the model prediction score after performing all feature changes. Although all this information could be described in a textual counterfactual, as mentioned earlier, the graphical representation is beneficial [45] to understanding and interpretation.

It's interesting to note how the information conveyed here differs from Figure 2. The Greedy chart tells us that "Sex" is the first best change one can make, despite the fact it does not increase the score by a lot initially. The CounterShapley chart tells us that, on average "Sex" has the biggest effect on the prediction score when averaged over every possible combination of features. The Greedy chart shows that "Salary" has the biggest impact, but only if "Sex" and "City" were changed beforehand. The CounterShapley chart tells us that, when considering every possible combination of feature changes, "Salary" only accounts for 23.2% of the overall change in prediction score. Thus, the large impact of "Salary" in the Greedy chart is only because the previous feature changes allowed it to make as big of an impact. Indeed, the CounterShapley value of "Sex" is large (40.6%), indicating that "Salary" impacts the prediction score a lot on the Greedy chart, but only because "Sex" is included in the set of feature changes. "Sex" does not have a large contribution on its own, but it allows other features to make a big impact.
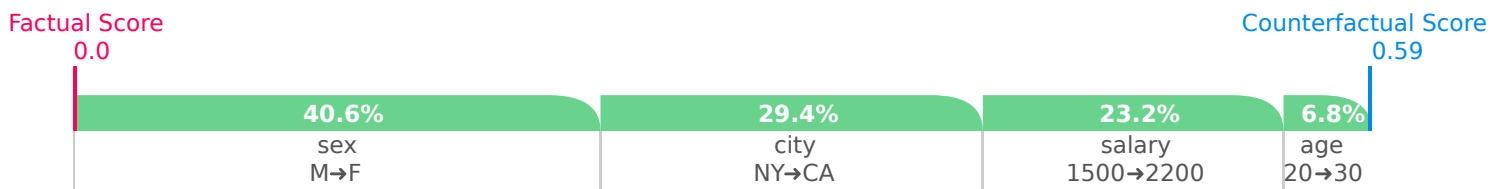


Figure 3: Counterfactual explanation visualization for *Explanation 2* using CounterShapley chart. It highlights the factual and counterfactual scores and, in the middle, the contribution to this prediction modification for each feature change in the counterfactual explanation.

## 5.3. Constellation Chart

Counterfactual explanations tend to have sparse feature changes, as we discussed in previous sections. The Constellation chart makes good use of this characteristic to display how each feature change combination influences the counterfactual prediction score. This chart has the same complexity as calculating the CounterShapley values, since it iterates over all feature combinations ($2^K$). Contrarily to both previous charts, it also has a more complex representation.

However, this type of chart can be useful to practitioners to have a complete view of how each counterfactual feature combination affects prediction scoring, allowing an extensive analysis of feature changes' influences. Moreover, it can also serve as a debugging tool for counterfactual explanations, since it is possible to verify if any subset of changes crosses the decision threshold, essentially finding a subset of the counterfactual that's also a valid counterfactual (see property 3 of the definition of counterfactuals in Section 3).

It has a similar application in describing importance values to counterfactual changes as the CounterShapley chart, but instead of summarizing them in single values, it detangles all possible feature's change association impacts over the model's prediction score. In Figure 4 we take the same *Explanation 2* counterfactual to generate a Constellation chart.
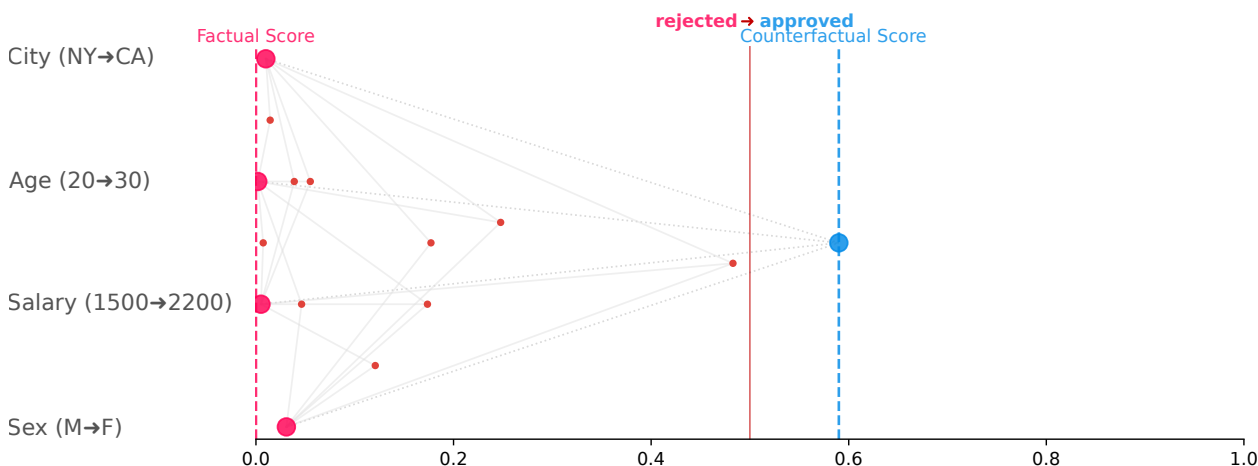


Figure 4: Counterfactual explanation visualization for *Explanation 2* using the Constallation chart.

The Constellation chart representation in Figure 4 shows how each single feature change ("Age", "Salary", "City", and "Sex") and their combinations affect the prediction score. On the left, we have the four counterfactual changes listed. The first pink dashed line (on the left) shows the factual score according to the $x$-axis. Next to the right, the four larger pink dots represent a single feature change. For example, the bottom large pink dot indicates the prediction score (according to the $x$-axis) when the feature "Sex" is changed from "M" to "F". Similarly, the other large pink dots correspond to the feature changes at the same height. The smaller pink dots represent the combinations of different feature changes. For instance, the first small pink dot (from the bottom to the top) indicates the prediction score when both features "Sex" and "Salary" are modified (this can also be deducted from the lines connecting to the larger dots). Notice that the height of the combination dots is equal to the average height of the features combined, but this height has no explicit meaning other than visualisation purposes. The rightmost point represents the prediction score when all feature changes are applied, as indicated by the lines connecting to all large dots on the left.

While this view looks more complex than the other previous methods, it also targets a more specialized public that is interested in understanding the underlying meaning of CounterShapley values and features' effects. For example, this chart already shows how the exclusion of "Age" still yields an instance that's very close to the decision threshold. Note for example how: "Sex" is indeed the first best feature one can change (lower left big pink dot), conform with the Greedy chart; and, on average, most small pink dots that are linked with "Sex" have indeed higher prediction scores than those not linked with "Sex", explaining its large CounterShapley value as seen in the CounterShapley chart. Finally, all small pink dots that are connected to both "City" and "Sex" yield the highest values. This is also congruent with the fact that these have the biggest CounterShapley values, as seen in the CounterShapley chart.

14

Moreover, this chart can help in other technical tasks, such as model improvement, since the multiple feature associations can reveal unexpected effects that may require further investigation. This will be shown in the next section.

## 6. Empirical Experiments

In the previous sections, we have introduced two main algorithms to assess the feature importance values of counterfactual explanations: CounterShapley and a Greedy approach. These counterfactual feature importance values now make it possible to directly compare counterfactual explanations with other feature importance methods, such as LIME [15] and SHAP [39]; but also to white-box models, such as the weights of a logistic regression.

In the following sections, we generate various counterfactual explanations using NICE [46], but the results work equally well for other counterfactual methods such as DiCE [43], ALIBIC [44], or SEDC [23]. The CFI values are compared to those of SHAP and LIME for various test cases. We show different experiments with multiple data, models, and methods that characterize how our CFI methods work, highlighting the differences with LIME and SHAP. In the first two sections, we use artificial data to demonstrate the characteristics of our methods, while the other three subsequent sections use real-world data, showing its applicability. We also explain cases in which the plots proposed here embed critical value for enhanced explanations and "debugging" counterfactual explanations. All experiments presented here are fully reproducible and can be found in the "Experiments" branch of our GitHub repository (see below).

Note that these comparisons are not quality assessments of these methods, as there is no consensus on what a "correct" explanation should be [47]. However, empirical testing is useful in providing a better hint at what each explanation approach tries to achieve.

Our work uses Python to make all calculations regarding the counterfactual scores, and the Matplotlib package [48] to create the graphical elements. Our implementation is open-source and independent from a specific algorithm, which makes it possible to integrate it with any counterfactual generation algorithm. The repository can be found on GitHub https://github.com/ADMAntwerp/CounterPlots. All empirical analyses use Python and open-source packages such as Scikit-Learn [49], Pandas [50], and Seaborn [51].

### 6.1. Simple Model Experiments

We use three simple, explainable (white-box) machine learning models (logistic regression, decision tree, and K-nearest neighbors) together with artificially generated data, and explain the predictions with three importance attribution methods: LIME, SHAP, and CounterShapley. Since we need a counterfactual generator for the CounterShapley calculation, we opt for NICE [46] to create the counterfactuals. We can then infer if the results given by explanation methods follow an expected pattern or not, by using the intrinsic explainability of these simpler models: feature weights of the logistic regression, decision nodes for the decision tree, and the overall data distribution for KNN.

For the logistic regression model, we explain 1,000 instances from a 5-feature synthetically generated dataset with two classes and a fixed random state (42) created with Scikit-Learn. Table 2 shows the logistic regression weights for each feature. Since individual explanations may not provide enough information to compare each method, we analyze the CFIs by aggregating the explanations that contain non-zero values (Table 3) and by plotting them in histogram charts (Figure 5).

Table 2: Logistic regression weights for the dataset's experiment features.

| Features | 1st | 2nd | 3rd | 4th | 5th |
|---|---|---|---|---|---|
| LR Weight | -0.59 | 0.67 | 0.78 | -0.34 | -0.59 |

The first main observation we can take from Figure 5 is the fact that the CounterShapley CFI method generally yields only positive importance scores, while the other methods include negative scores. Negative scores simply represent that an increase in these features leads to a probability decrease that an instance belongs to class 1. It highlights the substantially divergent interpretation from the CounterShapley values, which measure the impact of the feature changes on the class probability, knowing that the sum of all changes should be positive. A negative

Table 3: Share, in percentage, of the explanations with non-zero importance score for the logistic regression classification model.

| Features | 1st | 2nd | 3rd | 4th | 5th |
|---|---|---|---|---|---|
| CounterShapley | 58.29 | 53.30 | 75.04 | 25.49 | 53.83 |
| SHAP | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 |
| LIME | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 |

feature value here would mean that, not only is there at least one coalition including this feature that *decreases* the probability of a class switch, but these also outweigh the overall effect of the remaining coalitions that *increase* the class probability. Either there are more coalitions that have a negative effect on the class probability; or the negative effects are, on average, stronger than the positive effects. As long as the coalition including all feature changes (i.e. the counterfactual itself) yields a class flip, such feature changes can occur in counterfactual explanations. Cases like these are rare, but not impossible.

Second, although the average importance scores for SHAP and LIME are similar, their distribution is considerably divergent. SHAP appears to favor small values, as seen by the peak around 0, while LIME has an undefined distribution, generally assigning high-importance values to features,

Finally, let's compare the logistic regression weights (in Table 2) to the share of explanations that each feature is included (with an importance score different than zero). We see that the feature with the highest weight (3rd) is the most frequent in CounterShapley values; on the other hand, the other methods always assign scores for all features. While, again, this does not represent an advantage, it shows that CounterShapley is sparser given that it is based on counterfactual explanations.

For the decision tree, we made a very simple model with a depth of 2 and a total of 3 decision nodes using the same artificially generated data that we used for the logistic regression experiments. Figure 6 shows a representation of the model, where we can observe that only 3 (out of the 5 features) are present in the model: the first, second, and fifth. Figure 7 then shows the distribution of the importance scores obtained by the feature scoring methods.

Table 4: Share, in percentage, of explanations instances in which importance scores are different from zero for the decision tree model.

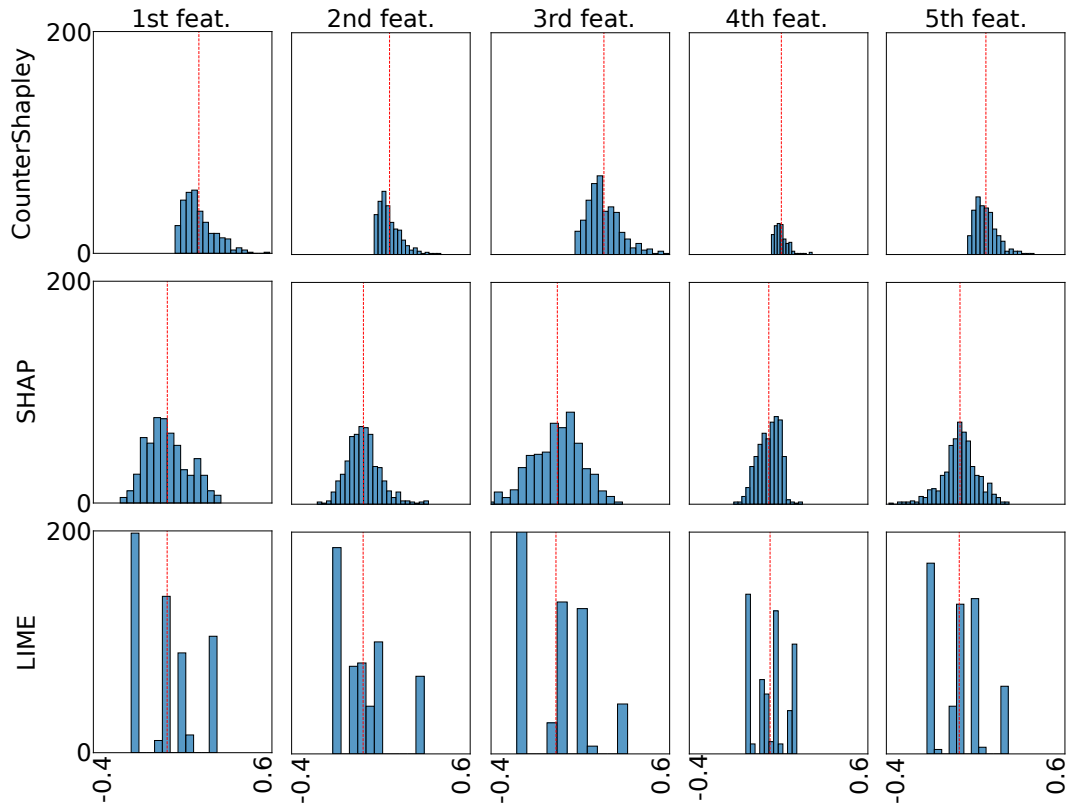| Features | 1st | 2nd | 3rd | 4th | 5th |
|---|---|---|---|---|---|
| CounterShapley | 44.18 | 20.00 | 0.0 | 0.0 | 44.18 |
| SHAP | 100.0 | 100.0 | 0.0 | 0.0 | 100.0 |
| LIME | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 |

Figure 5: Histogram with importance scores of all points that were classified as 0 in the logistic regression classification model, each column represents one feature while each row represents the CounterShapley, SHAP, and LIME distributions, respectively. The red dotted line indicates the distribution's average.

Again, we see a similar trend regarding the negative values assigned to importance scores for SHAP and LIME, which, similarly, means that increasing their respective features decreases the probability of class 1. Moreover, including the analysis of Table 4, we see that the 2nd and 3rd features have 0 importance for CounterShapley and SHAP, which is understandable (and expected) since these two features are not present in the decision tree nodes. Although LIME assigns importance scores to the 2nd and 3rd features, they are relatively low values compared to the other scores. Regarding the absolute importance values, it is clear on CounterShapley gives high importance to the 5th feature since it is the first splitting node and for the 1st feature since it splits more instances than the 2nd feature.

Finally, the last explainable model was a KNN classifier fitted in a 2-features dataset depicted in Figure 8. This dataset is also artificially created with Scikit-Learn, with two informative features and classes, a fixed random state (equal to 42) to guarantee reproducibility, and all other parameters as default. The results of the average importance scores are shown in Figure 9.

Table 5: Share of importance scores different from 0 of the KNN classifier instances using three different methods.

| Features | 1st (x) | 2nd (y) |
|---|---|---|
| CounterShapley | 62.26 | 81.61 |
| SHAP | 100.0 | 100.0 |
| LIME | 100.0 | 100.0 |

All methods assign a higher (absolute) importance to the 2nd ($y$-axis) feature. This is understandable since, by analyzing the data distribution in Figure 8, we can observe that the $y$-axis makes better segregation between
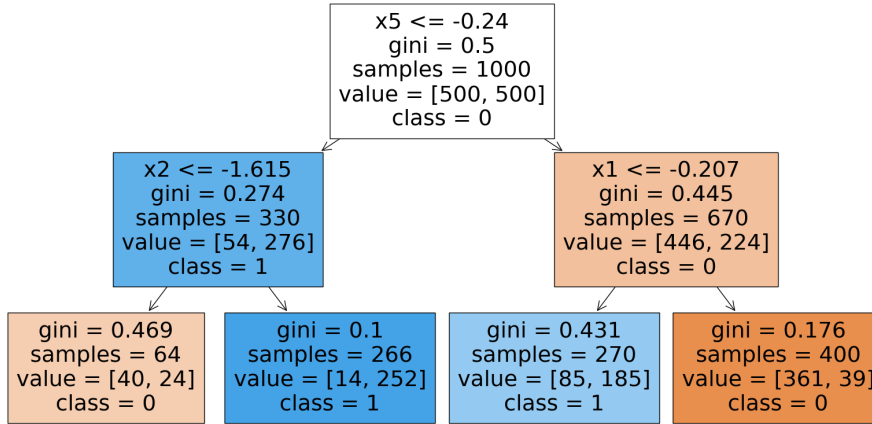
Figure 6: Simple decision tree model used for classification, in the first row of each decision node we show the condition, following, it presents the gini coefficient, number of total samples, the share of instances in each class (0 and 1, respectively) and, finally, the major class. For the node leaves, there's no condition statement and the majority class is the model's prediction class.
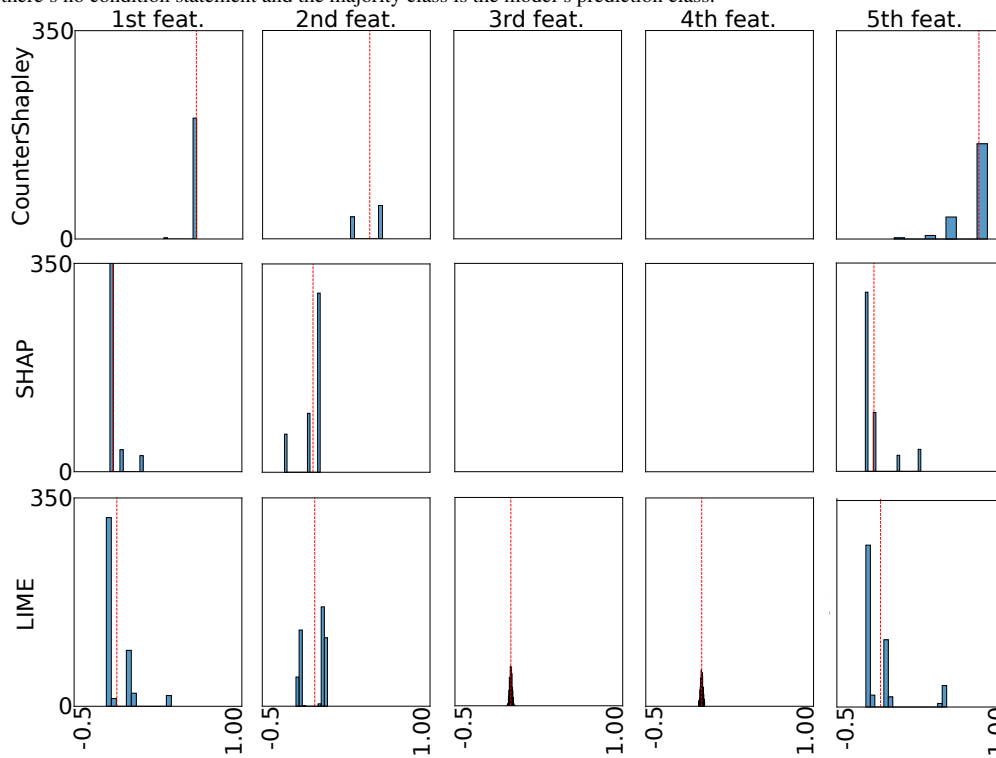


Figure 7: Histogram with importance scores of all points that were classified as 0 in the decision tree classification model, each column represents one feature while each row represents the CounterShapley, SHAP, and LIME distributions, respectively. The red dotted line indicates the distribution's average. For the 3rd and 4th features in the CounterShapley and SHAP methods, the histograms are empty because all importance scores are equal to zero.
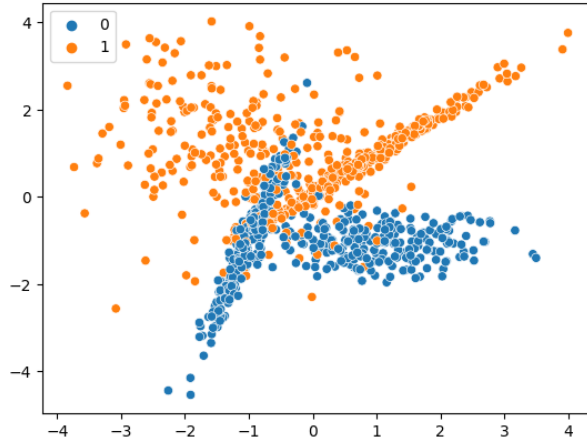
18

Figure 8: Instances classes of the dataset used to the KNN classifier. Blue points have class 0, while the orange points are classified as 1.

classes than the $x$-axis. Similarly to the previous cases, in Table 5, we see again the sparsity of the CounterShapley explanations and also that the most important feature (the 2nd) is more frequent in the explanations.

Considering all three experiments with the simple models described above, we see that CounterShapley assigns importance scores that are substantially different from SHAP and LIME. One major difference observed in all experiments is the fact that CounterShapley values focus on the score changes needed to flip the classification - given the already mentioned decision-driven nature. This characteristic can make the score interpretation easier since both LIME and SHAP assign importance scores to every class (in the case of binary classification, both 0 and 1 classes). However, we highlight that, depending on the research question, multiple importance scores for classes may be a desirable outcome, with the practitioner being on a task to consider each case in particular. The experiments also present, in general, the intrinsic model's justifications (weights, decision nodes, and data distribution) correspond to the average model's explanations. Consequently, all methods are aligned with the main features responsible for the classification. Nevertheless, as thoughtfully discussed in the previous sections, the numerical scoring values and intrinsic meaning of CounterShapley values are not the same as SHAP and LIME.

### 6.2. Simple Data Experiments

For the next experiments, we evaluate how CounterShapley, SHAP, and LIME assign importance scores to feature changes in a complex black box model based on random forest classifiers. However, if we use equally complex data, we wouldn't be able to say if the numerical results correspond to what is expected from the model and data. Therefore, we use simple datasets whose patterns can be clearly grasped by analyzing their class distributions. These data have multiple patterns and follow a similar concept as the experiments performed by Robnik-Sikonja *et al.* [52].

Our first experiment evaluates the impact of adding noisy training data by making a very simple dataset consisting of 3 informative features and 3 random features with no direct connection to the class labels. This dataset was artificially generated using the Scikit-Learn package, for a classification task with two labels and a fixed random state (42). Figure 10 shows the distribution of importance score values for informative features (1st, 2nd, and 3rd) and random features (4th, 5th, and 6th). Additionally, Table 6 shows the percentage of features that had non-zero importance scores.

The importance scores show that all three methods assign more relative importance to the informative features if compared to the random ones. For this specific dataset, all three methods also show a similar pattern having the 1st feature as the most important, followed by the 2nd and 3rd. One may question why the random features are not equal to zero, especially because in a previous experiment, in which the model only used a subset of features, both CounterShapley and SHAP assigned exactly 0 importance scores to unused features. The reason is that, differently from the latter case, the random features in this experiment influence the prediction scoring. We also see the characteristics mentioned in the previous experiments, like CounterShapley having positive scores, sparsity with the most important features being more frequent, and average scores different from SHAP and LIME, still happening in this case.
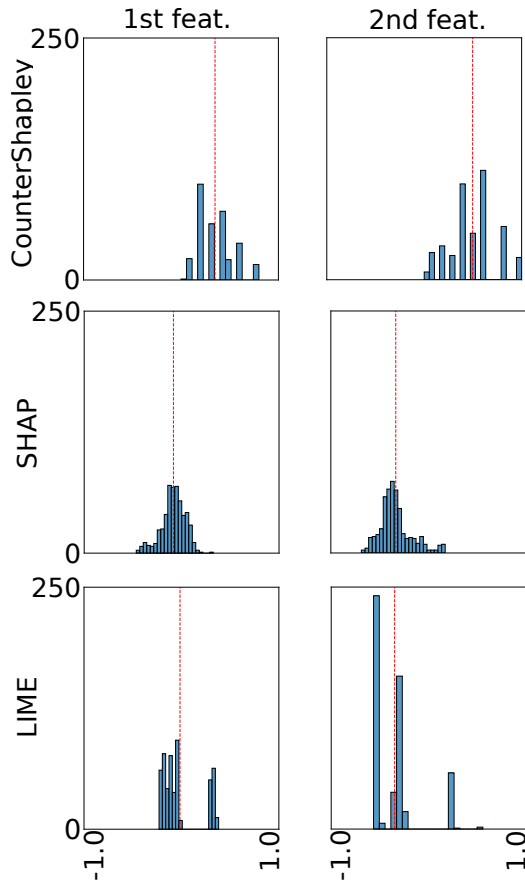
19

Figure 9: Histograms with importance score distribution of 0 classified instances in the KNN model for the 1st and 2nd features (columns) and different methods (rows).

In the following simple data experiments, we use only 2 features, and classes synthetically generated with Scikit-Learn to follow a certain geometric pattern that can be detected by simple visual analysis. The first experiment divides the data points into four quadrants. The separation between the upper and lower quadrants has different thresholds, as shown in Figure 11. This figure also shows the $x$-axis importance score using CounterShapley, SHAP, and LIME methods. Since we use a binary class and the random forest classifier assigns classes with a high degree of confidence, the $y$-axis importance charts are very similar but with inverse scores, with clear regions being darker and vice-versa.

The charts in Figure 11 show that SHAP and LIME have very similar behavior scoring importance; for the $y$ threshold equal to 0.5, they score about the same score for all points in any region. This can be justified since we have equally distributed classes over the 4 quadrants: therefore, overall, the $x$ and $y$ axis contribute equally to the output prediction. When the threshold is changed, the upper and lower quadrants have different scores for the $x$ feature. This is also a reasonable result because the region with increased height has a lower importance for $y$ since a variation in that direction has a reduced effect on probability while the $x$-axis is unchanged. However, the pattern observed for CounterShapley is substantially different if compared to the other two methods. As previously discussed, counterfactual explanations show the minimal feature change to modify the model's classification. Therefore, with a score based on this concept, the explanations are divided into regions where the $x$ feature is closer to the decision threshold (hereafter having a higher importance score) and regions in which the $y$ feature is closer to the decision threshold (consequently, $x$ has a lower importance score). So for the same data and model, CounterShapley gives different scores if compared to SHAP and LIME. All methods have reasonable justifications for their scoring strategy. However, we emphasize that CounterShapley seems more reasonable if we want to score counterfactual explanations since it uses the same concepts.

20

Table 6: Share, in percentage, of explanation instances which includes non-zero importance scores for the three correlated (1st, 2nd, 3rd, highlighted in bold) and three random features (4th, 5th, 6th, unrelated to the label class).

| Features | 1st | 2nd | 3rd | 4th | 5th | 6th |
|---|---|---|---|---|---|---|
| CounterShapley | 83.83 | 35.53 | 31.34 | 4.99 | 6.9 | 3.19 |
| SHAP | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 |
| LIME | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 |



Figure 10: Histograms with importance score distribution of 0 classified instances for the three correlated features (1st, 2nd, and 3rd) and the randomly generated features (4th, 5th, and 6th).

For the next experiment, the dataset has a pattern with two behaviors. For features up to $x = 0.5$ the classification has a linear ($x = y$) threshold, while for values higher than 0.5, it has a constant behavior ($y = 0.5$). Figure 12 illustrates this dataset class distribution and the importance scores for the $x$ (top) and $y$ (bottom) features. There we can see the complementary effect of importance scores, given that we have a binary classification and a predictor that assigns scores with a high degree of confidence.

The CounterShapley CFI method clearly separates the two regions in the chart, for the linear behavior region ($x < 0.5$) does not show a clear preference for either feature while, for the constant region, the $y$ feature has evident higher importance values. This behavior is justified because, in the linear region, the $x$ or $y$ feature distance to the classification threshold is about the same for most instances, while in the constant region, $y$ is usually the closest path to the model's decision threshold. Unlike the previous cases, SHAP and LIME had substantially different patterns; however, both generally give more importance to the $y$ feature. Nevertheless, their importance has a clearly more subtle relationship with the dataset class distribution if compared to CounterShapley importance scores.

Figure 13 shows the importance scores for a circular pattern in the dataset. In this case, the 0 class points are concentrated in the center, while the 1 class is around the circular region formed by the 0 class instances.

The analysis of importance values shows that LIME assigns about the same importance for the $x$ and $y$ features in all instances. This may happen given the sampling strategy, which randomly selects nearby points which end up normalizing the value given the circular geometry. Interestingly, SHAP and CounterShapley have some similarities where the middle top and bottom regions assign a higher relevance to the $y$ feature while the lateral regions assign more importance to the $x$ feature. However, the internal region of SHAP does not look to have this pattern. The justification
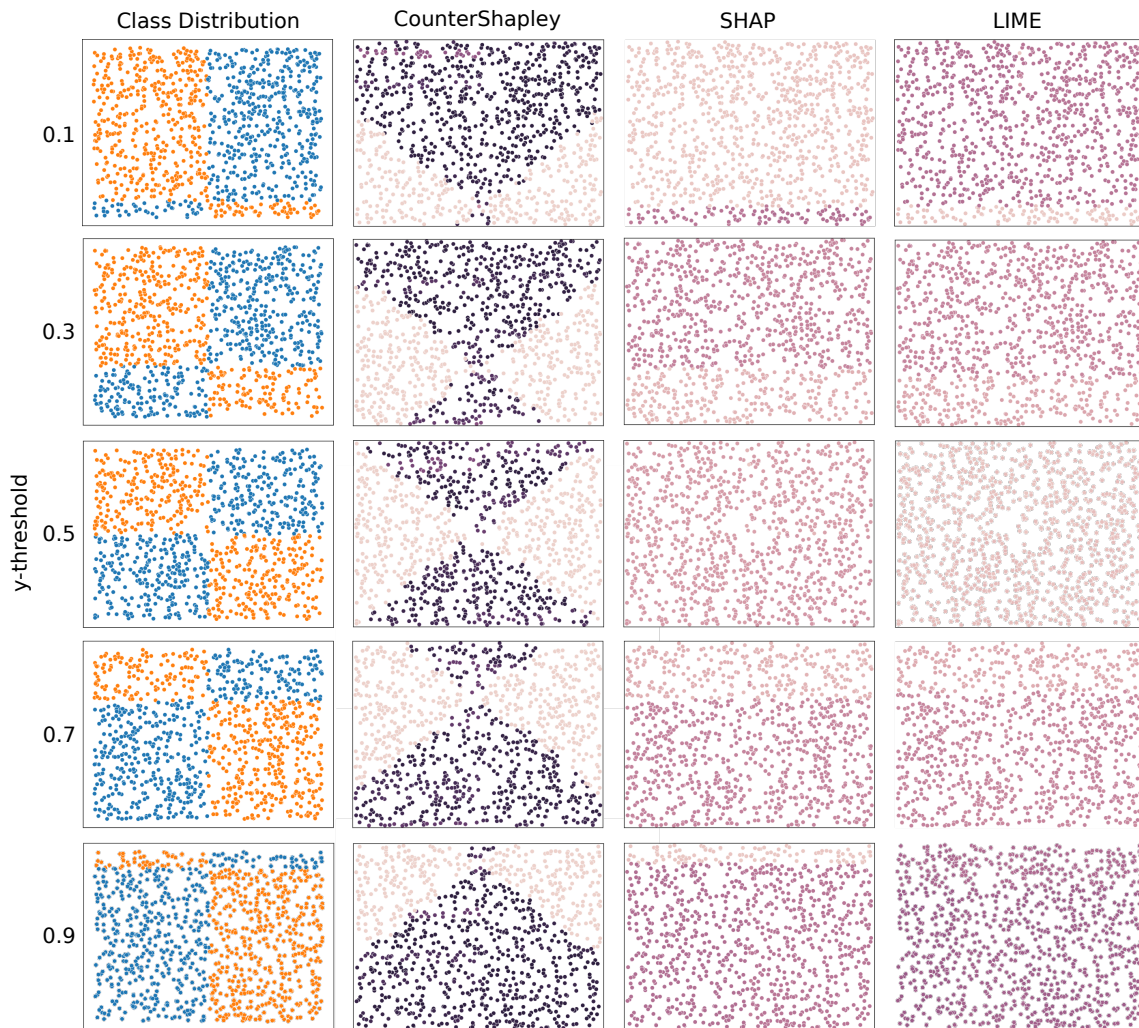
21

Figure 11: First column shows the class distribution (0-class blue, 1-class orange) of each instance in the dataset. The following columns show in a color scale (darker means higher values) the $x$ importance for each instance using CounterShapley, SHAP, and LIME.

Figure 12: Top chart shows the data distribution with class 0 (blue) and class 1 (orange) instances. The second and third row charts show importance weights for the $x$ and $y$ features, respectively.

for such a pattern in CounterShapley may be explained by the sparsity-optimized strategy that the counterfactual generator uses. Therefore, in lateral regions, the most optimal feature change to reach the decision threshold is modifying the $x$ feature, while it is $y$ otherwise.

All experiments with simple data above show again that CounterShapley values give considerably different importance values for feature changes, following a scoring pattern more consistent with the idea behind counterfactual explanations, minimal sparse changes to change prediction classification. Moreover, we can notice the capacity that CounterShapley values have to react differently to multiple patterns in the same dataset, where these regions are less evident if we use SHAP and LIME. This characteristic is intrinsically linked to the fact that counterfactual explanations are based on only two points (factual and counterfactual). In contrast, SHAP and LIME have a normalizing effect since they consider multiple points over the dataset. This difference does not necessarily represents an advantage of CounterShapley values because all methods present reasonable justifications for their scores. Therefore, practitioners must consider all these characteristics and their objectives before selecting an importance scoring method.

### 6.3. Replication Experiments

One of the main advantages of the methods presented in this paper is their broad compatibility with any counterfactual generation algorithm. The compatibility is not only in terms of theoretical requirements since they only need the factual, counterfactual, and model prediction functions, but also given our algorithmic framework that can be easily implemented using Python and popular data science packages. In this section, we show how our methods can be applied to three different counterfactual generation algorithms: NICE [46], DiCE [43], and ALIBIC [44]. We provide a snippet code of how they are implemented, highlighting the simplicity, and we perform some experiments with data and models which they use in their respective papers. All these experiments are publicly available on GitHub in the Experiments branch.

NICE [46] is a fast and reliable counterfactual generation algorithm that finds optimal explanations. It uses dataset labeled instances to find a first optimal instance with a different class and then perform multiple optimization iterations to improve the counterfactual with respect to a specific optimization objective (which can be sparsity, proximity, or
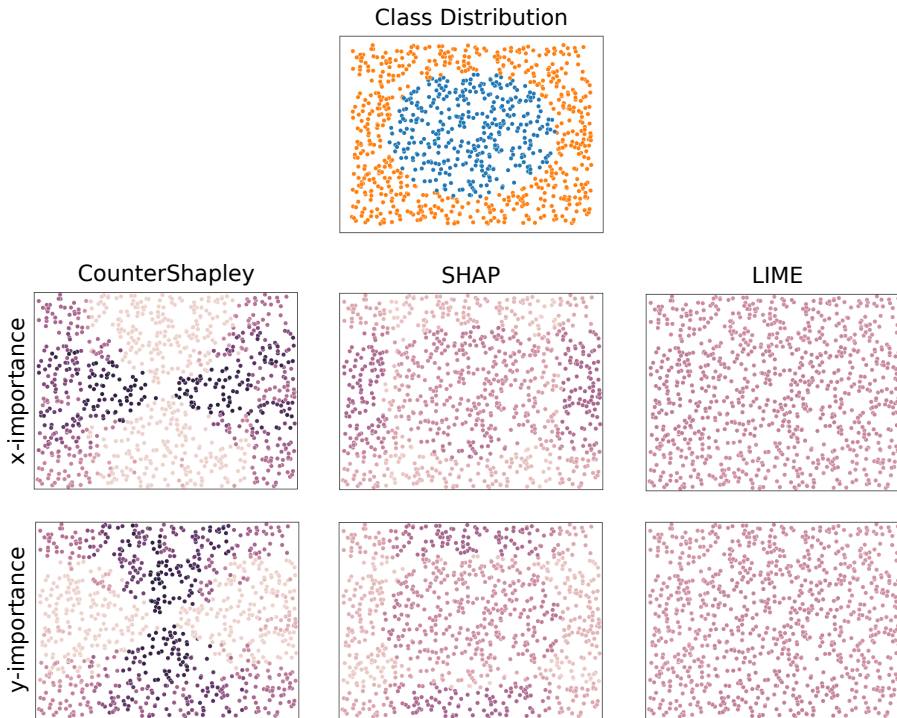
Figure 13: Top chart shows the data distribution with class 0 (blue) and class 1 (orange) instances. The second and third row charts show importance weights for the $x$ and $y$ features, respectively.

plausibility). Appendix Algorithm 4 shows how easily the NICE official package[4] can be adapted to generate all analyses proposed by the methods presented in this article.

Our method allows the creation of an object (showed in the Appendix Algorithm 4) that works exactly like NICE, since it is an extended class, but we change the output of the "explain" function to not being a simple list of the counterfactual points. Rather, the function's output returns "CounterPlot" objects, which are able to give the CounterShapley values for each feature and plot any of the three charts presented here (Greedy, CounterShapley, and Constellation). Figure 14 shows an example of the result obtained from an instance of the Wisconsin Breast Cancer dataset, as used in the NICE paper.

We also use the DiCE's [43] counterfactual generator official Python package[5]. This counterfactual generator can work independently of a specific model (model agnostic) and allows to generate a set of diverse counterfactual explanations. Appendix Algorithm 5 shows how a function can take the objects generated by the DiCE counterfactual generation and transform them into the object that allows all analysis presented in this paper.

Note that, in DiCE's case, the algorithmic framework is more complex since the package is responsible for the dataset encoding, and the resulting counterfactual is inside a custom object. Nevertheless, our algorithm can still be adapted and processes all the diverse counterfactual instances produced by DiCE. As a sample example, we use the same dataset DiCE uses in their package official documentation, the Adult dataset, and Figure 15 shows the results. In this case, we evaluate the features needed to change the prediction from low-income (class 0) to high-income (class 1).

ALIBIC [44] is also adapted using our framework using Appendix Algorithm 6. This counterfactual generator performs a gradient descent using a complex loss function that includes terms for distance, sparsity, and similarity to the dataset manifold. Similarly to what we have done with DiCE, we use the package official documentation[6] to get a

---

[4]https://github.com/DBrughmans/NICE

[5]https://github.com/interpretml/DiCE

[6]https://docs.seldon.io/projects/alibi/en/stable/examples/cfproto_housing.html
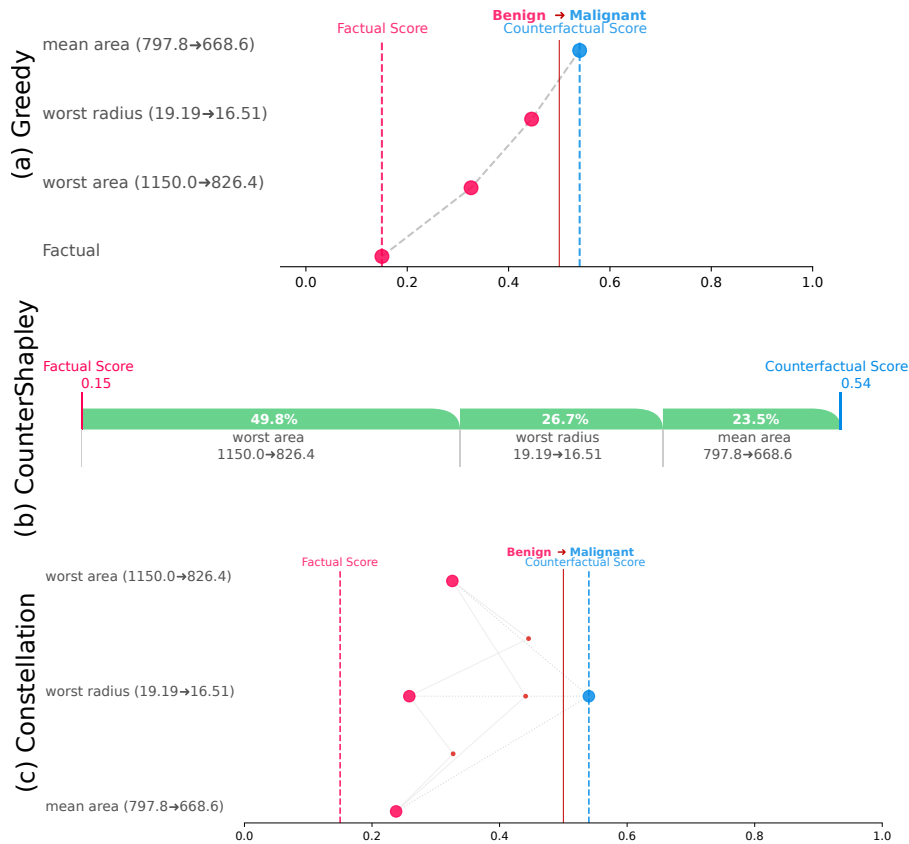
Figure 14: Greedy, CounterShapley, and Constellation charts generated using a NICE counterfactual explanation from the Wisconsin Brest Cancer dataset.

dataset, generate a model and counterfactual explanation. In this case, we use the California Housing dataset and an artificial neural network as the machine learning model. ALIBIC generates an object that includes multiple attributes, such as the counterfactual point, classes, and features' ranges.

ALIBIC's adaptation to generate our analysis is considerably simpler than that of DiCE, as DiCE uses multiple built-in methods from their framework, such as the dataset treatment and model adapter. For ALIBIC (and NICE), we use simpler strategies that directly use the factual point and the prediction model. Figure 16 shows an instance explanation example with the three chart types.

With the examples shown above, we showcase that not only the theoretical concepts described in this article are independent of a specific counterfactual generator, but also our algorithmic framework works with any type of generator. Therefore, this represents a substantial contribution to the counterfactual explanation area since it can readily embed a higher explanation value to novel and past algorithms.

### 6.4. Missclassification Experiments

One of the most valuable uses of explanation methods, especially those that focus on local/instance-level explanations, is the investigation of why a model assigned a wrong classification to a given instance. The explanation can reveal patterns or features that can be further investigated, leading to correction or improvements.

Counterfactual explanations, as already mentioned, have multiple characteristics (simple, sparse, decision-oriented, etc.) that make them a valuable method for investigating wrong classification instances. However, the isonomic view of feature changes' importance to the model's prediction may not give enough information for a complete understanding of the underlying causes of the misclassification. Therefore, this experiment shows an example of how the
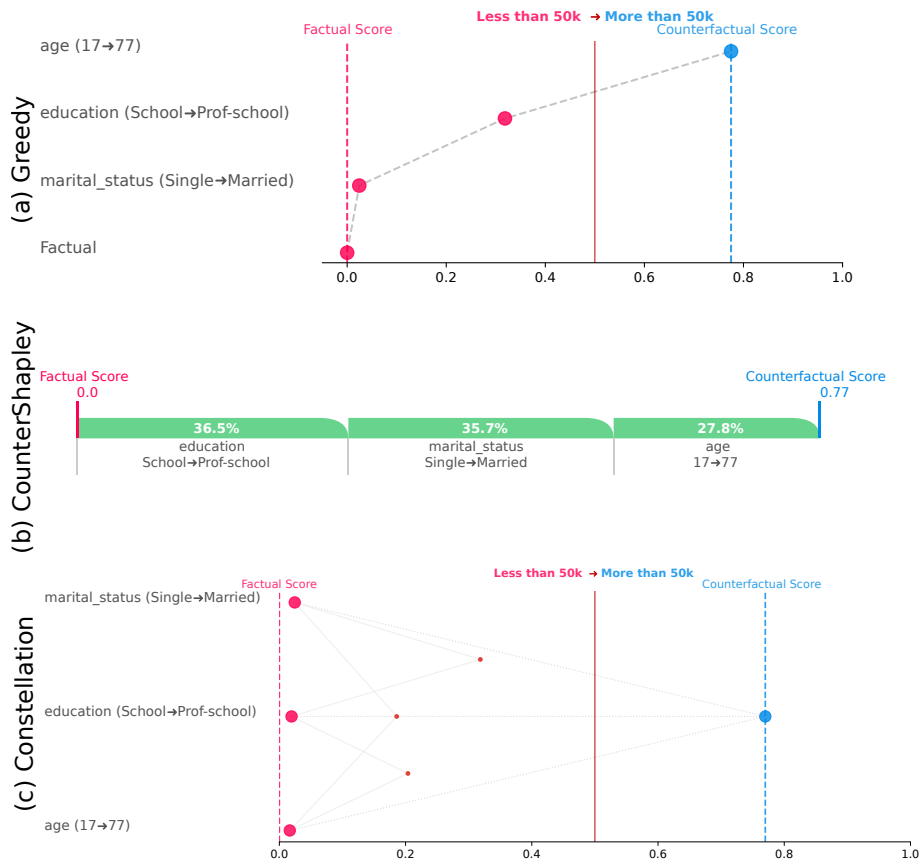
Figure 15: Feature change importance charts (Greedy, CounterShapley, and Constellation) for the Adult dataset using DiCE counterfactual generator in a sample instance.

analysis presented in this article can help to unveil more information about the model scoring and decision. To illustrate this application, we use the Bankruptcy Prediction dataset [53], which includes bankruptcy data from the Taiwan Economic Journal from 1999 to 2009. As a model, we use the Scikit-Learn random forest algorithm with default parameters and a fixed random seed (42).

Figure 17 shows the feature change importance for an instance that is incorrectly classified as a healthy company but got bankrupt. The counterfactual explanation and instances can show interesting insights about this misclassification, such as: a lower total asset growth rate would positively contribute towards a correct classification. This makes sense since slower growth generally is seen as a pessimistic indicator of a company's prospects. As this was the most important feature change, contributing to more than 50% of CounterShapley value, measures such as the acquisition of more data related to asset growth and feature engineering could potentially help the model achieve better performance. Moreover, the two other features also give us an interesting perspective on the misclassified instance. The second most important feature change (according to CounterShapley values) is a higher gross profit to sales share, and the third, with a similar importance value, is a higher net value per share. At first sight, these changes look unexpected since both rises are associated with positive outcomes. However, it's important to notice these two changes alone are not sufficient to flip the classification and must be associated with the first, most important feature change. This analysis enhances the informative value of counterfactual explanations, allowing us to investigate the pattern formed by these feature changes: maybe the association of higher shares and profit with lower asset growth indicates a purely speculative/momentaneous gain which is not sustainable in the long run. These unexpected associations can also mean that the model is not behaving as wanted, making spurious associations caused by overfitting or another problem in the dataset. In summary, the counterfactual explanations associated with the CFI methods and charts give
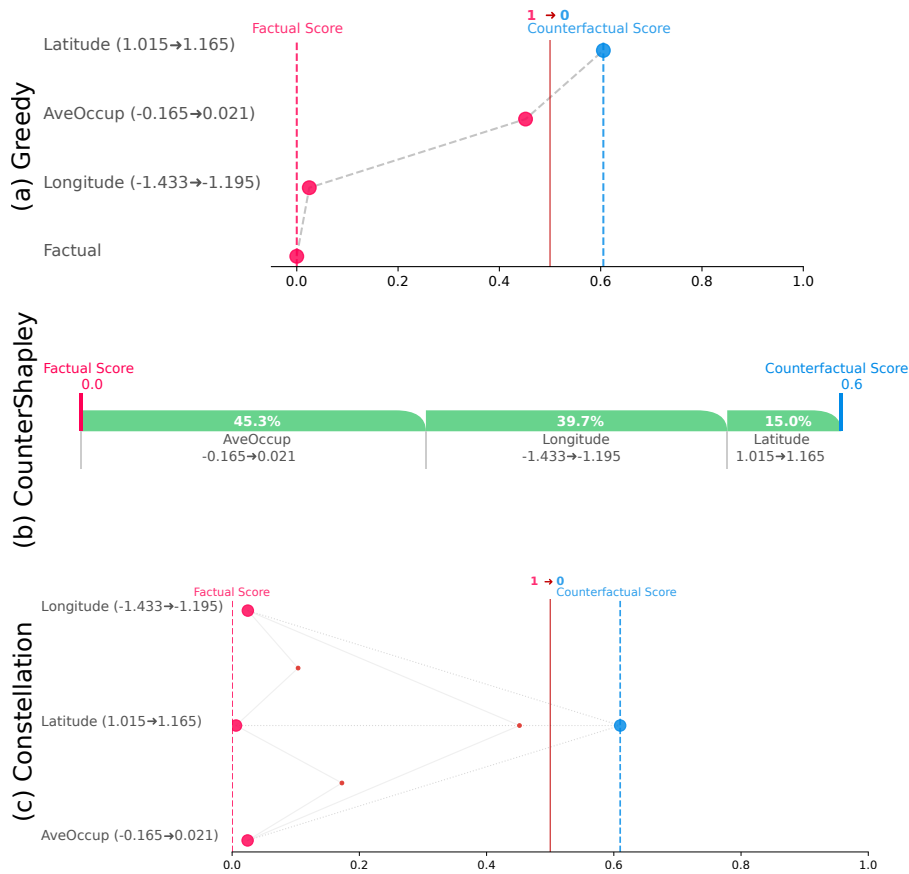
Figure 16: ALIBIC feature change importance charts, Greedy, CounterShapley, and Constellation using an instance from the California housing dataset.

us extra knowledge of how changes affect the prediction score and classification result, which can lead to better model understanding and improvements.

### 6.5. Special Cases

Finally, in this last section of the experiments, we show two special cases in which the analysis presented in this article could contribute to a better understanding of the model and counterfactual result evaluation. First, we talk about counterfactual feature changes which have negative CounterShapley values, revealing non-linear behaviors between features. Then, we present a way to identify counterfactual explanations that have more feature modifications than needed.

#### 6.5.1. Negative Contributions

Complex models can lead to highly non-linear behaviors in which features that usually have a certain behavior (for example, increase the probability of a class) completely change when a specific set of features are present. Simple counterfactual explanations cannot detect such behavior since they only disclose the feature modifications needed to make a change in the classification prediction. However, as shown in Figure 18, our analysis can detect features that, although they have a negative contribution over the prediction scoring, are necessary for a class change.

Our analysis uses the UCI's Wine Quality dataset[7] and a random forest classifier built with Scikit-Learn and default parameters. In Figure 18, we show one instance example that has a feature (*f3*) with a negative CounterShapley value

---

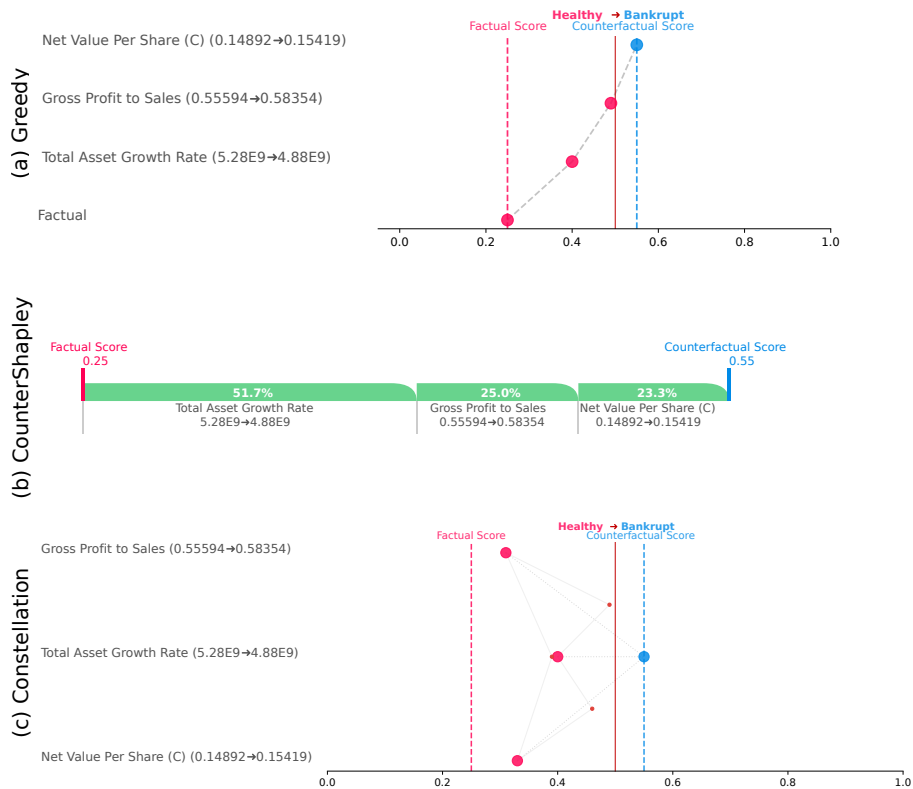[7]Dataset page: https://archive.ics.uci.edu/ml/datasets/wine+quality

Figure 17: NICE counterfactual explanation of an instance from the Bankruptcy Prediction dataset, which was incorrectly classified by the model as a healthy company (but it got bankrupt instead).

as seen in its respective CounterShapley chart. This can be further understood with the Constellation chart, which shows that *f3* decreases the inverse class probability if changed alone (bottom clear blue dot) or in association with other features. However, as can also be seen in the Greedy and Constellation chart, this feature is fundamental for the class change since, if absent, it is not sufficient to modify the model's predicted class.

### 6.5.2. Counterfactual subset

The definition of a counterfactual explanation includes the requirement of having an irreducible set of valid modifications (as described in Equation 8). Although the interpretation of valid modifications is debatable [22], being able to identify subsets of changes that form a counterfactual explanation may be useful for a better understanding of the prediction model and for the assessment of the counterfactual explanation.

The usual simple counterfactual results are not capable of showing this since they only show the feature changes needed for a class modification. Therefore, the analysis presented in this work may be a valuable tool for evaluating the counterfactual result. Figure 19 shows an example of how a counterfactual with class changing subsets would look in Greedy, CounterShapley, and Constellation charts.

In our example figure, the Greedy and CounterShapley charts indicate that the counterfactual explanation has more features than needed to change the classification change. The Constellation Chart may be the most appropriate analysis since it clearly shows which subsets of change lead to a counterfactual class. For this example, we can see the change in *occupation* feature is not needed to modify the prediction class.
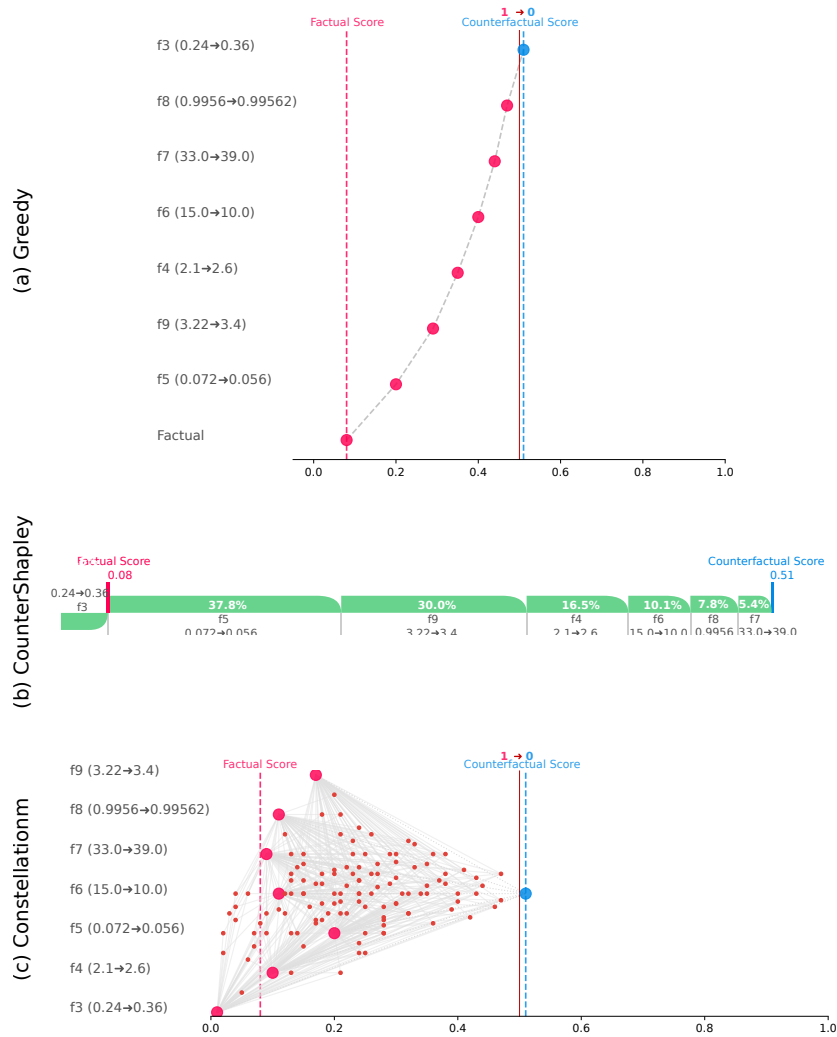
Figure 18: Analysis made over a counterfactual explanation generated by NICE that has negative CounterShapley value. Although, for the Greedy chart, this behavior cannot be directly observed, in CounterShapley and Constellation chart is possible to see that *f3* (and some subsets containing it) generate prediction scores lower than the factual score.

# 7. Conclusion

Machine learning predictions can be effectively explained with counterfactual explanations. Its simple and decision-driven nature makes it easier for even laypersons to get a sufficient understanding of why automated decisions are made. However, up to now, the representation of this type of explanation was mainly based on textual communication which, although has its own benefits, lacks the benefits of graphical elements. Moreover, the isonomic view of the counterfactual modified features may not give a complete picture of why a model has certain decisions since it does not explain the relationship between features and prediction scoring. This paper proposes two CFI methods as a solution to those challenges, Greedy and CounterShapley, as being approaches that can embed relative importance to counterfactual features (and their corresponding value modifications) and three chart types that highlight different aspects of the counterfactual explanation. Hence, the Greedy CFI method illustrates how a specific path, following the steepest increases of scores, links to the counterfactual changes, and CounterShapley fills the knowledge gap of how the prediction scores are affected by each changed feature in overall. Those CFI methods can provide better insights into how the model works, consequently augmenting the informative value of the explanations. Regarding the plotting strategies, CounterShapley charts use the feature change importance values to create simple yet descriptive graphs.
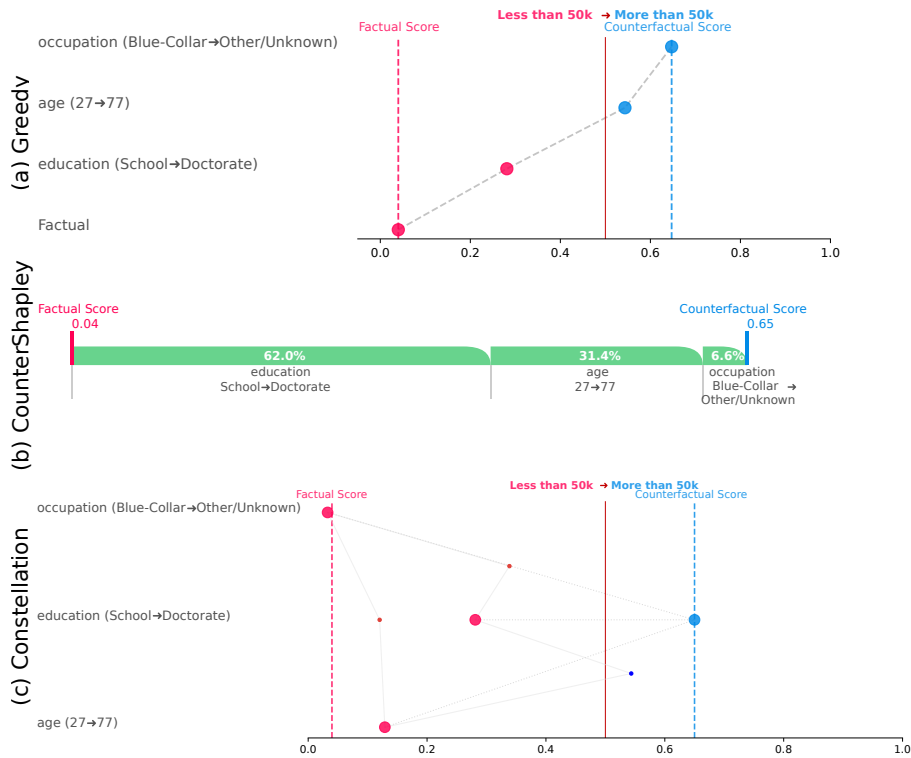
29

Figure 19: Example of a counterfactual with a subset that also flips the prediction class. We use the Adult dataset and the DiCE counterfactual generator for this case.

For the cases in which computational burden is an issue, we propose the Greedy Counterfactual chart to quickly provide a visual representation that describes the greediest strategy to achieve the counterfactual prediction score. The third chart type, the Constellation chart, takes advantage of the sparse nature of counterfactual explanations to show the possible feature associations and their impact on the prediction scoring. Despite the usefulness of these methods in augmenting the informative value of counterfactual explanations, they are not adequate for situations where the characteristics of counterfactuals are not desirable, for example, if an importance score is required for every model's feature. Finally, all methods and charts here described are implemented in an open-source package that can be easily integrated into any counterfactual generator, bringing immediate impact considering the multiple counterfactual generation algorithms available. As for future perspective, we envision the application of counterfactual-based scores and their visual representations in real-world cases, assessing the user reception in terms of understandability and preference over other explainability approaches.

# References

[1] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al., Language models are few-shot learners, Advances in neural information processing systems 33 (2020) 1877–1901.

[2] S. Ronanki, G. E. Henter, Z. Wu, S. King, A template-based approach for speech synthesis intonation generation using lstms., in: INTER-SPEECH, 2016, pp. 2463–2467.

[3] P. Neekhara, S. Hussain, S. Dubnov, F. Koushanfar, J. McAuley, Expressive neural voice cloning, in: Asian Conference on Machine Learning, PMLR, 2021, pp. 252–267.

[4] S. H. Yang, M. Chung, Self-imitating feedback generation using gan for computer-assisted pronunciation training, arXiv preprint arXiv:1904.09407 (2019).

[5] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins, et al., Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai, Information fusion 58 (2020) 82–115.

[6] C. Meske, E. Bunde, J. Schneider, M. Gersch, Explainable artificial intelligence: objectives, stakeholders, and future research opportunities, Information Systems Management 39 (2022) 53–63.

[7] S. Wachter, B. Mittelstadt, C. Russell, Counterfactual explanations without opening the black box: Automated decisions and the gdpr, Harv. JL & Tech. 31 (2017) 841.

[8] S. Goethals, D. Martens, T. Calders, Precof: Counterfactual explanations for fairness (2022).

[9] D. Martens, Data Science Ethics: Concepts, Techniques, and Cautionary Tales, Oxford University Press, 2022. doi:10.1093/oso/9780192847263.001.0001.

[10] T. Speith, A review of taxonomies of explainable artificial intelligence (xai) methods, in: 2022 ACM Conference on Fairness, Accountability, and Transparency, 2022, pp. 2239–2250.

[11] L. Breiman, Random forests, Machine learning 45 (2001) 5–32.

[12] M. Craven, J. Shavlik, Extracting tree-structured representations of trained networks, advances, neural information processing systems 8, 1996.

[13] D. Martens, B. Baesens, T. Van Gestel, Decompositional rule extraction from support vector machines by active learning, IEEE Transactions on Knowledge and Data Engineering 21 (2009) 178–191.

[14] T. Hastie, R. Tibshirani, J. Friedman, The Elements of Statistical Learning, Springer Series in Statistics, Springer New York Inc., New York, NY, USA, 2001.

[15] M. T. Ribeiro, S. Singh, C. Guestrin, "why should I trust you?": Explaining the predictions of any classifier, in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016, 2016, pp. 1135–1144.

[16] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, Advances in neural information processing systems 30 (2017).

[17] S. Verma, J. Dickerson, K. Hines, Counterfactual explanations for machine learning: A review, arXiv preprint arXiv:2010.10596 (2020).

[18] L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae, P. Prettenhofer, A. Gramfort, J. Grobler, R. Layton, J. VanderPlas, A. Joly, B. Holt, G. Varoquaux, API design for machine learning software: experiences from the scikit-learn project, in: ECML PKDD Workshop: Languages for Data Mining and Machine Learning, 2013, pp. 108–122.

[19] L. Shapley, Quota solutions op n-person games1, Edited by Emil Artin and Marston Morse (1953) 343.

[20] I. Stepin, J. M. Alonso, A. Catala, M. Pereira-Fariña, A survey of contrastive and counterfactual explanation generation methods for explainable artificial intelligence, IEEE Access 9 (2021) 11974–12001.

[21] D. Martens, F. Provost, Explaining documents' classifications, Center for Digital Economy Research (2011).

[22] R. M. B. de Oliveira, D. Martens, A framework and benchmarking study for counterfactual generating methods on tabular data, Applied Sciences 11 (2021) 7274.

[23] C. Fernández-Loría, F. Provost, X. Han, Explaining data-driven decisions made by ai systems: The counterfactual approach, arXiv preprint arXiv:2001.07417 (2020).

[24] R. M. Byrne, Precis of the rational imagination: How people create alternatives to reality, Behavioral and Brain Sciences 30 (2007) 439–453.

[25] P. Lipton, Contrastive explanation, Royal Institute of Philosophy Supplements 27 (1990) 247–266.

[26] T. Vermeire, D. Brughmans, S. Goethals, R. M. B. de Oliveira, D. Martens, Explainable image classification with evidence counterfactual, Pattern Analysis and Applications (2022) 1–21.

[27] Y. Ramon, T. Vermeire, O. Toubia, D. Martens, T. Evgeniou, Understanding consumer preferences for explanations generated by xai algorithms, arXiv preprint arXiv:2107.02624 (2021).

[28] D. Martens, F. Provost, Explaining data-driven document classifications, MIS quarterly 38 (2014) 73–100.

[29] F. Hohman, A. Head, R. Caruana, R. DeLine, S. M. Drucker, Gamut: A design probe to understand how data scientists understand machine learning models, in: Proceedings of the 2019 CHI conference on human factors in computing systems, 2019, pp. 1–13.

[30] J. Wexler, M. Pushkarna, T. Bolukbasi, M. Wattenberg, F. Viégas, J. Wilson, The what-if tool: Interactive probing of machine learning models, IEEE transactions on visualization and computer graphics 26 (2019) 56–65.

[31] F. Cheng, Y. Ming, H. Qu, Dece: Decision explorer with counterfactual explanations for machine learning models, IEEE Transactions on Visualization and Computer Graphics 27 (2020) 1438–1447.

[32] O. Gomez, S. Holter, J. Yuan, E. Bertini, Vice: visual counterfactual explanations for machine learning models, in: Proceedings of the 25th International Conference on Intelligent User Interfaces, 2020, pp. 531–535.

[33] P. P. Angelov, E. A. Soares, R. Jiang, N. I. Arnold, P. M. Atkinson, Explainable artificial intelligence: an analytical review, Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery 11 (2021) e1424.

[34] E. Winter, The shapley value, Handbook of game theory with economic applications 3 (2002) 2025–2054.

[35] E. Strumbelj, I. Kononenko, An efficient explanation of individual classifications using game theory, The Journal of Machine Learning Research 11 (2010) 1–18.

[36] M. R. Zafar, N. Khan, Deterministic local interpretable model-agnostic explanations for stable explainability, Machine Learning and Knowledge Extraction 3 (2021) 525–541.

[37] G. Visani, E. Bagli, F. Chesani, A. Poluzzi, D. Capuzzo, Statistical stability indices for lime: Obtaining reliable explanations for machine learning models, Journal of the Operational Research Society 73 (2022) 91–101.

[38] M. R. Zafar, N. M. Khan, Dlime: A deterministic local interpretable model-agnostic explanations approach for computer-aided diagnosis systems, arXiv preprint arXiv:1906.10263 (2019).

[39] E. Štrumbelj, I. Kononenko, Explaining prediction models and individual predictions with feature contributions, Knowledge and information systems 41 (2014) 647–665.

[40] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, S.-I. Lee, From local explanations to global understanding with explainable AI for trees, Nature Machine Intelligence 2 (2020) 56–67.

[41] S. Ghalebikesabi, L. Ter-Minassian, K. DiazOrdaz, C. C. Holmes, On locality of local explanation models, Advances in Neural Information Processing Systems 34 (2021) 18395–18407.

[42] J. Moore, N. Hammerla, C. Watkins, Explaining deep learning models with constrained adversarial examples, in: Pacific Rim international conference on artificial intelligence, Springer, 2019, pp. 43–56.

[43] R. K. Mothilal, A. Sharma, C. Tan, Explaining machine learning classifiers through diverse counterfactual explanations, in: Proceedings of the 2020 conference on fairness, accountability, and transparency, 2020, pp. 607–617.

[44] A. V. Looveren, J. Klaise, Interpretable counterfactual explanations guided by prototypes, in: Joint European Conference on Machine Learning and Knowledge Discovery in Databases, Springer, 2021, pp. 650–665.

[45] S. Ainsworth, A. Th Loizou, The effects of self-explaining when learning with text or diagrams, Cognitive science 27 (2003) 669–681.

[46] D. Brughmans, P. Leyman, D. Martens, Nice: an algorithm for nearest instance counterfactual explanations, arXiv preprint arXiv:2104.07411 (2021).

[47] G. Vilone, L. Longo, Notions of explainability and evaluation approaches for explainable artificial intelligence, Information Fusion 76 (2021) 89–106.

[48] J. D. Hunter, Matplotlib: A 2d graphics environment, Computing in Science & Engineering 9 (2007) 90–95.

[49] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, Scikit-learn: Machine learning in Python, Journal of Machine Learning Research 12 (2011) 2825–2830.

[50] T. pandas development team, pandas-dev/pandas: Pandas, 2020. URL: https://doi.org/10.5281/zenodo.3509134. doi:10.5281/zenodo.3509134.

[51] M. L. Waskom, seaborn: statistical data visualization, Journal of Open Source Software 6 (2021) 3021.

[52] M. Robnik-Šikonja, I. Kononenko, Explaining classifications for individual instances, IEEE Transactions on Knowledge and Data Engineering 20 (2008) 589–600.

[53] D. Liang, C.-C. Lu, C.-F. Tsai, G.-A. Shih, Financial ratios and corporate governance indicators in bankruptcy prediction: A comprehensive study, European journal of operational research 252 (2016) 561–572.

# Appendices

---

**Algorithm 4:** NICE's Python package adaptation for integrating the article's methods

---

```python
from nice import NICE
from counterplots import CreatePlot

class CounterPlotNice(NICE):

    def __init__(self, predict_fn, *args, **kwargs):
        super().__init__(predict_fn, *args, **kwargs)
        self.predict_fn = predict_fn

    def explain(self, X, target_class='other'):
        explanations = super().explain(X, target_class)

        out_exp = []

        for i in range(len(X)):
            out_exp.append(
                CreatePlot(factual=X[i],
                cf=explanations[i],
                model_pred=self.predict_fn)
            )

        return out_exp
```

---

**Algorithm 5:** Function which takes DiCE counterfactual results and generates this article's analysis and plots

```python
import json
import pandas as pd
from dice_ml import Dice
from counterplots import CreatePlot

# exp and dice_exp is generated using DiCE package, following official documentation

def dice_counterplots(dice_exp, exp):

    cf_data = json.loads(dice_exp.to_json())
    factual = cf_data['test_data'][0][0][:-1]
    feature_names = exp.data_interface.feature_names
    df_structure = exp.data_interface.data_df[:0].loc[:, feature_names]
    data_types = df_structure.dtypes.apply(lambda x: x.name).to_dict()

    def adjust_types(x):
        for i in range(x.shape[1]):
            if 'int' in list(data_types.values())[i]:
                x[:, i] = int(float(x[:, i]))
        return x

    def model_pred(x):
        scores = exp.predict_fn(df_structure.append(pd.DataFrame(x, columns=feature_names))).
                                                    numpy()
        return np.concatenate((1 - scores, scores), axis=1)

    out_exp = []

    for raw_cf in cf_data['cfs_list'][0]:
        cf = adjust_types(np.array([raw_cf[:-1]]))[0]

        out_exp.append(
            CreatePlot(
                factual=np.array(factual),
                cf=np.array(cf),
                model_pred=model_pred,
                feature_names=feature_names)
            )

    return out_exp
```

**Algorithm 6:** Function which takes DiCE counterfactual results and generates this article's analysis and plots

```python
from alibi.explainers import CounterfactualProto
from counterplots import CreatePlot

# factual and model are obtained from the dataset treatment/training
# and exp is created with the instance and model data using ALIBIC package

def CounterPlotAlibic(factual, exp, model):
    cf = exp.cf['X'][0]

    out_exp = CreatePlot(
        factual=factual,
        cf=cf,
        model_pred=model
    )

    return out_exp
```

## A. CounterShapley calculation example

If we take as an example a counterfactual **c** and a factual **x**, such that the counterfactual differs from the factual by 5 feature values:

$$|\delta| = K = 5 \tag{13}$$

Let us iterate through all the possibilities in which these features can be changed *one at a time*, while keeping track of the impact of each change to the prediction score. This is one of three ways in which a Shapley value can be calculated. The example presented here follows the definition in Equation 12:

$$\varphi_i(\mathcal{M}) = \sum_{V \subseteq \delta \setminus \{i\}} \frac{|V|! \, (K - |V| - 1)!}{K!} \left( \mathcal{M} \left( V \cup \{i\} \right) - \mathcal{M} \left( V \right) \right) \tag{12}$$

Let us start by considering the case where we change two (unspecified) feature values, and then change feature $i$:

$$[\times \ \times \ i \ \times \ \times] \tag{14}$$

Since $i$ is on the third spot in this example (we change two features, and then we change feature $i$), we have 4 remaining features we could fill in the first two spots marked as $\times$. This means that the sum over all coalitions will have $_4C_2 = 6$ elements for the case where $i$ is in position 3. This is an important result, as it equals the number of model predictions that will have to be made. This means that there are 6 coalitions $V$ of size 2 that do not have $i$. This condition appears in the subscript of the summation sign of the Shapley definition in Equation 12.

Note that, given some coalition of two features that are to be changed before $i$, it does not matter in which order these first two features are changed: both orderings will yield the same marginal contribution for $i$ in position 3. As we are iterating all order-independent coalitions of size 2, this iteration step is the only step a particular coalition will be considered for the first two places. To calculate its contribution compared to all possible configurations where feature $i$ is changed as third (including those where another order of the same coalition is considered), we must multiply whatever contribution $i$ in place 3 yields in this specific configuration by the amount of possible orderings of the precursory coalition: $N_{orders} = |V|! = 2! = 2$. This is done to properly weigh this contribution in the total average over all possible coalitions, and corresponds to the weight of $|V|!$ in Equation 12.

When we now iterate through all possible subsets of size 2 that do not have $i$, we can essentially calculate the marginal contribution of feature $i$ when it is changed as a third feature. This is equivalent to the summation sign in Equation 12, but where the summation would be limited to include only coalitions of size 2. Note that it now also does not matter which features will be changed later on. The features that are still to be changed have no influence on the model output difference yielded by changing $i$ as the third feature, which is the only contribution we are counting for now. The number of iteration steps where $i$ is changed as a third feature change, given that the first two $\times$ are already filled in, can be counted by considering how many features are left to be filled in on the places of the last two $\times$, which equals $(K - |V| - 1)! = (5 - 2 - 1)! = 2! = 2$. Multiplying this with $|V|!$ gives us the total amount of coalitions where $i$ is changed as a third feature (see again Equation 12), or equivalently, the amount of coalitions of size 2 that do not contain $i$.

We have now only considered the case where feature $i$ was changed as third. To gain a complete understanding of the influence of feature $i$, we must also consider the cases where it was changed first, second, fourth, and last. This is equivalent to placing $i$ on each $\times$ spot in Equation 14, or to considering coalitions of every possible size that do not have $i$. It is also equivalent to the summation sign in Equation 12, but now not limited anymore to a fixed coalition size. We can extend the previous logic to all coalition sizes, ranging from 1 to 5.

The only thing that remains is to weigh this sum of marginal contributions with the total amount of possible feature permutations by dividing by $K!$. In doing so, we have the average marginal feature contribution of feature $i$, and we have explained the last factor of Equation 12.

Given a set $\delta$ of size $K$, the total amount of terms Equation 12 will have to add together to calculate the Shapley value for a single feature $i$ equals:

$$N_{coalitions,i} = \sum_{|V|=0}^{K-1} \binom{K-1}{V} = 2^{K-1} - 1 \tag{15}$$