

# Counterfactual explanation

Chenyuan Zhao

Technische Universität München

**Abstract.** Used in classification problems. Counterfactual explanation treat the ML model as a black box, omitting the inner algorithm, instead, focus on the correlation between input and output. This explanation asks how the classification result will change if we alter the input.

**Keywords:** First keyword · Second keyword · Another keyword.

## 1 Introduction

People tend to ask themselves "what-if" questions to dream about another better results when unfortunate happens. These "what-if" arguments are usually a (slightly) different world in which a slight change in action may amend the current outcome. Counterfactual explanations (CF) are mostly applied for binary automated decision in social sensitive regions (such as disease diagnose, loan approval, school admission...), to audit possible bias and artifacts. [2] has pointed out three usages of a explanation of machine decisions: (a.) to understand the decision, (b.) to contest a decision, and (c.) to obtain instructions for better outcomes in the future. Counterfactual explanation, though ignoring the working principle of the model, satisfies all these demands.

Without "opening the black box", a CF algorithm generates a positive classified user-case based on a originally negative one (or vice versa). Providing that the algorithm suggests a higher GPA score or more deposit in bank account, the user understands the reason of the rejection, but also knows the measure to adopt for the desired outcome. However, if the algorithm suggests a change in the race or gender, which implies that the decision is made through a discrimination factor, the user is expected to contest the unfair decision.

In practical, there are still several features to taken into consideration. Sparsity, to ensure an actionable advice for the user, the generation should prefer as few changes in items as possible, changing one item vastly is better than changing several slightly. Diversity, to provide a list of choices, the algorithm should offer multiple CF cases varying in the changed item and in different extend. And finally, interpretability, the suggested CF case should be a possible case in the real life. Later, it will be leveraged that sparsity is contradictable with diversity and interpretability, optimizing sparsity usually causes downgrade in other 2 features.

## 2 Generation of a counterfactual explanation

### 2.1 by loss function

CF explanation doesn't require to understand the internal work principle of a model, instead focusing on which kind of minimal perturbation in the data subject will leads to a different classification result. Therefore, the search of a an appropriate CF example is boiled down to an optimizing problem with at least two requirements: (1.) classified as the counter class, and (2.) still similar to the original data with only necessary modifications. [2] has proposed the following formulation, which becomes the basis of the following research:

$$\mathbf{c} = \arg \min_{\mathbf{c}} \text{trgtloss}(f(\mathbf{c}), y) + d(\mathbf{x}, \mathbf{c}) \quad (1)$$

where  $\mathbf{x}$  is the original data,  $y$  is the target data class,  $\mathbf{c}$  is the generated counterfactual data, and  $f$  is the model so that  $f(\mathbf{c})$  is the new prediction of the model, i.e. the new class. The first part of the formula encourages a different class, while the second part penalizes large distances away from the original data.

To facilitate equation 1 one need to define it in detail, namely how to define "loss" and "distance", which has a significant impact on final results. Another main contribution of [2] is to define distance as  $L_1$  norm divided by MAD (median absolute deviation):

$$\text{dist} = \sum_{k=1}^K \frac{|\mathbf{x} - \mathbf{c}|}{MAD_k} \quad (2)$$

where MAD is defined for every feature  $k$  over the whole points set  $P$ :

$$MAD_k = \text{median}_{i \in P} (|X_{i,k} - \text{median}_{j \in P} (X_{j,k})|) \quad (3)$$

The  $L_1$  norm in equation 2 prone to generate zero entries, which ensures a sparse optimizing result. Normalising the norm is important as well, otherwise big range data would have heavy weights. Here MAD turns out to outperform standard deviation, because it cooperates with  $L_1$  norm better, and generates a even more sparse result.

Sometimes data may contain categorical features (e.g. occupation, gender...), it is counter-intuitive to define "distance" for these features. A simple matching distance is used, where distance one is assigned to different values and zero if the value remains unchanged.

$$\text{dist}_{cat}(\mathbf{c}, \mathbf{x}) = \sum_{k=1}^{K_{cat}} \mathbf{I}(c^k \neq x^k) \quad (4)$$

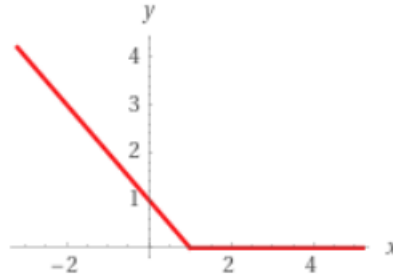
Combining equation 2 and 4 weighted by the number of continuous and categorical features, the most common used distance metric is shown as following:

$$d(\mathbf{x}, \mathbf{c}) = \frac{1}{K_{con}} \sum_{k=1}^{K_{con}} \frac{|\mathbf{x} - \mathbf{c}|}{MAD_k} + \frac{1}{K_{cat}} \sum_{k=1}^{K_{cat}} \mathbf{I}(c^k \neq x^k) \quad (5)$$

The target loss term is defined as  $L_2$  norm,  $(f(\mathbf{c}) - y)^2$ . Without loss of generality, in this article we choose 0 as the original class label, and 1 for the desired label. Note that the prediction of the model  $f(\mathbf{c})$  is a continuous value between 0 and 1. [1] argues that a valid counterfactual suffices if it bypasses the classification threshold of the model (typically 0.5). However, the choice of  $L_2$  norm always prefer the extreme value 1. The underlying issue is that counterfactuals, unless predicted exactly as 1, still receives penalty even if it is classified correctly. If the raw prediction of sigmoid/softmax layer is available ( $\cdot$ ), a better choice is the hinge-loss:

$$\text{hinge\_trgtloss} = \max(0, 1 - \text{logit}(f(\mathbf{c}))) \quad (6)$$

where  $\text{logit}(f(\mathbf{c}))$  is the final activation before entering a sigmoid/softmax output



**Fig. 1.** the hinge loss penalize wrong classification heavily, correct one near the boundary slightly and has no effect above a certain threshold

layer.

**diversity term**

**interpretability term**

## 2.2 by data evolution

## References

1. Mothilal, R.K., Sharma, A., Tan, C.: Explaining machine learning classifiers through diverse counterfactual explanations. In: Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency. pp. 607–617 (2020)
2. Wachter, S., Mittelstadt, B., Russell, C.: Counterfactual explanations without opening the black box: Automated decisions and the gdpr. Harv. JL & Tech. **31**, 841 (2017)