# Learning impartial policies for sequential counterfactual explanations using Deep Reinforcement Learning

Emmanouil Panagiotou[1,2][0000−0001−9134−9387] and Eirini Ntoutsi[2][0000−0001−5729−1003]

[1] Freie Universität Berlin, Department of Mathematics & Computer Science, Berlin, Germany
`emmanouil.panagiotou@fu-berlin.de`
[2] Universität der Bundeswehr München, Faculty for Informatik, Munich, Germany

**Abstract.** In the field of explainable Artificial Intelligence (XAI), sequential counterfactual (SCF) examples are often used to alter the decision of a trained classifier by implementing a sequence of modifications to the input instance. Although certain test-time algorithms aim to optimize for each new instance individually, recently Reinforcement Learning (RL) methods have been proposed that seek to learn policies for discovering SCFs, thereby enhancing scalability. As is typical in RL, the formulation of the RL problem, including the specification of state space, actions, and rewards, can often be ambiguous. In this work, we identify shortcomings in existing methods that can result in policies with undesired properties, such as a bias towards specific actions. We propose to use the output probabilities of the classifier to create a more informative reward, to mitigate this effect.

**Keywords:** Sequential counterfactuals · Reinforcement Learning · Model-agnostic

## 1 Introduction

Predictive Machine Learning (ML) models have been used in various fields to make predictions by generalizing knowledge learned from data. Yet, in recent years, generating explanations for model decisions has become increasingly important, shifting interest to the field of explainable AI (XAI) [7]. A popular category of XAI methods are the so-called "counterfactual explanations" that offer a "what-if" analysis of model decisions by identifying the minimal set of changes or interventions needed to alter the outcome of the model for a given instance, allowing users to understand the factors that influenced the model's decision. Common desiderata for counterfactual explanations, as highlighted by [4], include *proximity* (distance) to the original instance, *sparsity* in terms of changes made, and the *plausibility* of the counterfactual instances w.r.t. the particular problem domain.

Recently, a specific type of counterfactual explanations, called sequential counterfactuals (SCFs) has been proposed [14,16,17,13,15,5] that, rather than making instantaneous or simultaneous modifications to the input, propose a sequence of alterations, each building upon the previous one, until the classifier's decision changes. Such approaches allow for considering the consequences of certain changes to other attributes, for instance, increasing the age of an individual might require an increase in their level of education to maintain consistency or plausibility. SCF generation algorithms look for the optimal order of changes to an input instance, to alter the decision of a model.

In the generation of SCFs, certain methods employ search algorithms that need to initiate optimization for each new instance [14,8], leading to efficiency problems. On the other hand, feed-forward methods [18,3] leverage Reinforcement Learning (RL) to learn a scalable policy for generating SCFs. The problem is formulated as a Markov Decision Process (MDP), where *states* present the original input instance and its alternations, and *actions* correspond to all allowed changes to input features (discrete or continuous). The immediate *reward* for taking an action depends on the output of the black-box model and on the distance of the altered instance to its original state (proximity), to ensure minimal feature changes (action sparsity).

As is typical in RL, the problem formulation plays a crucial role in shaping the policy that the agent learns. In this study, we focus on analyzing certain modeling choices, specifically those related to the action space and the shaping of the reward function, which is influenced by the output of the classifier. We observe that rewarding the RL agent based solely on the sparse binary decision of the classifier can result in problematic policies that only modify features highly correlated with the classifier's output. While these methods demonstrate satisfactory performance in terms of action sparsity, satisfiability, and distance, they prove to be ineffective in practice as they learn to take identical actions for any given input, leading to what we call *feature over-utilization*. In extreme cases, this phenomenon leads to altering a single feature for any given input, thereby reducing the practicality and diversity of generated counterfactual explanations. We present an action entropy metric as a means to quantify the extent of feature over-utilization. Moreover, we propose to overcome this issue by taking continuous actions and implementing a denser reward function when the classifier's output probability is available instead of solely relying on the class labels.

The rest of this paper is organized as follows: related work is discussed in Section 2. In Section 3 we introduce our proposed method. Experimental results are provided in Section 4. Conclusions and future work are discussed in Section 5.

## 2   Related Work

In recent years, many XAI techniques have been proposed to solve the problem of finding SCFs for a pre-trained model. Model-agnostic methods are more generalizable, as they do not require differentiability or information on the internal

architecture of the models, but only use the output for decision-making. The majority are search algorithms that re-optimize for each new given instance, for example by searching in the neighborhood [2,1], or by population-based approaches [14,6]. These algorithms start from a given instance and navigate through trial and error to find a CF, often optimizing for multiple objectives. In our setting, a RL policy has to be learned that can be applied as a feed-forward method after training, on any given instance.

In [18] a Deep Q-Network (DQN) is used to predict discrete feature changes, which are defined as constant steps. Although the reward function incorporates classifier probabilities whenever feasible, its effect on the SCF generation is not discussed. Besides, the use of fixed increments presents issues by assuming all discrete features to be ordinal and restricting changes in continuous features to discrete steps only. This is further detailed in Section 3.

To overcome this problem, [3] proposes to use a parametrized DQN (P-DQN) [19] network, which combines a DQN with a policy gradient method (DDPG) [11], that can be described as the continuous counterpart of a DQN. Doing that, P-DQN can choose among all features using the DQN, and predict a more precise continuous feature change using the DDPG network. Nonetheless, the reward depends only on the class labels, which leads to over-utilization of some features as presented in our experiments.

## 3   Overcoming feature over-utilization through informative reward design

In this section, we address the issue of action over-utilization that arises when employing sparse rewards in RL techniques for generating SCFs. To overcome this issue, we propose the adoption of a denser reward function that employs the probabilities obtained from the black-box classifier.

### 3.1   Limitations of existing approaches

**DQN:** As in [18], a DQN takes discrete steps of ($\pm 0.05$) for continuous features that are in the range of $[-1, 1]$, and ($\pm 1$) for discrete features. This formulation is problematic because i) some discrete features can be nominal, e.g. Occupations don't necessarily have a specific ordering, ii) for continuous features, the step size is discrete and constant, e.g. working hours $\pm 30 min$, and iii) the constant step size of $\pm 0.05$ does not take feature interdependencies into account, for example, it can translate to $\pm 30 min$ for working hours, but $\pm 2500 EUR$ for capital gain, depending on the normalization method. The paper mentions that a dense reward using the black-box probabilities is applied when available. Furthermore, results report a very low action sparsity (close to 1) and distance (0.04), which is desired when generating SCFs, as a lesser number of changes, closer to the original instance, are more interpretable. However, this also means that for all input instances, CFs are found after changing a single feature just once, as the distance is almost the step size of 0.05, leading to feature over-utilization and

scarcity of diverse options. Such an unvarying single-step policy does not provide useful SCF explanations, and potentially other XAI approaches such as global feature importance methods [10,12] are more fitting in this case. Due to these many disadvantages, we do not follow this approach.

**P-DQN:** As proposed in [3] the P-DQN [19] model is used to perform continuous changes sequentially. Taking an action consists of a high-level feature choice $k$ (discrete) and the new value of that feature $u_k$ (continuous). This assumption considers all feature changes $u_k$ to be continuous, which necessitates discretization when handling discrete features. However, P-DQN remains more accurate and adaptable compared to the DQN method, as it enables continuous changes for continuous features without the restriction of discrete steps. Analytically, the action space for P-DQN is defined as:

$$A = \left\{ (k, u_k) \mid u_k \in [-1, 1] \quad k \in [K] \right\} \tag{1}$$

Where $[K] = \{1, 2, ..., K\}$ are the different features. Each time a feature is changed using an action, the timestep $t$ is incremented. A state $s_t = (x_t, b_t)$ is defined as a tuple of some instance $x_t$ concatenated with a binary indicator vector of previously used actions $b_t$. When taking an action $a_t = (k, u_k)$ the new state $s_{t+1} = (x_{t+1}, b_{t+1})$ is created by updating the value of the chosen feature $k$ as $u_k$, i.e. $x_{t+1}[k] = u_k$, and incrementing the indicator vector $b_{t+1}[k] \mathrel{+}= 1$. Finally, the reward function used in this method [3] is defined as:

$$R_t^{bin} = \begin{cases} Pen & \text{if } s_{t+1} = \text{failure,} \\ Pos & \text{if } s_{t+1} = \text{success,} \\ 0 & else. \end{cases} - \beta * \delta(s_{t+1}, s_0) \tag{2}$$

Where $Pen$ is a negative penalty if a failure terminal state is reached, i.e. violating a feature constraint or re-using the same action, $Pos$ is a large positive reward if a success terminal state is reached (counterfactual found), $\delta()$ is a distance measure (the $l_2$ norm in our case), and $\beta$ is a parameter controlling the weight of the distance measure. It is evident that this reward function is sparse since it only has access to the binary class labels and not the probabilities. Therefore we refer to this reward function as $R^{bin}$. In Section 4, we demonstrate that the learned policy using this reward is suboptimal in terms of action variety.

### 3.2  P-DQN & classifier probability rewards

To construct a more dense reward, we assume that the classifier $f(\mathcal{X}) = \mathcal{Y}$ makes a binary decision $\mathcal{Y}$ given an input $\mathcal{X}$ based on a probability measure. We denote the probability of an instance belonging to the target class as $P(\mathcal{X}) = P(f(\mathcal{X}) = \text{target})$. This way we can redefine the reward function for a non-terminal state as follows:

$$R_t^{prob} = \alpha * \left( P(s_{t+1}) - P(s_0) \right) - \beta * \delta(s_{t+1}, s_0) \tag{3}$$

Where $P(s_{t+1}) - P(s_0)$ is the difference in target class probability between the next state and the initial state and $\alpha, \beta$ are parameters controlling the balance between the difference of probability and distance respectively. This is beneficial because the RL agent has immediate feedback on each step regarding the direction taken in relation to the target class. In contrast, when using the aforementioned binary reward ($R^{bin}$), the agent must explore the space without receiving any positive feedback until a successful state is reached (found CF). This additional information on every step results in a more dense reward function, that not only discourages large changes but also guides towards the target class. We denote this reward function with $R^{prob}$ since it incorporates the utilization of target probabilities.

## 4    Experiments & Results

In this section, we compare the P-DQN model's performance when using the binary reward $R^{bin}$ and when employing the target class probabilities reward $R^{prob}$.

### 4.1    Datasets

Motivated by previous approaches we use four tabular datasets in our comparative study. Specifically, the Adult Income, Credit Approval, German Credit, and German Risk [9]. Each dataset has a binary output label and is comprised of mixed features. Some datasets also have immutable features (e.g. Race) that are considered when defining the states but are not taken into account in the action space.

**Table 1.** Characteristics of the four public datasets used.

| Dataset | N. instances | Features | | |
|---|---|---|---|---|
| | | N. of Continuous | N. of Discrete | N. of Immutable |
| Adult Income | 48842 | 4 | 3 | 4 |
| Credit Approval | 690 | 5 | 10 | 0 |
| German Credit | 1000 | 5 | 11 | 4 |
| German Risk | 1000 | 3 | 4 | 2 |

### 4.2    Evaluation measures

Based on previous works on SCF generation [3,18,14], we use the metrics of satisfiability, distance (i.e. proximity), and action sparsity to evaluate the quality of the generated CF examples. Additionally, we measure the target probability of the found CF and define a novel entropy metric to evaluate the diversity of the actions taken by the policy. When an agent has learned to modify only a

specific feature for any given instance, the entropy metric will be low. Conversely, policies with higher entropy exhibit a tendency to modify a diverse range of features depending on the input provided. All metrics are analytically defined hereafter:

- CF satisfiability $= \frac{\text{N. of counterfactuals found}}{\text{N. of test instances}}$, measures the percentage of counterfactuals found, i.e. ended in a success state.
- Action sparsity $=$ N. of features changed to reach a CF (averaged over all test instances), no repeating actions allowed.
- $l_2$ norm distance of the CF to the original state (averaged over all test instances)
- entropy $\mathbb{H} = -\frac{\sum_t \sum_k p_t^k \log_2(p_t^k)}{\log_2(K^2)}$ where $p_t^k = \frac{\text{N. of times feature } k \text{ was changed at step } t}{\text{N. of total feature changes}}$

We calculate the entropy of an agent by assigning *probabilities* $p_t^k$ for every feature $k$ to change at each step $t$. We accomplish this by counting how many times the agent changed a feature at each step and dividing by the total number of changes. We normalize dividing by the maximum possible entropy $H_{max} = \log_2(K^2)$, where $K$ corresponds to the total number of mutable features, which is equal to the maximum timesteps, considering that repeating actions are prohibited.

### 4.3   Experimental setup

Regarding the reward functions of Eq. 2, 3, the distance metric used is $\delta = l_2$ norm, and the hyper-parameters $(\alpha, \beta) = (10, 1)$ are set empirically to give more weight to the change of probability, given that distance values can be larger. We set the large positive reward for a found CF as $Pos = \alpha * thr = 5$ where $thr = 0.5$ is the probability threshold of the classifier, and the large negative penalty for a failure state, $Pen = -2 * Pos = -10$. We use 80% of the data for training across 40.000 environment episodes and 20% for evaluation. Each experiment is repeated five times with different random seeds. Experimental results are reported as the average values and standard deviations over all seed initializations.

### 4.4   Results

We compare the P-DQN model using the binary reward used in [3], denoted as $R^{bin}$, to our proposed reward $R^{prob}$ that we define in Section 3. Results are shown in Table 2.

It is evident that the original approach (P-DQN $R^{bin}$) is mostly performing better in terms of satisfiability, sparsity, and distance. This indicates that it finds SCFs more consistently (satisfiability) for most test instances while taking fewer steps to achieve that (action sparsity) and with close proximity to the original instance. Nevertheless, the entropy of the policy learned by this method is low for all datasets, i.e. under 63%, compared to our approach where the policy has higher entropy over 68%.

It should be emphasized that the objectives of action sparsity and entropy can be contradicting. In particular, a policy that consistently alters a single feature strongly correlated with the output target class may exhibit a desired low sparsity metric but also underperforms with a significantly low entropy value. This occurs for example for the Adult Income dataset where P-DQN $R^{bin}$ performs very well at action sparsity but with a huge loss in entropy.

We can further visualize this issue by plotting the Sankey diagrams for both policies in Figure 1. Sankey diagrams are used to depict the flow of information between sets of values called nodes.
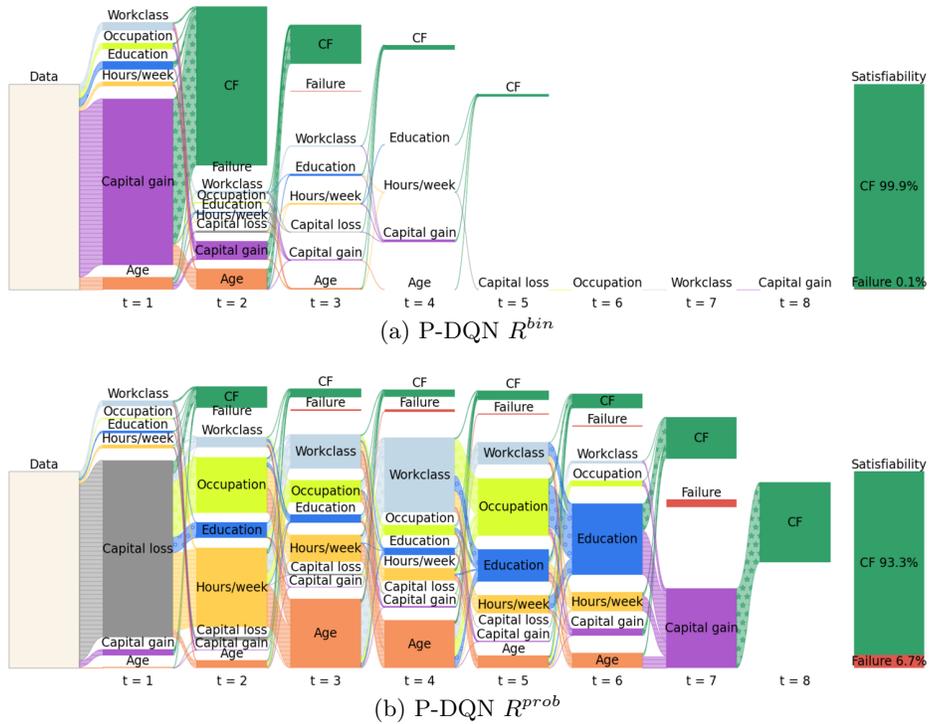


(a) P-DQN $R^{bin}$



(b) P-DQN $R^{prob}$

**Fig. 1.** Sankey diagrams for (a) binary reward and (b) our target probability reward.

In our case, the x-axis corresponds to the distinct timesteps of the agent during the evaluation episodes. On the y-axis, the first node represents all initial instances of the test data. For each test instance and at each step, the trained RL agent can alter a single feature, causing the flow to split into different nodes. This process continues until a CF is found (green), or a failure state is reached (red). The final node represents the percentage of CFs found, i.e. the satisfiability. We can see that our approach finds SCFs by using a much greater variety of actions, due to the denser reward function that incorporates the target probabilities.

**Table 2.** Comparison of both methods on all evaluation metrics and datasets. The results are presented as the average (with standard deviation) over five random initializations for each experiment.

| Dataset | Method | Satisfiability ↑ | Sparsity ↓ | distance $\delta$ ↓ | $P()$ ↑ | entropy $\mathbb{H}$ ↑ |
|---|---|---|---|---|---|---|
| Adult Income | P-DQN $R^{bin}$ | **1.00** (0.00) | **1.18** (0.06) | **1.25** (0.04) | **0.75** (0.02) | 0.29 (0.06) |
| | P-DQN $R^{prob}$ | 0.92 (0.04) | 5.35 (0.09) | 2.69 (0.23) | 0.66 (0.03) | **0.82** (0.03) |
| Credit Approval | P-DQN $R^{bin}$ | 0.80 (0.40) | **1.77** (0.15) | **1.25** (0.65) | 0.55 (0.21) | 0.45 (0.07) |
| | P-DQN $R^{prob}$ | **0.93** (0.02) | 8.16 (2.04) | 3.94 (0.74) | **0.65** (0.05) | **0.84** (0.01) |
| German Credit | P-DQN $R^{bin}$ | 0.97 (0.05) | **1.46** (0.05) | **1.39** (0.04) | **0.70** (0.02) | 0.58 (0.02) |
| | P-DQN $R^{prob}$ | **0.97** (0.01) | 1.71 (0.16) | 1.47 (0.06) | 0.69 (0.02) | **0.68** (0.05) |
| German Risk | P-DQN $R^{bin}$ | **0.97** (0.03) | **1.81** (0.37) | **1.52** (0.11) | **0.63** (0.02) | 0.67 (0.10) |
| | P-DQN $R^{prob}$ | 0.81 (0.33) | 2.48 (0.47) | 1.57 (0.71) | 0.56 (0.09) | **0.74** (0.14) |

## 5   Conclusion and Future Work

In this study, we focus on the problem of learning a policy to generate Sequential Counterfactuals (SCF) using Reinforcement Learning (RL). We identify issues in related methods that lead to policies over-utilizing features, thus becoming ineffective in practice by repeatedly taking identical actions for any input. To address this issue, we propose taking continuous actions and implementing a denser reward function that considers the classifier's output probabilities instead of relying only on class labels. By introducing an action entropy metric, we quantify the extent of feature over-utilization and demonstrate the effectiveness of our approach in mitigating this problem. These findings highlight the importance of carefully designing the problem formulation, reward function, and modeling choices, particularly when generating SCFs.

For future research, we will focus on generating SCFs for mixed features without transforming them into a common data type. We aim to construct RL agents that work in mixed action spaces and explore distance metrics that take mixed feature interdependencies into account.

## 6   Acknowledgements

## References

1. Brughmans, D., Leyman, P., Martens, D.: Nice: an algorithm for nearest instance counterfactual explanations. Data Mining and Knowledge Discovery pp. 1–39 (2023)

2. Cai, Y., Zimek, A., Ntoutsi, E.: Xproax-local explanations for text classification with progressive neighborhood approximation. In: 2021 IEEE 8th International Conference on Data Science and Advanced Analytics (DSAA). pp. 1–10. IEEE (2021)
3. Chen, Z., Silvestri, F., Wang, J., Zhu, H., Ahn, H., Tolomei, G.: Relax: Reinforcement learning agent explainer for arbitrary predictive models. In: Proceedings of the 31st ACM International Conference on Information & Knowledge Management. pp. 252–261 (2022)
4. Dandl, S., Molnar, C., Binder, M., Bischl, B.: Multi-objective counterfactual explanations. In: International Conference on Parallel Problem Solving from Nature. pp. 448–469. Springer (2020)
5. Dandl, S., Molnar, C., Binder, M., Bischl, B.: Multi-objective counterfactual explanations. In: Parallel Problem Solving from Nature–PPSN XVI: 16th International Conference, PPSN 2020, Leiden, The Netherlands, September 5-9, 2020, Proceedings, Part I. pp. 448–469. Springer (2020)
6. Dandl, S., Molnar, C., Binder, M., Bischl, B.: Multi-objective counterfactual explanations. In: International Conference on Parallel Problem Solving from Nature. pp. 448–469. Springer (2020)
7. Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., Pedreschi, D.: A survey of methods for explaining black box models. ACM computing surveys (CSUR) **51**(5), 1–42 (2018)
8. Keane, M.T., Smyth, B.: Good counterfactuals and where to find them: A case-based technique for generating counterfactuals for explainable ai (xai). In: Case-Based Reasoning Research and Development: 28th International Conference, IC-CBR 2020, Salamanca, Spain, June 8–12, 2020, Proceedings 28. pp. 163–178. Springer (2020)
9. Kelly, M., Longjohn, R., Nottingham, K.: The uci machine learning repository. https://archive.ics.uci.edu (2023)
10. Kommiya Mothilal, R., Mahajan, D., Tan, C., Sharma, A.: Towards unifying feature attribution and counterfactual explanations: Different means to the same end. In: Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society. pp. 652–663 (2021)
11. Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., Wierstra, D.: Continuous control with deep reinforcement learning. arXiv preprint arXiv:1509.02971 (2015)
12. Lundberg, S.M., Lee, S.I.: A unified approach to interpreting model predictions. Advances in neural information processing systems **30** (2017)
13. Mahajan, D., Tan, C., Sharma, A.: Preserving causal constraints in counterfactual explanations for machine learning classifiers. arXiv preprint arXiv:1912.03277 (2019)
14. Naumann, P., Ntoutsi, E.: Consequence-aware sequential counterfactual generation. In: Machine Learning and Knowledge Discovery in Databases. Research Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part II 21. pp. 682–698. Springer (2021)
15. Pawelczyk, M., Broelemann, K., Kasneci, G.: Learning model-agnostic counterfactual explanations for tabular data. In: Proceedings of The Web Conference 2020. pp. 3126–3132 (2020)
16. Ramakrishnan, G., Lee, Y.C., Albarghouthi, A.: Synthesizing action sequences for modifying model decisions. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 5462–5469 (2020)

17. Tsirtsis, S., De, A., Rodriguez, M.: Counterfactual explanations in sequential decision making under uncertainty. Advances in Neural Information Processing Systems **34**, 30127–30139 (2021)
18. Verma, S., Hines, K.E., Dickerson, J.P.: Amortized generation of sequential algorithmic recourses for black-box models. In: AAAI Conference on Artificial Intelligence (2021)
19. Xiong, J., Wang, Q., Yang, Z., Sun, P., Han, L., Zheng, Y., Fu, H., Zhang, T., Liu, J., Liu, H.: Parametrized deep q-networks learning: Reinforcement learning with discrete-continuous hybrid action space. arXiv preprint arXiv:1810.06394 (2018)