

0.0.1 What are Counterfactuals?

Counterfactuals have various definitions. However, their semantic meaning refers to “a conditional whose antecedent is false”[1]. A simpler definition from Starr states that counterfactual modality concerns itself with *what is not, but could or would have been*. Both definitions are related to linguistics and philosophy. Within AI and the mathematical framework various formal definitions can be found in the causal inference[2] literature. Here, ^{CITE} describes a counterfactual as “[**Causal inference definition**]”. What binds all of these definitions is the notion of causality within “what if” scenarios.

However, for this paper, we will use the understanding established within the eXplainable AI (XAI) context. Within XAI, counterfactuals act as a prediction which “describes the smallest change to the feature values that changes the prediction to a predefined output” according to Molnar[3, ch.9.3]. Note that XAI mainly concerns itself with the explanation of *models*, which are always subject to inductive biases and therefore, inherently subjective. The idea behind counterfactuals as explanatory tool¹ is simple. We understand the outcome of a model, if we know *what* outcome would occur *if* we changed its input. For instance, let's declare a sequence 1 as *ABCDEF^G*. Then a counterfactual *ABCDEX^Z* would tell us that **F** (probably) caused **G** in sequence 1. As counterfactuals only address explanations of one model result and not the model as a whole, they are called *local* explanations[3, ch.9.3]. According to Molnar *Valid* counterfactuals satisfy **four** criteria[3, ch.9.3]: [**FIX Indentation!**]

Similarity: A counterfactual should be similar to the original instance. If the counterfactual to sequence 1 was *AACDEX^Z* we would already have difficulties to discern whether B or F or both caused G at the end of sequence 1. Hence, we want to be able to easily compare the counterfactual with the original. We can archive this by either minimizing their mutual distance.

Precision: A counterfactual should produce the desired outcome if possible. This characteristic is ingrained in Molnar's definition. However, as the model might not be persuaded to change its prediction, we relax this condition. We say that we want to increase the likelihood of the outcome as much as possible. If the counterfactual *ABCDEX^Z* ends with Z but this sequence is highly unrealistic, we cannot be certain of our conclusion for sequence 1. Therefore, we want the outcome's likelihood

¹There are other explanatory techniques in XAI like *feature importances* but counterfactuals are considered the most human-understandable

to be at least higher under the counterfactual than under the factual instance.

Sparcity: In line with the notion of similarity, we want to change the original instance only minimally. Multiple changes impede the understanding of causal relationships in a sequence.

Feasibility: Each counterfactual should be feasible. As an example, a sequence *ABCDE1G* would not be feasible if numerals are not allowed. **Typically we can use data to ensure this property. However, the open-world assumption makes this criterion quite hard to fulfill. We have to expect longer sequences that have not been encountered, yet.**

All four criteria allow us to assess the validity of each generated counterfactual and thus, help us to define an evaluation metric for each individual counterfactual. However, we also seek to optimise certain qualities on the population level of the counterfactual candidates.

Diversity: We typically desire multiple diverse counterfactuals. One counterfactual might not be enough to understand the causal relationships in a sequence. In the example above, we might have a clue that F causes G, but what if G is not only caused by F? If we are able to find counterfactuals *VBCDEFH* and *ABCDEXZ* but all other configurations lead to G, then we know positions 1 and 6 cause G.

Realism: For a real world application, we still have to evaluate their *reasonability* within the applied domain. This is a characteristic that can only be evaluated by a domain expert.

We refer to both sets of validity criterions as *individual validity* and *population validity*. However, to remain concise, we will use *validity* to refer to the individual criterions only. We will explicitly mention *population validity* if we refer to criterions that concern the population.

0.0.2 The Challenges of Counterfactual Sequence Generation

The current literature surrounding counterfactuals exposes a number of challenges when dealing with counterfactuals.

The most important disadvantage of counterfactuals is the Rashomon Effect[3, ch.9.3]. If all of the counterfactuals are valid, but contradict each other, we have to decide which of the *truths* are worth considering.

This decision reveals the next challenge of evaluation [CITE](#) . Although, the criteria can support us with the decision, it remains an open question *how* to evaluate counterfactuals. Every automated measure comes with implicit assumptions and they cannot guarantee a realistic explanations. We still need domain experts to assess their *validity*.

The generation of counterfactual sequences contribute to both former challenges, due to the combinatorial expansion of the solution space. This problem is common for counterfactual sentence generation and has been addressed within the Natural Language Processing (NLP) [CITE](#) . However, as process mining data not only consist of discrete objects like *words*, but also event and case features, the problem remains a daunting task. So far, little work has gone into the generation of multivariate counterfactual sequences like process instances [CITE](#) .