

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. A prominent figure within the causal inference discipline is Pearl et al., who postulates that a “*kind of statement – an ‘if’ statement in which the ‘if’ portion is untrue or unrealized – is known as a counterfactual*” [4]. What binds all of these definitions is the notion of causality within *what-if* scenarios.

For this paper, we 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, p. 212]. 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 \mathbf{G}* . Then a counterfactual *ABCDEX \mathbf{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, p. 212]. According to Molnar *Valid* counterfactuals satisfy **four** criteria [3, p. 212]:

Similarity: A counterfactual should be similar to the original instance. If the counterfactual to sequence 1 was *AACDEX \mathbf{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.

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. In other words, impossible values are not allowed. As an example, a sequence *ABCDE1 \mathbf{G}* would not be feasible if numerals are not allowed. Typically we can use data to ensure this property. However, the *open-world assumption* impedes

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

this solution. With *open-world*, we mean that processes may change and introduce behaviour that has not been measured before. Especially for long and cyclical sequences, we have to expect previously unseen sequences.

Likelihood: 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 *ABCDE**XZ*** 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 to be at least higher under the counterfactual than under the factual instance.

All four criteria allow us to assess the viability 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 *ABCDE**XZ*** 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 viability criterions as *individual viability* and *population viability*. However, to remain concise, we will use *viability* to refer to the individual criterions only. We will explicitly mention *population viability* if we refer to criterions that concern the population.