

# Alterfactual Explanations - The Relevance of Irrelevance for Explaining AI Systems

Silvan Mertes<sup>1\*</sup>, Christina Karle, Tobias Huber<sup>1</sup>, Katharina Weitz<sup>1</sup>,  
Ruben Schlagowski<sup>1</sup>, Elisabeth André<sup>1</sup>

<sup>1</sup>Chair for Human-Centered Artificial Intelligence, Augsburg University  
{first, second}@uni-a.de

## Abstract

Explanation mechanisms from the field of Counterfactual Thinking are a widely-used paradigm for Explainable Artificial Intelligence (XAI), as they follow a natural way of reasoning that humans are familiar with. However, all common approaches from this field are based on communicating information about features or characteristics that are especially important for an AI's decision. We argue that in order to fully understand a decision, not only knowledge about relevant features is needed, but that the awareness of irrelevant information also highly contributes to the creation of a user's mental model of an AI system. Therefore, we introduce a new way of explaining AI systems. Our approach, which we call Alterfactual Explanations, is based on showing an alternative reality where irrelevant features of an AI's input are altered. By doing so, the user directly sees which characteristics of the input data can change arbitrarily without influencing the AI's decision. We evaluate our approach in an extensive user study, revealing that it is able to significantly contribute to the participants' understanding of an AI. We show that alterfactual explanations are suited to convey an understanding of different aspects of the AI's reasoning than established counterfactual explanation methods.

## 1 Introduction

With the steady advance of Artificial Intelligence (AI), and the resulting introduction of AI-based applications into everyday life, more and more people are being directly confronted with decisions made by AI algorithms [Stone *et al.*, 2016]. As the field of AI advances, so does the need to make such decisions explainable and transparent. The development and evaluation of *Explainable AI* (XAI) methods is important not only to provide end users with explanations that increase acceptance and trust in AI-based methods, but also to empower researchers and developers with insights to improve their algorithms.

\*Contact Author

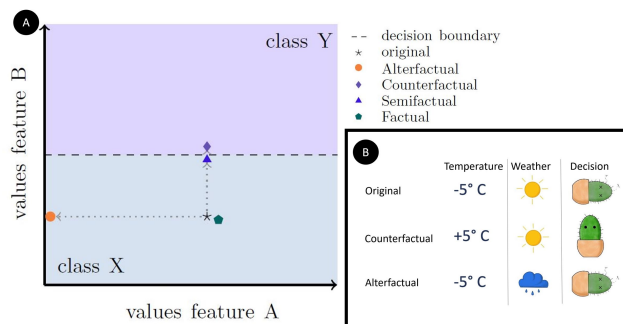


Figure 1: (A) Conceptual comparison of factual, counter-, semi-, and alterfactual explanations. The diagram shows the original input which is to be explained below the decision boundary belonging to class X. A factual explanation could be the nearest neighbor, located anywhere around the original input. A semifactual explanation would be located in minimal distance directly next to the decision boundary, but still below it. A counterfactual explanation would be above it in the region of class Y, but barely so. An alterfactual explanation would move in parallel to the decision boundary, indicating which feature values would not modify the model's decision. Note that this diagram is highly simplified - normally, there are more than two features, the decision boundary is more complex, etc. (B) Examples of a counterfactual and an alterfactual explanation. Input features to a fictional decision system to be explained are *temperature* and *weather*, whereas the former is relevant and the latter is irrelevant to the AI's decision on whether a cactus survives or not.

The need for XAI methods has prompted the research community to develop a bewildering number of different approaches to unraveling the black boxes of AI models. A considerable part of these approaches is based on telling the user of the XAI system in various ways *which* features of the input data are important for a decision (often called *Feature Attribution*) [Arrieta *et al.*, 2020]. Other methods, which are closer to human habits of explanation, are based on the paradigm of *Counterfactual Thinking* [Miller, 2019]. Procedures that follow this guiding principle try answering the question of *What if...?* by showing an alternative reality and the corresponding decision of the AI. Here, in contrast to feature attribution mechanisms, not only the importance of the various features is emphasized. Rather, it is conveyed, even if only indirectly, *why* features are relevant. Prominent exam-

ples of these explanatory mechanisms are *Counterfactual Explanations* and *Semifactual Explanations* [Kenny and Keane, 2020]. Counterfactual explanations show a version of the input data that is altered just enough to change an AI’s decision. By doing so, the user is shown not only *which* features are relevant to the decision, but more importantly, *how* they would need to be changed to result in a different decision of the AI. Semifactual explanations follow a similar principle, but they modify the relevant features of the input data to an extent that the AI’s decision does not change just yet.

All of these methods have in common that they focus on the *important* features. However, awareness of irrelevant features can also contribute substantially to the complete understanding of a decision domain, as knowledge of the important features for the AI does not necessarily imply knowledge of the unimportant ones.

For example, if we want to investigate whether an AI system is subject to some bias regarding its predictions, we often want to know explicitly whether a particular feature is completely irrelevant to a classifier. As a concrete example, consider an AI system that assesses a person’s creditworthiness based on various characteristics, and we want to study that system regarding its fairness. If that system was completely fair, a counterfactual explanation would be of the form: *If your income was higher, you would be creditworthy*. However, this explanation does not exclude the possibility that your skin color also influenced the AI’s decision. It only shows that the income had a high impact on the AI. An explanation confined to the irrelevant features, on the other hand, might say *No matter what your skin color is, the decision would not change*. In this case, direct communication of irrelevant features ascertains, that the system is fair with regards to skin color. Conventional counterfactual thinking explanation paradigms do not provide this information directly.

To address this issue, this paper introduces and evaluates a novel explanatory paradigm. We call explanations that follow this paradigm *Alterfactual Explanations*. This principle is based on showing the user of the XAI system an alternative reality that leads to the exact same decision of the AI, but where only irrelevant features change. All relevant features of the input data, on the other hand, remain the same. As this type of explanation conveys completely different information than common methods, we investigate whether the mental model that users have of the explained AI system is also formed in a different way, or can even be improved. We show that the communication of features unimportant to the decision contributes significantly to the understanding and formation of a mental model of AI systems.

For this purpose, Section 2 gives an overview of related methods and approaches. In Section 3, the principle of alterfactual explanations is explained in further detail. In Section 4 we lay the foundation for further research regarding alterfactual explanations by evaluating the potential of our approach in a user study. In this study, we used an imaginary AI to compare how the most widely used counterfactual thinking paradigm, namely counterfactual explanations, compares to our new concept of alterfactual explanations in terms of mental model creation and explanation satisfaction. As we regard our approach not as a substitution, but rather as supplement

to counterfactual explanations, we further explore the combination of both explanation methods in our user study. Section 5 reports the results of our study. We discuss our findings in Section 6, before we conclude our work and give an outlook into future research in Section 7.

## 2 Related Work

As the approach presented in this paper can be counted to the class of XAI methods that work by inducing counterfactual thinking processes, it is important to gain an understanding of how common methods from this field work. Therefore, this section gives an overview on related concepts that also try to answer the question: *What if...?*

### Factual Explanations

Factual explanations are the traditional way of explaining by example, and often provide a similar instance from the underlying data set (adapted or not) for the input data point that is to be explained [Keane *et al.*, 2021b]. Other approaches do not choose an instance from the dataset, but generate new ones [Guidotti *et al.*, 2019]. The idea behind factual explanations is that similar data instances lead to similar decisions, and the awareness of those similarities leads to a better understanding of the model. Thus, they aim to answer the question *If the input would look more like this Factual Explanation, what would the decision be?*.

### Contrastive Explanations

Contrastive explanations in the context of XAI answer the question of why a decision occurred relative to some other decision that could also have happened for a specific input [Stepin *et al.*, 2021]. The decision that happened is often called *fact*, while *foil* refers to another possible decision it is being contrasted with [Miller, 2019]. Contrastive explanations are usually structured as *Decision X as compared to decision Y occurred because features  $f_1 \dots f_n$  are present and features  $f'_1 \dots f'_n$  are absent* [Verma *et al.*, 2020]. They, therefore, highlight the required minimally present and absent features to achieve the given decision.

### Counterfactual Explanations

Counterfactual explanations are often conflated with contrastive explanations, but do actually state which changes to a given input would be necessary to achieve a contrastive output, i.e. the foil, by providing an example [Stepin *et al.*, 2021]. Counterfactual explanations are a common method humans naturally use when attempting to explain something and answer the question of *Why not ...?* [Miller, 2019; Byrne, 2019]. In XAI, counterfactual explanations are usually used for classification tasks [Verma *et al.*, 2020]. Counterfactual explanations should be minimal, which means they should change as little in the original input as possible to cross the decision boundary of the model between the fact and foil class [Keane *et al.*, 2021b; Miller, 2021]. In the context of generating counterfactuals explanations, the foil class will be referred to as target class. Since counterfactual explanations are a way for a user to understand what they would need to change in order to achieve a more favorable outcome for themselves, many researchers have emphasized that

counterfactual explanations should be actionable and feasible, i.e., should provide a user with an example that is achievable and realistic in real life [Barocas *et al.*, 2020; Ustun *et al.*, 2019]. In certain scenarios, modern approaches for generating counterfactual explanations have shown significant advantages over feature attribution mechanisms in terms of mental model creation and explanation satisfaction [Mertes *et al.*, 2022]. Wachter *et al.* [2017] name multiple advantages of counterfactual explanations, such as being able to detect biases in a model, providing insight without attempting to explain the complicated inner state of the model, and often being efficient to compute.

### Semifactual Explanations

Similar to counterfactual explanations, semifactual explanations are an explanation type humans commonly use. They follow the pattern of *Even if X, still P*, which means that even if the input was changed in a certain way, the prediction of the model would still not change to the foil [McCloy and Byrne, 2002]. In an XAI context, this means that an example, based on the original input, is provided that modifies the input in such a way that moves it toward the decision boundary of the model, but stops just before crossing it [Kenny and Keane, 2020]. Similar to counterfactual explanations, semifactual explanations can be used to guide a user’s future action, possibly in a way to deter them from moving toward the decision boundary [Keane *et al.*, 2021b].

## 3 The Concept of Alterfactual Explanations

The basic idea of alterfactual explanations introduced in this paper is to strengthen the user’s mental model of an AI by showing irrelevant attributes of a predicted instance. Hereby, we understand irrelevance as the property that the corresponding feature, regardless of its value, does not contribute in any way to the decision of the AI model. When looking at models that are making decisions by mapping some sort of input data  $x \in X$  to output data  $y \in Y$ , the so-called *decision boundary* describes the region in  $X$  which contains data points where the corresponding  $y$  that is calculated by the model is ambiguous, i.e., lies just between different instances of  $Y$ . Thus, irrelevant features can be thought of as features that do not contribute to a data point’s distance to the decision boundary.

On the other hand, the information that is carried out by an explanation should be communicated as clearly as possible. As the information that is contained in an alterfactual explanation consists of the *irrelevance* of certain features, it should somehow be emphasized that these features can take *any* possible value. If we would change the respective features only to a small amount, the irrelevance is not clearly demonstrated to the user. Therefore, we argue that an alterfactual explanation should change the affected features to the maximum amount possible. By doing so, we communicate that the feature, *even if it is changed as much as it can change*, still does not influence the decision.

We take those two considerations as the base for the definition of an alterfactual explanation:

Let  $d : X \times X \rightarrow \mathbb{R}$  be a distance metric on the input space  $X$ . An *alterfactual explanation* for a

model  $M$  is an altered version  $a \in X$  of an original input data point  $x \in X$ , that maximizes the distance  $d(x, a)$  whereas the distance to the decision boundary  $B \subset X$  and the prediction of the model do not change:  $d(x, B) = d(a, B)$  and  $M(x) = M(a)$

Thus, the main difference between an alterfactual explanation and a counterfactual or semifactual explanation becomes clear: While the latter methods alter features resulting in a decreased distance to the decision boundary, the former method tries to keep that distance fixed. Further, while counterfactual explanations as well as semifactual explanations try to keep the overall change to the original input minimal [Keane *et al.*, 2021a; Kenny and Keane, 2020], alterfactual explanations do exactly the opposite, which is depicted in Figure 1A. Figure 1B illustrates the difference between counterfactual and alterfactual explanations using a simple example.

## 4 User Study

In order to validate if our approach of focusing only on irrelevant features for explaining an AI system helps users to form correct mental models of the system, we performed an online user study. Prior to the real study, a pilot study (n=14) was conducted to find out whether subjects could cope with the tasks.

### 4.1 Hypotheses

The hypotheses that we addressed in our user study are as follows:

1. Mental Model Creation
  - (a) Alterfactual explanations lead to a more correct mental model of the AI than no explanations.
  - (b) Alterfactual explanations lead to similarly good mental models as counterfactual explanations.
  - (c) The combination of alterfactual and counterfactual explanations outperform both alterfactual as well as counterfactual explanations in terms of mental model creation.
2. Explanation Satisfaction
  - (a) Alterfactual explanations lead to a similarly good explanation satisfaction as counterfactual explanations.
  - (b) The combination of alterfactual and counterfactual explanations outperform both alterfactual as well as counterfactual explanations in terms of explanation satisfaction.

### 4.2 Methodology

In order to test the hypotheses stated above, an online user study was conducted. We used a between-subject design with four conditions:

- **Alterfactual condition.** Participants in that condition were presented with original input features to an AI as well as alterfactual explanations.

- **Counterfactual condition.** Participants in that condition were presented with the original features as well as the counterfactual explanations.
- **Combination condition.** Participants in that condition were presented with the original features as well as both the alterfactual and the counterfactual explanations.
- **No Explanation condition.** Participants in that condition were presented only with the original features. No explanation was shown.

Between-subject was chosen, mainly because we wanted to avoid order effects and mitigate the risk of fatigue. In the study, the participants were presented with an imaginary AI. The participants were told that the AI decides if hypothetical historical documents are forged or not. This specific scenario was chosen as it is not present in most people’s everyday life, ensuring that the mental model of the AI that the participants develop is predominantly induced by the explanations that they are presented with during the study and do not stem from prior knowledge of the domain. The AI gets different inputs to work with. We designed the imaginary AI so that it follows a set of rules (unknown to the participants), where each input feature has a specific relevance to the AI. Those features are as follows:

- **Parchment Color.** The documents can be either of *light*, *medium* or *dark* parchment.
- **Word Count.** A single integer in the range [1, 500].
- **Year of Creation.** The documents were created sometime between 200 BC and 200 AD.

The rules which the fictional AI uses to decide if a document is forged are:

- A document is forged if the word count is equal to or below 50.
- A document is forged if the word count is between 51 and 150 *and* the parchment color is light or medium.
- In all other cases, the document is considered to be authentic

Therefore, one attribute is always relevant (word count), one is relevant only in some cases (parchment color), and one is always irrelevant (year of creation). After answering some questions about their demographic background, the participants were given some general information about the data and AI used in the experiment. They were told that some historical documents had been found, and some of them had already been identified as forgeries. Furthermore, they were told that an AI had been trained to detect forgeries based on a short description of the documents containing the three attributes mentioned above. The three attributes were shown along with their value ranges. An exemplary input to the fictional AI was displayed in a table. Additionally, we explained which explanation type the participant was going to be shown during the study, and how the explanation type works. The participants were provided with example explanations that could be revealed by clicking a button. After using that button, the explanations were shown next to the original input. An example explanation is shown in Figure 2. Following this introduction,

the participants were given two example inputs and corresponding explanations in order to familiarise themselves with the document descriptors and the mechanism to reveal the explanations. After that, each participant was quizzed about the information that was given up to that point. By doing so, we could exclude subjects who did not conscientiously participate in the study. After the quiz, a short training phase followed. In this phase, the participants were shown four exemplary document descriptors. Explanations for the AI’s decisions, as well as the decisions itself were shown as well. The training phase was conducted to give the participants another chance to get comfortable with the explanation type and the domain itself. Subsequently, the study itself started. It was divided into three parts: For assessing the participants’ mental model of the AI, we used (i) a prediction task and (ii) a questionnaire about the AI’s rule set. To assess the participants’ explanation satisfaction, we used (iii) an explanation satisfaction questionnaire.

### Mental Model Creation (i): Prediction Task

The goal of the prediction task was to detect how well the participants could anticipate the classifier’s decisions, which provides a quick window into how well they *understood* the AI [Hoffman *et al.*, 2018]. To this end, eight example inputs with explanations were shown in a random order. Four examples were classified as *forged* by the AI, whereas four examples were predicted as being *authentic*. As proposed by Hoffman *et al.* [2018], the decision of the AI was *not* shown, but had to be predicted by the participants. The idea of such a prediction task is that a good explanation should help to build a correct mental model of the AI, allowing to understand its decision process to an extent that those decisions can be predicted by the user. Additionally to the prediction of the AI’s decision, participants had to choose how confident they were in their prediction on a 7-point Likert scale (0 = not at all confident, 6 = very confident). Furthermore, they had to justify their prediction in a free text form. Participants that were in the *No Explanation* condition did not see any explanations but had to rely on the original input data for their predictions. For every single prediction task, explanations had to be revealed by pressing the *Explain* button. By doing so, we were able to track if the participants really looked at the explanations.

### Mental Model Creation (ii): Understanding Questionnaire

To assess if the participants developed a correct mental model of the AI’s decision process, for each feature (i.e., parchment color, word count, year of creation), they were explicitly asked how much they agreed that it was relevant to the AI’s decision on a 5-point Likert scale (0 = strongly disagree, 4 = strongly agree) after completing all predictions of the Prediction Task. Thus, while the Prediction Task can be seen as implicit measurement of the mental model’s correctness, our Understanding Questionnaire directly measures if participants understood the relevance of different features.

### Explanation Satisfaction

In order to validate hypotheses 2a and 2b, we used the Explanation Satisfaction Scale proposed by Hoffman *et al.* [2018]

which consists of seven items, rated on a 5-point Likert scale (0 = strongly disagree, 5 = strongly agree).

Finally, the participants had the possibility to give free text feedback. The whole study was built using the *oTree* framework by Chen *et al.* [2016].

### 4.3 Participants

113 Participants between 24 and 71 years ( $M = 41.2$ ,  $SD = 10.2$ ) were recruited via Amazon MTurk. 62 of them were male, 48 female, 1 non-binary, and 2 preferred not to answer this question. Only participants with an *MTurk Masters Qualification* were allowed to participate, and subjects that did not pass the quiz were excluded from the study to minimize bias due to unconscious participants. The participants were randomly separated in the four conditions. Subjects of the three explanation conditions that did not look at a single explanation during the whole study were moved to the *No Explanation* condition for evaluation. Participants got paid a base reward of 5.00\$ and another 0.50\$ for each right prediction in the Prediction Task. By communicating that bonus payment before participation, we wanted to further motivate the participants to stay focused on the study. Only 5.3% of the participants had no experience with AI. Most of the participants (86.7%) have heard from AI in the media. In general, 79.7% of the participants were expecting a positive or extremely positive impact of AI systems in the future.

## 5 Results

### 5.1 Mental Model Creation

To investigate the impact of the four different experimental conditions<sup>1</sup> on the (1) understanding and (2) prediction accuracy, we conducted a MANOVA. We found a significant difference, Pillai's Trace = 0.13,  $F(6,218) = 2.52$ ,  $p = .022$ .

The following ANOVA revealed that only the understanding of the participants showed significant differences between the conditions:

- *Understanding*:  $F(3,109) = 3.90$ ,  $p = .011$ .
- *Prediction Accuracy*:  $F(3,110) = 2.63$ ,  $p = .217$ .

As displayed in Figure 3, the post-hoc t-tests showed that the participants' understanding was significantly better in the *Alterfactual* condition compared to all other conditions. The effect size  $d$  is calculated according to Cohen [2013]<sup>2</sup>:

- **alterfactual vs. counterfactual**:  $t(109) = 2.58$ ,  $p = .011$ ,  $d = 0.89$  (large effect).
- **alterfactual vs. combination**:  $t(109) = 3.11$ ,  $p = .002$ ,  $d = 1.24$  (large effect).
- **alterfactual vs. no explanation**:  $t(109) = 2.86$ ,  $p = .005$ ,  $d = 0.82$  (large effect).

The results indicate that alterfactual explanations help participants understand the relevant features more correctly than

<sup>1</sup>(*Alterfactual* condition, *Counterfactual* condition, *Combination* condition, *No Explanation* condition)

<sup>2</sup>Interpretation of the effect size is:  $d < .5$  : small effect;  $d = 0.5-0.8$  : medium effect;  $d > 0.8$  : large effect

in all other conditions. Interestingly, the combination of alterfactual and counterfactual explanations leads to a worse performance and understanding by the participants (see Figure 3).

Therefore, hypothesis 1a holds, because alterfactual explanations outperformed the *No Explanation* condition as well as the *Combination* condition. Hypotheses 1b and 1c have to be rejected because alterfactuals explanations also outperformed counterfactual explanations as well as the combination of both explanation types in the context of mental model creation.

Wondering about the results, especially about the fact that the *No Explanation* condition outperformed the *Combination* condition, we took a closer look, which of the features (i.e., word count, parchment color, year of creation) the participants did or did not understand in each condition. For this, we compared the amount of the correct features between the group, using a MANOVA. We found a significant difference, Pillai's Trace = 0.26,  $F(9,327) = 3.49$ ,  $p < .001$ .

The following ANOVA revealed that only the feature *parchment color* showed significant differences between the conditions:

- *Parchment color*:  $F(3,109) = 10.49$ ,  $p < .001$ .
- *Word count*:  $F(3,109) = 0.03$ ,  $p = .099$ .
- *Creation year*:  $F(3,109) = 1.22$ ,  $p = .305$ .

As displayed in Figure 3, the post-hoc t-tests showed that the participants' correct understanding of the relevance of the parchment color feature were significant better in the *Alterfactual* condition, compared to all other conditions:

- **alterfactual vs. counterfactual**:  $t(109) = 5.21$ ,  $p < .001$ ,  $d = 1.80$  (large effect).
- **alterfactual vs. combination**:  $t(109) = 4.34$ ,  $p < .001$ ,  $d = 1.72$  (large effect).
- **alterfactual vs. no explanation**:  $t(109) = 4.24$ ,  $p < .001$ ,  $d = 1.21$  (large effect).

### 5.2 Explanation Satisfaction

The ANOVA revealed that there were no significant differences between the three explanation conditions,  $F(2,42) = 1.57$ ,  $p = .219$ , indicating that participants felt not specific satisfied by one of the explanation conditions.

Therefore, hypothesis 2b has to be rejected, as the combination of alterfactual and counterfactual explanations does not lead to a higher explanation satisfaction of the participants. Nevertheless, hypothesis 2a holds since the alterfactual explanations do not differ significantly compared to counterfactual explanations.

## 6 Discussion

The results of our user study show novel insights into the explanatory performance of the different XAI approaches.

**Alterfactual Explanations Support Global Understanding of Users.** First of all, although not significantly differing from the other conditions in the Prediction Task, subjects that

	Descriptor	Alterfactual	Counterfactual
Word Count	40	40	51
Parchment Color	dark	light	dark
Creation Year	150 BC	200 AD	150 BC

Figure 2: A sample document descriptor with explanations. In the *Combination* condition, both an alter- and a counterfactual explanation were shown. Subjects in the *Alterfactual* and *Counterfactual* conditions did not see the respective other explanation type. Subjects in the *No Explanation* condition did not see an explanation at all, but only the original document descriptor.

were provided with alterfactual explanations performed significantly better in the Understanding Task than all other participants. This indicates that direct communication of information about irrelevant features does indeed offer benefits. Contrary to our original assumption, the alterfactual explanations outperformed even the more traditional counterfactual explanations. Different from the Prediction Task, the Understanding Task directly surveys the users’ mental models regarding the relevance of the input features. Thus, we argue that alterfactual explanations work better when it comes to the communication of how important different features are for a decision in general, although they do not convey a better understanding of which exact decision will be made when presented with a concrete input sample compared to counterfactual explanations. This suggests that alterfactual explanation could find application in scenarios where a global understanding of the AI system is important. Our investigation of the participants’ feature-specific understanding strengthens this assumption: The alterfactual explanations’ better performance in the Understanding Task mainly stems from users presented with alterfactual explanations having a significantly better understanding of the importance of the *parchment* feature. As that feature was relevant in some cases and irrelevant in others, understanding its relevance highly depends on global understanding of the model. However, future studies have to be conducted to assess the capability of alterfactual explanations to induce a global understanding of an AI’s decision process in a broader scope.

**Too many Explanations can Overstrain Users.** It seems very surprising that the combination of alterfactual and counterfactual explanations performs poorly, although they contain more information than any other condition. We assume that this stems from the fact that more information comes with higher demands on the users’ attentiveness. We argue that the participants in the *Combination* condition were simply overwhelmed by the wealth of information. This finding is in line with *cognitive load* research that emphasized the fact that too much information can overwhelm users [Sweller *et al.*, 1998]. Future research has to find ways to communicate this vast amount of information without overburdening

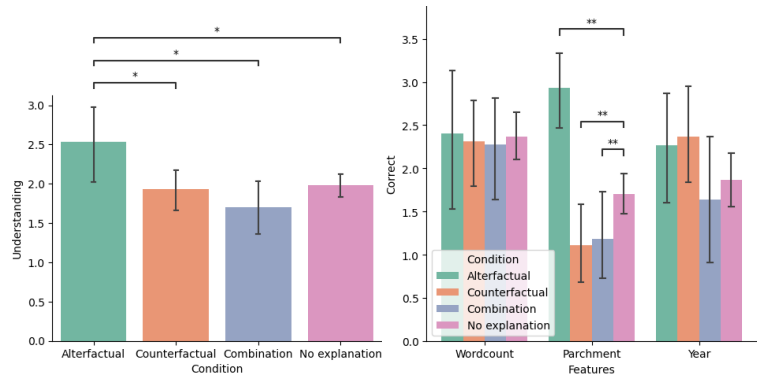


Figure 3: Impact of the four experimental conditions on the understanding of the relevant features of the AI. Alterfactual explanations outperformed all other conditions in helping participants to understand the relevant features of the AI system. Best viewed in color. Error bars represent the 95% CI. \* $p < .05$ , \*\* $p < .001$ .

users.

**Alterfactual Explanations are Equally Satisfying as Counterfactual Explanations.** Furthermore, we found no significant differences regarding Explanation Satisfaction between the three conditions that were presented with some kind of explanation. We argue again that the combination of counterfactual and alterfactual explanations could have overwhelmed the user also in this regards. However, we see that alterfactual explanations lead to similarly good Explanation Satisfaction as the traditional counterfactual explanations, making them a viable approach for real-world XAI scenarios.

## 7 Conclusion & Outlook

In this work, we presented a new XAI paradigm that we call *Alterfactual Explanations*. Our approach is based on only communicating information about features that are irrelevant to an AI’s decision. A user study that we conducted showed that alterfactual explanations show huge potential for the field of XAI. In an Understanding Task measuring the capabilities of users to tell which features of an example input to an AI are important for its decision, alterfactual explanations significantly outperformed the more traditional counterfactual explanations as well as the combination of alterfactual and counterfactual explanations. Surprisingly, combining counterfactual and alterfactual explanations did not result in more correct mental models. We showed that alterfactual explanations lead to a similar good Explanation Satisfaction as counterfactual explanations. As alterfactual and counterfactual explanations convey a different kind of information, future research has to investigate how the combination of the two can leverage the best of both worlds to create even better explanations without overwhelming the user. Further, our study did not include a comparison to Semifactual Explanations. As the concept of those is also based on showing an alternative reality in which the decision does not change, it is likely that users get confused by the differences between Alterfactual and Semifactual explanations. Therefore, the advantages and disadvantages of those two concepts have to be evaluated in further research.

## References

- Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bénéttot, Siham Tabik, Alberto Barabado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*, 58:82–115, 2020.
- Solon Barocas, Andrew D Selbst, and Manish Raghavan. The hidden assumptions behind counterfactual explanations and principal reasons. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 80–89, 2020.
- Ruth MJ Byrne. Counterfactuals in explainable artificial intelligence (xai): Evidence from human reasoning. In *IJCAI*, pages 6276–6282, 2019.
- Daniel L Chen, Martin Schonger, and Chris Wickens. otree—an open-source platform for laboratory, online, and field experiments. *Journal of Behavioral and Experimental Finance*, 9:88–97, 2016.
- Jacob Cohen. *Statistical power analysis for the behavioral sciences*. Academic press, 2013.
- Riccardo Guidotti, Anna Monreale, Fosca Giannotti, Dino Pedreschi, Salvatore Ruggieri, and Franco Turini. Factual and counterfactual explanations for black box decision making. *IEEE Intelligent Systems*, 34(6):14–23, 2019.
- Robert R Hoffman, Shane T Mueller, Gary Klein, and Jordan Litman. Metrics for explainable ai: Challenges and prospects. *arXiv preprint arXiv:1812.04608*, 2018.
- Mark T Keane, Eoin M Kenny, Eoin Delaney, and Barry Smyth. If only we had better counterfactual explanations: Five key deficits to rectify in the evaluation of counterfactual xai techniques. *arXiv preprint arXiv:2103.01035*, 2021.
- Mark T Keane, Eoin M Kenny, Mohammed Temraz, Derek Greene, and Barry Smyth. Twin systems for deepcbr: A menagerie of deep learning and case-based reasoning pairings for explanation and data augmentation. *arXiv preprint arXiv:2104.14461*, 2021.
- Eoin M Kenny and Mark T Keane. On generating plausible counterfactual and semi-factual explanations for deep learning. *arXiv preprint arXiv:2009.06399*, 2020.
- Rachel McCloy and Ruth MJ Byrne. Semifactual “even if” thinking. *Thinking & Reasoning*, 8(1):41–67, 2002.
- Silvan Mertes, Tobias Huber, Katharina Weitz, Alexander Heimerl, and Elisabeth André. Ganterfactual—counterfactual explanations for medical non-experts using generative adversarial learning. *Frontiers in artificial intelligence*, 5, 2022.
- Tim Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artificial intelligence*, 267:1–38, 2019.
- Tim Miller. Contrastive explanation: A structural-model approach. *The Knowledge Engineering Review*, 36, 2021.
- Ilia Stepin, Jose M Alonso, Alejandro Catala, and Martín Pereira-Fariña. A survey of contrastive and counterfactual explanation generation methods for explainable artificial intelligence. *IEEE Access*, 9:11974–12001, 2021.
- Peter Stone, Rodney Brooks, Erik Brynjolfsson, Ryan Calo, Oren Etzioni, Greg Hager, Julia Hirschberg, Shivaram Kalyanakrishnan, Ece Kamar, Sarit Kraus, et al. One hundred year study on artificial intelligence: Report of the 2015-2016 study panel. Technical report, Technical report, Stanford University, 2016.
- John Sweller, Jeroen JG Van Merriënboer, and Fred GWC Paas. Cognitive architecture and instructional design. *Educational psychology review*, 10(3):251–296, 1998.
- Berk Ustun, Alexander Spangher, and Yang Liu. Actionable recourse in linear classification. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 10–19, 2019.
- Sahil Verma, John Dickerson, and Keegan Hines. Counterfactual explanations for machine learning: A review. *arXiv preprint arXiv:2010.10596*, 2020.
- Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening the black box: Automated decisions and the gdpr. *Harv. JL & Tech.*, 31:841, 2017.