

On counterfactuals in general and the structural theory of counterfactuals in particular

Isabelle Mayerhofer



Supervisors:



Abstract

Counterfactual reasoning is a key aspect of human rationality. Being able to reflect about alternative outcomes of the past, present and future had different choices been made, or circumstances been different, is a universal human ability. Classifying counterfactuals however comes with many challenges, as their interpretation tends to be subjective and context-dependent. It has even been argued they are a case of failing of classical logic. Judea Pearl is among the researchers assuming there must be a common underlying mental representation of counterfactual thinking that can be algorithmitized. Pearl offers what he calls a structural approach to the classification of counterfactuals. This paper gives an overview of the main challenges that lie within the classification of counterfactuals before then offering an introduction to Pearl's structural theory of counterfactuals and an insight of possible applications within the realm of artificial intelligence.

1 Introduction and terminology

1. If Oswald did not kill Kennedy, someone else did. [3]
2. If Oswald had not killed Kennedy, someone else would have. [3]

These two sentences on the assassination of the 35th U.S. president John F. Kennedy serve as the classic working example for reasoning about uncertainty in the realm of counterfactuality [19]. Picking up on this line, they will also be used in the following paper. Both sentences are conditionals, the first sentence is in the indicative form and the second in the subjunctive. Focus of this paper is the subjunctive conditional which is generally understood to be the equivalent of a *counterfactual* [1]. The terms *counterfactual* and *subjunctive conditional* will be used interchangeably in the following and sentence 2 will serve as main working example. A counterfactual consists of an *antecedent* and a *consequent*. The antecedent is the part of the sentence following the word 'if'. It specifies

an event contrary to a given real-world observation [4]. In the example the antecedent is 'Oswald had not killed Kennedy' implying that Oswald was not the person who shot the former U.S. president. The antecedent is hence contrary to the given historical fact about Kennedy's death ¹. The consequent is the part of a conditional sentence following the antecedent. In the example it is hence 'someone else would have'. It specifies the event expected to hold in an alternative reality where the antecedent is true [4]. Hence, in alternative possible world where Oswald did indeed not kill Kennedy, this counterfactual's consequent implies that someone else than Oswald was the murderer. While human beings tend to universally accept the truth of sentences in the indicative as given in example 1, they equally as universally do not accept the subjunctive conditional as true [19]. Given the historical fact that Kennedy was assassinated, someone has to be the culprit, be it Oswald or Kennedy which makes for the acceptance of the sentence in the indicative. However, taking the counterfactual literally and assuming it is a fact that Oswald would have been killed by someone else if had not been for Kennedy, is not accepted as a general truth [19]. Instead, the counterfactual 'If Oswald had not killed Kennedy, someone else would have' opens doors to imagining infinite alternative worlds where Oswald was not killed, be it a world in which Kennedy was indeed shot by someone else or one in which Kennedy might even still be alive.

Especially in situations involving planning, humans tend to apply counterfactual reasons and implement lessons learned from the past [1]. Finding explanations for why humans are able to make the distinction between indicative and subjunctive conditionals has been of interest for research in the fields of psychology, philosophy and computer science for many years. Research in this area has however remained widely subjective and open to an author's individual interpretation [19]. Judea Pearl, philosopher and computer scientist, aims to capture the human universal feature of counterfactual reasoning in a way that can be algorithmitized. He does so by providing a structural theory of counterfactuals that enables to teach machines how to master counterfactual reasoning. Pearl's theory offers a computation of counterfactuals that is widely recognized and applied to various areas of artificial intelligence.

In the following paper, in a first part, the key problems with classifying counterfactuals shall be addressed and illustrated. In the second and main part, the scientific work of Judea Pearl on the structural theory of counterfactuals shall be introduced. The paper concludes with giving insight to areas of application of counterfactuals in the field of artificial intelligence.

2 Related work

Scientific work in the realm of counterfactuals is of interest across the domains of especially philosophy, linguistics, the social sciences and artificial intelligence. As of relevance for this paper, some prominent authors among the domains of philosophy and artificial intelligence shall be named: Stalnaker [21] was among those attempting to classify coun-

¹Even though there remains some controversy until the very day.

terfactuals with means of classical logic. Lewis attempted to define counterfactuals within the possible world's theory, an attempt which shall briefly be outlined in section 3.3 [15], [14]. Besides and often in cooperation with Pearl further researchers focus on a structural theory of counterfactuals. These are, among others, Halpern [11], Galles [10] and Briggs [6]. However, Pearl seems to be the most widely cited and most influential, especially when it come to applications within the field of artificial intelligence.

3 Classifying Counterfactuals: Challenges

3.1 Non-truthfunctionality and context-sensitivity

Counterfactuals are said to be a "failing of classical [truth-functional propositional] logic" [1] by many. Propositional logic is the basis for all other logical systems [5]. A proposition is the minimal unit of a propositional formula. Propositions can be assigned truth values. Connecting the propositions via boolean operators, as e.g. negation, conjunction or disjunction, results in a propositional formula. The latter again has a truth value of its own determined by the rules governing the connection of propositions and boolean operators [5]. Herein lies the problem when it comes to classifying counterfactuals: one of the rules determining a proposition's truth value is known as the *material conditional*, also known as *material implication*. Figure 1 gives the rules for connecting antecedent p and consequent q via the boolean operator *implication*. Generally, the most striking quality of the material conditional is that it is only false in cases where p is true and q is false.

p	q	$p \rightarrow q$
T	T	T
T	F	F
F	T	T
F	F	T

Figure 1:
Material im-
plication: truth
table [5]

The problem coming up with the classification of counterfactuals lies, however, elsewhere. As the following two subjunctive conditionals show, the material conditional cannot capture structures like them:

1. If Janis Joplin were alive today, she would drive a Mercedes-Benz. [1]
2. If Janis Joplin were alive today, she would metabolize food. [1]

In both sentences the antecedents and consequents are false: Janis Joplin died in 1970 and hence cannot metabolize food. Further, as an opposer of consumerism, she would likely not drive a Mercedes Bence [1]. Taken as a whole, counterfactual 1 is false and 2 is true. The material conditional does not capture a case like 1, where a false antecedent and false consequent lead to a false propositional formula. The material conditional, which is a central element within classical propositional logic, is hence not fulfilled. While sentence 2 is abiding to the rules of the material conditional, neither the latter nor any other boolean operator are able to capture sentence like 1 and 2 at once, which is why the classification of counterfactuals within classical logic is often said to be impossible [1].

While there is no truth-functional connective for counterfactuals, a possible notation for counterfactuals was given by Stalnaker 1968 [21]: It is $\phi > \psi$, i.e. 'If it had been the case that ϕ , then it would have been the case that ψ '. The operator $>$ is also called the *corner* [21].

Scepticism as to whether counterfactuals at all carry truth or false values led to the assumption of counterfactuals' truth-conditions being context-sensitive, determined by the speaker's individual and subjective perception of the truth of the antecedent [15], [16], [1]. Consider the paper's working example 'If Oswald had not killed Kennedy, someone else would have.' This would perhaps evaluate to true when the speaker takes into account a bad political climate. Pearl addresses this specific example of context-dependency as 'public resentment theory' [19].

3.2 Non-monotonicity of counterfactual antecedents

Classical logic further has the characteristic of being *monotonic*. This means that adding information to a given context does not lead to a change in truth values. In *non-monotonic* reasoning the contrary is the case, it is oftentimes not possible to make the same derivations after adding further information than before having done so [13]. Counterfactuals are a case of non-monotonicity which again makes it difficult to capture them by the means of classical logic. The following three sentences, based on [1], show how the truth value of a counterfactual can flip, when adding more information to the antecedent. Given $I > \neg H$ in sentence 1, I being the antecedent and $\neg H$ the consequent, $\neg H$'s truth value flips to H when adding further information U in 2. The truth value flips again in 3 when even further information ($I \wedge U \wedge T$) is added to the antecedent.

1. If I had gone on a holiday, the project would not have been in jeopardy. $I > \neg H$
2. If I had gone on a holiday and you had too, the project would have been in jeopardy.
 $(I \wedge U) > H$
3. If I had gone on a holiday and you had too but another colleague had worked through the night, the project would not have been in jeopardy. $(I \wedge U \wedge T) > \neg H$

In the context of the *possible worlds theory*, the two semantic analyses of counterfactuals *strict conditional analysis* and *similarity analysis* capture the counterfactuals' non-truthfunctionality and non-monotonicity of their antecedents. This is enabled, put in a nutshell, by focussing on cases when antecedent and counterfactual are true and hence also the consequent is true [1], [16]. Giving a further detailed account is beyond the scope of this paper, however as Judea Pearl constantly contrasts his structural theory of counterfactuals with the possible world's theory, with special focus on the accounts of David Lewis, a brief outline thereof shall be given.

3.3 Possible worlds and Lewis's closest-worlds semantics

A possible world, put simply, is an alternative way a world could be now, in the future, or in the past [1]. In his theory on counterfactuals Lewis takes possible worlds into account [15]. According to him possible worlds are ordered by a measure of similarity, a proposal he transfers to the evaluation of counterfactuals. Lewis denotes a counterfactual as ' B if it were A ', A is hence the antecedent and B the consequent. According to Lewis those worlds most similar to the actual world and in which the antecedents are true, are worlds in which also the consequent is true [1]. Pearl, takes equations as measurement for the evaluation of counterfactuals and not similarity.

4 Structural theory of counterfactuals

"Any theory of counterfactuals, be it of the possible worlds or "truth functional" variety should be deemed incomplete, until it is algorithmitized in sufficient details to allow a robot to correctly evaluate sentences on which humans agree." - Judea Pearl, 2011. [19]

In a nutshell, Pearl's basic thesis is that counterfactuals are generated by symbolic operations on a model. This model represents an agent's knowledge about the world and the agent's belief about functional relationships within it. A counterfactual is accepted when the truth of the antecedent is established, which is done via slightly modifying the model. This procedure shall in the following be explained in as much detail as the scope of this paper allows it, beginning with the model.

4.1 Causal Model

While a model in general is an abstracted representation of an aspect of reality, Pearl implements *causal* models in his structural theory of counterfactuals. As Pearl [18] states, a special kind of language is needed in order to formalize counterfactuals; it has to enable to mirror the distinction between relationships in the world that are static and such that are transitory. Causal models enable this via illustrating causal relationships. While causal relationships are sentences implying an action, like e.g. ' A will be true if B is done' and simple causal relationships as ' A might cause B ', they can also encode counterfactual relationships [18]. An example for the latter would be that ' A would have been different were it not for B ' [18]. A formal definition of the causal model is as follows [17]:

Definition 4.1. (Causal Model)

A causal model is a tripel $M = \langle U, V, F \rangle$

- $U = \{U_1, U_2, \dots, U_n\}$
- $V = \{V_1, V_2, \dots, V_n\}$, determined by $U \cup V$
- $F = \{f_1, f_2, \dots, f_n\}$

U is a set of *exogenous* variables that are determined by external circumstances beyond the model's encoding [19]. An exogenous variable can be paraphrased as an *independent* variable and thought of as *cause* [20]. V are *endogenous* variables, also *dependent* variables and can again be seen as *effect* of the exogenous variables [20]. These variables are determined by $U \cup V$. F is a set of functions where each function f_i maps a value to a variable v_i . If an exogenous variable is instantiated, all of the variables in V are instantiated as well, with a unique value determined by U 's instantiation. The variables are the basic units [12] of the model and are *propositional* variables [19]. In Pearl's structural causal model, these are sentences on an agent's belief about the world.

In Pearl's work causal models take the form of *directed graphs*. While the graph's nodes are the endogenous variables, the relationships between them are represented by the graph's edges. The latter are equipped by an arrow indicating the direction of causality. In case the arrow points from X to Y , X is said to be Y 's *parent* and the latter is hence the *child*. Hence, a directed graph consists of a set of ordered pairs of variables in V [12].

One of the advantages of the causal models over standard propositional logic is that causal models are able to encode information that is not only static but can be determined by external changes. The latter are contained in the model via explicitly representing mechanisms that are altered when it comes to such external changes [18]. The reason for Pearl's theory on counterfactuals being called *structural* lies within its origins tracing back to structural equation models from the 1940s [19]. A structural equation model can be thought of as a causal model which contains structural equations about how every child node is causally dependent on its parent node [12]. The equations are equivalent to functions in F [17]. Hence, the terms *causal model* and *structural* causal model are used interchangeably. The aim of applying the structural causal model is to enable that a counterfactual is accepted by a machine in a way comparable to humans' acceptance of counterfactuals.

While a very basic idea of how counterfactuals are encoded in a structural causal model was already given, a specification shall now follow: The sentence ' Y would be y had X been x in situation $U = u$ ' is what Pearl refers to as the *basic counterfactual entity* [19]. The formal notation is $Y_x(u) = y$. A paraphrase of the basic counterfactual entity is 'If X had been x in situation $U = u$, Y would be y '. This paraphrase shows more clearly that " X had been x " is the antecedent. When the antecedent is true, acceptance of the counterfactual is achieved. Computation of the counterfactual is hence done via "surgically modifying" [19] the model. This 'surgery', staying in line with Pearl's vocabulary, is also referred to as *intervention* [12]. By intervening in the model's structure, a variable's value is set, creating a new causal structure w.r.t. this variable but without changing causal processes relevant for other variables [12]. The surgery is performed via manipulation of the subjunctive 'had X been x ' within the basic counterfactual entity. The goal of the surgery is to establish the truth of the antecedent. The surgery hence consists of instantiating x . Figure 2 serves as illustration of the surgery performed on the structural causal model on the left, performed via a do-operator (shown on the arrow in the middle), leading to a modified submodel M_x . Source of figure 2 is [2].

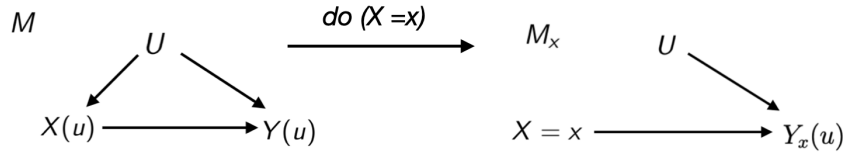


Figure 2: Causal model M and its submodel M_x resulting from action $\text{do}(X = x)$

Ultimately, Pearl gives the formal definition of a counterfactual $Y_x(u)$ in a model M as being "the solution for Y in the surgically modified submodel M_x " [19]. This leads to the definition of a submodel, given as follows in definition 4.2 [17].

Definition 4.2. (Submodel) With M given as causal model in the form of a directed graph, X being a set of variables in V and x a specific realization of X , a submodel M_x of M is the triple:

$$M_x = \langle U, V, F_x \rangle$$

where

$$F_x = \{f_i : V_i \notin X\} \cup \{X = x\}$$

M_x is hence the *effect of action* of $\text{do}(X = x)$ on M [17].

F_x again is the result of deleting all functions f_i from F that are in correspondance to values in X . The deleted functions f_i are replaced by the set $X = x$ [17].

Pearl further generalizes his structural causal model of counterfactuals to a *probabilistic causal model*, which consists of a causal model and a probabilistic function over the set of exogenic variable U , $P(u)$, it is hence a pair $\langle M, P(u) \rangle$ [17]. As of Pearl there are three steps to computing the probability of a counterfactual sentence $P(Y_x = y|e)$, e being any propositional evidence. The steps are *abduction*, *action* and *prediction*. In the *abduction* step, $P(u)$ is updated in order to achieve $P(u|e)$, which is, translated into a temporal metaphor, equivalent to explaining a past event via evidence provided by the present [19]. In the *action* step, variables in X are replaced by $X = x$. This, again in temporal metaphors, corresponds to a time travel in which this action is actually performed. In the *prediction* step the future Y is predicted based on a new version of the past and the new condition $X = x$ and the probability of $Y = y$ via the modified model is computed [19].

Going back in time and bending history corresponds to humans' counterfactual thinking; considering an alternative reality that would have emerged if the past had been different. Different outcomes are determined by decisions as well as external circumstances that lie beyond one's control [4]. In the following it shall be illustrated how Pearl applies his theory to the example on Kennedy's assassination [19].

4.2 Structural Theory of Counterfactuals: an example

Pearl evaluates the counterfactual 'If Oswald had not killed Kennedy, someone else would have' via firstly assuming it to be a fact that Oswald was the person who killed Kennedy. Then Pearl goes on a time travel, as he describes it, and reruns the historical events in an imagined setting where Oswald, contrarily to the historical fact, not killed Kennedy. Pearl formalizes this in three steps, given in figure 4 [19]. In each step a directed graph mirrors the evaluating of the counterfactual. The graphs will be addressed via M . In order to more easily comprehend Pearl's approach it is helpful to keep in mind that the graphs are representations of an agent's belief about the world and its relationships within. Each of the steps shall be explained separately, beginning with the graphs' nodes represented by the following propositional variables:

- OS : Oswald killed Kennedy.
- SE : Someone else killed Kennedy.
- KD : Kennedy is dead.
- M_{OS} : Motivations and enabling conditions for Oswald to be the killer.
- M_{SE} : Motivations and enabling conditions for s.o. else to be the killer.

The edges' arrows give the relationships between the nodes.

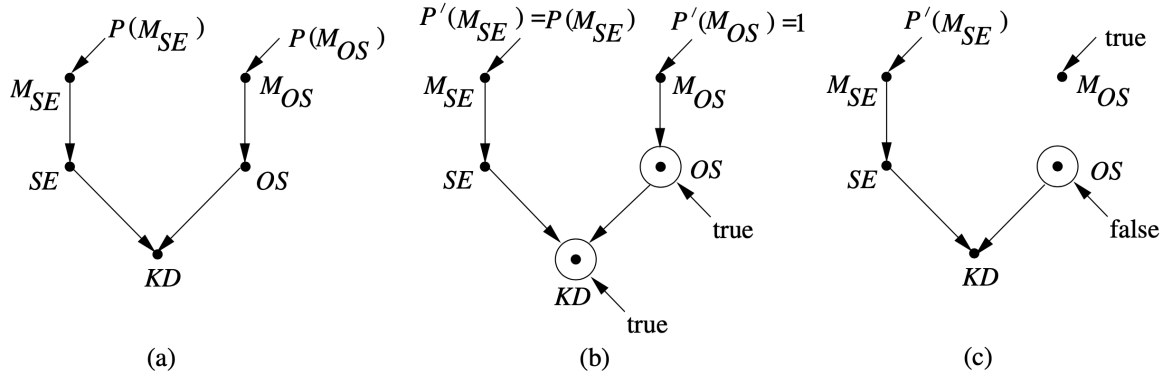


Figure 3: Causal graphs showing steps to evaluating the counterfactual, source: [19]

To begin with graph a): M depicts an agent's general belief state before having the information that Oswald killed Kennedy. The two root variables M_{SE} and M_{OS} have the respective prior probabilities $P(M_{SE})$ and $P(M_{OS})$. As can be seen in a) causal reasoning about Kennedy's death (KD) can be represented by two possible paths. One is that an unidentified person, 'SE', has a motivation to kill Kennedy, which results in the person killing Kennedy and ultimately Kennedy being dead. The other is that specifically

Oswald has a motivation to kill the former president which leads to the actual crime and Kennedy's death.

When an agent learns about the fact that Kennedy was killed by Oswald, the root variables' prior probabilities are updated ²: $P'(M_{SE}) = P(M_{SE}|KD, OS) = P(M_{SE})$ $P'(M_{OS}) = P(M_{OS}|KD, OS) = 1$. This leads to graph b) where the propositional variables OS and KD are now assigned the truth values 'true'. b) is hence a representation of an agent's belief state of a world in which Kennedy's killer is Oswald.

The transition between b) and c) is what corresponds to figure 3 above. c) shows the result of going back in time and assuming Oswald did not kill Kennedy. This is where the surgery comes into play: the truth values of KD and OS and the link between M_{OS} and OS are deleted and OS is set to *false*.

Now, as no longer captured within Pearl's illustration given in figure 3, the posterior probabilities are calculated via $P'(M_{OS})$ and $P'(M_{SE})$ and the newly set fact $OS = false$ which leads to $P'(SE) = P(M_{SE})$.

In words, the antecedent 'Oswald had not killed Kennedy' has hence been satisfied, and it is accepted that there is a possibility that someone else killed Kennedy. This corresponds to the human universal perception of counterfactuals as explained in the introductory section.

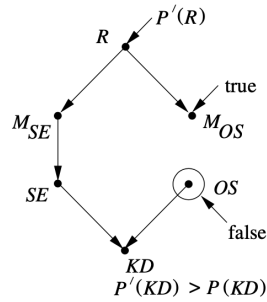


Figure 4:
Resentment
theory [19].

Pearl further takes into account the possible influence of external circumstances on the evaluation of a counterfactual, via referring to the 'public resentment' theory, already briefly addressed in section 3.1. Figure 4, [19], gives the submodel incorporating it. R is the factor of public resentment. A bad political climate where Kennedy was given the blame leads to a higher probability of Kennedy being killed by assassination, if not by Kennedy, then by someone else [19]. This probability is higher than in cases where the resentment theory is not considered, see $P'(KD) > P(KD)$ in figure 4.

Pearl's structural theory of counterfactuals has hence shown how it is possible to model human counterfactual thinking in a way that can be algorithmitized.

5 Summary and areas of application

This paper has provided an introduction to reasoning about uncertainty with regards to counterfactuals. An account on the difficulties that come along with the classification of counterfactuals has been given. The structural theory by Judea Pearl, offering an approach to representing and computing counterfactuals, has been introduced. While Pearl himself states there is no proof whatsoever that his account on counterfactuals is an adequate representation of the mental processes that are active when humans engage in

²The formula in this chapter are not given in figure 3 but are also from [19]

counterfactual reasoning, his account did prove to be successfully applicable to real-world AI applications [17].

An example are counterfactual algorithms for risk minimization in cases of *epistemic* uncertainty. Epistemic uncertainty is the term for uncertainty related to missing data [8]. Counterfactual algorithms can be of use in various online settings, where algorithms make predictions based on log data from user interaction. These predictions can e.g. be used for user-specific add placement. Log data can never cover all possible options a user might take and is hence never complete. Counterfactual algorithms for risk minimization, put simply, hence include possible alternative user input in their training process [22], [9]. This enables the algorithm to make predictions on possible alternative user input.

A further area of application can be found in so-called *explainable* artificial intelligence. Training neural networks on massive amounts of data can provide information to make predictions and decisions e.g. on a person's solvency or future health status. The neural nets' underlying calculations are however at levels of complexity that can be impossible for humans to comprehend. Counterfactuals can be incorporated in explainable AI in order to provide models of the decisions made by the AI algorithm, providing insight why decision A was made instead of decision B [7].

Further examples of applications are policy analysis and retrospective reasoning. As humans use counterfactual reasoning for planning purposes, especially this aspect has proved to be successful when applying to artificial intelligence when it comes to agents and policy-making. Here Pearl's structural theory can find direct application [4].

A main aspect that makes counterfactuals so crucial for reasoning about uncertainty is further that for many settings there is an infinite amount of possible outcomes, given the fact that an infinite amount of alternative choices could have or can be taken. Consider the mentioned example of logging user data, as e.g. obtained from queries used as input for a search machine. The queries in regards to their semantics, amount and timing are deeply idiosyncratic. Hence although making predictions based on user input is possible there always remains uncertainty about different potential outcomes.

Thinking in terms of counterfactuals is also linked to free will and imagination [1]. A person who knows that a choice A will lead to a different outcome as a choice B can imagine the effects of these different choices and decide which choice to make. A famous quote by Albert Einstein is that phantasy is more important than knowledge as there are limits to knowledge while there are none to imagination. If there are no limits to imagination and imagination is linked to free will, it might be put into question if it is really desirable to enable machines to master counterfactual reasoning.

References

- [1] <https://plato.stanford.edu/entries/counterfactuals/>. Last visited July 23, 2021.
- [2] <http://media.nips.cc/Conferences/2013/nips-dec2013-pearl-bareinboim-tutorial-full.pdf>. Pearl, Judea. "Causes and counterfactuals: concepts, principles and tools." NIPS 2013 Tutorial. Last visited July 23, 2021.
- [3] Ernest W Adams. *Logic of conditionals*. Reidel Dordrecht, 1975.
- [4] Alexander Balke and Judea Pearl. Counterfactuals and policy analysis in structural models. *arXiv preprint arXiv:1302.4929*, 2013.
- [5] Mordechai Ben-Ari. *Propositional Logic: Formulas, Models, Tableaux*, pages 7–47. Springer London, London, 2012.
- [6] Rachael Briggs. Interventionist counterfactuals. *Philosophical studies*, 160(1):139–166, 2012.
- [7] Ruth MJ Byrne. Counterfactuals in explainable artificial intelligence (xai): Evidence from human reasoning. In *IJCAI*, pages 6276–6282, 2019.
- [8] William R Clements, Bastien Van Delft, Benoît-Marie Robaglia, Reda Bahi Slaoui, and Sébastien Toth. Estimating risk and uncertainty in deep reinforcement learning. *arXiv preprint arXiv:1905.09638*, 2019.
- [9] Andrew Forney, Judea Pearl, and Elias Bareinboim. Counterfactual data-fusion for online reinforcement learners. In *International Conference on Machine Learning*, pages 1156–1164. PMLR, 2017.
- [10] David Galles and Judea Pearl. An axiomatic characterization of causal counterfactuals. *Foundations of Science*, 3(1):151–182, 1998.
- [11] Joseph Y Halpern. Axiomatizing causal reasoning. *Journal of Artificial Intelligence Research*, 12:317–337, 2000.
- [12] Christopher Hitchcock. Causal Models. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Summer 2020 edition, 2020.
- [13] Paul Krause and Dominic Clark. *Non-monotonic Logic*, pages 143–196. Springer Netherlands, Dordrecht, 1993.
- [14] David Lewis. Causation. *The journal of philosophy*, 70(17):556–567, 1974.
- [15] David Lewis. *Counterfactuals*. John Wiley & Sons, 2013.

- [16] EJ Lowe. The truth about counterfactuals. *The Philosophical Quarterly* (1950-), 45(178):41–59, 1995.
- [17] Judea Pearl. Probabilities of causation: three counterfactual interpretations and their identification. *Synthese*, 121(1):93–149, 1999.
- [18] Judea Pearl. *Causality*. Cambridge university press, 2009.
- [19] Judea Pearl. The algorithmization of counterfactuals. *Annals of Mathematics and Artificial Intelligence*, 61(1):29–39, 2011.
- [20] Neil J Smelser, Paul B Baltes, et al. *International encyclopedia of the social & behavioral sciences*, volume 11. Elsevier Amsterdam, 2001.
- [21] Robert C Stalnaker. A theory of conditionals. In *Ids*, pages 41–55. Springer, 1968.
- [22] Adith Swaminathan and Thorsten Joachims. Counterfactual risk minimization: Learning from logged bandit feedback. In *International Conference on Machine Learning*, pages 814–823. PMLR, 2015.