

An interactive demonstration of counterfactual truth conditions

Bachelor Thesis

Andreas Paul Bruno Lönne

`loenne@campus.tu-berlin.de`

Technische Universität Berlin

discourse Degree program: Bachelor Informatik / Computer Science

Abstract

In this thesis, I address the scarcity of online resources showcasing counterfactual truth conditions in an intuitive and digestible manner. To this end I (i) formulate a semantic satisfiability game for counterfactual sentences; (ii) prove its correctness; (iii) prove that it always halts after a finite number of moves; (iv) develop a browserbased web-application, that makes the semantic game of counterfactuals playable.

1 Introduction

In this thesis I make the attempt to create an application that is able to convey Lewis' counterfactual truth conditions by way of a semantic two-player game and inspire players to learn more about counterfactuals.

Background and literature

To this end i defined a semantic game of counterfactuals and implemented a browserbased demonstration game.

The document is laid out as follows:

First i will explain counterfactuals.

First i will begin by introducing the counterfactual logic I am basing the semantic game of counterfactuals on. Then i will give game-theoretical definitions and formulate two versions of the semantic game. After that ...

2 Counterfactuals

Counterfactuals are statements about what might or would have been the case, if things took place differently than they did. One may think to themselves "If I had not forgotten about my appointment, I would have been punctual". Or one may wonder "If Alexander the Great had not died at the age of 32 and attacked europe, would the Romans have defeated him?". Such *counterfactual thought*—that is the thought of alternate outcomes—is essential for reasoning, deduction and cognitive function. [Byr16] Due to its abundance in human thought, the ability to imagine alternate realities seems trivial to most. But making rigorous statements or claims about them is difficult. This is because communicating a complete

and consistent account of the state of affairs of an alternate reality—similar in complexity to ours—is difficult, if not impossible. One may attempt to circumvent this issue by giving the state of affairs of an alternate reality as a deviation from the state of affairs of reality. But consider this. Take our previous example about Alexander the Great and assume that the imagined alternate reality is identical to our reality, except that Alexander the Great did not die at the age of 32 and attacked Europe. Now suppose we know about reality, that Alexander's troops remained outside of Europe. Then this should also be the case in the alternate reality we attempted to describe. If we are to assume that an army cannot be in two places at the same time and Alexander could not have attacked Europe without his army, then Alexander's troops could not have remained outside of Europe and attacked Europe at the same time. We find, that simply deviating from our own reality in a few concrete ways may produce internally inconsistent alternate realities, which cannot be alternate realities, because they are not ways the world could have been.

To avoid this problem we forgo describing alternate realities altogether. We call a complete and consistent way the world is or could have been, a possible world. And agree that we refer to a possible world, most similar to ours, where our stipulation is true. So our example is to be read as "In a world most similar to our own, where Alexander the Great did not die at the age of 32 and attacked Europe, the Romans would have defeated him". In this way—although we may not know the state of affairs of a possible world—we are able to assign definite truth values to our counterfactual sentence for any possible state of affairs at that world. However we need to note, that this approach assumes the notion of comparative similarity between possible worlds. Which means that given any 3 worlds w, v_1, v_2 , with respect to their overall similarity from the standpoint of w , either

- w is more similar to v_1 , than to v_2 ,
- w is more similar to v_2 , than to v_1 ,
- or w is equally similar to v_1 and v_2 .

While the notion of an aggregate overall similarity between possible worlds appears justified at first glance, it has been subject of contention. [Mor10]

2.1 Lewis' counterfactual operators

With these introductory thoughts out of the way, let us talk in greater detail about the counterfactual operators themselves. Lewis introduces the counterfactual would $\Box \rightarrow$ and counterfactual might $\Diamond \rightarrow$ operators as binary modal operators. [Lew73] When we write $\varphi \Box \rightarrow \psi$, we call the formula φ the antecedent and the formula ψ the consequent. We may informally rewrite one of our former examples as "I did not forget about my appointment $\Box \rightarrow$ I was punctual" and read it the following way.

Read $\varphi \Box \rightarrow \psi$ as "If it were the case that φ , then it would be the case that ψ ", and read $\varphi \Diamond \rightarrow \psi$ as "If it were the case that φ , then it might be the case that ψ ".

Lewis defines the truth conditions of his operators with respect to a system of spheres, that is defined as follows.

Let $\$$ be an assignment to each possible world i of a set $\$i$ of sets of possible worlds. Then $\$$ is called a (centered) system of spheres, and the members of each $\$i$ are called spheres around i , if and only if, for each world i , the following conditions hold.

- (C) $\$i$ is centered on i ; that is, the set $\{i\}$ having i as its only member belongs to $\$i$.
- (1) $\$i$ is nested; that is, whenever S and T belong to $\$i$, either S is included in T or T is included in S .
 - (2) $\$i$ is closed under unions; that is, whenever \mathcal{S} is a subset of $\$i$ and $\bigcup \mathcal{S}$ is the set of all worlds j such that j belongs to some member of \mathcal{S} , $\bigcup \mathcal{S}$ belongs to $\$i$.
 - (3) $\$i$ is closed under (nonempty) intersections; that is, whenever \mathcal{S} is a nonempty subset of $\$i$ and $\bigcap \mathcal{S}$ is the set of all worlds j such that j belongs to every member of \mathcal{S} , $\bigcap \mathcal{S}$ belongs to $\$i$.

Fig. 1: Lewis' (centered) system of spheres

Illustrate example system of spheres and explain $\$$

A world where φ holds is called a φ -world. And his counterfactual truth conditions are:

$\varphi \Diamond \rightarrow \psi$ is true at a world i (according to a system of spheres $\$$) if and only if both

- (1) some φ -world belongs to some sphere S in $\$i$, and
- (2) every sphere S in $\$i$ that contains at least one φ -world contains at least one world where $\varphi \wedge \psi$ holds.

Fig. 2: Lewis' counterfactual might truth conditions

$\varphi \Box \rightarrow \psi$ is true at a world i (according to a system of spheres $\$$) if and only if either

- (1) no φ -world belongs to any sphere S in $\$i$, or
- (2) some sphere S in $\$i$ does contain at least one φ -world, and $\varphi \rightarrow \psi$ holds at every world in S .

Fig. 3: Lewis' counterfactual would truth conditions

3 The semantic game of counterfactuals

two-player game about proving counterfactual sentences

I will give two formulations

- the first making the limit assumption
- the second not doing so

3.1 Counterfactual logic

This section is concerned with defining counterfactual logic. First it provides a brief definition of well-formed counterfactual formulas, then introduces possible world semantics and concludes by defining the truth conditions of counterfactual logic.

3.1.1 Counterfactual formulas

We call $\Phi = \{\varphi, \psi, \dots\}$ the set of all well-formed counterfactual formulas.

Definition 1 (Well-formed counterfactual formula). Given an infinite set of atomic formula symbols $Atoms = \{x, y, \dots\}$ and an alphabet $A = \{\perp, \top, \neg, \vee, \wedge, \Diamond, \Box, \Diamond\rightarrow, \Box\rightarrow\} \cup Atoms$, the structure of every well-formed counterfactual formula is expressed through the following Backus-Naur form.

$$\varphi, \psi ::= \perp \mid \top \mid x \mid \neg\varphi \mid \varphi \vee \psi \mid \varphi \wedge \psi \mid \Diamond\varphi \mid \Box\varphi \mid \varphi \Box\rightarrow \psi \mid \varphi \Diamond\rightarrow \psi \quad (1)$$

3.1.2 Possible worlds semantics

In order to evaluate the non-truth-functional connectives of counterfactual logic, we define a variation of the Kripke structure in [Kri63]. To evaluate Lewis' counterfactual truth conditions [Lew73], we also introduce a notion of similarity between possible worlds to the accessibility relation.

Definition 2 (counterfactual kripke structure). A *counterfactual kripke structure* is an ordered triple (W, \rightsquigarrow, F) , where $W = \{w, v, \dots\}$ is the set of all possible worlds, $\rightsquigarrow: W \times \mathbb{R} \times W$ is the similarity relation, and $F: W \rightarrow 2^{Atoms}$ is an assignment of each world, to a set of atomic formulas.

Let us explain further the members of our counterfactual kripke structure.

The set of all possible worlds W represents all complete and self-consistent ways reality could have been. F in turn, describes the state of affairs at each world. It assigns to each world a set of atomic propositions, that are the case at it. The similarity relations purpose is twofold. It serves as an accessibility relation, restricting accessibility between worlds and also carries information about comparative similarity between worlds. In particular this means for any two worlds w and v , that iff $w \rightsquigarrow^r v$, then v is called accessible from w . $w \rightsquigarrow^r v$ also means, that r is a measure of similarity between w and v , from the standpoint of w . We can suppose another world v' , with $w \rightsquigarrow^{r'} v'$ and compare r and r' . If $r < r'$, we can conclude, that v is more similar to w than v' , from the standpoint of w . Additionally we note, that for our purposes \rightsquigarrow will always contain a tuple $(w, 0, w)$ for each world $w \in W$. This is because we assume analogously to Kripke, that worlds are self-similar and are accessible from themselves. Finally we introduce $W_w = \{w' \mid w \rightsquigarrow^r w'\}$ as a shorthand for the set of all accessible worlds, from a world w .

3.1.3 Truth conditions of counterfactual logic

Given a counterfactual kripke structure $S = (W, \sim, F)$ and a world w ,

$S, w \not\models \perp$.

$S, w \models \top$.

$S, w \models x$ iff $x \in F(w)$.

$S, w \models \neg\varphi$ iff $S, w \not\models \varphi$.

$S, w \models \varphi \vee \psi$ iff ($S, w \models \varphi$ or $S, w \models \psi$).

$S, w \models \varphi \wedge \psi$ iff ($S, w \models \varphi$ and $S, w \models \psi$).

$S, w \models \Diamond\varphi$ iff a world $w' \in W_w$ exists, such that $S, w' \models \varphi$.

$S, w \models \Box\varphi$ iff for every world $w' \in W_w$, it is true that $S, w' \models \varphi$.

$S, w \models \varphi \Diamond\rightarrow \psi$, iff both

- (1) a world $w' \in W_w$ exists, such that $S, w' \models \varphi$, and
- (2) for every world w'' , for which an r'' exists, such that $w \xrightarrow{r''} w''$ and $S, w'' \models \varphi \wedge \neg\psi$ are true, a world w^* and an r^* exist, such that $r^* \leq r''$ and $w \xrightarrow{r^*} w^*$ and $S, w^* \models \varphi \wedge \psi$.

$S, w \models \varphi \Box\rightarrow \psi$, iff either

- (1) no world $w' \in W_w$ exists, such that $S, w' \models \varphi$, or
- (2) a world w' and an r exist, such that $w \xrightarrow{r} w'$ and $S, w' \models \varphi$ and for each world w^* , for which an r^* exists, such that $r^* \leq r$ and $w \xrightarrow{r^*} w^*$, it is true that $S, w^* \models \psi \vee \neg\varphi$.

3.2 The semantic game of counterfactuals

This subsection offers an overview and introduction to the semantic game of counterfactuals. The semantic game of counterfactuals is a sequential two-player satisfiability game, wherein a *defender* d tries to prove a counterfactual formula and an *attacker* a tries to disprove it. The defender begins the game as the *active player*, with a counterfactual formula at a possible world. As the active player makes moves, the counterfactual formula is resolved step by step, until a player resolves a \top and wins or cannot make any move and loses. Throughout the game, moves may also change the active player and possible world.

3.2.1 Game state

Game states form the backbone of the semantic game of counterfactuals, by describing its game positions. Each game state describes the game at a discrete point in time, inbetween moves.

Definition 3 (Game state). A game state is a tuple that can take any one of the forms

- (1) $(\varphi, w)_p$
- (2) $(\varphi, w, e)_p$
- (3) $(\varphi, w, w', r)_p$

where $\varphi \in \Phi$; $w, w' \in W$; $r \in \mathbb{R}$ and $p, e \in \{a, d\}$.

Of these types of game states (1) is the most common. It carries a counterfactual formula, the current world and is indexed with the active player. It is employed to resolve atomic formulas x, y, \dots and the propositional symbols $\top, \neg, \vee, \diamond$; while also serving as an initial and final state type for other resolutions. (2) is a type of game state, that represents a player temporarily becoming active and choosing a state transition. e holds the previous active player, which will resume being active after the next state transition. Lastly game state type (3) is an extension of (1), describing a game state where a sphere of accessibility has been chosen. w' is a *sphere-delimiting world*, lying exactly on the surface of said sphere of accessibility with radius r . We will look at this in further detail in the next subsection. We call the set of all game states G .

3.2.2 Moves

When the active player transforms the current game state into another and progresses the game, we call that a move. Formally we represent this in the following way.

Definition 4 (Move). A move m is a tuple $m \in G \times G$.

The first element of the tuple is the initial game state, while the second one is the resulting game state. We define the set of moves \rightarrow specific to the semantic game of counterfactuals.

Definition 5 (Moves). \rightarrow is a set of moves, such that for each world w , player $p \dots$, the following is true. And for each $p, np \in \{a, d\} \setminus \{p\}$.

$$\begin{aligned}
(x, w)_p &\xrightarrow{[x \in F(w)]} (\top, w)_p & (2) \\
(\neg\varphi, w)_p &\rightarrow (\varphi, w)_{np} & (3) \\
(\varphi \vee \psi, w)_p &\rightarrow (\varphi, w)_p & (4) \\
(\varphi \vee \psi, w)_p &\rightarrow (\psi, w)_p & (5) \\
(\varphi \wedge \psi, w)_p &\rightarrow (\varphi \wedge \psi, w, p)_{np} & (6) \\
(\varphi \wedge \psi, w, e)_p &\rightarrow (\varphi, w)_e & (7) \\
(\varphi \wedge \psi, w, e)_p &\rightarrow (\psi, w)_e & (8) \\
(\Diamond\varphi, w)_p &\xrightarrow{[w \prec w']} (\varphi, w')_p & (9) \\
(\Box\varphi, w)_p &\rightarrow (\Box\varphi, w, p)_{np} & (10) \\
(\Box\varphi, w, e)_p &\xrightarrow{[w \prec^r w']} (\varphi, w')_e & (11) \\
(\varphi \Diamond\rightarrow \psi, w)_p &\xrightarrow{[w \prec^r w']} (\varphi \Diamond\rightarrow \psi, w, w', r)_{np} & (12) \\
(\varphi \Diamond\rightarrow \psi, w, w', r)_p &\rightarrow (\varphi \wedge \psi, w')_{np} & (13) \\
(\varphi \Diamond\rightarrow \psi, w, w', r)_p &\xrightarrow{[w \prec^s w^*, r^* < r]} (\neg\varphi, w^*)_{np} & (14) \\
(\varphi \Box\rightarrow \psi, w)_p &\rightarrow (\varphi \Box\rightarrow \psi, w, p)_{np} & (15) \\
(\varphi \Box\rightarrow \psi, w, e)_p &\xrightarrow{[w \prec^r w']} (\varphi \Box\rightarrow \psi, w, w', r)_e & (16) \\
(\varphi \Box\rightarrow \psi, w, w', r)_p &\rightarrow (\neg\varphi \vee \psi, w')_p & (17) \\
(\varphi \Box\rightarrow \psi, w, w', r)_p &\xrightarrow{[w \prec^s w^*, r^* < r]} (\varphi, w^*)_p & (18)
\end{aligned}$$

Fig. 4: Moves of the semantic game of counterfactuals

3.2.3 The transition system

The moves, that can be played in a given game state, are determined by a transition system L . It contains transitions between game states, which we call moves. To arrive at our definition for L , let us first define game states and transitions themselves.

Definition 6 (Transition system). A transition system is a tuple (G, \rightarrow) , where G is the set of all game states and \rightarrow is a transition relation.

3.3 Limit Assumption

- explain limit assumption
- justify limit assumption (finite graph satisfies it anyways, nicer interactive presentation)
- alternative semantic game rules

3.4 Alternative game rules

$$(x, w)_p \xrightarrow{[x \in F(w)]} (\top, w)_p \quad (19)$$

$$(\neg\varphi, w)_p \rightarrow (\varphi, w)_{np} \quad (20)$$

$$(\varphi \vee \psi, w)_p \rightarrow (\varphi, w)_p \quad (21)$$

$$(\varphi \vee \psi, w)_p \rightarrow (\psi, w)_p \quad (22)$$

$$(\varphi \wedge \psi, w)_p \rightarrow (\varphi \wedge \psi, w, p)_{np} \quad (23)$$

$$(\varphi \wedge \psi, w, e)_p \rightarrow (\varphi, w)_e \quad (24)$$

$$(\varphi \wedge \psi, w, e)_p \rightarrow (\psi, w)_e \quad (25)$$

$$(\Diamond\varphi, w)_p \xrightarrow{[w \rightsquigarrow w']} (\varphi, w')_p \quad (26)$$

$$(\Box\varphi, w)_p \rightarrow (\Box\varphi, w, p)_{np} \quad (27)$$

$$(\Box\varphi, w, e)_p \xrightarrow{[w \rightsquigarrow w']} (\varphi, w')_e \quad (28)$$

//TODO: Check truth conditions

$$(\varphi \Diamond\rightarrow \psi, w)_p \rightarrow (\varphi \Diamond\rightarrow \psi, w, p)_{np} \quad (29)$$

$$(\varphi \Diamond\rightarrow \psi, w, e)_p \rightarrow (vac, \varphi \Diamond\rightarrow \psi, w)_e \quad (30)$$

$$(vac, \varphi \Diamond\rightarrow \psi, w)_p \xrightarrow{[w \rightsquigarrow w']} (\varphi, w')_{np} \quad (31)$$

$$(\varphi \Diamond\rightarrow \psi, w, e)_p \xrightarrow{[w \rightsquigarrow w']} (\varphi \Diamond\rightarrow \psi, w, w', r)_e \quad (32)$$

$$(\varphi \Diamond\rightarrow \psi, w, w', r)_p \rightarrow (\neg\varphi \vee \psi, w')_p \quad (33)$$

$$(\varphi \Diamond\rightarrow \psi, w, w', r)_p \xrightarrow{[w \rightsquigarrow w^*, r^* \leq r]} (\varphi \wedge \psi, w^*)_p \quad (34)$$

//TODO: Accurate truth conditions

$$(\varphi \Box\rightarrow \psi, w)_p \rightarrow (\Box\neg\varphi, w)_p \quad (35)$$

$$(\varphi \Box\rightarrow \psi, w)_p \xrightarrow{[w \rightsquigarrow w']} (\varphi \Box\rightarrow \psi, w, w', r)_{np} \quad (36)$$

$$(\varphi \Box\rightarrow \psi, w, w', r)_p \rightarrow (\varphi, w')_{np} \quad (37)$$

$$(\varphi \Box\rightarrow \psi, w, w', r)_p \xrightarrow{[w \rightsquigarrow w^*, r^* \leq r]} (\neg\varphi \vee \psi, w^*)_{np} \quad (38)$$

Fig. 5: alternate rules of the semantic game of counterfactuals

3.5 Playing the game

Definition 7 (Semantic game of counterfactuals). The semantic game of counterfactuals is an ordered tuple (I, E, L) , where $I \in G$ is the initial game state; $E \subseteq G$ is the set of ending game states; L is a transition system for game states.

TODO: Add ending game states - Rules

(stuck player loses)

- game is like a conversation between players
- show something through their interplay
- show truth of counterfactual formula
- Game produces a play
- game states (formula, current world, player)
- game outcomes (winning/losing)
- moves
- play (sequence of moves)
- winning play
- game tuple (state, moves)
- subgame
- strategy (mapping from game states to moves (positional game, thus state-based is okay, otherwise partial mapping from partial plays to moves))

3.6 Correctness proof

3.7 Termination proof

4 The educational game

4.1 Correspondance between semantic game and computer game

4.2 Artificial intelligence

Algorithm 1 Adjusted MiniMax Algorithm

Input: s (game state)

```
function MINIMAX( $s$ )                                ▶ Calculate best move
   $Q \leftarrow [s]$                                 ▶ Setup queue, expanded node and score lists
   $E \leftarrow []$ 
   $SC \leftarrow []$ 
  while  $Q$  not empty do
     $current \leftarrow \text{POP}(Q)$ 
     $e \leftarrow \text{false}$ 
     $states \leftarrow \text{NEXTSTATES}(s)$ 
    Let  $n$ 
    if  $c$  not in  $E$  then
      PUSH( $E$ ,  $current$ )
       $e \leftarrow \text{true}$ 
      for each state in  $states$  do
        if state not in  $E$  then
          if state in  $Q$  then
            REMOVEEACH( $Q$ , state)
          end if
          PUSH( $Q$ , state)
        end if
      end for
    end if
    if  $e = \text{false}$  then // TODO:
    end if
  end while
end function
```

NEXT calculates the next game states and returns them

5 Results

5.1 Player feedback

References

[Byr16] Ruth MJ Byrne. Counterfactual thought. *Annual review of psychology*, 67(1):135–157, 2016.

- [Kri63] Saul A. Kripke. Semantical considerations on modal logic. *Acta Philosophica Fennica*, 16:83–94, 1963.
- [Lew73] David K. Lewis. *Counterfactuals*. Basil Blackwell, Oxford, 1973.
- [Mor10] Michael Morreau. It simply does not add up: Trouble with overall similarity. *The journal of philosophy*, 107(9):469–490, 2010.