

SUBJUNCTIVE CONDITIONALS AND STRUCTURAL EQUATIONS

02.27.12

WILLIAM STARR

Phil 6710 / Ling 6634 – Spring 2012

1 Review

Goodman’s Puzzle A counterfactual $A > B$ is true in w just in case B follows all the facts about w which are co-tenable with A .

- ‘Co-tenability’ isn’t logical consistency. Counterfactually assuming A definitely requires giving up $\neg A$, but it may require giving up things that followed by law from $\neg A$. Also, assuming A may require adding things that follow from it and the laws.
- But what are laws? Further, doesn’t this require a total theory of how to revise one’s beliefs when you learn that A is false.

The Lewis-Stalnaker Strategy A counterfactual $A > B$ is true in w just in case B is true in all of the A -worlds most **similar** to w . We can say enough about similarity to capture the logic of subjunctive conditionals without proposing a solution to Goodman’s problem.

- Similarity is formally modeled with a *selection function* f :
 - f takes a world w and a proposition p , and returns the set of p -worlds most similar to w , e.g. $f(w_0, p) = \{w_1, w_3, w_9\}$.
- By placing minimal constraints on which worlds count as most similar, Lewis and Stalnaker are able to get a plausible (though debated) logic for subjunctive conditionals.
- Some of these constraints are motivated by the intuitive concept of similarity itself
- Some are imposed just to get the right logic

Lewis-Stalnaker Theory (Stalnaker 1968; Stalnaker & Thomason 1970; Lewis 1973)

- $\phi > \psi$ is true at w just in case all of the ϕ -worlds most **similar** to w are ψ -worlds
 - Most similar according to the selection function f
 - f takes a proposition p and a world w and returns the p -worlds most similar to w
- $\llbracket \phi > \psi \rrbracket^f = \{w \mid f(w, \llbracket \phi \rrbracket^f) \subseteq \llbracket \psi \rrbracket^f\}$

(Momentarily making ‘Limit Assumption’: there are most similar worlds)

- Logic of $>$ is determined by constraints on f (where $p, q \subseteq W$ and $w \in W$):
 - (a) $f(w, p) \subseteq p$ **success**
 - (b) $f(w, p) = \{w\}$, if $w \in p$ **strong centering**
 - (c) $f(w, p) \subseteq q$ & $f(w, q) \subseteq p \implies f(w, p) = f(w, q)$ **uniformity**
 - (d) $f(w, p)$ contains *at most* one world **uniqueness**

Stalnaker Constraints (a)-(d)

Limited Lewis Constraints (a)-(c)

- ‘Limited Lewis’: Lewis if he had accepted the Limit Assumption

1.1 The Challenges

- My editorial on context sensitivity:
 - A proposition is conveyed $w/A > B$ only when a particular f is filled in
 - ▶ There are literally thousands of f ’s meeting the formal constraints
 - Stalnaker and Lewis regard this filling in of f as a standard case of context sensitivity
 - ▶ Context sensitivity is the process of using information mutually available in the utterance context to interpret an utterance
 - ▶ More precisely: interpretations of utterances that would fail to communicate anything without using information mutually available in the context
 - We must all therefore have the means for getting f ’s from available information
 - Lewis and Stalnaker: ordinary concept of similarity serves this role
 - **Issue 1:** filling in f seems to require solving Goodman’s puzzles
 - ▶ Maybe there is so much vagueness in f , this is an instance of the more general puzzle of how communication with vague language works
 - **Issue 2:** the facts that determine similarity make uttering the subjunctive redundant
 - ▶ Rather than taking f as fixed and communicating something on the basis of it, subjunctives seem to inform us about f
 - ▶ This doesn’t mesh with the standard model of context sensitivity
 - **Issue 3:** various examples demonstrate that it is not our intuitive concept of similarity that is put to use in determining the truth value of subjunctive conditionals

Email: will.starr@cornell.edu.

URL: <http://williamstarr.net>.

- Fine (1975: 452):
 - (1) If Nixon had pressed the button there would have been a nuclear holocaust
 $B > H$
 - Plausibly, (1) is true, or can at least be supposed to be.
 - Suppose further that there never will be a nuclear holocaust.
 - For every $B \wedge H$ -world, there will be a closer $B \wedge \neg H$ -world
 - ▶ In this world a small change prevents the holocaust, such as a malfunction in the electrical detonation system
 - The idea: surely a world where Nixon presses the button and a malfunction prevents nuclear holocaust is more like our own than one where there is a nuclear holocaust!
 - So it would seem that the Lewis-Stalnaker theory predicts (1) to be false!
- Tichý (1976: 271):
 - (2) a. Invariably, if it is raining, Jones wears his hat
 - b. If it is not raining, Jones wears his hat at random
 - c. Today, it is raining and so Jones is wearing his hat
 - d. But, even if it hadn't been raining, Jones would have been wearing his hat
 - Given (2a-c), (2d) seems clearly incoherent/false/bad
 - Why is this a counterexample to the Stalnaker-Lewis theory?
 - ▶ In the actual world $w_{@}$, Jones is wearing his hat.
 - ▶ So in the non-raining-worlds most similar to $w_{@}$, Jones is wearing his hat
 - ▶ But then Stalnaker-Lewis predict that (2d) is true!
- Lewis (1979) articulated a system of weights for similarity which was, among other things, supposed to address these counterexamples:
 - (1) It is of the first importance to avoid big, widespread, diverse violations of law. ('big miracles')
 - (2) It is of the second importance to maximize the spatio-temporal region throughout which perfect match of particular fact prevails.
 - Maximize the time periods over which the worlds match exactly in matters of fact
 - (3) It is of the third importance to avoid even small, localized, simple violations of law. ('little miracles')
 - (4) It is of little or no importance to secure approximate similarity of particular fact, even in matters that concern us greatly.
- This system works for the two counterexamples by not counting particular matters of fact towards similarity

- As Lewis (1979: 466-7) acknowledged, this gives up the idea that intuitive similarity is involved in evaluating subjunctive conditionals
 - But this is a problem!
- Further the system of weights doesn't cover a simple variant of Tichý's case:
 - (3) a. Before Jones opens the curtain to see what the weather is like, he flips a coin
 - b. If it's not raining and the coin comes up heads, he wears his hat
 - c. If it's not raining and the coin comes up tails, he doesn't wear his hat
 - d. Invariably, if it is raining he wears his hat
 - e. Today, the coin came up heads and it is raining, so Jones is wearing his hat
 - f. But, even if it hadn't been raining, Jones would have been wearing his hat

(Veltman 2005: 164)

 - Here, when you counterfactually suppose that it isn't raining, to *do* keep fixed the particular fact that Jones is wearing his hat
 - So particular matters of fact are sometimes kept fixed!
- A simpler case that illustrates the same phenomenon:
 - (4) [Suppose there is a circuit such that the light is on exactly when both switches are in the same position (up or not up). At the moment switch one is down, switch two is up and the lamp is out.] If switch one had been up, the lamp would have been on. (Lifschitz)
 - We keep fixed the fact that switch two is up!
- Lewis (1979: 472) notes that in cases like these, things come out differently and 'would like to know why'
- Diagnosis by Veltman (2005: 164):

Similarity of particular fact is important, but only for facts that do not **depend** on other facts. Facts stand and fall together. In making a counterfactual assumption, we are prepared to give up everything that depends on something that we must give up to maintain consistency. But, we want to keep in as many independent facts as we can.

 - In Tichý's case, when we counterfactually suppose that it isn't raining we don't keep fixed the fact that he is wearing his hat because his wearing his hat depended on the fact that it is raining
 - In the variant, the outcome of the coin flip was independent of raining, so we keep it fixed when we suppose that it wasn't raining
 - ▶ But then it follows, because of (3b), that Jones would be wearing his hat
- Veltman (2005) proposes an analysis of counterfactuals that is sensitive to **dependence**
 - A version of premise semantics, couched in a situation semantics
 - Situation semantics: worlds are made up of facts
 - Spells out what it is for one fact to depend on another

- But, as Schulz (2007:101) notes, the Veltman (2005) theory does not make the right prediction for (4)
 - Evidencing that the Veltman (2005) analysis of dependence is not quite right
- Schulz (2011, 2007: §5.6) and Hiddleston (2005) develop analyses of subjunctive conditionals inspired by the analysis given in Pearl (1995, 2000: Ch.7) based on the causal models pioneered by Pearl (1993) and Spirtes *et al.* (1993)
 - It offers great promise in correctly capturing our cases while also capturing the logic of subjunctive conditionals
- Causal models provide a powerful tool for reasoning about causal dependence
 - But not all subjunctive conditionals concern causal dependencies
 - So we'll want to revisit the question of whether the analysis can extend to other kinds of dependence
 - ▶ Schaffer "Structural Equations for Metaphysical Dependence" MS. Rutgers, certainly thinks it can
- We'll pursue this strategy too, but note where it is taking us
 - Back to directly answering Goodman's puzzles, rather than attempting an 'abstract analysis'!
- The structural equations approach is also quite nice because there's a lot of explicit work on how it can be used to formulate a theory of explanation: Woodward (2003), Halpern & Pearl (2005a,b)

2 Causal Models and Structural Equations

- Why does a computer scientist care about causation?
 - When a robot with some knowledge about its environment performs an action that changes that environment, how does the robot decide to update its knowledge to reflect this change? (McCarthy & Hayes 1969)
 - Striking the match tends to lead to fire, but not in some circumstances, like when it's wet, or when the striking surface is too smooth, or when there's no oxygen, or when the robot is confusing a twig for a match and so on.
 - This is the *frame problem* (Shanahan 2009)
 - Pearl's example:
 - ▶ *Input:*
 - (1) Suitcases open iff both locks are open
 - (2) The right lock is open
 - (3) The suitcase is closed

▶ *Query:*

- (1) What would happen if we open the left lock?
 - (2) The right lock might get closed
- A lot like our cases!
 - In standard propositional logics, atomic sentences are assigned independent truth-values
 - A valuation v , is a simple function from atomics to truth-values
 - Pearl's causal models give up this assumption
 - The truth-value of an atomic, say D , can **depend** on the truth-value of others, say A and B
 - These dependencies are *functional*
 - If all that D depends on are A and B , then D ' truth-value is uniquely determined by A and B
 - Actually, Pearl's most important contributions have come in making these kind of models probabilistic
 - But we're not going to go into that
 - Following Simon (1953), Pearl thinks of dependencies as invariant 'causal mechanisms'
 - Like Simon, he represents them out as a set of 'structural equations'
 - D 's truth depends on both A and C being true, $D := A \wedge C$
 - Or D 's truth depends on one of them being true, $D := A \vee C$
 - You can picture the models underlying these equations as directed graphs

2.1 Pearl's Prisoner Example

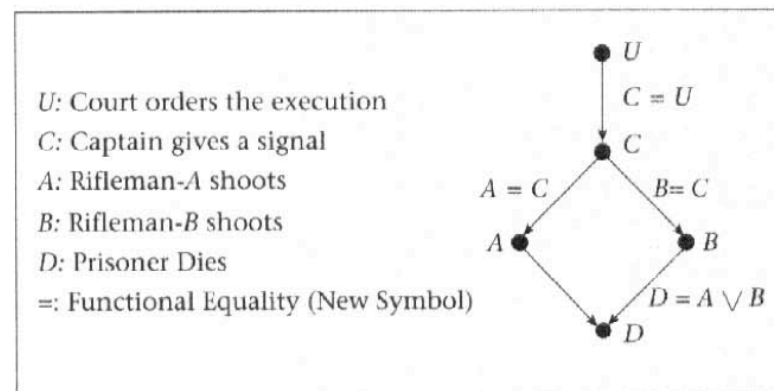


Figure 4. Causal Models at Work (the Impatient Firing Squad).

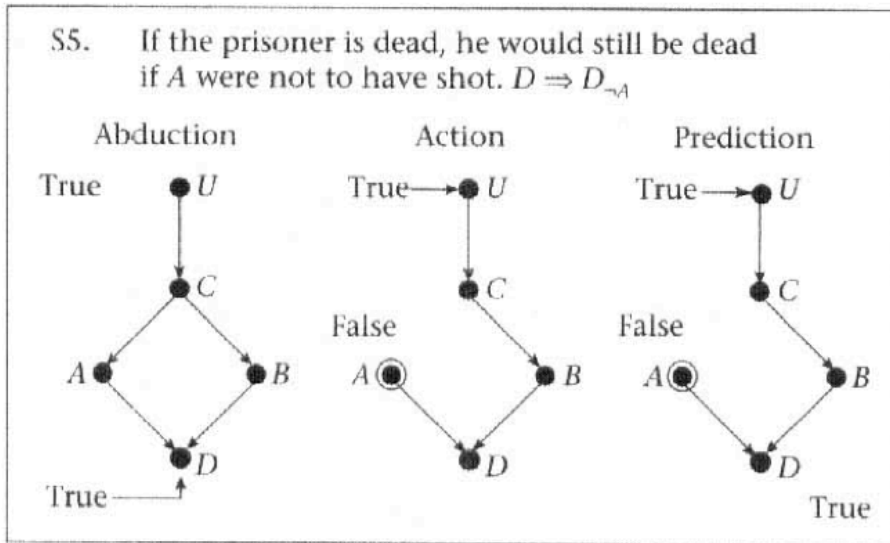


Figure 7. Three Steps to Computing Counterfactuals.

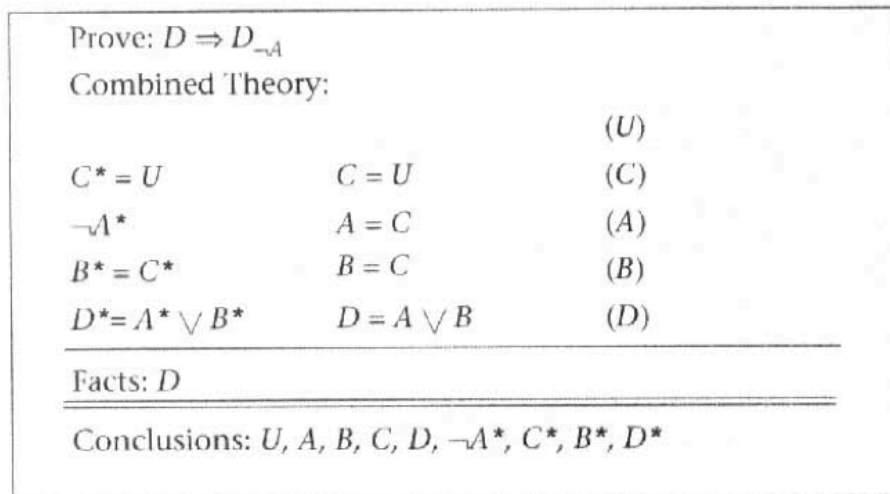


Figure 8. Symbolic Evaluation of Counterfactuals.

- The account of the light example works pretty much the same!
 - Instead of keeping fixed that B shot, we keep fixed that switch two is up

3 A Logic of Structural Equations: the structure of possible worlds

- Pearl's approach is promising, but it has pretty severe limitations
 - Exogenous atomic sentences (ones with no arrows going into them) cannot be manipulated by 'intervention': if the order hadn't been given, prisoner would be alive!
 - What happens with logically complex antecedents?
 - If switch one alone controlled the light, then the light would be on
 - Not all subjunctive conditionals are about causal connections
 - If this cup were red, it wouldn't be blue
- Starting point: classical possible worlds are valuations
 - Fix the truth values of each atomic sentence
 - Picture each atomic as a dot, which is black if false, white if true.

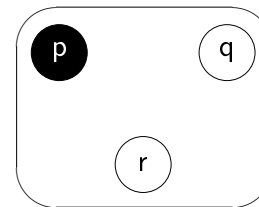


Fig. 1. Classical possible world w

$$\begin{aligned} w(p) &= 0 \\ w(q) &= 1 \\ w(r) &= 1 \end{aligned}$$

Fig. 2. System of equations for w

- Now depart from the classical picture (Starr 2012):
 - The dependencies between facts endow worlds with a structure

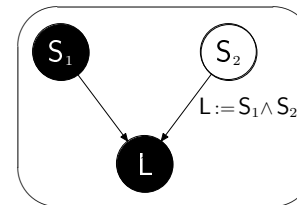


Fig. 3. A structured possible world w

$$\begin{aligned} w(S_1) &= 0 & (5) \\ w(S_2) &= 1 & (6) \\ w(L) &= w(S_1) \cdot w(S_2) & (7) \\ &= 0 \end{aligned}$$

Fig. 4. Equations for w

- \neg, \wedge and \vee all have arithmetic counterparts operating on 1 and 0

\neg	\wedge	\vee
$1 - x$	$x \cdot y$	$(x + y) - (x \cdot y)$

- To evaluate the counterfactual $S_1 > L$, create world w_{S_1}
 - Step 1: intervention
 - ▶ Eliminate old assignment for S_1 , line (8)
 - ▶ Make S_1 1, line (9)
 - Step 2: projection
 - ▶ Apply equation (11) to solve for L
 - ▶ New result: $w_{S_1}(L) = 1!$

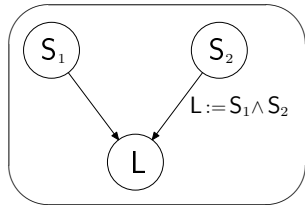


Fig. 5. The New World w_{S_1}

$$\begin{aligned}
 w(S_1) &\neq 1 & (8) \\
 w_{S_1}(S_1) &= 1 & (9) \\
 w_{S_1}(S_2) &= w(S_2) = 1 & (10) \\
 w_{S_1}(L) &= w(L) = w_{S_1}(S_1) \cdot w_{S_1}(S_2) & (11) \\
 &= 1 & (12)
 \end{aligned}$$

Fig. 6. Equations for w_{S_1}

- To see how this works better, consider a slightly modified scenario:
 - Switch 1 turns on a servo that controls switch 2: $S_2 := S_1$
 - Switch 2 turns on the light: $L := S_2$
 - Currently, switch 1 is up, so 2 is up and the light is on
- (13) If switch 2 were up, the light would be off
- This comes out true, because after setting S_2 to 1, the equation connecting it will make L come out as 1 too
 - Lesson: you only keep fixed facts which are not determined by facts you are counterfactually giving up
- These interventions are quite like Lewis' miracles
 - We go to a world exactly like w except w 's mechanisms have been broken to allow a particular fact to hold
- So this account captures Fine's Nixon case nicely.

3.1 What are Worlds?

- The idea spelled out for w :
 - Each independent atomic mapping: $\{\langle S_2, 1 \rangle, \langle S_1, 0 \rangle\}$
 - ▶ **Situation:** a part of the atomic mapping of a world (any set of pairs of atomics and truth-values)

- ▶ A partial function from atomics to truth-values
- Pair L with a dependency function d that determines the truth of L from a pairing of S_1, S_2 with truth-values:

S_1	S_2	L	
1	1	1	◦ d maps L and $\{\langle S_1, 1 \rangle, \langle S_2, 1 \rangle\}$ to 1
1	0	0	◦ d maps L and $\{\langle S_1, 1 \rangle, \langle S_2, 0 \rangle\}$ to 0
1	0	0	◦ d maps L and $\{\langle S_1, 0 \rangle, \langle S_2, 1 \rangle\}$ to 0
1	1	0	◦ d maps L and $\{\langle S_1, 0 \rangle, \langle S_2, 0 \rangle\}$ to 0

- We can then define an abuse of notation:
 - $w(L) = t \iff \exists s \subseteq w \ \& \ d(L, s) = t$
 - So $w(L) = 0$
- So, a world is a function from some independent atomics to truth-values, together with a dependency function for each dependent atomic
- $w = \{\langle S_2, 1 \rangle, \langle S_1, 0 \rangle, d\}$
 - $d = \{ \langle \langle L, \{\langle S_1, 1 \rangle, \langle S_2, 1 \rangle\} \rangle, 1 \rangle, \langle \langle L, \{\langle S_1, 1 \rangle, \langle S_2, 0 \rangle\} \rangle, 0 \rangle, \langle \langle L, \{\langle S_1, 0 \rangle, \langle S_2, 1 \rangle\} \rangle, 0 \rangle, \langle \langle L, \{\langle S_1, 0 \rangle, \langle S_2, 0 \rangle\} \rangle, 0 \rangle \}$
- Dependencies are functions from pairs of atomics and situations to truth-values

3.2 Extending the Analysis

- We'd like to define the notation w_ϕ for any non-subjunctive ϕ
- Let's first try just with compounds of atomics:
 - w_A is the world exactly like w except that it assigns A to 1
 - $w_{\neg A}$ is the world exactly like w except that it assigns A to 0
 - $w_{A \wedge B}$ is the world exactly like w except that it assigns A and B to 1
 - $w_{A \vee B}$ is the world exactly like w except that it assigns ?????????
- An idea: allow w_ϕ to be a set of worlds
- Second try:
 - w_A are the worlds exactly like w except that they assign A to 1
 - $w_{\neg A}$ are the worlds exactly like w except that they assign A to 0
 - $w_{A \wedge B}$ are the worlds exactly like w except that they assign A and B to 1
 - $w_{A \vee B}$ are the worlds exactly like w except that they either assign A to 1 or B to 1

- Attempt to generalize:
 - $w_{-\phi}$ are the worlds exactly like w except that they assign....????
 - ▶ Currently, w doesn't assign non-atomics like ϕ truth-values
- One slightly ugly solution:
 - For the atomics A_1, \dots, A_n in ϕ :
 - ▶ Eliminate $\langle A_1, x_1 \rangle, \dots, \langle A_n, x_n \rangle$ from w , call it w^*
 - ▶ Take all $\langle A_1, t_1 \rangle, \dots, \langle A_n, t_n \rangle$, which when added to w^* create a world which makes ϕ true
 - ▶ Each such world makes it into the set w_ϕ

Dependency Semantics for Subjunctives

- $\llbracket \phi > \psi \rrbracket = \{w \mid w_\phi \subseteq \llbracket \psi \rrbracket\}$
- $\phi > \psi$ is true iff either ψ is independent of ϕ and true, or else ϕ is sufficient for bringing about ψ when holding fixed all those facts that do not depend upon ϕ .¹

3.3 Remaining Issues

- Consider a case where one switch controls a light:
 - The switch is up and the light on
 - But if the light had been off, then if you had flipped the switch up, the light would have come on
- What about backtrackers:
 - If the light had been on, the two switches would have to have been up
- Comparison with Briggs (to appear)?

A The Logic of Structural Equations

- SYNTAX:
 - **Atomic formulas:** A, B, C, \dots , **Connectives:** $\neg, \wedge, \vee, \rightarrow, >$
 - ▶ $>$ can only embed on the right, e.g. $A > (B > C)$
- SEMANTICS:
 - **Situations:** for some $A \subseteq \mathcal{At}$, $s : A \mapsto \{1, 0\}$
 - ▶ S is the set of all situations
 - ▶ $s = \{\langle A_0, t_0 \rangle, \dots, \langle A_n, t_n \rangle, \dots\}$

¹ This intuitive paraphrase is from Cumming (2009:1).

- ▶ A potentially partial function from atomics to truth-values
- **Dependencies:** for some $A \in \mathcal{At}$, $d : \{A\} \mapsto (S \mapsto \{1, 0\})$
 - ▶ D is the set of all dependencies
 - ▶ d takes an atomic and then a situation and delivers the atomic's truth-value in that situation
 - ▶ $d = \{\langle A, \langle \{\langle B_0, t_0 \rangle, \dots, \langle B_n, t_n \rangle\}, t \rangle \rangle\}$
- **Worlds:** $w = s \cup \{d_0, \dots, d_n\}$, for some $s \in S$ iff:
 - ▶ $s \neq \emptyset$
 - ▷ At least one fact is independently settled
 - ▶ $\text{dom } s \cup \text{dom } d_0 \cup \dots \cup \text{dom } d_n = \mathcal{At}$
 - ▷ Every atomic gets assigned a truth-value
 - ▶ $\text{dom } s \cap \text{dom } d_0 \cap \dots \cap \text{dom } d_n = \emptyset$
 - ▷ No atomic gets assigned truth-values by both s and d_0, \dots, d_n
 - ▶ If $d \in w$, $A \in \text{dom } d$, $d' \in w$, $B \in \text{dom } d'$ and $A \in \text{dom } (\text{dom } d(A))$ then $B \notin \text{dom } s$
 - ▷ Dependence is acyclical: if A depends, in part, on B then B cannot depend at all on A.
 - ▶ $\text{dom } d_0 \cap \dots \cap \text{dom } d_n = \emptyset$
 - ▷ Dependencies are total: no two dependencies settle the same facts

References

- BRIGGS, R (to appear). 'Interventionist Counterfactuals.' *Philosophical Studies*. URL [http://www.rachaelbriggs.net/Rachael_Briggs/CV_\(with_online_papers\)_files/CM_Counterfac_6.pdf](http://www.rachaelbriggs.net/Rachael_Briggs/CV_(with_online_papers)_files/CM_Counterfac_6.pdf).
- CUMMING, S (2009). 'On What Counterfactuals Depend.' Ms. UCLA.
- FINE, K (1975). 'Review of Lewis' *Counterfactuals*.' *Mind*, **84**: 451–8.
- HALPERN, J & PEARL, J (2005a). 'Causes and Explanations: A Structural-Model Approach. Part I: Causes.' *British Journal for Philosophy of Science*, **56**.
- HALPERN, J & PEARL, J (2005b). 'Causes and Explanations: A Structural-Model Approach. Part II: Explanations.' *British Journal for Philosophy of Science*, **56**: 889–911.
- HIDDLESTON, E (2005). 'A Causal Theory of Conditionals.' *Noûs*, **39**(4): 632–657.
- LEWIS, DK (1973). *Counterfactuals*. Cambridge, Massachusetts: Harvard University Press.
- LEWIS, DK (1979). 'Counterfactual Dependence and Time's Arrow.' *Noûs*, **13**: 455–476.
- MCCARTHY, J & HAYES, PJ (1969). 'Some Philosophical Problems from the Standpoint of Artificial Intelligence.' In B MELTZER & D MICHIE (eds.), *Machine Intelligence 4*, 463–502. Edinburgh: Edinburgh University Press.

-
- PEARL, J (1993). ‘Graphical Models, Causality and Intervention.’ *Statistical Science*, **8(3)**: 266–273.
- PEARL, J (1995). ‘Causation, Action, and Counterfactuals.’ In A GAMMERMAN (ed.), *Computational Learning and Probabilistic Learning*. New York: John Wiley and Sons.
- PEARL, J (2000). *Causality: Models, Reasoning, and Inference*. Cambridge, England: Cambridge University Press.
- SCHULZ, K (2007). *Minimal Models in Semantics and Pragmatics: free choice, exhaustivity, and conditionals*. Ph.D. thesis, University of Amsterdam: Institute for Logic, Language and Information, Amsterdam. URL <http://www.illc.uva.nl/Publications/Dissertations/DS-2007-04.text.pdf>.
- SCHULZ, K (2011). “‘If you’d wiggled A, then B would’ve changed”.’ *Synthese*, **179**: 239–251. 10.1007/s11229-010-9780-9, URL <http://dx.doi.org/10.1007/s11229-010-9780-9>.
- SHANAHAN, M (2009). ‘The Frame Problem.’ In EN ZALTA (ed.), *The Stanford Encyclopedia of Philosophy*, winter 2009 edn. URL <http://plato.stanford.edu/archives/win2009/entries/frame-problem/>.
- SIMON, HA (1953). ‘Causal Ordering and Identifiability.’ In WC HOOD & TC KOOPMANS (eds.), *Studies in Econometric Method*, 49–74. New York: Wiley.
- SPIRITES, P, GLYMOUR, C & SCHEINES, R (1993). *Causation, Prediction, and Search*. Berlin: Springer-Verlag.
- STALNAKER, RC (1968). ‘A Theory of Conditionals.’ In N RESCHER (ed.), *Studies in Logical Theory*, 98–112. Oxford: Basil Blackwell Publishers.
- STALNAKER, RC & THOMASON, RH (1970). ‘A Semantic Analysis of Conditional Logic.’ *Theoria*, **36**: 23–42.
- STARR, WB (2012). ‘The Structure of Possible Worlds.’ Ms. Cornell University.
- TICHÝ, P (1976). ‘A Counterexample to the Stalnaker-Lewis Analysis of Counterfactuals.’ *Philosophical Studies*, **29**: 271–273.
- VELTMAN, F (2005). ‘Making Counterfactual Assumptions.’ *Journal of Semantics*, **22**: 159–180. URL <http://staff.science.uva.nl/~veltman/papers/FVeltman-mca.pdf>.
- WOODWARD, J (2003). *Making Things Happen: A Theory of Causal Explanation*. Oxford: Oxford University Press.