

Counterfactuals in Critical Thinking with Application to Morality

Luís Moniz Pereira¹ and Ari Saptawijaya²

¹ NOVA Lab. for Computer Science and Informatics, Universidade Nova de Lisboa, Portugal
lmp@fct.unl.pt

² Faculty of Computer Science, Universitas Indonesia, Indonesia.
saptawijaya@cs.ui.ac.id

Abstract. Counterfactuals are conjectures about what would have happened, had an alternative event occurred. It provides lessons for the future by virtue of contemplating alternatives; it permits thought debugging; it supports a justification why different alternatives would have been worse or not better. Typical expressions are: “If only I were taller . . .”, “I could have been a winner . . .”, “I would have passed, were it not for . . .”, “Even if . . . the same would follow”. Counterfactuals have been well studied in Linguistics, Philosophy, Physics, Ethics, Psychology, Anthropology, and Computation, but not much within Critical Thinking. The purpose of this study is to illustrate counterfactual thinking, through logic program abduction and updating, and inspired by Pearl’s structural theory of counterfactuals, with an original application to morality, a common concern for critical thinking. In summary, we show counterfactual reasoning to be quite useful for critical thinking, namely about moral issues.

Keywords: Critical Thinking, Counterfactual Reasoning, Abduction, Morality.

1 Counterfactual Reasoning

Counterfactual literally means contrary to the facts. Counterfactual reasoning involves thoughts on what could have happened, had some matter – action, outcome, etc. – been different in the past. Counterfactual thinking covers everyday experiences, like regret: “If only I had told her I love her!”, “I should have studied harder”; or guilt responsibility, blame, causation: “If only I had said something sooner, then I could have prevented the accident”. The general form is: “If the *Antecedent* had been true, then the *Consequent* would have been true”.

Counterfactuals have been well studied in Linguistics, Philosophy, Physics, Ethics, Psychology, Anthropology, and Computation [2–4, 6, 8, 9, 14, 16, 17, 19, 22, 23, 29, 35], but oddly not much within Critical Thinking. However, people often think how things that matter to them might have turned out differently [15]. Researchers from psychology have asked: “Why do people have such a strong tendency to generate counterfactuals?”; “What functions does counterfactual thinking serve?”; “What are the determinants of counterfactual thinking?”; “What are its adaptive and psychological consequences?”. Human’s ability for the mental time travel required by counterfactual thinking relies on

their use of episodic memory. Without this memory humans would be unable to form a stable concept of self along time, consider what might counterfactually have happened instead, and hence human cultures would not have been able to consider evolution paths that took into account past alternatives.

In this paper, counterfactual reasoning is enacted using a three-step logic evaluation procedure [28], inspired by the structure-based approach of [22], viz.

1. **Abduction:** to explain past circumstances in the presence of observed evidence, i.e., use the given evidence to determine the unchanging external background circumstances;
2. **Action:** to adjust the logical causal model to comply with the antecedent of the counterfactual, i.e., to impose the truth of the antecedent's hypotheses by means of a forced intervention on the model; and
3. **Prediction:** to predict if the counterfactual's consequent deductively follows, subsequently to steps 1 and 2, i.e., to compute the truth-value of the consequent in the modified intervened model.

The approach is realised by means of logic program abduction and updating. Abduction chooses from available hypotheses (the set A of abducibles) the exogenous variables that constitute the situation's background – i.e., those abducibles or their negations, that best explain the observed given evidence O . An abduced explanation, E , is a subset of A that finds the specific values for exogenous variables, which lend an explanatory support to all currently observed evidence. Note that the abduction procedure guarantees the abduced explanation to be consistent, i.e., disallows both abducible a and its negation a^* to hold in explanation E .³ Subsequent to abduction, updating modifies those rules to be updated and fixes the initially abduced exogenous background context of the counterfactual statement. That is, updates the knowledge base with some preferred explanation to the current observations and, additionally, the updating also permits causal intervention on the causal knowledge model, namely by means of hypothetical updates to the rules, achieved via reserved predicate *make* (illustrated in examples below), so as to render the knowledge base consistently compliant with the antecedent of the counterfactual.

Consider an example [3]: *Lightning hits a forest, and a devastating forest fire breaks out. The forest was dry, after a long hot summer.* Let us add more causes for forest fire, i.e., there are two possible alternative causes: *storm* – presuming the lightning – or *barbecue*. The model of this example consists in a set of abducibles

$$A = \{storm, barbecue, storm^*, barbecue^*\}$$

and program P :

$$\begin{aligned} fire &\leftarrow barbecue, dry_leaves. \\ fire &\leftarrow barbecue^*, lightning, dry_leaves. \\ lightning &\leftarrow storm. \\ dry_leaves. \end{aligned}$$

³ In the sequel, starred atoms stand for their negations.

Take counterfactual statement: *If only there had not been lightning, then the forest fire would not have occurred.*

Step 1: Given the observation $O = \{lightning, fire\}$, abduce its explanations E (a subset of A). Note that the observations assure us that both the antecedent and the consequent of the counterfactual were *factually* false. Two possible explanations for O : $E_1 = \{storm, barbecue^*\}$ and $E_2 = \{storm, barbecue\}$. Say E_1 is preferred for consideration, on a criterion of simplicity. Then fix its abduced background context for the counterfactual: i.e., update program P with E_1 .

Step 2: Update program P , via an automated transformation, to get a new program T by adding:

```
make(lightning*)           % Intervention: If there had not been lightning...
lightning ← make(lightning). % Note that lightning or otherwise are
lightning* ← make(lightning*). % now available only by intervention.
```

where $make/1$ represents an explicit intervention on the model, by forcing its argument true. It corresponds to Pearl's $do/1$ operator. Because the intervention must be made explicit, an implicit default representation would not be adequate. e.g. $make(lightning^*)$ is explicitly imposing on the model that there was no lightning.

Plus, for irrelevancy and consistency, the transformation deletes:

$$lightning \leftarrow storm.$$

Step 3: Verify if the conclusion “the forest fire would not have occurred” is true. Since *fire* is not provable, ‘*not fire*’ holds in the semantics of T for explanation $E_1 = \{storm, barbecue^*\}$ with intervention $make(lightning^*)$. The counterfactual is valid.

2 Counterfactuals in Morality

Typically, people think critically about what they should or should not have done when they examine decisions in moral situations. It is therefore natural for them to engage in counterfactual thoughts of alternatives in such settings. Counterfactual thinking has been investigated in the context of moral reasoning, notably by psychology experimental studies [3], e.g., to understand the kind of critical counterfactual alternatives people tend to think of in contemplating moral behaviours, and the influence of counterfactual thoughts in moral judgment [15, 30].

Morality and normality judgments typically correlate. Normality infuses morality with causation and blame judgments. The importance of control, namely the possibility of intervention, is highlighted in theories of blame that presume someone responsible only if they had some control over the outcome [37]. The explicit controlled interventions expressed by the counterfactual premises enable to interfere with normality, and hence with blame and cause judgments.

As argued by [6], the function of counterfactual thinking is not just limited to the evaluation process, but occurs also in the reflection one. Through evaluation, counterfactuals help correct wrong behaviour in the past, thus guiding future moral decisions.

Moreover, counterfactually thinking about guilt or shame is useful to prevent their future arising, a process of self-cleansing or self-debugging [20]. Reflection, on the other hand, permits momentary experiential simulation of possible alternatives, thereby allowing careful consideration before a decision is made, and to subsequently justify it.

The investigation in this paper pertains to how moral issues can innovatively be expressed with counterfactual reasoning by resorting to the aforementioned approach. In particular, its application for examining viewpoints on moral permissibility is scrutinized, exemplified by classic moral dilemmas from the literature on the Doctrine of Double Effect (DDE) [18], and the Doctrine of Triple Effect (DTE) [12].

DDE is often invoked to explain the permissibility of an action that causes a harm, by distinguishing whether this harm is a mere side effect of bringing about a good result, or if this harm is rather the actual means to bringing about the same good end [18]. In [11], DDE has been utilized to explain the consistency of judgments, shared by subjects from demographically diverse populations, on a series of moral dilemmas.

Counterfactuals may provide a general way to examine DDE in dilemmas, e.g., the classic trolley problem [7], by distinguishing between cause and side effect of performing an action to achieve a goal. This distinction between causes and side effects may explain the permissibility of an action in accordance with DDE. That is, if some morally wrong effect E happens to be a cause for a goal G that one wants to achieve by performing an action A , and E is not a mere side effect of A , then performing A is impermissible. The counterfactual form below, in a setting where action A is performed to achieve goal G , expresses this:

“If *not E* had been true, then *not G* would have been true.”

The evaluation of this counterfactual form identifies permissibility of action A from its effect E , by identifying whether the latter is a necessary cause for goal G or a mere side effect of action A . That is, if the counterfactual proves valid, then E is instrumental as a cause of G , and not a mere side effect of action A . Since E is morally wrong, achieving G that way, by means of A , is impermissible; otherwise, not.

Note that the evaluation of counterfactuals in this application is considered from the perspective of agents who perform the action, rather than from others’ (e.g., observers). Moreover, the emphasis on causation in this application focuses on agents deliberate actions, rather than on causation and counterfactuals in general, cf. [4, 22].

In the next examples, the aforementioned general counterfactual method is illustrated by taking off-the-shelf military morality cases [33]. Consider “Terror Bombing”, *teb* for short, which means: Bombing a civilian target during a war, thus killing many civilians, in order to terrorise the enemy, and thereby getting them to end the war. DDE affirms *teb* impermissible.

On the other hand, “Tactical bombing” (*tab*) means: Bombing a military target, which will effectively end the war, but with the foreseen consequence of killing the same large number of civilians nearby. DDE affirms *tab* permissible.

Modeling Terror Bombing. Take set of abducibles $A = \{teb, teb^*\}$ and program P :

$$\begin{array}{ll} end_war \leftarrow terror_civilians. & terror_civilians \leftarrow kill_civilian. \\ kill_civilian \leftarrow target_civilian. & target_civilian \leftarrow teb. \end{array}$$

Counterfactual: *If civilians had not been killed, the war would not have ended.*

The evaluation follows.

Step 1: Observations $O = \{kill_civilian, end_war\}$ with explanation $E = \{teb\}$.

Step 2: Produce program T from P :

```
make(kill_civilians*)           % Intervention: If civilians had not been killed...
kill_civilians ← make(kill_civilians). % Killing civilians or otherwise
kill_civilians* ← make(kill_civilians*). % is now available only by intervention.
```

Simply deleting *kill_civilians* and adding *kill_civilians**, without employing *make/1*, would throw away the structural information about the intervention, which would hinder the program providing detailed explanations about intervention.

Plus, for irrelevancy and consistency, delete:

$$kill_civilian \leftarrow target_civilian.$$

Step 3: The counterfactual is valid since conclusion “the war would not have ended” is true. Indeed, ‘*not end_war*’ holds in the semantics of updated T , added with the abducted, and adopted, unchanging background fact E . Hence, the morally wrong action *kill_civilians* is an instrument to achieve the goal *end_war*. It is a cause of *end_war* by performing *teb*, and not a mere side effect of *teb*. Therefore, *teb* is DDE morally impermissible.

Modeling Tactical Bombing. Take set of abducibles $A = \{tab, tab^*\}$ and program P :

```
end_war ← target_military.
kill_civilian ← tab.      target_military ← tab.
```

The counterfactual is the same as above. The evaluation follows.

Step 1: Observations $O = \{kill_civilian, end_war\}$ with explanation $E = \{tab\}$.

Step 2: Produce T from P , obtaining same T as in the terror bombing’s model.

And, for irrelevancy and consistency, now delete:

$$kill_civilian \leftarrow tab.$$

Step 3: The counterfactual is not valid, since its conclusion “the war would not have ended” is false. Indeed, *end_war* holds in the semantics of updated T plus E . Hence, the morally wrong *kill_civilian* is a just side effect of achieving the goal *end_war*. Therefore, *tab* is DDE morally permissible.

A more complex scenario can challenge this application of counterfactuals, to distinguish moral permissibility according to DDE vs. DTE. DTE [12] refines DDE particularly on the notion about harming someone as an intended means to harm the person, or instead harming the person only because it is a causal happenstance towards some goal. That is, DTE distinguishes further between doing an action with the intended goal of a harm effect to occur, and doing an action even though a harming effect will instrumentally occur. The latter is a new category of action, which is not accounted for in DDE. Though DTE also classifies the former as impermissible, it is more tolerant to the

latter (the third effect), i.e., it treats as permissible those actions performed just because instrumental harm will occur.

Kamm proposed DTE to accommodate a variant of the trolley problem, viz., the Loop Case [34]:

A trolley is headed toward five people walking on the track, and they will not be able to get off the track in time. The trolley can be redirected onto a side track, which loops back towards the five. A fat man sits on this looping side track, whose body will by itself stop the trolley. Is it morally permissible to divert the trolley to the looping side track, thereby hitting the man and killing him, but saving the five?

This case strikes most moral philosophers that diverting the trolley is permissible [21]. Referring to a psychology study [11], 56% of its respondents judged that diverting the trolley in this case is also permissible. To this end, DTE may provide the justification of its permissibility [12]. Nonetheless, DDE views diverting the trolley in the Loop case as impermissible.

Modeling Loop Case. Take set of abducibles $A = \{divert, divert^*\}$ and program P , where *save*, *divert*, *hit*, *tst*, *mst* stand for *save the five*, *divert the trolley*, *man hit by the trolley*, *train on the side track* and *man on the side track*, respectively:

save \leftarrow *hit*. *hit* \leftarrow *tst, mst*. *tst* \leftarrow *divert*. *mst*.

Counterfactual: *If the man had not been hit by the trolley, the five people would not have been saved.* The evaluation follows.

Step 1: Observations $O = \{hit, save\}$ with explanation $E = \{divert\}$.

Step 2: Produce program T from P :

```
make(hit*)           % Intervention: If the man had not been hit by the trolley...
hit ← make(hit).     % The man being hit by the trolley or otherwise
hit* ← make(hit*).   % is now available only by intervention.
```

And, for irrelevancy and consistency, now delete:

hit \leftarrow *tst, mst*.

Step 3: The counterfactual is valid, since its conclusion “the five people would not have been saved” is true. Indeed, *not save* holds in the semantics of updated T plus E . Hence, *hit*, as a consequence of action *divert*, is instrumental as a cause of goal *save*. Therefore, *divert* is DDE morally impermissible.

DTE considers diverting the trolley as permissible, since the man is already on the side track, without any deliberate action performed in order to place him there. In the above program, we have the fact *mst* ready, without abducing any ancillary action. The validity of the counterfactual “*if the man had not been on the side track, then he would not have been hit by the trolley*”, which can easily be verified, ensures that the unfortunate event of the man being hit by the trolley is indeed the consequence of the man being on the side track. The lack of deliberate action (say, by pushing the man – *push* for short) in order to place him on the side track, and whether the absence of this action still causes the unfortunate event (the third effect) is captured by the

counterfactual “*if the man had not been pushed, then he would not have been hit by the trolley*”. This counterfactual is not valid, because the new observation $O = \{push, hit\}$ has no explanation: *push* is not in the set of abducibles A , and moreover there is no fact *push* either. This means that even without this hypothetical but unexplained deliberate action of pushing, the man would still have been hit by the trolley (just because he is already on the side track). In summary, though *hit* is a consequence of *div* and instrumental in achieving *save*, no deliberate action is required to cause *mst*, in order for *hit* to occur. Hence *divert* is DTE morally permissible.

In order to further distinguish moral permissibility with respect to DDE and DTE, we also consider a variant of the Loop case, viz., the *Loop-Push* case – see also the Extra Push case in [12]. Differently from the Loop case, in this Loop-Push case the looping side track is initially empty, and besides the diverting action, an ancillary action of pushing a fat man in order to place him on the side track is additionally performed.

Modeling Loop-Push Case. Take set of abducibles $A = \{divert, push, divert^*, push^*\}$ and program P :

$$save \leftarrow hit. \quad hit \leftarrow tst, mst. \quad tst \leftarrow divert. \quad mst \leftarrow push.$$

Recall the counterfactuals considered in the discussion of DDE and DTE of the Loop case:

- “*If the man had not been hit by the trolley, the five people would not have been saved.*” The same observation $O = \{hit, save\}$ provides an extended explanation $E = \{divert, push\}$. That is, the pushing action needs to be abduced for having the man on the side track, so the trolley can be stopped by hitting him. The same intervention $make(hit^*)$ is applied to the same transform T , resulting in a valid counterfactual: *not save* holds in the semantics of updated T plus E .
- “*If the man had not been pushed, then he would not have been hit by the trolley.*” The relevant observation is $O = \{push, hit\}$, explained by $E = \{divert, push\}$. Whereas this counterfactual is not valid in DTE of the Loop case, it is valid in the Loop-Push case. Given rule $push^* \leftarrow make(push^*)$ in the transform T and intervention $make(push^*)$, we verify that *not hit* holds in the semantics of updated T plus E .

From the validity of these two counterfactuals it can be inferred that, given the diverting action, the ancillary action of pushing the man onto the side track causes him to be hit by the trolley, which in turn causes the five to be saved. In the Loop-Push, DTE agrees with DDE that such a deliberate action (pushing) performed in order to bring about harm (the man hit by the trolley), even for the purpose of a good or greater end (to save the five), is likewise impermissible.

3 Conclusions and Further Work

Computational morality [1, 36] is a burgeoning field that emerges from the need of imbuing autonomous agents with the capacity of moral decision making to enable them to function in an ethically responsible manner via their own ethical decisions. It has attracted the artificial intelligence community, and brought together perspectives from

various fields: philosophy, anthropology, cognitive science, neuroscience, and evolutionary biology. The overall result of this interdisciplinary research is not just important for equipping agents with some capacity for making moral judgments, but also to help better understand morality, via the creation and testing of computational models of ethical theories.

This paper presented a formulation of counterfactuals evaluation by means of logic program abduction and updating. The approach corresponds to the three- step process in Pearl's structural theory, despite omitting probability to concentrate on a naturalised logic. Furthermore, counterfactual reasoning has been shown quite useful for critical thinking, namely about moral issues, where (non-probabilistic) moral reasoning about permissibility is examined by employing this logic program approach to distinguish between causes and the side effects that are the result of agents actions to achieve goals.

In Pearl's theory, intervention is realised by superficial revision, i.e., by imposing the desired value to the intervened node and cutting it from its parent nodes. This is also the case in the approach presented here, achieved by hypothetical updates via the reserved predicate *make*. Other subtle ways of intervention may involve deep revision, realisable with logic programs (cf. [25]), and minimal revision (cf. [5]).

Logic program abduction was used in [13] and [26] to model moral reasoning in various scenarios of the trolley problem, both from DDE and DTE viewpoints, sans counterfactuals. Abducibles are used to represent decisions, where impermissible actions are ruled out using an integrity constraint, and a posteriori preferences are eventually enacted to come up with a moral decision from the remaining alternatives of action. Subsequent work [10] refines it with uncertainty of actions and consequences in several scenarios of the trolley problem by resorting to probabilistic logic programming P-log [2].

Side effects in abduction have been investigated in [24] through the concept of inspection points; the latter are construed in a procedure by 'meta-abducting' a specific abducible, *abduced(a)*, whose function is only checking that its corresponding abducible *a* is indeed already abduced elsewhere. Therefore, the consequence of the action that triggers this 'meta-abducting' is merely a side effect. Indeed, inspection points may be employed to distinguish a cause from a mere side effect, and thus may provide an alternative or supplement to counterfactuals employed for the same purpose.

Counterfactuals may as well be suitable to address moral justification, via 'compound counterfactuals': *Had I known what I know today, then if I were to have done otherwise, something preferred would have followed*. Such counterfactuals, by imagining alternatives with worse effect – the so-called *downward counterfactuals* [16] – may provide justification for what was done due to lack of the current knowledge. This is accomplished by evaluating what would have followed if the intent had been otherwise, other things (including present knowledge) being equal. It may justify that what would have followed is no morally better than the actual ensued consequence. We are currently investigating the application of counterfactuals to justify an exception for an action to be permissible [27, 31], which may lead to agents' argumentation following contractualism of [32].

Acknowledgements

AS acknowledges the support from Fundação para a Ciência e a Tecnologia (FCT/MEC) Portugal, grant SFRH/BD/72795/2010. LMP acknowledges the support from FCT/MEC NOVA LINCS PEst UID/CEC/04516/2013. We thank Emmanuelle-Anna Dietz for the fruitful discussions.

References

1. M. Anderson and S. L. Anderson, editors. *Machine Ethics*. Cambridge University Press, New York, NY, 2011.
2. C. Baral and M. Hunsaker. Using the probabilistic logic programming language P-log for causal and counterfactual reasoning and non-naive conditioning. In *Procs. 20th International Joint Conference on Artificial Intelligence (IJCAI)*, 2007.
3. R. M. J. Byrne. *The Rational Imagination: How People Create Alternatives to Reality*. MIT Press, Cambridge, MA, 2007.
4. J. Collins, N. Hall, and L. A. Paul, editors. *Causation and Counterfactuals*. MIT Press, Cambridge, MA, 2004.
5. E-A. Dietz, S. Hölldobler, and L. M. Pereira. On conditionals. In *Procs. Global Conference on Artificial Intelligence (GCAI 2015)*, 2015.
6. K. Epstude and N. J. Roese. The functional theory of counterfactual thinking. *Personality and Social Psychology Review*, 12(2):168–192, 2008.
7. P. Foot. The problem of abortion and the doctrine of double effect. *Oxford Review*, 5:5–15, 1967.
8. M. L. Ginsberg. Counterfactuals. *Artificial Intelligence*, 30(1):35–79, 1986.
9. J. Y. Halpern and C. Hitchcock. Graded causation and defaults. *British Journal for the Philosophy of Science*, 66:413–457, 2015.
10. T. A. Han, A. Saptawijaya, and L. M. Pereira. Moral reasoning under uncertainty. In *Procs. 18th International Conference on Logic for Programming, Artificial Intelligence and Reasoning (LPAR)*, volume 7180 of LNCS, pages 212–227. Springer, 2012.
11. M. Hauser, F. Cushman, L. Young, R. K. Jin, and J. Mikhail. A dissociation between moral judgments and justifications. *Mind and Language*, 22(1):1–21, 2007.
12. F. M. Kamm. *Intricate Ethics: Rights, Responsibilities, and Permissible Harm*. Oxford University Press, Oxford, UK, 2006.
13. R. Kowalski. *Computational Logic and Human Thinking: How to be Artificially Intelligent*. Cambridge University Press, New York, NY, 2011.
14. D. Lewis. *Counterfactuals*. Harvard University Press, Cambridge, MA, 1973.
15. D. R. Mandel, D. J. Hilton, and P. Catellani, editors. *The Psychology of Counterfactual Thinking*. Routledge, New York, NY, 2005.
16. K. D. Markman, I. Gavanski, S. J. Sherman, and M. N. McMullen. The mental simulation of better and worse possible worlds. *Journal of Experimental Social Psychology*, 29:87–109, 1993.
17. R. McCloy and R. M. J. Byrne. Counterfactual thinking about controllable events. *Memory and Cognition*, 28:1071–1078, 2000.
18. A. McIntyre. Doctrine of double effect. In E. N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Center for the Study of Language and Information, Stanford University, Fall 2011 edition, 2004. <http://plato.stanford.edu/archives/fall2011/entries/double-effect/>.

19. S. Migliore, G. Curcio, F. Mancini, and S. F. Cappa. Counterfactual thinking in moral judgment: an experimental study. *Frontiers in Psychology*, 5:451, 2014.
20. P. M. Niedenthal, J. P. Tangney, and I. Gavanski. “If Only I weren’t” Versus “If Only I Hadn’t”: Distinguishing Shame and Guilt in Counterfactual Thinking. *Journal of Personality and Social Psychology*, 67(4):585–595.
21. M. Otsuka. Double effect, triple effect and the trolley problem: Squaring the circle in looping cases. *Utilitas*, 20(1):92–110, 2008.
22. J. Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, Cambridge, MA, 2009.
23. L. M. Pereira, J. N. Aparício, and J. J. Alferes. Counterfactual reasoning based on revising assumptions. In *Procs. International Symposium on Logic Programming (ILPS 1991)*, pages 566–577. MIT Press, 1991.
24. L. M. Pereira, P. Dell’Acqua, A. M. Pinto, and G. Lopes. Inspecting and preferring abductive models. In K. Nakamatsu and L. C. Jain, editors, *The Handbook on Reasoning-Based Intelligent Systems*, pages 243–274. World Scientific Publishers, 2013.
25. L. M. Pereira, E-A. Dietz, and S. Hölldobler. Abductive framework for counterfactual reasoning in logic programming. Draft, Available from <http://centria.di.fct.unl.pt/~lmp/publications/online-papers/counterfactuals.pdf>, 2015.
26. L. M. Pereira and A. Saptawijaya. Modelling Morality with Prospective Logic. In M. Anderson and S. L. Anderson, editors, *Machine Ethics*, pages 398–421. Cambridge U. P., 2011.
27. L. M. Pereira and A. Saptawijaya. Abduction and beyond in logic programming with application to morality. Accepted at *Frontiers of Abduction, a special issue of IfCoLog Journal of Logics and their Applications*. Available from (preprint): <http://goo.gl/yhmZzy>, 2015.
28. L. M. Pereira and A. Saptawijaya. Counterfactuals, logic programming and agent morality. In R. Urbaniak and G. Payette, editors, *Logic, Argumentation & Reasoning*. Springer, 2016 (forthcoming).
29. N. J. Roese. Counterfactual thinking. *Psychological Bulletin*, 121(1):133–148, 1997.
30. N. J. Roese and J. M. Olson, editors. *What Might Have Been: The Social Psychology of Counterfactual Thinking*. Psychology Press, Hove, UK, 2009.
31. A. Saptawijaya and L. M. Pereira. Logic programming applied to machine ethics. In *Procs. 17th Portuguese International Conference on Artificial Intelligence (EPIA)*, volume 9273 of *LNAI*. Springer, 2015.
32. T. M. Scanlon. *What We Owe to Each Other*. Harvard University Press, Cambridge, MA, 1998.
33. T. M. Scanlon. *Moral Dimensions: Permissibility, Meaning, Blame*. Harvard University Press, Cambridge, MA, 2008.
34. J. J. Thomson. The trolley problem. *The Yale Law Journal*, 279:1395–1415, 1985.
35. J. Vennekens, M. Bruynooghe, and M. Denecker. Embracing events in causal modeling: Interventions and counterfactuals in CP-logic. In *JELIA 2010*, volume 6341 of *LNCS*, pages 313–325. Springer, 2010.
36. W. Wallach and C. Allen. *Moral Machines: Teaching Robots Right from Wrong*. Oxford University Press, Oxford, UK, 2009.
37. B. Weiner. *Judgments of Responsibility: A Foundation for a Theory of Social Conduct*. The Guilford Press, New York, NY, 1995.