

ON THE POSSIBILITY OF A GENERAL PURGE OF SELF-REFERENCE

LUCAS ROSENBLATT

Universidad de Buenos Aires – CONICET

l_rosenblatt@hotmail.com

Abstract

My aim in this paper is to gather some evident in favor of the view that a general purge of self-reference is possible. I do this by considering a modal-epistemic version of the Liar Paradox introduced by Roy Cook. Using yabloesque techniques, I show that it is possible to transform this circular paradoxical construction (and other constructions as well) into an infinitary construction lacking any sort of circularity. Moreover, contrary to Cook's approach, I think that this can be done without using any controversial multi-modal rules, i.e., the usual rules from normal epistemic and modal logic are enough to show the paradoxicality of the infinitary construction.

KEY WORDS: Yablo's Paradox; Self-reference; Infinitary Logic; Modal Logic.

Resumen

Mi objetivo en este trabajo es ofrecer cierta evidencia a favor de la tesis según la cual una purga general de la autorreferencia es posible. Hago esto considerando una versión modal-epistémica de la Paradoja del Mentiroso introducida por Roy Cook. Usando técnicas yablescas, muestro que es posible transformar esta construcción paradójica circular (y también otras construcciones) en una construcción infinitaria que carece de cualquier forma de circularidad. Más aún, en contra de la propuesta de Cook, muestro que esto puede hacerse sin utilizar ninguna regla multimodal controversial, esto es, las reglas usuales de la lógica modal y la lógica epistémica son suficientes para mostrar la paradojicidad de la construcción infinitaria.

PALABRAS CLAVE: Paradoja de Yablo; Autorreferencia; Lógica infinitaria; Lógica modal.

Some philosophers¹ have seen in Yablo's paradox the starting point of a general method for purging every self-referential paradox of its self-referentiality². Sorensen, for example, says that "[t]he simplicity of Yablo's paradox invites the conjecture that all [of the self-referential] paradoxes can be purged of self-reference" (1998, p. 150). If such a method were available, there would be a compelling argument in favor of the view

¹ In particular, Sorensen (1998) and Schlenker (2007).

² I will use 'self-referential' and 'circular' interchangeably throughout the paper.

that paradoxicality must not (maybe never) be blamed on circularity. Cook (2013) rightly notices that there are two ways to make Sorensen's idea more precise:

Weak Purge: Given any self-referential construction Σ in some language L , there is some (possibly distinct) language L^* such that L^* contains a non-self-referential Yabloesque analogue of Σ .

Strong Purge: Given any self-referential construction Σ in some language L , there is, in L itself, a non-self-referential Yabloesque analogue of Σ .

In his book Cook presents a purposed-build infinitary language L_P in which certain paradoxes can be represented. This language only contains one type of formula: (possibly infinite) conjunctions of sentences of the form ' S_n is false', where S_n is the name of a sentence of L_P . Cook uses a function δ to provide the denotation of every sentence name. For example, the Liar sentence can be expressed by $F(S_1)$, where $\delta(S_1) = F(S_1)$. And the Yablo sequence can be expressed as "the unwinding" of the Liar sentence³:

$$\begin{aligned}\delta(S_{\langle 1,1 \rangle}) &= F(S_{\langle 2,1 \rangle}) \wedge F(S_{\langle 3,1 \rangle}) \wedge F(S_{\langle 4,1 \rangle}) \wedge \dots \\ \delta(S_{\langle 2,1 \rangle}) &= F(S_{\langle 3,1 \rangle}) \wedge F(S_{\langle 4,1 \rangle}) \wedge F(S_{\langle 5,1 \rangle}) \wedge \dots \\ \delta(S_{\langle 3,1 \rangle}) &= F(S_{\langle 4,1 \rangle}) \wedge F(S_{\langle 5,1 \rangle}) \wedge F(S_{\langle 6,1 \rangle}) \wedge \dots \\ &\dots\end{aligned}$$

Cook thinks that if the strong version of the purge is to have any chance of success, the language in which to carry out the unwindings must be at least as strong as the language of arithmetic. This gives us a sort of dilemma. On the one hand, since L_P is specifically designed to model certain paradoxes, it can only be used to argue in favor of the weak version of the purge. On the other hand, if the language of arithmetic is used, then it can be shown that Yablo's sequence is circular (provided that circularity is understood as being a fixed point (of a certain sort))⁴. Cook's conclusion is that unwindings are not useful to carry out the strong version of the purge.

³ The pairs are ordered by the usual lexicographical order. The reason for using pairs (or n -tuples, depending on the example) in the unwinding has to do with the possibility of unwinding a sequence of sentences. The reader interested in the technical details for obtaining the unwinding of an L_P sentence (or sequence of sentences) can see Cook (forthcoming).

⁴ Cook (forthcoming) has put forward a couple of important distinctions concerning the different types of fixed points (strong/weak, sentential/predicative). Since I do not think (see below) that fixed points provide an adequate definition of circularity, I will not concern myself with these subtleties here.

I will claim that this is a false dilemma. Even if L_P is not a good candidate to carry out the purge, it is possible to use a richer infinitary language (or maybe even arithmetic itself⁵) to argue in favor of the strong version of the purge. Of course, the success of the purge will depend on our ability to provide a definition of circularity suitable for the language. Presenting such a definition is beyond the scope of this paper, but it is interesting to note that Cook does not consider the possibility of vindicating the strong version of the purge not by finding some method different from the one based on unwindings but by finding a different definition of circularity.

With this conceptual background in mind, let me illustrate how the same unwinding operation performed on the Liar can be performed on non-semantic circular constructions. To do things properly I am going to work with an infinitary propositional language $L_{\square KT(x)}$ ⁶. The vocabulary of $L_{\square KT(x)}$ is the following: a class C of sentence names $\{S_0, S_1, S_2, \dots\}$, a set of contingent sentences $\{p_0, p_1, p_2, \dots\}$ ⁷, *falsum* (\perp), negation (\neg), conjunction (\wedge), the material conditional (\supset), the truth predicate (T), the possibility operator (\diamond), and the knowledge operator (K). There is a denotation function $\delta: C \rightarrow \{\text{sentences of } L_{\square KT(x)}\}$ that assigns a semantic value (in particular, a sentence of $L_{\square KT(x)}$) to each of the sentence names in C .

The non-semantic example Cook mentions is the Modal Knower⁸, the sentence that says of itself that it is not knowable. In my present framework this can be formalized as $\delta(S_1) = \neg \diamond KT(S_1)$

The unwinding of this sentence is:

$$\begin{aligned} \delta(S_{\langle 1,1 \rangle}) &= \neg \diamond KT(S_{\langle 2,1 \rangle}) \wedge \neg \diamond KT(S_{\langle 3,1 \rangle}) \wedge \neg \diamond KT(S_{\langle 4,1 \rangle}) \wedge \dots \\ \delta(S_{\langle 2,1 \rangle}) &= \neg \diamond KT(S_{\langle 3,1 \rangle}) \wedge \neg \diamond KT(S_{\langle 4,1 \rangle}) \wedge \neg \diamond KT(S_{\langle 5,1 \rangle}) \wedge \dots \\ \delta(S_{\langle 3,1 \rangle}) &= \neg \diamond KT(S_{\langle 4,1 \rangle}) \wedge \neg \diamond KT(S_{\langle 5,1 \rangle}) \wedge \neg \diamond KT(S_{\langle 6,1 \rangle}) \wedge \dots \\ &\dots \end{aligned}$$

⁵ It would have to be second-order arithmetic, since Yablo's sequence is satisfiable in first-order arithmetic.

⁶ I assume that $L_{\square KT(x)}$ does not count as a purposed-build language. It is just an extension of some infinitary language $L_{(\kappa, \omega)}$.

⁷ We need contingent sentences because otherwise truth and necessity will collapse.

⁸ Actually Cook talks about the Knower, but he formalizes it as $\delta(S_1) = UK(S_1)$, and says that 'UK(x)' should be read as 'x is not knowable'. So it appears that he is thinking about a modal version of the Knower. Even if that is not the case, the modal version is interesting in its own right. (The reason for using a version without negation is that the proof showing that the unwinding of the Modal Knower does not contain fixed points depends on the language failing to contain a primitive negation operator. However, since I will not give an account of circularity in terms of fixed points, negation will be introduced as a primitive symbol).

To prove the paradoxicality of these constructions Cook suggests using the following two multi-modal rules⁹:

Rule K1: if $T(S_m) \vdash T(S_n)$ and $T(S_m) \vdash \neg \Diamond KT(S_n)$, then $\vdash \neg \Diamond KT(S_m)$

Rule K2: if $T(S_m) \vdash \perp$, then $\vdash \neg \Diamond KT(S_m)$

Another important aspect of the system is that given a particular denotation function δ , we can extend the deductive system by incorporating rules for it:

δ -*elim* if $\delta(S_n)$, then infer $T(S_n)$

δ -*intro* if $T(S_n)$, then infer $\delta(S_n)$

So the consistency of the extended system will depend on the denotation function being used. Interestingly, these rules are enough to obtain a contradiction from the Modal Knower and from its unwinding.

Another issue is whether the same sort of argument is available for circular constructions involving belief and believability. In particular, it would be interesting to see what would happen if there were a belief operator¹⁰ in the language and a couple of rules for it. If a general purge of self-reference is possible, there should be a way of proving the paradoxicality of the Believer and the Modal Believer (the sentence that says of itself that it is not rationally believed and the sentence that says of itself that it is not rationally believable, respectively) in the extended system. However, things are problematic for belief. In particular, if we mimic Cook's system for knowability we would have the following two rules:

Rule B1: if $T(S_m) \vdash T(S_n)$ and $T(S_m) \vdash \neg \Diamond BT(S_n)$, then $\vdash \neg \Diamond BT(S_m)$

Rule B2: if $T(S_m) \vdash \perp$, then $\vdash \neg \Diamond BT(S_m)$

Unfortunately, these rules are much more controversial than the ones for knowability. B2, for example, implies there cannot be an agent and a time such that the agent rationally believes at that time something that logically leads to a contradiction. This strikes me as too strong.

Be that as it may, this does not imply that there is no other way in which the purge of self-reference could be carried out. In

⁹ Cook uses a sequent calculus. Here I will work with a Fitch-style natural deduction system.

¹⁰ I am dealing with rational belief, so the intended reading of B is 'at some time someone rationally believes'.

particular, it is not difficult to see that if some reasonable principles concerning \Box and K (and B) are accepted, the paradoxicality of several (all?) modal-epistemic paradoxes and their unwindings is provable without resorting to multi-modal rules. The same holds of the usual paradoxes involving just one modality, such as the Liar, the Knower and the Believer.

Let us assume that we have the following unimodal principles and rules for \Box and K ¹¹:

(\Box -fact)	$\Box\Phi \rightarrow \Phi$	(K -fact)	$K\Phi \rightarrow \Phi$
(\Box -nec)	If $\vdash \Phi$, then $\vdash \Box\Phi$	(K -nec)	If $\vdash \Phi$, then $\vdash K\Phi$
(\Box -dist)	$\Box(\Phi \rightarrow \Psi) \rightarrow (\Box\Phi \rightarrow \Box\Psi)$	(K -dist)	$K(\Phi \rightarrow \Psi) \rightarrow (K\Phi \rightarrow K\Psi)$

The following is a simple proof of the paradoxicality of the Modal Knower in which only these unimodal rules are used. Proof: Assume $KT(S_1)$, where $\delta(S_1) = \neg\Diamond KT(S_1)$. By (K -fact), we can infer $T(S_1)$. Using δ -intro, we obtain $\neg\Diamond KT(S_1)$, which is equivalent to $\Box\neg KT(S_1)$. Using (\Box -fact) we reach a contradiction between $\neg KT(S_1)$ and $KT(S_1)$. By *reductio*, we have $\neg\neg KT(S_1)$. Applying (\Box -nec) yields $\Box\neg\neg KT(S_1)$, from which we infer $\neg\Diamond\neg\neg KT(S_1)$. By δ -elim we obtain $T(S_1)$. At this point we can apply (K -nec), which gives us $KT(S_1)$. And this contradicts $\neg\neg KT(S_1)$.

Next we show that if the modal-epistemic logic for K and \vdash has the unimodal rules specified above, the paradoxicality of the unwinding of the Modal Knower is provable as well. Proof:

1)	$T(S_{\langle 1,1 \rangle})$	Assumption
2)	$\wedge\{\neg\Diamond KT(S_{\langle n,1 \rangle}) \text{ for every } n > 1\}$	δ -intro
3)	$\neg\Diamond KT(S_{\langle 2,1 \rangle})$	\wedge -elim
4)	$\neg\Diamond KT(S_{\langle 3,1 \rangle})$	\wedge -elim
....	
w)	$\wedge\{\neg\Diamond KT(S_{\langle n,1 \rangle}) \text{ for every } n > 2\}$	\wedge -intro
w+1)	$T(S_{\langle 2,1 \rangle})$	δ -elim
w+2)	$T(S_{\langle 1,1 \rangle}) \supset T(S_{\langle 2,1 \rangle})$	\supset -intro
w+3)	$K(T(S_{\langle 1,1 \rangle}) \supset T(S_{\langle 2,1 \rangle}))$	(K -nec)
w+4)	$KT(S_{\langle 1,1 \rangle}) \supset KT(S_{\langle 2,1 \rangle})$	(K -dist)
w+5)	$T(S_{\langle 1,1 \rangle})$	Assumption
w+6)	$\wedge\{\neg\Diamond KT(S_{\langle n,1 \rangle}) \text{ for every } n > 1\}$	δ -intro

¹¹ For belief I assume (B -dist), (B -nec) plus the following principle that plausibly holds of an idealized notion of belief:

(B -semifact) $B\neg B\Phi \rightarrow \neg B\Phi$

$\omega+7)$	$\neg\Diamond\text{KT}(S_{\langle 2,1 \rangle})$	\wedge -elim
$\omega+8)$	$\Box\neg\text{KT}(S_{\langle 2,1 \rangle})$	$\Box\neg =_{\text{df}} \neg\Diamond$
$\omega+9)$	$\neg\text{KT}(S_{\langle 2,1 \rangle})$	$(\Box\text{-fact})$
$\omega+10)$	$\text{T}(S_{\langle 1,1 \rangle}) \supset \neg\text{KT}(S_{\langle 2,1 \rangle})$	\supset -intro
$\omega+11)$	$\text{KT}(S_{\langle 1,1 \rangle}) \supset \text{T}(S_{\langle 1,1 \rangle})$	(K-fact)
$\omega+12)$	$\text{KT}(S_{\langle 1,1 \rangle}) \supset \neg\text{KT}(S_{\langle 2,1 \rangle})$	Transitivity of \supset
$\omega+13)$	$\neg\text{KT}(S_{\langle 1,1 \rangle})$	\neg -intro
$\omega+14)$	$\Box\neg\text{KT}(S_{\langle 1,1 \rangle})$	$(\Box\text{-nec})$
$\omega+15)$	$\neg\Diamond\text{KT}(S_{\langle 1,1 \rangle})$	$\Box\neg =_{\text{df}} \neg\Diamond$
....	
$2\omega+15)$	$\neg\Diamond\text{KT}(S_{\langle 2,1 \rangle})$	$\Box\neg =_{\text{df}} \neg\Diamond$
....	
$3\omega+15)$	$\neg\Diamond\text{KT}(S_{\langle 3,1 \rangle})$	$\Box\neg =_{\text{df}} \neg\Diamond$
....	
$\omega^2)$	$\wedge\{\neg\Diamond\text{KT}(S_{\langle n,1 \rangle}) \text{ for every } n > 1\}$	\wedge -intro
$\omega^2+1)$	$\text{T}(S_{\langle 1,1 \rangle})$	δ -elim
$\omega^2+2)$	$\text{KT}(S_{\langle 1,1 \rangle})$	(K-nec)
$\omega^2+3)$	$\Box\neg\text{KT}(S_{\langle 1,1 \rangle}) \supset \neg\text{KT}(S_{\langle 1,1 \rangle})$	$(\Box\text{-fact})$
$\omega^2+4)$	$\neg\Box\neg\text{KT}(S_{\langle 1,1 \rangle})$	modus tollens
$\omega^2+5)$	$\Diamond\text{KT}(S_{\langle 1,1 \rangle})$	$\Box =_{\text{df}} \neg\Diamond\neg$
$\omega^2+6)$	\perp	\perp -intro ¹²

Let me finish by saying what remains to be done. If Sorensen's idea is to be fully vindicated by means of the strategy presented above, some conditions must be set up on what counts as an adequate unwinding. A very reasonable demand is that the unwinding should have the same semantic status as the construction being unwinded. This is what Cook calls the "Covariation criterion". Informally, we are demanding that a construction C should be semantically similar to its unwinding. This requires, *at least*, that a construction C be paradoxical/non-paradoxical if and only if its unwinding is paradoxical/non-paradoxical. I believe that a lot of work has to be done in order to show that this criterion is satisfied, especially for a language like $L_{\Box\text{KT}(x)}$: we need an adequate interpretation for an *infinitary* language that has a truth *predicate* and operators for necessity, knowledge, and maybe other notions too. Also, we need to come up with a way of interpreting the language such that every construction

¹² There are analogous proofs for other modal-epistemic paradoxes (and their unwindings), such as the $\Diamond\text{K}$ -version of Curry's sentence, the Knower, the Believer, the Modal Believer, etc. Obtaining an inconsistency from the ones involving belief (and their unwindings) requires (B-semifact).

is semantically similar (in a relevant way) to its unwinding. This is not to say that the project is not viable, I am just pointing out what remains to be done if the strong version of the purge is to be vindicated.

A second reasonable demand is that no sentence in the unwinding should be circular. We could call this demand the “Failure of Circularity Criterion”. I have already mentioned that there are several ways to flesh this out. If circularity is to be coded as being a fixed point (of some sort), then the idea is that no sentence in the unwinding should be a fixed point (of that sort). If circularity is explained in some other way, then the criterion must be characterized differently.

One positive aspect of infinitary languages (lacking quantifiers) is that it is not hard to provide an acceptable definition of circularity for its sentences (since the difficulties are usually associated with quantified statements)¹³. Additionally, it is important to notice that I need something slightly weaker than a definition. It would be enough to provide a sufficient condition for non-circularity (or a necessary condition for circularity) and to show that the infinitary unwindings always satisfy that condition.

References

- Cook, R. (forthcoming), *The Yablo Paradox: An Essay on Circularity*, Oxford, Oxford University Press.
- Schlenker, P. (2007), “The Elimination of Self-Reference: Generalized Yablo Series and the Theory of Truth”, *Journal of Philosophical Logic*, 36, pp. 251-307.
- Sorensen, R. (1998), “Yablo’s Paradox and Kindred Infinite Liars”, *Mind*, 107, pp. 137-155.

¹³ For instance, since $L_{\omega\text{-KT}(x)}$ has no quantifiers, we could use something like Picollo’s notion of m-circularity to give a full account of circularity in this language (see Picollo’s contribution to this volume).