

*Logika w zastosowaniach kognitywistycznych*

**Paradoks wszechwiedzy logicznej**  
*(logical omniscience paradox)*  
**i wybrane metody jego unikania**

(notatki do wykładów)

Andrzej Wiśniewski  
Andrzej.Wisniewski@amu.edu.pl

*wersja beta 1.1*

Niech  $\otimes$  będzie metazmienną, której wartościami są operatory wiedzy lub operatory przekonań. Niech  $A, B, \dots$  będą metajęzykowymi zmiennymi przebiegającymi zbiór formuł języka (zdaniowej) logiki epistemicznej. Symbol  $\vdash$  czytamy "jest tezą". Rozważmy zależności o schematach:

- LO1:** *jeśli*  $\vdash A$ , *to*  $\vdash \otimes A$
- LO2:** *jeśli*  $\vdash A \rightarrow B$ , *to*  $\vdash \otimes A \rightarrow \otimes B$
- LO3:** *jeśli*  $\vdash A \leftrightarrow B$ , *to*  $\vdash \otimes A \leftrightarrow \otimes B$
- LO4:**  $\vdash \otimes A \rightarrow \otimes(A \vee B)$
- LO5:**  $\vdash (\otimes A \vee \otimes B) \rightarrow \otimes(A \vee B)$
- LO6:**  $\vdash \neg(\otimes A \wedge \otimes \neg A)$
- LO7:**  $\vdash \otimes(A \rightarrow B) \rightarrow (\otimes A \rightarrow \otimes B)$
- LO8:**  $\vdash \otimes(A \leftrightarrow B) \rightarrow (\otimes A \leftrightarrow \otimes B)$
- LO9:**  $\vdash \otimes(A \wedge B) \rightarrow \otimes A \wedge \otimes B$
- LO10:**  $\vdash \otimes A \wedge \otimes B \rightarrow \otimes(A \wedge B)$

Zależności te obowiązują we wszystkich rozważanych dotąd logikach wiedzy i/lub przekonañ. Z drugiej strony, co najmniej niektóre z tych zależności można interpretować jako przypisujące podmiotowi epistemicznemu wszechwiedzę logiczną (*logical omniscience*).

Dlaczego?

Rozważmy zależność:

**LO1:** *jeśli*  $\vdash A$ , *to*  $\vdash \otimes A$

Niech  $\otimes$  oznacza "podmiot *a* wie, że".

Zależność **LO1** implikuje wówczas m.in., że podmiot *a* zna wszystkie tezy *KRZ* (ściślej: że dla każdej tezy *KRZ*, *A*, jest tak, że podmiot *a* wie, że *A*).

A teraz przypuśćmy, że  $\otimes$  oznacza "podmiot *a* jest (całkowicie) przekonany, że *A*".

Wówczas, na mocy **LO1**, każda teza *KRZ* jest (całkowitym) przekonaniem podmiotu *a*.

Można zapytać, czy nie są to zbyt daleko idąca idealizacje.

Teraz rozważmy zależność:

**LO2:** *jeśli*  $\vdash A \rightarrow B$ , *to*  $\vdash \otimes A \rightarrow \otimes B$

Jak pamiętamy, formuła postaci  $A \rightarrow B$  (języka *KRZ*) jest tezą *KRZ* wtw formuła  $B$  wynika logicznie na gruncie *KRZ* z formuły  $A$ .

Przypuśćmy, że podmiot  $a$  wie/jest przekonany, że  $A$ . Wówczas, na mocy zależności **LO2**, do wiedzy/przekonań podmiotu  $a$  należą wszystkie sądy wynikające logicznie z  $A$ . W tym także te, które wprawdzie wynikają logicznie z  $A$ , ale są na tyle "odległymi konsekwencjami" sądu  $A$ , iż rozważany podmiot nie wie/nie jest przekonany/nie jest świadomy tego, że wynikają one z  $A$ .

Również tutaj można mieć uzasadnioną wątpliwość, czy nie jest to zbyt daleko idąca idealizacja.

Z kolei zależność:

**LO3:** *jeśli*  $\vdash A \leftrightarrow B$ , *to*  $\vdash \otimes A \leftrightarrow \otimes B$

zdaje się orzekać o podmiocie  $a$  co następuje: gdy  $A$  oraz  $B$  są logicznie równoważne, to podmiot  $a$  wie/ jest przekonany, że  $A$  zawsze i tylko, gdy podmiot  $a$  wie/ jest przekonany, że  $B$ .

Jest oczywiste, że jest to daleko idąca idealizacja.

O zależnościach **LO4-LO10** mówiliśmy na konwersatorium, stąd też pominię tutaj ich analizę.

Zauważmy teraz, że obowiązywanie w rozważanych logikach epistemicznych zależności **LO1**, **LO2** i **LO3** jest prostym następstwem faktu, iż logiki te budujemy na bazie normalnych modalnych logik zdaniowych.

Z kolei, patrząc od strony semantycznej, gdy budujemy semantyki dla (zdaniowych) logik epistemicznych tak, jak robiliśmy to dotąd, zależności **LO1**, **LO2** i **LO3** obowiązywać muszą. Aby ich obowiązywanie uchylić, musimy zatem zmodyfikować semantykę.

Do powstających w logikach epistemicznych paradoksów wszechwiedzy logicznej można podejść dwojako.

Po pierwsze, można twierdzić, że w istocie wskazują one na przedmiot rozważań danej logiki epistemicznej, logiki wiedzy i/lub przekonań. Otóż przedmiotem tym nie są tylko odpowiednie nastawienia sądzeniowe (*propositional attitudes*) żywione aktualnie przez podmiot – nawet wyidealizowany – lecz również coś więcej: zobowiązania (*commitments*) podmiotu do żywienia tych nastawień.

Co to konkretnie znaczy? Cóż, zapraszam na wykład :)

Po drugie, można uznać, że występowanie paradoksów wszechwiedzy logicznej w poszczególnych logikach epistemicznych jest – jednak! – istotną wadą tych logik, a następnie starać się budować logiki epistemiczne wolne od tych paradoksów.

Wymaga to jednak zastosowania jakichś "nowych" środków semantycznych i/lub odejścia od normalnych modalnych logik zdaniowych.

# **Modele Rantali**

Fiński logik Veikko Rantala podał w 1982 r. pewną semantyczną metodę "blokowania" paradoksów wszechwiedzy logicznej – i zarazem semantyczną metodę budowy logik epistemicznych wolnych od tych paradoksów.

Mówiąc najogólniej, główny pomysł Rantali polega na wzięciu pod uwagę obok światów możliwych (*possible worlds*) również tzw. niemożliwych światów możliwych (*impossible possible worlds*).

Pojęcie niemożliwego świata możliwego (*sic!*) było już, nawiasem mówiąc, używane wcześniej, gdy budowano semantyki typu Kripkego dla logik istotnie słabszych od **K**.

Intuicja leżąca u podstaw pojęcia niemożliwego świata możliwego jest następująca: jest to taki świat, w którym wszystko jest możliwe – nawet sprzeczność! – ale nic nie jest konieczne.



Dla uproszczenia załóżmy, że rozważamy jednopodmiotową zdaniową logikę przekonań, w której jedynym operatorem epistemicznym jest **B** (*believes*).

Niech  $W$  będzie niepustym zbiorem, natomiast  $W^\#$  będzie podzbiorem zbioru  $W$ . Intuicyjnie rzecz biorąc,  $W^\#$  to zbiór niemożliwych światów możliwych, różnica  $W \setminus W^\#$  to zbiór "możliwych" światów możliwych, natomiast  $W$  – to zbiór "możliwych i niemożliwych" światów możliwych.

Zasadnicza idea jest teraz następująca: w światach "niemożliwych" wszystkie formuły są traktowane jak "atomy". W konsekwencji wartość logiczna formuły złożonej w świecie "niemożliwym" nie zależy od wartości logicznych jej formuł składowych w tym świecie. Powoduje to, że – przykładowo - formuła o schemacie  $A \wedge \neg A$  może być prawdziwa w świecie "niemożliwym"  $w$ ; w świecie takim może być prawdziwa implikacja  $A \rightarrow B$ , chociaż jest w nim prawdziwy jej poprzednik i fałszywy jest następnik *etc.*

**Model Rantali** to czwórka uporządkowana:

$$(\$) \quad \langle W, R, v, W^\# \rangle,$$

gdzie  $W$  i  $W^\#$  są rozumiane jak wyżej, a ponadto:

- $v$  jest funkcją wartościowania o następujących własnościach:
  - $v$  przyporządkowuje każdej parze  $\langle p_i, w \rangle$ , gdzie  $p_i$  jest zmienną zdaniową, a  $w \in W \setminus W^\#$  (tj.  $w$  jest "możliwym" światem możliwym), dokładnie jedną z wartości logicznych:  $1, 0$ ;
  - $v$  przyporządkowuje każdej parze  $\langle A, w \rangle$ , gdzie  $A$  jest (dowolną) formułą, a  $w \in W^\#$  (tj.  $w$  jest niemożliwym światem możliwym) dokładnie jedną z wartości logicznych:  $1, 0$ .
- $R \subseteq W \times W$  jest relacją (epistemicznej) alternatywności taką, że:
  - $R$  jest seryjna w zbiorze  $W$ ,
  - $R$  jest przechodnia i euklidesowa w zbiorze  $W \setminus W^\#$  (tj. w zbiorze "możliwych" światów możliwych rozważanego modelu),

- o a ponadto  $R$  spełnia pewne dalsze warunki o których powiemy za chwilę.

Pojęcie **prawdziwości formuły  $A$  w świecie  $w$  modelu Rantali**  $M = \langle W, R, v, W^\# \rangle$ , symbolicznie  $M \models^w A$ , określa się tak:

- jeśli  $w \in W^\#$ , to:  $M \models^w A$  wtw  $v(A, w) = 1$ ;
- jeśli  $w \in W \setminus W^\#$ , to:
  - (1)  $M \models^w p_i$  wtw  $v(p_i, w) = 1$ ;
  - (2)  $M \models^w \neg A$  wtw  $M \text{ non } \models^w A$ ;
  - (3)  $M \models^w (A \wedge B)$  wtw  $M \models^w A$  oraz  $M \models^w B$ ;
  - (4)  $M \models^w (A \vee B)$  wtw  $M \models^w A$  lub  $M \models^w B$ ;
  - (5)  $M \models^w (A \rightarrow B)$  wtw  $M \text{ non } \models^w A$  lub  $M \models^w B$ ;
  - (6)  $M \models^w (A \leftrightarrow B)$  wtw  $M \models^w A$  zawsze i tylko, gdy  $M \models^{w^*} A$ ;
  - (7)  $M \models^w \mathbf{B}A$  wtw dla każdego  $w^* \in W$  takiego, że  $wRw^*$ :  $M \models^{w^*} A$ .

Wspomniane wcześniej dodatkowe warunki nakładane na relację alternatywności  $R$  są następujące:

- ❖ dla każdego  $w \in W \setminus W^\#$ , dla każdego  $w^* \in W^\#$ : jeśli  $wRw^*$  oraz  $v(\neg B \neg A, w^*) = 1$ , to dla pewnego  $w^{**} \in W$  takiego, że  $wRw^{**}$  zachodzi  $M \models^{w^{**}} A$ .
- ❖ dla każdego  $w \in W \setminus W^\#$ , dla każdego  $w^* \in W^\#$ : jeśli  $wRw^*$  oraz  $v(BA, w^*) = 1$ , to dla każdego  $w^{**} \in W$  takiego, że  $wRw^{**}$  zachodzi  $M \models^{w^{**}} A$ .

**Komentarz:** zapraszam na wykład :)

Prawdziwość formuły w modelu Rantali postaci  $\langle W, R, v, W^* \rangle$  jest definiowana (nieco) niestandardowo: formuła  $A$  jest **prawdziwa w modelu Rantali**  $\langle W, R, v, W^\# \rangle$  wtw dla każdego  $w \in W \setminus W^\#$  jest tak, że  $M \models^w A$

-- tj. gdy  $A$  jest prawdziwa w każdym "możliwym" świecie możliwym rozważanego modelu.

Teraz możemy rozważyć klasę formuł, które **są prawdziwe w każdym modelu Rantali**. Będą to prawa logiki epistemicznej (tutaj: logiki przekonań) wyznaczonej przez daną klasę modeli Rantali.

Otóż do klasy tej należą m.in. formuły:

$$4_B: \quad \mathbf{B}p \rightarrow \mathbf{B}\mathbf{B}p$$

$$5_B: \quad \neg\mathbf{B}p \rightarrow \mathbf{B}\neg\mathbf{B}p$$

Dla nas jednak istotniejsze jest to, które formuły *nie* są prawami logiki epistemicznej wyznaczonej przez analizowaną klasę modeli Rantali.

Rozważmy dla przykładu formułę:

$$\mathbf{B}p \rightarrow \mathbf{B}(p \vee q)$$

Weźmy model Rantali  $\mathbf{M} = \langle \mathbf{W}, \mathbf{R}, \mathbf{v}, \mathbf{W}^\# \rangle$  spełniający następujące warunki:

- $w \in \mathbf{W} \setminus \mathbf{W}^\#, w' \in \mathbf{W}^\#$  oraz  $wRw'$ ;
- $\mathbf{v}(p, w^*) = 1$  dla każdego  $w^* \in \mathbf{W}$  takiego, że  $wRw^*$ ;
- $\mathbf{v}(p \vee q, w') = 0$ .

Wówczas mamy  $M \models^w \mathbf{B}p$  oraz  $M \text{ non} \models^w \mathbf{B}(p \vee q)$ , czyli:

$$M \text{ non} \models^w (\mathbf{B}p \rightarrow \mathbf{B}(p \vee q)).$$

Skoro  $w \in W \setminus W^\#$ , znaczy to, że formuła:

$$\mathbf{B}p \rightarrow \mathbf{B}(p \vee q)$$

nie jest prawdziwa w pewnym modelu Rantali. Tak więc możemy powiedzieć, że zależność o schemacie:

$$\mathbf{LO4}: \quad \vdash \otimes A \rightarrow \otimes(A \vee B)$$

nie obowiązuje w logice epistemicznej wyznaczonej przez (rozważaną) klasę modeli Rantali.

W podobny sposób można pokazać, że w logice tej nie obowiązuje żadna z zależności o schematach **LO5-LO10**. Pozostawiam to Państwu jako ćwiczenie :)

Rozważmy teraz przypadek zależności o schemacie:

**LO1:** *jeśli*  $\vdash A$ , *to*  $\vdash \otimes A$

Niech  $A$  ma postać  $p \rightarrow p$ . Bierzemy model Rantali  $\mathbf{M} = \langle W, R, v, W^\# \rangle$  i świat  $w' \in W^\#$  taki, że  $v(p \rightarrow p, w') = \mathbf{0}$ . Jest oczywiste, że formuła  $p \rightarrow p$  jest prawdziwa w każdym "możliwym" świecie możliwym modelu  $\mathbf{M}$ , tj. że  $\mathbf{M} \models^w (p \rightarrow p)$  dla każdego  $w \in W \setminus W^\#$ . Jest niemniej oczywiste, że formuła  $p \rightarrow p$ , podobnie jak każdy inny aksjomat rachunkowozdaniowy, jest prawdą w każdym modelu Rantali, czyli wszystkie aksjomaty rachunkowozdaniowe są tezami logiki epistemicznej wyznaczonej przez (rozważane) modele Rantali. Przypuśćmy teraz, że  $w$  jest dowolnym ale ustalonym światem z  $W \setminus W^\#$  oraz że  $wRw'$ . Dostajemy:

$$\mathbf{M} \text{ non } \models^w \mathbf{B}(p \rightarrow p)$$

czyli kontrprzykład dla zależności o schemacie **LO1**.

W podobny sposób można pokazać, że w logice epistemicznej wyznaczonej przez (rozważane) modele Rantali nie obowiązują zależności o schematach **LO2** i **LO3**. I to pozostawiam Państwu jako ćwiczenie :)

# **Modele sitowe**



Wśród formuł wynikających logicznie z danej formuły mamy zarówno takie, których podmiot jest świadomy (*aware*), jak i takie, których podmiot świadomy nie jest. Kolejna semantyczna metoda "blokowania" paradoksów wszechwiedzy logicznej nawiązuje do tej intuicji, zaliczając do wiedzy/ przekonań podmiotu tylko te konsekwencje logiczne aktualnie żywionych przekonań/ posiadanej wiedzy, których podmiot jest świadomy.

Metoda ta została zaproponowana przez Fagina i Halperna w 1988 r.

Podobnie jak poprzednio, dla uproszczenia rozważamy jednopodmiotową logikę przekonań z operatorem **B** (*believes*), ale bez operatora dualnego wobec **B**. Do języka wprowadzamy natomiast nowy operator **Ś**; formułę o schemacie  $\check{S}A$  czytamy "podmiot jest świadomy tego, że *A*".

**Modelem sitowym** nazywamy czwórkę uporządkowaną:

$$(\$ \$) \quad \langle W, R, v, f \rangle$$

gdzie  $W$  jest niepustym zbiorem ("światów możliwych"),  $R \subseteq W \times W$  jest relacją epistemicznej alternatywności, która jest seryjna, przechodnia i euklidesowa w zbiorze  $W$ , a  $v$  jest funkcją wartościowania, przyporządkowującą każdej parze postaci  $\langle p_i, w \rangle$ , gdzie  $p_i$  jest zmienną zdaniową oraz  $w \in W$ , dokładnie jedną z wartości logicznych:  $1$  lub  $0$ , natomiast  $f$  jest funkcją o argumentach w zbiorze  $W$  i wartościach będących zbiorami formuł (intuicyjnie: funkcją, przyporządkowującą światom zbior tych wszystkich formuł, których podmiot jest świadomy w tych światach), spełniającą następujący warunek:

- ❖ dla każdego  $w \in W$ : elementami  $f(w)$  są wszystkie podstawienia następujących formuł:
  - $\mathbf{B}p \rightarrow \neg \mathbf{B}\neg p$
  - $\mathbf{B}p \rightarrow \mathbf{B}\mathbf{B}p$
  - $\neg \mathbf{B}p \rightarrow \mathbf{B}\neg \mathbf{B}p$ .

**Komentarz:** oczywisty, ale i tak zapraszam na wykład :)

**Prawdziwość** formuły  $A$  w świecie  $w$  modelu sitowego  $M = \langle W, R, v, f \rangle$ , symbolicznie:  $M \models^w A$ , definiujemy następująco:

- (1)  $M \models^w p_i$  wtw  $v(p_i, w) = 1$ ;
- (2)  $M \models^w \neg A$  wtw  $M \text{ non } \models^w A$ ;
- (3)  $M \models^w (A \wedge B)$  wtw  $M \models^w A$  oraz  $M \models^w B$ ;
- (4)  $M \models^w (A \vee B)$  wtw  $M \models^w A$  lub  $M \models^w B$ ;
- (5)  $M \models^w (A \rightarrow B)$  wtw  $M \text{ non } \models^w A$  lub  $M \models^w B$ ;
- (6)  $M \models^w (A \leftrightarrow B)$  wtw  $M \models^w A$  zawsze i tylko, gdy  $M \models^w B$ ;
- (7)  $M \models^w \acute{S}A$  wtw  $A \in f(w)$ ;
- (8)  $M \models^w \mathbf{B}A$  wtw  $A \in f(w)$  oraz dla każdego  $w^* \in W$  takiego, że  $wRw^*$ :  $M \models^{w^*} A$ .

**Prawdziwość w modelu** sitowym to prawdziwość w każdym świecie tego modelu. **Tezami** logiki epistemicznej wyznaczonej przez (rozważaną) klasę modeli sitowych są formuły prawdziwe w każdym modelu sitowym.

Jak łatwo zauważyć, tezami analizowanej logiki są m.in. wszystkie aksjomaty rachunkowozdaniowe oraz wszystkie podstawienia formuł **D<sub>B</sub>**, **4<sub>B</sub>** i **5<sub>B</sub>**. W logice tej *nie zachodzą* natomiast zależności o schematach **LO1-LO5** i **LO7-LO10**. Dlaczego? Cóż, zapraszam na wykład :)

\*\*\*

Poza przedstawionymi wyżej, istnieje rzecz jasna wiele innych metod "blokowania" paradoksu wszechwiedzy logicznej – oraz wiele logik epistemicznych, w których nie występują co najmniej niektóre formy tego paradoksu. Logiki te są jednak zwykle bardziej skomplikowane – syntaktycznie i/lub semantycznie – niż logiki epistemiczne przedstawiane na tym cyklu wykładów.

Informacje na temat wielu z takich logik mogą Państwo zaczerpnąć np. z artykułu przeglądowego Gocheta i Gribomonta "Epistemic Logic".