

Deception in Epistemic Causal Logic

Chiaki Sakama

Wakayama University Japan

August 19, 2021

What is deception?

To cause to believe what is false (OED 1989)

What kind of logic is needed for formal account of deception?

To **cause** to believe what is false (OED 1989)

- a Logic needs to express **causality** between a deceptive act and its effect.

What kind of logic is needed for formal account of deception?

To cause to believe what is false (OED 1989)

- a Logic needs to express **causality** between a deceptive act and its effect.
- a Logic needs to express **belief states** of agents.

What kind of logic is needed for formal account of deception?

To cause to believe what is false (OED 1989)

- a logic needs to express causality between a deceptive act and its effect.
- a logic needs to express belief states of agents.

For this purpose, we use

epistemic causal logic

= causal logic [Giunchiglia, et al., AIJ 2004]

+ belief modality

Epistemic Causal Logic: Definitions

causal rule

A **causal rule** is of the form:

$$\phi \Rightarrow \psi \quad (\phi, \psi : \text{propositional formula})$$

meaning “ ψ is caused if ϕ is true.”

A **(causal) theory** is a finite set of causal rules.

model

Given a theory T and an interpretation I , define

$$T^I = \{ \psi \mid (\phi \Rightarrow \psi) \in T \text{ for some } \phi \text{ and } I \models \phi \}.$$

I is a **model** of T if I is the unique model of T^I .

If every model of T satisfies a formula F , written $T \models F$.

If T has no model, written $T \models \perp$.

Epistemic Causal Logic: Axioms

$U_{ab}^t\phi$: an agent a utters a sentence ϕ to an agent b at time t .

$B_a^t\phi$: an agent a believes a sentence ϕ at time t

Axioms for utterance and beliefs

(axioms of utterance): $U_{ab}^t\phi \Rightarrow U_{ab}^t\phi$ and $\neg U_{ab}^t\phi \Rightarrow \neg U_{ab}^t\phi$.

(axioms of belief): $B_a^t\phi \Rightarrow B_a^t\phi$ and $\neg B_a^t\phi \Rightarrow \neg B_a^t\phi$.

$B_a^t\phi \equiv B_a^t\psi$ if $\phi \equiv \psi$.

$B_a^t(\phi \wedge \psi) \equiv B_a^t\phi \wedge B_a^t\psi$.

(axioms of inertia): $B_a^t\phi \wedge B_a^{t+1}\phi \Rightarrow B_a^{t+1}\phi$

$\neg B_a^t\phi \wedge \neg B_a^{t+1}\phi \Rightarrow \neg B_a^{t+1}\phi$.

(axiom of truth): $B_a^t\top$ for any t .

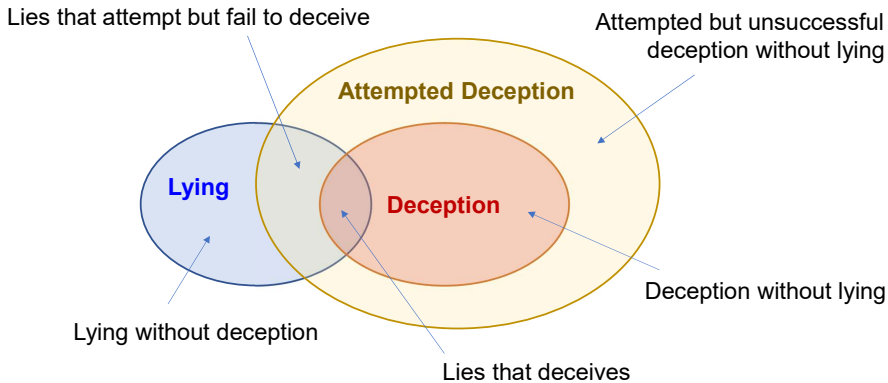
(axiom of rationality): $\neg B_a^t\perp$ for any t if a is *rational*.

(axiom of credibility): $U_{ab}^t\phi \Rightarrow B_b^{t+1}\phi$ if b is *credulous*.

(axiom of reflection): $U_{ab}^t\phi \Rightarrow B_b^{t+1}B_a^t\phi$ if b is *reflective*.

Lying, deception and attempted deception

(Carson, T.L. "Lying and Deception: Theory and Practice", 2010)



Deception by Lying

Lying (a : speaker, b : hearer, ϕ : sentence)

$$\text{LIE}_{ab}^t(\phi) \stackrel{\text{def}}{=} B_a^t \neg \phi \wedge U_{ab}^t \phi$$

(a lies to b if a utters a believed-false sentence ϕ to b at t)

Deception by Lying

$$\text{DBL}_{ab}^{t+1}(\phi) \stackrel{\text{def}}{=} \neg \phi \wedge (\text{LIE}_{ab}^t(\phi) \Rightarrow B_b^{t+1} \phi).$$

(a lies to b at t on a false sentence ϕ , which causes b 's believing ϕ at the next time step $t + 1$)

Note

- In lying, a speaker a believes ϕ but the actual falsity of ϕ is not required.
- In DBL, the actual falsity of ϕ is required.

Properties of DBL (1)

DBL does not happen if a sentence ϕ is true

$$\phi \wedge \text{DBL}_{ab}^{t+1}(\phi) \models \perp$$

Lying on a false sentence succeeds to deceive if a hearer is credulous, i.e., $U_{ab}^t \phi \Rightarrow B_b^{t+1} \phi$

$$\neg \phi \wedge \text{LIE}_{ab}^t(\phi) \models B_b^{t+1} \phi$$

DBL on the valid sentence always fails

$$\text{DBL}_{ab}^{t+1}(\top) \models \perp$$

DBL on the contradictory sentence fails if a hearer is rational, i.e., $\neg B_b^t \perp$

$$\text{LIE}_{ab}^t(\perp) \wedge \text{DBL}_{ab}^{t+1}(\perp) \models \perp \quad \text{if } b \text{ is rational.}$$

Properties of DBL (2)

DBL fails if a rational hearer believes the contrary

$$B_b^t \neg \phi \wedge \text{LIE}_{ab}^t(\phi) \wedge \text{DBL}_{ab}^{t+1}(\phi) \models \perp \text{ if } b \text{ is rational.}$$

If a rational hearer is credulous, DBL succeeds even if the hearer believes the contrary

$$\neg \phi \wedge B_b^t \neg \phi \wedge \text{LIE}_{ab}^t(\phi) \models B_b^{t+1} \phi$$

if b is credulous and rational.

A hearer b does not believe that a speaker a is lying if b is rational, reflective (i.e., $U_{ab}^t \phi \Rightarrow B_b^{t+1} B_a^t \phi$) and believes that a is also rational.

$$\text{LIE}_{ab}^t(\phi) \wedge B_b^{t+1}(\text{LIE}_{ab}^t(\phi)) \models \perp \text{ if } b \text{ is rational, reflective, and believes that a speaker } a \text{ is rational.}$$

Deception without Lying

Deception by Bluffing

$$\text{DBB}_{ab}^{t+1}(\phi) \stackrel{\text{def}}{=} \neg\phi \wedge (\text{BLUF}_{ab}^t(\phi) \Rightarrow B_b^{t+1}\phi)$$

where $\text{BLUF}_{ab}^t(\phi) \stackrel{\text{def}}{=} \neg B_a^t\phi \wedge \neg B_a^t\neg\phi \wedge U_{ab}^t\phi$

Deception by Truthful Telling

$$\text{DBT}_{ab}^{t+1}(\phi) \stackrel{\text{def}}{=} \neg\phi \wedge (\text{TRT}_{ab}^t(\phi) \Rightarrow B_b^{t+1}\phi)$$

where $\text{TRT}_{ab}^t(\phi) \stackrel{\text{def}}{=} B_a^t\phi \wedge U_{ab}^t\phi$

Deception by Omission (or withholding information)

$$\text{DBO}_{ab}^{t+1}(\phi) \stackrel{\text{def}}{=} \phi \wedge (\text{WI}_{ab}^t(\phi) \Rightarrow \neg B_b^{t+1}\phi)$$

where $\text{WI}_{ab}^t(\phi) \stackrel{\text{def}}{=} B_a^t\phi \wedge \neg U_{ab}^t\phi$

Intentional Deception

Intentional DBL and DBB

- I-DBL $_{ab}^{t+1}(\phi) \stackrel{def}{=} \neg\phi \wedge (\text{LIE}_{ab}^t(\phi) \wedge B_a^t B_b^{t+1}\phi \Rightarrow B_b^{t+1}\phi)$
- I-DBB $_{ab}^{t+1}(\phi) \stackrel{def}{=} \neg\phi \wedge (\text{BLUF}_{ab}^t(\phi) \wedge B_a^t B_b^{t+1}\phi \Rightarrow B_b^{t+1}\phi)$

By $B_a^t B_b^{t+1}\phi$, a speaker a believes that a hearer b will believe the false sentence ϕ in the next time step.

Intentional Deception

Intentional DBL and DBB

- I-DBL $_{ab}^{t+1}(\phi) \stackrel{def}{=} \neg\phi \wedge (\text{LIE}_{ab}^t(\phi) \wedge B_a^t B_b^{t+1} \phi \Rightarrow B_b^{t+1} \phi)$
- I-DBB $_{ab}^{t+1}(\phi) \stackrel{def}{=} \neg\phi \wedge (\text{BLUF}_{ab}^t(\phi) \wedge B_a^t B_b^{t+1} \phi \Rightarrow B_b^{t+1} \phi)$

Intentional DBT

$$\text{I-DBT}_{ab}^{t+1}(\phi, \psi) \stackrel{def}{=} \neg\psi \wedge (\text{TRT}_{ab}^t(\phi) \wedge B_a^t (B_b^{t+1} \phi \supset B_b^{t+1} \psi) \wedge B_a^t \neg\psi \Rightarrow B_b^{t+1} \psi)$$

A speaker a truthfully tells ϕ while a believes that a hearer b 's believing ϕ leads to b 's believing another false sentence ψ in the next time step.

Intentional Deception

Intentional DBL and DBB

- I-DBL $_{ab}^{t+1}(\phi) \stackrel{def}{=} \neg\phi \wedge (\text{LIE}_{ab}^t(\phi) \wedge B_a^t B_b^{t+1}\phi \Rightarrow B_b^{t+1}\phi)$
- I-DBB $_{ab}^{t+1}(\phi) \stackrel{def}{=} \neg\phi \wedge (\text{BLUF}_{ab}^t(\phi) \wedge B_a^t B_b^{t+1}\phi \Rightarrow B_b^{t+1}\phi)$

Intentional DBT

$$\text{I-DBT}_{ab}^{t+1}(\phi, \psi) \stackrel{def}{=} \neg\psi \wedge (\text{TRT}_{ab}^t(\phi) \wedge B_a^t (B_b^{t+1}\phi \supset B_b^{t+1}\psi) \wedge B_a^t \neg\psi \Rightarrow B_b^{t+1}\psi)$$

Intentional DBO

$$\text{I-DBO}_{ab}^{t+1}(\phi) \stackrel{def}{=} \phi \wedge (\text{WI}_{ab}^t(\phi) \wedge B_a^t \neg B_b^t \phi \Rightarrow \neg B_b^{t+1}\phi)$$

A speaker a withholds ϕ while believing b 's ignorance of ϕ , which causes b 's disbelieving ϕ (or prevents b from believing ϕ) in the next time step.

Indirect Deception

- $\text{IN-DBL}_{ac}(\phi) \stackrel{def}{=} (\text{I-})\text{DBL}_{ab}^{t+1}(\phi) \wedge \text{DBT}_{bc}^{t+2}(\phi)$
(a 's lying on ϕ results in b 's believing a false sentence ϕ , and then b 's truthful telling on ϕ results in c 's believing ϕ)
- $\text{IN-DBB}_{ac}(\phi) \stackrel{def}{=} (\text{I-})\text{DBB}_{ab}^{t+1}(\phi) \wedge \text{DBT}_{bc}^{t+2}(\phi)$
- $\text{IN-DBT}_{ac}(\phi) \stackrel{def}{=} \text{DBT}_{ab}^{t+1}(\phi) \wedge \text{DBT}_{bc}^{t+2}(\phi)$
- $\text{IN-DBO}_{ac}(\phi) \stackrel{def}{=} (\text{I-})\text{DBO}_{ab}^{t+1}(\phi) \wedge \neg U_{bc}^{t+1}\phi \Rightarrow \neg B_c^{t+2}\phi$
(a 's withholding ϕ results in b 's disbelieving a true sentence ϕ . Then b does not inform c of ϕ , which results in c 's disbelieving ϕ)
- $\text{IN-I-DBT}_{ac}(\phi, \psi) \stackrel{def}{=} \text{I-DBT}_{ab}^{t+1}(\phi, \psi) \wedge \text{DBT}_{bc}^{t+2}(\psi)$

Self-Deception

Self-deception by lying produces contradictory belief

$$\text{LIE}_{aa}^t(\phi) \wedge (\text{I-})\text{DBL}_{aa}^{t+1}(\phi) \models B_a^{t+1} \perp$$

$B_a^t \neg \phi$ in $\text{LIE}_{aa}^t(\phi)$ implies $B_a^{t+1} \neg \phi$ by the axioms of inertia.
 $\text{DBL}_{aa}^{t+1}(\phi)$ implies $B_a^{t+1} \phi$. Then, $B_a^{t+1} \neg \phi \wedge B_a^{t+1} \phi \equiv B_a^{t+1} \perp$

If a rational agent is credulous, self-DBL does not involve contradictory belief

$$\text{LIE}_{aa}^t(\phi) \wedge (\text{I-})\text{DBL}_{aa}^{t+1}(\phi) \not\models B_a^{t+1} \perp$$

if a is credulous and rational

A credulous agent revises its belief from $B_a^t \neg \phi$ to $B_a^{t+1} \phi$.
 $B_a^{t+1} \phi$ implies $\neg B_a^{t+1} \neg \phi$ by the axiom of rationality. Then
the axioms of inertia do not produce $B_a^{t+1} \neg \phi$ from $B_a^t \neg \phi$.

Self-Deception

Self-deception by bluffing does not produce contradictory belief

$$\text{BLUF}_{aa}^t(\phi) \wedge (\text{I-})\text{DBB}_{aa}^{t+1}(\phi) \models \neg B_a^t \phi \wedge B_a^{t+1} \phi$$

Self-deception by omission does not produce contradictory belief

$$\text{WI}_{aa}^t(\phi) \wedge (\text{I-})\text{DBO}_{aa}^{t+1}(\phi) \models B_a^t \phi \wedge \neg B_a^{t+1} \phi$$

(a person, who believes something true but does not refer to it, will *forget* it.)

Self-deception by truthful-telling does not contradict while inconsistency arises if accompanied by intention

- $\text{TRT}_{aa}^t(\phi) \wedge \text{DBT}_{aa}^{t+1}(\phi) \not\models B_a^{t+1} \perp$
- $\text{TRT}_{aa}^t(\phi) \wedge \text{I-DBT}_{aa}^{t+1}(\phi, \psi) \models B_a^{t+1} \perp$

Summary

- Different types of deception are formulated using epistemic causal logic.
- From the computational perspective, a causal theory handled in this study is translated into a logic program under the answer set semantics.
- The current framework is extended to handle more complicated cases by taking *a theory of mind* into consideration.