



Department of Philosophy, Tsinghua University

# A Logic for Desire Based on Causal Inference

Kaibo Xie & Jialiang Yan  
xiekaibozju@gmail.com  
yan-ji19@mails.tsinghua.edu.cn

8th International Conference on Logic, Rationality and  
Interaction (LORI-VIII)  
October 16, 2021

## ① Introduction

1.1 Motivation

1.2 Why the Desire is Not For Its Own Sake

1.3 Why the Causality is Needed

## ② The Logic ID

2.1 Causal Modelling

2.2 Language and Semantics of Logic ID

2.3 The Property of Logic ID

**Desire** is the primitive component that drives the agent's behavior. Theories of desire have been established in many disciplines. We study the reasoning about desire in order to predict and characterize tendency of the agent's action.



In this paper, we propose an interpretation of “desiring  $\phi$ ” based on all of the most relevant alternatives that satisfies  $\phi$  which are preferred by the agent. The relevance is evaluated by the causality. And then a logic based on causal inferences is given.

Having a desire means, as we explain, this desire can bring about what the agent most prefers at the moment.

There are two main feature capturing the desire in our proposal:

- What people desire is not for its own sake.
- Causality should be taken into consideration.

*Robin was in charge of an important project in the company recently, and the tremendous pressure made him sleepless every night. Robin is a healthy - conscious person and poor sleep will make him sluggish the next day which bothers Robin most. So having a good sleep is what he needs most at the moment. Therefore, on this sleepless night, Robin wants to take some sleep-helping pills to help him sleep well, even if these pills have some side effects.*

- $P$ : “pills are taken”
- $S$ : “Robin sleeps well”

Now assume that there are four possible worlds:

- $w_1: P \wedge \neg S$
- $w_2: \neg P \wedge S$
- $w_3: P \wedge S$
- $w_4: \neg P \wedge \neg S$

The preference structure can be characterized by the following preference order:

$$w_1 < w_4 < w_3 < w_2$$

*Robin's company offered him a chance to study in the Netherlands and he attached great importance to it. But his savings are small, so he doesn't want to pay for the tuition. Unfortunately, due to the epidemic, he cannot go abroad. Even after paying the tuition, he can only stay at home and take online lessons.*

(The original version of this example comes from Maria Aloni.)

- $P'$ : “paying the tuition”
- $S'$ : “studying in Netherlands”

Now assume that there are four possible worlds:

- $w'_1: P' \wedge \neg S'$
- $w'_2: \neg P' \wedge S'$
- $w'_3: P' \wedge S'$
- $w'_4: \neg P' \wedge \neg S'$

The preference structure can be characterized by the following preference order:

$$w'_1 < w'_4 < w'_3 < w'_2$$



To capture these two features of instrumental desire, we introduce a **desire-causality model** in our paper.

- Following the traditional approach of formalizing preference, e.g. [von Wright, 1972] and [van Benthem & Liu, 2007], we use a pre order over all the possible worlds to represent the preference structure of an agent.
- We adopt a causal model which makes use of the interventionist approach to causality from [Halpern, 2000] and [Pearl, 2002].

Halpern (2000) and Pearl (2002)

## Causal Variables

- $\mathcal{U} = \{U_1, \dots, U_m\}$  is a set of *exogenous* variables,
- $\mathcal{V} = \{V_1, \dots, V_n\}$  is a set of *endogenous* variables
- $\mathcal{R}(X)$  is the non-empty range of the variable  $X \in \mathcal{U} \cup \mathcal{V}$ .

## Desire-Causality Model

Let  $\mathcal{S} = \langle \mathcal{U}, \mathcal{V}, \mathcal{R} \rangle$ . A desire-causality model for  $\mathcal{S}$  is a tuple  $\langle \mathcal{F}, \mathcal{A}, \leq \rangle$ .

- **Structural Functions**  $\mathcal{F} = \{\mathcal{F}_{V_j} \mid V_j \in \mathcal{V}\}$ . For each endogenous variable  $V_j$ ,  $\mathcal{F}_{V_j}$  is a mapping from all assignments to  $\mathcal{U} \cup \mathcal{V} \setminus \{V_j\}$  to  $\mathcal{R}(V_j)$ .  $\mathcal{F}$  is assumed to be recursive.
- **Valuation Function**  $\mathcal{A}$  assigns to every  $X \in \mathcal{U} \cup \mathcal{V}$  a value  $\mathcal{A}(X) \in \mathcal{R}(X)$ .  $\mathcal{A}$  has to *comply with*  $\mathcal{F}_{V_j}$ .
- **Preference Ordering**  $\leq$  is a pre order over all possible assignments to  $\mathcal{U} \cup \mathcal{V}$ .

## Language of Logic ID

Formulas  $\phi$  of the language  $\mathcal{L}$  based on  $\mathcal{S}$

$$\phi ::= X=x \mid \neg\phi \mid \phi \wedge \phi \mid D(\vec{X}=\vec{x}) \mid (\vec{X}=\vec{x}) \Box \rightarrow \phi$$

The semantics of “desire  $X = x$ ”: after an intervention forcing  $X = x$  to be true, the world results from it will be more preferred than the current world.

## Intervention on desire-causality models

Let  $M = \langle \mathcal{F}, \mathcal{A}, \leq \rangle$  be a DC-model based on  $\mathcal{S}$ .

$M_{\vec{X}=\vec{x}} = \langle \mathcal{F}_{\vec{X}=\vec{x}}, \mathcal{A}_{\vec{X}=\vec{x}}^{\mathcal{F}}, \leq \rangle$  is the DC model resulting from an intervention setting the values of variables in  $X_1, \dots, X_n$  to  $x_1, \dots, x_n$ :

- $\mathcal{F}_{\vec{X}=\vec{x}}$  is as  $\mathcal{F}$  except that, for each endogenous variable  $X_i$  in  $\vec{X}$ , the function  $f_{X_i}$  is replaced by a *constant* function  $f'_{X_i}$  that returns the value  $x_i$  regardless of the values of all other variables.
- $\mathcal{A}_{\vec{X}=\vec{x}}^{\mathcal{F}}$  is the updated solution to the updated structural equations.

Let  $\langle \mathcal{F}, \mathcal{A}, \leq \rangle$  be a DC-model based on  $\mathcal{S}$

## Truth condition of formulas in $\mathcal{L}$

- $\langle \mathcal{F}, \mathcal{A}, \leq \rangle \models X=x$  iff  $\mathcal{A}(X) = x$
- $\langle \mathcal{F}, \mathcal{A}, \leq \rangle \models (\vec{X}=\vec{x}) \Box \rightarrow \phi$  iff  $\langle \mathcal{F}_{\vec{X}=\vec{x}}, \mathcal{A}_{\vec{X}=\vec{x}}^{\mathcal{F}}, \leq \rangle \models \phi$
- $\langle \mathcal{F}, \mathcal{A}, \leq \rangle \models D(\vec{X}=\vec{x})$  iff  $\mathcal{A}^{\mathcal{F}} < \mathcal{A}_{\vec{X}=\vec{x}}^{\mathcal{F}}$

$X$  stands for taking sleep-helping pills.  $Y$  stands for having a good sleep.

$X$  is an exogenous variable.

$\mathcal{F}_Y$ : the value of  $Y$  is equal to the value of  $X$

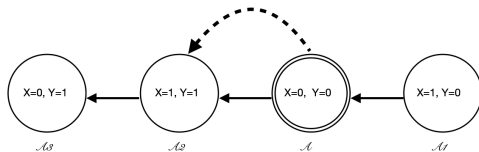


Figure: Robin's desire of pills

- $\mathcal{A}(X) = 0, \mathcal{A}(Y) = 0$
- According to  $\mathcal{F}_{X=1}$ ,  $\mathcal{A}_{X=1}(X) = 1, \mathcal{A}_{X=1}(Y) = 1$

$X$  stands for taking sleep-helping pills.  $Y$  stands for having a good sleep.

$X$  is an exogenous variable.

$\mathcal{F}_Y$ : the value of  $Y$  is equal to the value of  $X$

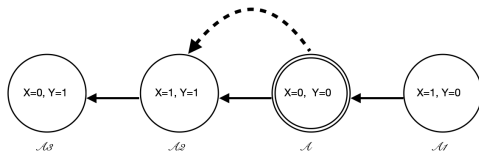


Figure: Robin's desire of pills

- $\mathcal{A}(X) = 0, \mathcal{A}(Y) = 0$
- According to  $\mathcal{F}_{X=1}$ ,  $\mathcal{A}_{X=1}(X) = 1, \mathcal{A}_{X=1}(Y) = 1$
- $\mathcal{A} < \mathcal{A}_{X=1}$ , therefore  $\langle \mathcal{F}, \mathcal{A}, \leq \rangle \models D(X = 1)$



$S$  stands for studying in Netherlands;  $T$  stands for paying tuition fee.  $S$  is an exogenous variable.

$\mathcal{F}_T$ : the value of  $T$  is equal to the value of  $S$

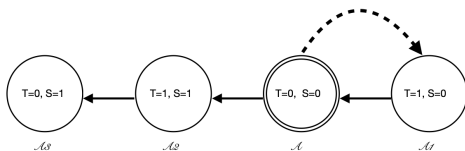


Figure: Robin's desire of tuition

- $\mathcal{A}(S) = 0, \mathcal{A}(T) = 0$ ;
- According to  $\mathcal{F}_{T=1}$ ,  $\mathcal{A}_{T=1}(T) = 1, \mathcal{A}_{T=1}(S) = 0$

$S$  stands for studying in Netherlands;  $T$  stands for paying tuition fee.  $S$  is an exogenous variable.

$\mathcal{F}_T$ : the value of  $T$  is equal to the value of  $S$

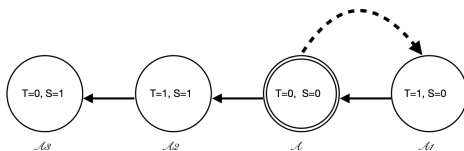


Figure: Robin's desire of tuition

- $\mathcal{A}(S) = 0, \mathcal{A}(T) = 0$ ;
- According to  $\mathcal{F}_{T=1}$ ,  $\mathcal{A}_{T=1}(T) = 1, \mathcal{A}_{T=1}(S) = 0$
- $\mathcal{A}_{S=1} < \mathcal{A}$ , therefore  $\langle \mathcal{F}, \mathcal{A}, \leq \rangle \models \neg D(T = 1)$

## Valid formulas of Logic ID

- $\vec{X}=\vec{x} \square \rightarrow D(\vec{Y}=\vec{y}) \rightarrow D(\vec{X}=\vec{x} \wedge \vec{Y}=\vec{y})$ , if  $\vec{X}$  and  $\vec{Y}$  are disjoint
- $(\vec{X}=\vec{x} \square \rightarrow \vec{Y}=\vec{y}) \wedge D(\vec{X}=\vec{x}) \rightarrow D(\vec{X}=\vec{x} \wedge \vec{Y}=\vec{y})$ , if  $\vec{X}$  and  $\vec{Y}$  are disjoint

## Axioms of Logic ID:

- (A modified version of) Axioms for counterfactuals based on causal models in Halpern (2000)
- Reduction axioms for counterfactuals
- Axiom for the interaction between desire and conditionals

## Axioms of Logic ID:

- (A modified version of) Axioms for counterfactuals based on causal models in Halpern (2000)
- Reduction axioms for counterfactuals
- Axiom for the interaction between desire and conditionals

|                |   |  |
|----------------|---|--|
| P              | $\varphi$   | for $\varphi$ an instance of a propositional tautology   |
| MP             | From $\varphi_1$ and $\varphi_1 \rightarrow \varphi_2$ infer $\varphi_2$  |  |
| Generalization | From $\phi$ infer $\vec{X} = \vec{x} \square \rightarrow \phi$  |  |
| A <sub>1</sub> | $\vec{X} = \vec{x} \square \rightarrow Y = y \rightarrow \neg \vec{X} = \vec{x} \square \rightarrow Y = y'$   | for $y, y' \in \mathcal{R}(Y)$ with $y \neq y'$  |
| A <sub>2</sub> | $\bigvee_{y \in \mathcal{R}(Y)} \vec{X} = \vec{x} \square \rightarrow Y = y$  |  |
| A <sub>3</sub> | $(\vec{X} = \vec{x} \square \rightarrow (Y = y) \wedge \vec{X} = \vec{x} \square \rightarrow (Z = z)) \rightarrow \vec{X} = \vec{x}, Y = y \square \rightarrow (Z = z)$                       |  |
| A <sub>4</sub> | $\vec{X} = \vec{x}, Y = y \square \rightarrow (Y = y)$  |  |
| A <sub>5</sub> | $(\vec{X} = \vec{x}, Y = y \square \rightarrow (Z = z) \wedge \vec{X} = \vec{x}, Z = z \square \rightarrow (Y = y)) \rightarrow \vec{X} = \vec{x} \square \rightarrow (Z = z)$                | for $Y \neq Z$   |
| A <sub>6</sub> | $(X_0 \rightsquigarrow X_1 \wedge \dots \wedge X_{k-1} \rightsquigarrow X_k) \rightarrow \neg(X_k \rightsquigarrow X_0)$  |  |
| A <sub>7</sub> | $U = u \leftrightarrow \vec{X} = \vec{x} \square \rightarrow U = u$   | for $U \in \mathcal{U}$ with $U \notin \vec{X}$  |
| A <sub>¬</sub> | $\vec{X} = \vec{x} \square \rightarrow \neg \varphi \leftrightarrow \neg \vec{X} = \vec{x} \square \rightarrow \varphi$   |  |
| A <sub>∧</sub> | $\vec{X} = \vec{x} \square \rightarrow (\varphi_1 \wedge \varphi_2) \leftrightarrow (\vec{X} = \vec{x} \square \rightarrow \varphi_1 \wedge \vec{X} = \vec{x} \square \rightarrow \varphi_2)$ |  |
| A <sub>∅</sub> | $\vec{X} = \vec{x} \square \rightarrow (\vec{Y} = \vec{y} \square \rightarrow \varphi) \leftrightarrow \vec{X}' = \vec{x}', \vec{Y} = \vec{y} \square \rightarrow \varphi$                    | with $\vec{X}' = \vec{x}'$ the<br>subassignment of $\vec{X} = \vec{x}$ to<br>$\vec{X}' := \vec{X} \setminus \vec{Y}$ |

---

|                       |  |  |
|-----------------------|--|--|
| <b>A<sub>8</sub></b>  | $\vec{X}=\vec{x} \rightarrow \neg D(\vec{X}=\vec{x})$  |  |
| <b>A<sub>9</sub></b>  | $(\vec{X}=\vec{x} \square \rightarrow Z=z) \rightarrow (\vec{X}=\vec{x} \square \rightarrow D(\vec{Y}=\vec{y})) \leftrightarrow (\vec{X}=\vec{x}, Z=z \square \rightarrow D(\vec{Y}=\vec{y}))$ | if $\vec{Y}=\vec{y}$ is a full assignment<br>of $\mathcal{U} \cup \mathcal{V}$ and $Z \notin \vec{X}$              |
| <b>A<sub>10</sub></b> | $(\vec{X}=\vec{x} \square \rightarrow \vec{Y}=\vec{y}) \rightarrow (D(\vec{X}=\vec{x}) \leftrightarrow D(\vec{X}=\vec{x} \wedge \vec{Y}'=\vec{y}'))$   | with $\vec{Y}'=\vec{y}'$ the<br>sub-assignment of $\vec{Y}=\vec{y}$ for<br>$\vec{Y}' := \vec{Y} \setminus \vec{X}$ |
| <b>A<sub>11</sub></b> | $(\vec{X}=\vec{x} \square \rightarrow D(\vec{Y}=\vec{y})) \wedge D(\vec{X}=\vec{x}) \rightarrow D(\vec{X}'=\vec{x}', \vec{Y}=\vec{y})$   | with $\vec{X}'=\vec{x}'$ the<br>sub-assignment of $\vec{X}=\vec{x}$ for<br>$\vec{X}' := \vec{X} \setminus \vec{Y}$ |

---



## Semantics of counterfactuals with complex antecedent

Briggs (2012):  $\chi \Box \rightarrow \phi$  holds iff  $\phi$  holds at every submodels generated by interventions on  $M$  corresponding to states that verify  $\chi$



## Semantics of counterfactuals with complex antecedent

Briggs (2012):  $\chi \Box \rightarrow \phi$  holds iff  $\phi$  holds at every submodels generated by interventions on  $M$  corresponding to states that verify  $\chi$

## Semantics of desire for complex propositions

Briggs (2012):  $D(\chi)$  holds iff every assignment in submodels generated by interventions on  $M$  corresponding to states that verify  $\chi$  is preferred.

$$(D(\phi_1 \vee \phi_2)) \leftrightarrow (D\phi_1 \wedge D\phi_2 \wedge (D(\phi_1 \wedge \phi_2)))$$

**THANK YOU!**