# Probabilistic Multi-agent Only-believing

**Qihui Feng**, Gerhard Lakemeyer
RWTH Aachen University

May 13, 2024

# Overview

# Introduction

Knowledge and Belief:

- $K(fair(\text{Coin}) \wedge \neg fair(\text{Die}))$
- $B(fair(\text{Coin}) : 0.8)$

Levesque proposed **only-knowing** to precisely capture the (non-)beliefs:

- $O(fair(\text{Coin})) \models \neg K(fair(\text{Die})) \wedge \neg K(\neg fair(\text{Die}))$
- $O(fair(\text{Coin})) \models \neg B(fair(\text{Die}) : r)$ for any $r \in [0, 1]$

Research on only-knowing:

- Probabilistic only-believing: The logic $\mathcal{OBL}$
- Projection reasoning: $O(KB_1) \rightarrow [action]O(KB_2)$
- Multi-agent only-knowing:
  - Previous works in both propositional and first-order cases
  - No first-order account faithfully follows Levesque's principle of only-knowing.

# Introduction (cont'd)

Levesque's notion of only-knowing: Given the only-knowing of the agent, any subjective formula will either be inferred or disproved.

- $O(\text{KB}) \models K\beta$ iff $\text{KB} \models \beta$; $\quad O(\text{KB}) \models \neg K\beta$ iff $\text{KB} \not\models \beta$.

Desiderata for multi-agent extension:

i) non-beliefs on irrelevant items:
$O_a(\neg \textit{fair}(\text{Coin}) \wedge K_b(\textit{fair}(\text{Coin}))) \models \neg K_a K_b(\textit{fair}(\text{Die}))$

ii) non-beliefs on mental states with deeper nesting
$O_a(\neg \textit{fair}(\text{Coin}) \wedge K_b(\textit{fair}(\text{Coin}))) \models \neg K_a K_b K_a(\textit{fair}(\text{Coin}))$

To model only-knowing up to all depths is semantically difficult.

To model only-knowing up to depth $k$?

- Modality $O_a^{(k)}$: agent $a$'s only-knowing(believing) up to depth $k$.

The new desiderata: for $K_a\beta$ with depth no more than $k$,

- Either $O_a^{(k)}\alpha \models K_a\beta$ or $O_a^{(k)}\alpha \models \neg K_a\beta$?
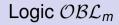
# Overview

# Logic $\mathcal{OBL}_m$

A first-order modal logic with equality and featured with:

- A finite set of agents, e.g. $Ag = \{a, b\}$;
- Modalities for belief and only-believing for each agent.

# Syntax

- Standard FO formulae;
- $B_i(\alpha\colon r)$: $\alpha$ is believed by agent $i$ with degree $r$ (written $K_i\alpha$ if $r = 1$).
- $O_i^{(k)}(\alpha\colon r)$: all that agent $i$ believes up to depth $k$ is $\alpha$ with degree $r$.
  - $O_a^{(1)}(\textit{fair}(\text{Coin}))$: agent $a$ knows $\textit{fair}(\text{Coin})$ and nothing else about the world.
  - $O_a^{(2)}(\textit{fair}(\text{Coin}))$: agent $a$ knows $\textit{fair}(\text{Coin})$ and nothing else about the world, and nothing about Bob's beliefs about the world.

# Semantics (Knowledge and Beliefs)

A model is a tuple $(w, e_a, e_b)$ with world $w$ and epistemic states $e_a$ and $e_b$.

A **world** $w \in \mathcal{W}$ is a set of ground atoms.

- $w, e_a, e_b \models \textit{fair}(\text{Coin})$ iff $\textit{fair}(\text{Coin}) \in w$

**Epistemic states** are defined inductively:

- **1-distribution** assigns each world a probability: $\mathcal{W} \to \mathbb{R}_{[0,1]}$;
- **1-epistemic state** is a set of 1-distributions.

### Example

$w_1 = \{\textit{fair}(\text{Coin})\}$ and $w_2 = \{\textit{fair}(\text{Die})\}$, $d^1(w) = \begin{cases} 0.5 & w \in \{w_1, w_2\} \\ 0 & \textit{otherwise}. \end{cases}$

Let $e_a = \{d^1\}$, then $w, e_a, e_b \models B_a(\textit{fair}(\text{Coin}) \colon 0.5)$

## Semantics (Nested Beliefs)

$\mathcal{E}^1$ denotes the set of all 1-epistemic states. For any $k > 1$,

- **$k$-distribution** $d^k : (\mathcal{W} \times \mathcal{E}^{k-1}) \to \mathbb{R}_{[0,1]}$
- a **$k$-epistemic state** is a set of $k$-distributions

### Example

Let $w_1 = \{fair(\text{Coin})\}$, $w_2 = \{fair(\text{Die})\}$, $w_3 = \{fair(\text{Coin}), fair(\text{Die})\}$.

$$\tilde{d}^1(w) = \begin{cases} 0.5 & w \in \{w_1, w_3\} \\ 0 & \textit{otherwise}. \end{cases} \qquad d^2(w, e_b^1) = \begin{cases} 0.3 & w = w_1, e_b^1 = \{\tilde{d}^1\} \\ 0.7 & w = w_2, e_b^1 = \{\tilde{d}^1\} \\ 0 & \textit{otherwise}. \end{cases}$$

Let $e_a = \{d^2\}$, then

- $w, e_a, e_b \models B_a(fair(\text{Coin}) \colon 0.3)$
- $w, e_a, e_b \models K_a(K_b(fair(\text{Coin})))$

# Semantics (Only-Believing)

Suppose that $e_a \in \mathcal{E}^2$ (2-epistemic state)

- $w, e_a, e_b \models O_a^{(2)}(\neg\textit{fair}(\text{Coin}) \wedge K_b(\textit{fair}(\text{Coin})))$ iff for any 2-distribution $d$,

$$d \in e_a \iff w, \{d\}, e_b \models K_a(\neg\textit{fair}(\text{Coin}) \wedge K_b(\textit{fair}(\text{Coin})))$$

i.e. there is a **unique** $e_a \in \mathcal{E}^2$ which satisfies $O_a^{(2)}(\neg\textit{fair}(\text{Coin}) \wedge K_b(\textit{fair}(\text{Coin})))$
Every $e_a \in \mathcal{E}^k$ can be **uniquely** "reduced" to an $e_a' \in \mathcal{E}^{k-1}$ s.t.

$$w, e_a, e_b \models \alpha \text{ iff } w, e_a', e_b \models \alpha \text{ for any } \alpha \text{ not deeper than k-1}$$

For $e_a \in \mathcal{E}^3$, $w, e_a, e_b \models O_a^{(2)}(\neg\textit{fair}(\text{Coin}) \wedge K_b(\textit{fair}(\text{Coin})))$ iff
$e_a$ reduce to $e_a' \in \mathcal{E}^2$ and $w, e_a', e_b \models O_a^{(2)}(\neg\textit{fair}(\text{Coin}) \wedge K_b(\textit{fair}(\text{Coin})))$

# Entailment and Validity

Compatibility:

- Formulae like $O_a^{(1)}(\neg \textit{fair}(\text{Coin}) \land \textit{\textbf{K}}_b(\textit{fair}(\text{Coin})))$ are illegal.
- *e compatible* with $\alpha$: the depth of $e$ is not less than the depth of $\alpha$

We say $\Sigma$ entails $\alpha$ (written $\Sigma \models \alpha$) iff:

- For each model $(w, e_a, e_b)$ compatible with $\Sigma, \alpha$, if $(w, e_a, e_b) \models \sigma$ for all $\sigma \in \Sigma$, then $(w, e_a, e_b) \models \alpha$

We say $\alpha$ is valid iff $\{\} \models \alpha$

# Overview

## Properties of Knowledge

$\mathcal{OBL}_m$ follows the KD45$_n$ properties. For agent $i \in Ag$,

- (Nec) If $\models \alpha$, then $\models K_i \alpha$
- (K) $\models K_i \alpha \wedge K_i(\alpha \supset \beta) \supset K_i \beta$
- (D) $\models K_i \alpha \supset \neg K_i \neg \alpha$
- (4) $\models K_i \alpha \supset K_i K_i \alpha$
- (5) $\models \neg K_i \alpha \supset K_i \neg K_i \alpha$
- $K_i \alpha \wedge \neg \alpha$ is satisfiable

Barcan formulae:

- $\models \forall x . K_i \alpha \supset K_i \forall x . \alpha$
- $\models \exists x . K_i \alpha \supset K_i \exists x . \alpha$

# Properties of Beliefs

The degree of belief follows the laws of probability:

- $\models B_i(\alpha \colon r) \supset \neg B_i(\alpha \colon r')$ for $r' \neq r$
- $\models B_i(\alpha \colon r) \supset B_i(\neg \alpha \colon 1 - r)$
- $\models B_i(\alpha \wedge \beta \colon r) \wedge B_i(\alpha \wedge \neg \beta \colon r') \supset B_i(\alpha \colon r + r')$

# Only-Believing

Modality $O_a^{(k)}$ precisely captures agent $a$'s beliefs and non-beliefs up to depth $k$.

- non-beliefs on irrelevant items:
$$O_a^{(2)}(\neg \textit{fair}(\text{Coin}) \wedge K_b(\textit{fair}(\text{Coin}))) \models \neg K_a K_b(\textit{fair}(\text{Die}))$$

- non-beliefs on deeper mental states:
$$O_a^{(2)}(\neg \textit{fair}(\text{Coin}) \wedge K_b(\textit{fair}(\text{Coin}))) \not\models \neg K_a K_b K_a(\neg \textit{fair}(\text{Coin}))$$
$$O_a^{(3)}(\neg \textit{fair}(\text{Coin}) \wedge K_b(\textit{fair}(\text{Coin}))) \models \neg K_a K_b K_a(\neg \textit{fair}(\text{Coin}))$$

For $i \in Ag$, given $i$-objective formulae $\alpha$ and $\beta$ s.t. the depth of $K_i(\beta)$ not greater than $k$, then $O_i^{(k)}(\alpha)$ entails either $K_i \beta$ or $\neg K_i \beta$.

- $O_i^{(k)}(\alpha) \models K_i \beta$ iff $\alpha \models \beta$;     $O_i^{(k)}(\alpha) \models \neg K_i \beta$ iff $\alpha \not\models \beta$

# Autoepistemic Reasoning

$\mathcal{OBL}_m$ can represent defaults about another agent's beliefs:

## Example

Let KB $= \{\neg \textit{fair}(\text{Coin})\}$,
$\delta = \forall r.\left(r \neq 0 \supset \neg \boldsymbol{B}_a(\neg \boldsymbol{K}_b(\textit{fair}(\text{Coin})): r)\right) \supset \boldsymbol{K}_b(\textit{fair}(\text{Coin}))$
Bob believes $\textit{fair}(\text{Coin})$ unless otherwise (Bob does not believes $\textit{fair}(\text{Coin})$) is believed (by Alice) with a non-zero degree

- $\boldsymbol{O}_a^{(2)}(\text{KB} \wedge \delta) \models \boldsymbol{K}_a \boldsymbol{K}_b(\textit{fair}(\text{Coin}))$
- $\boldsymbol{O}_a^{(2)}(\text{KB} \wedge \delta \wedge \boldsymbol{K}_b(\textit{fair}(\text{Coin}))) \models \boldsymbol{K}_a \boldsymbol{K}_b(\textit{fair}(\text{Coin}))$
- $\boldsymbol{O}_a^{(2)}(\text{KB} \wedge \delta \wedge \neg \boldsymbol{K}_b(\textit{fair}(\text{Coin}))) \models \neg \boldsymbol{K}_a \boldsymbol{K}_b(\textit{fair}(\text{Coin}))$

# Conclusion

In this work, we

- propose an logical account for multi-agent only-believing
- prove properties on beliefs and only-believing
- explore the capability of default reasoning about nested beliefs

For future work:

- extend to belief after actions ✓
- develop mechanisms for projection reasoning
- join common beliefs and only-believing

Thank you!