

DISHONEST REASONING BY ABDUCTION

Chiaki Sakama
Wakayama University

Background and Motivation

2

- People often behave **dishonestly** in daily life
- Few studies investigate inference mechanisms and computational methods of dishonest reasoning in AI
- This is a bit surprise because one of the goals of AI is to better understand human intelligence and to mechanize human reasoning

Contribution

3

- Exploring a computational logic for **dishonest reasoning**
- Formulating different types of dishonesty such as **lie, bullshit** and **withholding information**
- Characterizing dishonest reasoning in terms of **extended abduction**

Logic Programs with Disinformation

4

□ Logic Program

- A program \mathbf{K} consists of rules of the form:

$L_1 ; \dots ; L_l \leftarrow L_{l+1}, \dots, L_m, \text{not } L_{m+1}, \dots, \text{not } L_n$

where L_i is a literal and **not** is **negation as failure**

- The semantics of \mathbf{K} is given by its **answer sets**
- $\mathbf{K} \models L$ if L is included in every answer set of \mathbf{K} ;
- $\mathbf{K} \models \perp$ if \mathbf{K} has no answer set (or inconsistent)

□ Logic Program with Disinformation (or LPD)

- $\langle \mathbf{K}, \mathbf{D} \rangle$ with a program \mathbf{K} and a set \mathbf{D} of ground literals
- For any L in \mathbf{D} , either (a) $\mathbf{K} \models \neg L$ or (b) $\mathbf{K} \not\models L$ and $\mathbf{K} \not\models \neg L$

Deductive Dishonesty

5

- Misleading another agent to **deduce** wrong conclusion
- Two different types of deductive dishonesty
 - ▣ **offensive dishonesty**: behave dishonestly to have a positive (or wanted) outcome that would not be gained by telling the truth
 - ▣ **defensive dishonesty**: behave dishonestly to avoid a negative (or unwanted) outcome that would not be gained by telling the truth
- In each case, an agent can perform different categories of dishonest reasoning
 - ▣ **Lie**: to tell a believed-false sentence
 - ▣ **Bullshit**: to tell a sentence that is neither believed to be true nor believed to be false
 - ▣ **Withholding Information**: to fail to offer information that would help someone acquire true beliefs and/or correct false beliefs

Offensive Dishonesty

6

- $\langle \mathbf{K}, \mathbf{D} \rangle$: LPD, O^+ : a ground literal representing a positive outcome s.t. $\mathbf{K} \models O^+$
- Suppose a pair (\mathbf{I}, \mathbf{J}) of sets of ground literals satisfying:
 - $(\mathbf{K} \setminus \mathbf{J}) \cup \mathbf{I} \models O^+$
 - $(\mathbf{K} \setminus \mathbf{J}) \cup \mathbf{I} \not\models \perp$
 - $\mathbf{I} \subseteq \mathbf{D}$ and $\mathbf{J} \subseteq \mathbf{K}$
- Then, (\mathbf{I}, \mathbf{J}) is called
 - **lie** for O^+ if $\mathbf{I} \neq \emptyset$ and $\mathbf{K} \models \neg L$ for some $L \in \mathbf{I}$
 - **bullshit (or BS)** for O^+ if $\mathbf{I} \neq \emptyset$ and $\mathbf{K} \not\models \neg L$ for any $L \in \mathbf{I}$
 - **withholding information (or WI)** for O^+ if $\mathbf{I} = \emptyset$

Example

7

- A salesperson believes that a product will be sold if the quality is good. However, he believes that the quality is not good.
- The situation is represented by an LPD $\langle \mathbf{K}, \mathbf{D} \rangle$ where $\mathbf{K} = \{ \text{sold} \leftarrow \text{quality}, \neg \text{quality} \leftarrow . \}$ and $\mathbf{D} = \{ \text{quality} \}$
- To have a positive outcome $O^+ = \{ \text{sold} \}$, he introduces $\mathbf{I} = \{ \text{quality} \}$ to \mathbf{K} and eliminates $\mathbf{J} = \{ \neg \text{quality} \}$ from \mathbf{K} .
- As a result, $(\mathbf{K} \setminus \mathbf{J}) \cup \mathbf{I} \models O^+$. In this case, (\mathbf{I}, \mathbf{J}) is an offensive lie.

Defensive Dishonesty

8

- $\langle \mathbf{K}, \mathbf{D} \rangle$: LPD, O^- : a ground literal representing a negative outcome s.t. $\mathbf{K} \models O^-$
- Suppose a pair (\mathbf{I}, \mathbf{J}) of sets of ground literals satisfying:
 - $(\mathbf{K} \setminus \mathbf{J}) \cup \mathbf{I} \not\models O^-$
 - $(\mathbf{K} \setminus \mathbf{J}) \cup \mathbf{I} \not\models \perp$
 - $\mathbf{I} \subseteq \mathbf{D}$ and $\mathbf{J} \subseteq \mathbf{K}$
- Then, (\mathbf{I}, \mathbf{J}) is called
 - **lie** for O^- if $\mathbf{I} \neq \emptyset$ and $\mathbf{K} \models \neg L$ for some $L \in \mathbf{I}$
 - **bullshit (or BS)** for O^- if $\mathbf{I} \neq \emptyset$ and $\mathbf{K} \not\models \neg L$ for any $L \in \mathbf{I}$
 - **withholding information (or WI)** for O^- if $\mathbf{I} = \emptyset$

Example

9

- A salesperson believes that an order will be canceled if the quality is not good. However, **he has no information on the quality of the product.**
- The situation is represented by an LPD $\langle \mathbf{K}, \mathbf{D} \rangle$ where $\mathbf{K} = \{ \text{canceled} \leftarrow \text{not quality} \}$ and $\mathbf{D} = \{ \text{quality} \}$
- To avoid a negative outcome $O^- = \{ \text{canceled} \}$, he introduces $\mathbf{I} = \{ \text{quality} \}$ to \mathbf{K} .
- As a result, $\mathbf{K} \cup \mathbf{I} \neq O^-$. In this case, (\mathbf{I}, \emptyset) is **defensive BS.**

Abductive Dishonesty

10

- Interrupting another agent to **abduce** correct explanations
- Two different types of abductive dishonesty
 - ▣ **abductive dishonesty for positive evidences:** behave dishonestly to explain a positive evidence that is occurred
 - ▣ **abductive dishonesty for negative evidences:** behave dishonestly to explain a negative evidence that is not occurred
- In each case, an agent can perform different categories of dishonest reasoning – **Lie**, **BS** or **WI**
- A knowledge base of an agent includes a **secret set** of literals that he wants to conceal from another agent

Abductive Dishonesty for Positive Evidences

11

- $\langle \mathbf{K}, \mathbf{D} \rangle$: LPD, Σ : a secret set, E^+ : a ground literal representing a positive evidence s.t. $\mathbf{K} \models E^+$ and $\mathbf{K} \setminus \Sigma \not\models E^+$
- Suppose a pair (\mathbf{I}, \mathbf{J}) of sets of ground literals satisfying:
 - $(\mathbf{K} \setminus (\Sigma \cup \mathbf{J})) \cup \mathbf{I} \models E^+$
 - $(\mathbf{K} \setminus (\Sigma \cup \mathbf{J})) \cup \mathbf{I} \not\models \perp$
 - $\mathbf{I} \subseteq \mathbf{D}$ and $\mathbf{J} \subseteq \mathbf{K}$
- Then, (\mathbf{I}, \mathbf{J}) is called
 - **lie** for E^+ if $\mathbf{I} \neq \emptyset$ and $\mathbf{K} \models \neg L$ for some $L \in \mathbf{I}$
 - **bullshit (or BS)** for E^+ if $\mathbf{I} \neq \emptyset$ and $\mathbf{K} \not\models \neg L$ for any $L \in \mathbf{I}$
 - **withholding information (or WI)** for E^+ if $\mathbf{I} = \emptyset$

Example

12

- Sam is coming home late because he is cheating on his wife. Observing the late arrival, his wife might abduce his cheating. Sam does not want this abduction to take place, so he makes up another reason: **he did overtime at work**. He hopes this disinformation will stop her abduction.
- The situation is represented by an LPD $\langle \mathbf{K}, \mathbf{D} \rangle$ where
 $\mathbf{K} = \{ \text{late} \leftarrow \text{cheat}. \text{late} \leftarrow \text{overtime}. \text{cheat} \leftarrow. \neg \text{overtime} \leftarrow. \}$,
 $\mathbf{D} = \{ \text{overtime} \}$ and the secret set $\Sigma = \{ \text{cheat} \}$
- In face of the positive evidence $E^+ = \text{late}$, he introduces $\mathbf{I} = \{ \text{overtime} \}$ to \mathbf{K} and eliminates $\mathbf{J} = \{ \neg \text{overtime} \}$ from \mathbf{K} .
- As a result, $(\mathbf{K} \setminus (\Sigma \cup \mathbf{J})) \cup \mathbf{I} \models E^+$. In this case, (\mathbf{I}, \mathbf{J}) is an **abductive lie**.

Abductive Dishonesty for Negative Evidences

13

- $\langle \mathbf{K}, \mathbf{D} \rangle$: LPD, Σ : a secret set, E^- : a ground literal representing a negative evidence s.t. $\mathbf{K} \not\models E^-$ and $\mathbf{K} \setminus \Sigma \models E^-$
- Suppose a pair (\mathbf{I}, \mathbf{J}) of sets of ground literals satisfying:
 - $(\mathbf{K} \setminus (\Sigma \cup \mathbf{J})) \cup \mathbf{I} \not\models E^-$
 - $(\mathbf{K} \setminus (\Sigma \cup \mathbf{J})) \cup \mathbf{I} \not\models \perp$
 - $\mathbf{I} \subseteq \mathbf{D}$ and $\mathbf{J} \subseteq \mathbf{K}$
- Then, (\mathbf{I}, \mathbf{J}) is called
 - **lie** for E^- if $\mathbf{I} \neq \emptyset$ and $\mathbf{K} \models \neg L$ for some $L \in \mathbf{I}$
 - **bullshit (or BS)** for E^- if $\mathbf{I} \neq \emptyset$ and $\mathbf{K} \not\models \neg L$ for any $L \in \mathbf{I}$
 - **withholding information (or WI)** for E^- if $\mathbf{I} = \emptyset$

Example

14

- Sam and his wife promised to have a dinner at a restaurant. But Sam does not come to the restaurant on time because he is arguing with his girlfriend over the phone. Sam then excuses that **he mistook the time**.
- The situation is represented by an LPD $\langle \mathbf{K}, \mathbf{D} \rangle$ where
 $\mathbf{K} = \{ \text{on-time} \leftarrow \text{not call}, \text{remember} \leftarrow \text{call} \leftarrow \text{remember} \leftarrow \}$,
 $\mathbf{D} = \{ \text{remember} \}$ and the secret set $\Sigma = \{ \text{call} \}$
- In face of the negative evidence $E^- = \text{on-time}$, he eliminates $\mathbf{J} = \{ \text{remember} \}$ from \mathbf{K} .
- As a result, $\mathbf{K} \setminus (\Sigma \cup \mathbf{J}) \neq E^-$. In this case, (\emptyset, \mathbf{J}) is an **abductive WI**.

Preference between Dishonesties

15

- **Quantitative Measure**: Comparing the same type of dishonesties, **the smaller the better**
- Let (\mathbf{I}, \mathbf{J}) and $(\mathbf{I}', \mathbf{J}')$ be two lies/BS/WI for the same outcome/evidence. Then, (\mathbf{I}, \mathbf{J}) is **more or equally preferred to** $(\mathbf{I}', \mathbf{J}')$ (written $(\mathbf{I}, \mathbf{J}) \geq (\mathbf{I}', \mathbf{J}')$) if $\mathbf{I} \subseteq \mathbf{I}'$ and $\mathbf{J} \subseteq \mathbf{J}'$.
The most preferred one is called a **minimal dishonesty**.
- **Qualitative Measure**: Comparing different types of dishonesties,
 - ▣ **WI is preferable to BS and Lies**, since WI introduces no disinformation
 - ▣ **BS is preferable to Lies**, since BS is consistent with an agent's belief
- Let $(\mathbf{I}_1, \mathbf{J}_1)$, $(\mathbf{I}_2, \mathbf{J}_2)$ and $(\mathbf{I}_3, \mathbf{J}_3)$ be a lie, BS and WI for the same outcome/evidence, respectively. Then, $(\mathbf{I}_3, \mathbf{J}_3) > (\mathbf{I}_2, \mathbf{J}_2) > (\mathbf{I}_1, \mathbf{J}_1)$

Extended Abduction

[Inoue & Sakama, IJCAI-95]

16

- An **abductive program** is a pair $\langle \mathbf{K}, \mathbf{A} \rangle$ where \mathbf{K} is a logic program and \mathbf{A} is a set of ground literals representing hypotheses (called **abducibles**)
- Given a **positive observation** G^+ as a ground literal satisfying $\mathbf{K} \not\models G^+$, a pair (\mathbf{I}, \mathbf{J}) of sets of ground literals is an **explanation** of G^+ if (i) $(\mathbf{K} \setminus \mathbf{J}) \cup \mathbf{I} \models G^+$, (ii) $(\mathbf{K} \setminus \mathbf{J}) \cup \mathbf{I} \not\models \perp$, (iii) $\mathbf{I} \subseteq \mathbf{A} \setminus \mathbf{K}$ and $\mathbf{J} \subseteq \mathbf{A} \cap \mathbf{K}$
- Given a **negative observation** G^- as a ground literal satisfying $\mathbf{K} \models G^-$, a pair (\mathbf{I}, \mathbf{J}) of sets of ground literals is an **anti-explanation** of G^- if (i) $(\mathbf{K} \setminus \mathbf{J}) \cup \mathbf{I} \not\models G^-$, (ii) $(\mathbf{K} \setminus \mathbf{J}) \cup \mathbf{I} \not\models \perp$, (iii) $\mathbf{I} \subseteq \mathbf{A} \setminus \mathbf{K}$ and $\mathbf{J} \subseteq \mathbf{A} \cap \mathbf{K}$
- An (anti-)explanation (\mathbf{I}, \mathbf{J}) is **minimal** if $\mathbf{I}' \subseteq \mathbf{I}$ and $\mathbf{J}' \subseteq \mathbf{J}$ imply $\mathbf{I}' = \mathbf{I}$ and $\mathbf{J}' = \mathbf{J}$ for any (anti-)explanation $(\mathbf{I}', \mathbf{J}')$

Example

17

- Tweety is a bird and normally flies. One day an agent observes that Tweety does not fly. He then assumes that Tweety broke its wing.
- The situation is represented by an abductive program $\langle \mathbf{K}, \mathbf{A} \rangle$ where $\mathbf{K} = \{ \text{flies} \leftarrow \text{bird}, \text{not broken-wing. } \text{bird} \leftarrow . \}$, $\mathbf{A} = \{ \text{broken-wing} \}$
- In this case, the negative observation $G^- = \text{flies}$ has the anti-explanation $(\mathbf{I}, \mathbf{J}) = (\{ \text{broken-wing} \}, \emptyset)$ s.t. $\mathbf{K} \cup \mathbf{I} \neq G^-$
- The agent revises \mathbf{K} to $\mathbf{K}' = \mathbf{K} \cup \{ \text{broken-wing} \}$. After several days, the agent observes that Tweety flies as before. He then considers that the wound was healed.
- In this case, the positive observation $G^+ = \text{flies}$ has the explanation $(\mathbf{I}, \mathbf{J}) = (\emptyset, \{ \text{broken-wing} \})$ s.t. $\mathbf{K}' \setminus \mathbf{J} \models G^+$

Computing Dishonesties by Abduction

18

- There are structural similarities between deductive dishonesty and extended abduction
 - ▣ Viewing a positive outcome as a positive observation, **an offensive dishonesty (I,J)** for the outcome wrt $\langle K, D \rangle$ is identified with **an explanation** of the observation wrt $\langle K, L(K) \cup D \rangle$ where $L(K) = K \cap \text{Lit}$ and **Lit** is the set of all ground literals in the language
 - ▣ Viewing a negative outcome as a negative observation, **a defensive dishonesty (I,J)** for the outcome wrt $\langle K, D \rangle$ is identified with **an anti-explanation** of the observation wrt $\langle K, L(K) \cup D \rangle$
- Similar correspondences are observed between abductive dishonesty and extended abduction

Deductive Dishonesty vs. Extended Abduction

19

- $\langle \mathbf{K}, \mathbf{D} \rangle$: LPD, O^+ : positive outcome,
 O^- : negative outcome.
- ▣ (\mathbf{I}, \mathbf{J}) is a (minimal) offensive dishonesty for O^+ wrt $\langle \mathbf{K}, \mathbf{D} \rangle$ iff (\mathbf{I}, \mathbf{J}) is a (minimal) explanation of O^+ wrt $\langle \mathbf{K}, L(\mathbf{K}) \cup \mathbf{D} \rangle$
- ▣ (\mathbf{I}, \mathbf{J}) is a (minimal) defensive dishonesty for O^- wrt $\langle \mathbf{K}, \mathbf{D} \rangle$ iff (\mathbf{I}, \mathbf{J}) is a (minimal) anti-explanation of O^- wrt $\langle \mathbf{K}, L(\mathbf{K}) \cup \mathbf{D} \rangle$

Abductive Dishonesty vs. Extended Abduction

20

- $\langle \mathbf{K}, \mathbf{D} \rangle$: LPD, E^+ : positive evidence,
 E^- : negative evidence.
- ▣ (\mathbf{I}, \mathbf{J}) is a (minimal) abductive dishonesty for E^+ wrt
 $\langle \mathbf{K}, \mathbf{D} \rangle$ iff (\mathbf{I}, \mathbf{J}) is a (minimal) explanation of E^+ wrt
 $\langle \mathbf{K} \setminus \Sigma, L(\mathbf{K}) \cup \mathbf{D} \rangle$
- ▣ (\mathbf{I}, \mathbf{J}) is a (minimal) abductive dishonesty for E^- wrt
 $\langle \mathbf{K}, \mathbf{D} \rangle$ iff (\mathbf{I}, \mathbf{J}) is a (minimal) anti-explanation of E^-
wrt $\langle \mathbf{K} \setminus \Sigma, L(\mathbf{K}) \cup \mathbf{D} \rangle$

Computational Complexities

21

- The following 3 decision problems are considered. Given a propositional LPD $\langle \mathbf{K}, \mathbf{D} \rangle$ and a positive/negative outcome/evidence X ,
 - ▣ Does there **exist** a deductive/abductive dishonesty (\mathbf{I}, \mathbf{J}) for X ?
 - ▣ Is a literal L is **relevant** to some (minimal) deductive/abductive dishonesty for the outcome/evidence? (i.e., $L \in \mathbf{I} \cup \mathbf{J}$ for some (\mathbf{I}, \mathbf{J}) for X)
 - ▣ Is a literal L is **necessary** for every (minimal) deductive/abductive dishonesty for the outcome/evidence? (i.e., $L \in \mathbf{I} \cup \mathbf{J}$ for every (\mathbf{I}, \mathbf{J}) for X)

Summary of Complexity Results

22

	Deductive Dishonesty		Abductive Dishonesty	
	Positive Outcome	Negative Outcome	Positive Evidence	Negative Evidence
Existence	Σ^P_3	Σ^P_2	Σ^P_3	Σ^P_2
Relevance (minimal)	Σ^P_3 (Σ^P_4)	Σ^P_2 (Σ^P_3)	Σ^P_3 (Σ^P_4)	Σ^P_2 (Σ^P_3)
Necessity (minimal)	Π^P_3 (Π^P_3)	Π^P_2 (Π^P_2)	Π^P_3 (Π^P_3)	Π^P_2 (Π^P_2)

*Each entry represents completeness for the respective class.

Final Remark

- Extended abduction is computed using **answer set programming** (ASP) [Sakama & Inoue, TPLP 2003]
- (Correspondence between dishonest reasoning and extended abduction)
+ (computation of extended abduction in ASP)
⇒ (computation of dishonest reasoning in ASP).
(The method is provided in the paper).
- The logical framework of dishonest reasoning and its relationship to abduction do **not** depend on a particular logic.

Related Studies

by the author and his colleagues

24

- Chiaki Sakama, Martin Caminada and Andreas Herzig
A Logical Account of Lying.
12th European Conference on Logics in AI (JELIA),
LNAI 6341 , 2010.
- Chiaki Sakama and Martin Caminada
The Many Faces of Deception.
Thirty Years of Nonmonotonic Reasoning (NonMon@30),
Lexington, KY, USA, October 2010.
- Chiaki Sakama
Logical Definitions of Lying.
14th International Workshop on Trust in Agent Societies,
Taipei, Taiwan, May 2, 2011

Related Studies

by the author and his colleagues

25

- Ngoc-Hieu Nguyen, Tran Son, Enrico Pontelli and Chiaki Sakama
ASP-Prolog for Negotiation Among Dishonest Agents.
11th International Conference on Logic Programming and Nonmonotonic Reasoning (LPNMR), LNAI 6645, 2011.
- Chiaki Sakama, Tran Son and Enrico Pontelli
A Logical Formulation for Negotiation Among Dishonest Agents.
IJCAI-2011, Barcelona ([presented on Friday afternoon](#))
- Papers are available at the author's home page:
<http://www.wakayama-u.ac.jp/~sakama>