

人狼知能達成におけるエージェントの推論モデル

大澤博隆^{†1} 鳥海不二夫^{†2} 稲葉通将^{†3}
片上大輔^{†4} 梶原健吾^{†2} 篠田孝祐^{†5}

著者らは人工知能によって人狼というゲームを解くことを目標としている。人狼ゲームを解くためには、他者の意図を記述し、自己や他者の心的モデルを記述し、扱えるような手法が必要となる。本研究では BDI 論理を用いたエージェントの推論記述法を提案した。また、実際の人狼プレイログを用いて、このような BDI 論理でどのような記述が可能となるか検証した。

Agent's Reasoning Model for Achieving AIWolf

HIROTAKA OSAWA^{†1} FUJIO TORIUMI^{†2} MICHIMASA INABA^{†3}
DAISUKE KATAGAMI^{†4} KENGO KAJIWARA^{†2} KOSUKE SHINODA^{†5}

The authors are trying to solve the game of werewolf by artificial intelligence. It is required to describe and handle a mental model of others and a mental model of a player in others for constructing artificial intelligence for solving the game. The authors designed reasoning of agents for describe werewolf game based on BDI logics. We referring online playing log to verify how reasoning describe a werewolf game.

1. はじめに：ゲーム課題としての人狼の特徴

「ある村に、人間の姿に化けられる人食い人狼が現れた。人狼は人間と同じ姿をしており、昼間には区別がつかず、夜に村人たちをひとりずつ襲っていく。村の中に潜んだ人狼を暴き出すため、村人たちは互いを疑いつつ、毎夜一人を処刑していくことにした」

以上が、コミュニケーションゲームとして知られている『人狼』の大まかなカバーストーリーである。20世紀のソ連において、マフィアというゲームから派生したと思われる「人狼」は、言語のみを使う抽象化されたコミュニケーションゲームでありながら、全世界で楽しまれているゲームである。

人狼ゲームは通常の完全情報ゲームと異なり、情報が隠されている。また、事前に決定される情報は各人の役職のみであり、他の殆どの情報は共有されない。このため、各人の役職や各人の占い先、吊り先、疑惑対象、信頼対象、役職推理など、全ての情報は言語を用いて伝達される、という特徴がある。特に人狼では、実際にあった現象を確認する事実確認的発話だけでなく、行為遂行的発話と呼ばれる、発話自体が行為として認識される種類の発話を適切に行うことが、ゲーム内で勝利するために不可欠である[1]。

ゲーム課題として人狼を見た場合、コミュニケーション以外の客観的情報、勝敗決定要因がほとんど存在しないため、言語学的な課題の達成要件が直接勝率に影響する、と

いう特徴がある。ゲーム外から手に入る情報を極力削減し、人狼ゲームの標準化を促したのが、2004年に誕生した掲示板型のゲームサイト「人狼 BBS」である[2]。人狼 BBS ではカードゲームにおける昼と夜の区別を無くし、狼は人間同士の議論と並行して、狼同士の秘密の会話を行い、お互いに発話を修正することを可能とした。双方の陣営が対等な立場で制約にとらわれず、言語的な会話のみで勝負を行う論理的なコミュニケーションゲームの性質を備えることになった。また人狼 BBS では参加者は匿名の仮ログイン名を使って参加し、インタフェースとして老若男女織り交ぜた、共通のキャラクターセットの中からキャラクターを選んでプレイを行う。これにより、各プレイヤーは自身のゲーム外での情報を隠し、純粋な言語ゲームを行うことを達成している。本研究でエージェント課題として扱う人狼は、このような標準化された人狼ゲームが中心となる。

人狼では情報の確かさが立場により異なるため、最終的な尤度の推定はエージェントにより異なり、結論も個々で異なる。互いに異なる意見を持つエージェントたちは、最終的に自分たちの意見を集約し、合理的な決断を下す必要がある。ここでは推理よりも、説得や信頼といった要素が重要なキーワードとして浮かび上がってくる。人間も、人狼も、生き残りたければ信頼を勝ち取らねばならない。自分の意見を通すためには、自分の意見が信頼に足るものであることを他者に説明する必要がある。この段階では、他者のモデル化だけではなく、他者から見た自己のモデル化、

^{†1} 筑波大学
University of Tsukuba
^{†2} 東京大学
The University of Tokyo
^{†3} 広島市立大学
Hiroshima City University

^{†4} 東京工芸大学
Tokyo Polytechnic University
^{†5} 電気通信大学
The University of Electro-Communications

信頼出来る発話の演出が重要となってくる。

本研究では、このような推論を扱えるエージェント内部モデルとして、BDI 論理を用いたエージェントの思考方法を提案する。提案した記述方法を用いて、実際のプレイログを記述した。

2. オンラインゲームとしての人狼：人狼 BBS

ゲーム外から手に入る情報を極力削減し、人狼ゲームの標準化を促したのが、2004年に誕生した掲示板型のゲームサイト「人狼 BBS」である[2]。人狼 BBS ではカードゲームにおける昼と夜の区別を無くし、狼は人間同士の議論と並行して、狼同士の秘密の会話をし、お互いに発話を修正することを可能とした。双方の陣営が対等な立場で制約にとらわれず、言語的な会話のみで勝負を行う論理的なコミュニケーションゲームの性質を備えることになった。また人狼 BBS では参加者は匿名の仮ログイン名を使って参加し、インタフェースとして老若男女織り交ぜた、共通のキャラクターセットの中からキャラクターを選んでプレイを行う。これにより、各プレイヤーは自身のゲーム外での情報を隠し、純粋な言語ゲームを行うことを達成している。本研究でエージェント課題として扱う人狼は、このような標準化された人狼ゲームが中心となる。

3. BDI 論理によるエージェントの推論モデル

本研究ではエージェントの推論を様相論理の一形態である BDI 論理で記述される形式に書き換える。

3.1 BDI 論理のために用いる演算子の記述：主体と確率の追加

本研究では BDI 論理について、エージェントが想定する分岐した可能世界について扱えるように拡張した Rao の定義を使う[3]。BDI 論理では通常の論理式に加えて、表 1 のような BDI 演算子が追加される。また、それぞれの様相の主体と確信度を表すため、新出らの BDI 論理拡張を参考に、心的な主体と心的状態の確率を記述できるように拡張した[4]。

新出らの研究では確率的状況を X_e パラメータの拡張という形で表しており、ある場合から次時点への推移を確率的に表せるような拡張を行っており、BDI オペレータ自体の確率的な記述を行っていない。これに対し、筆者らは心的状態への確率的拡張を記述する。これは人狼において、各プレイヤーはある状況に至った過程を推理する場合、複数の原因が想定でき、また全ての原因を特定できるわけではないからである。 X_e パラメータによって確率を記述する場合、遷移前と遷移後の数が明示されており、遷移後の各状態の確率の合計が 1 になる必要がある。これに対し、心的状態を確率で記述する場合、状態の合計が 1 にならない場合が考えられる。このため、筆者らの方法は新出らの方法で記述できる範囲を記述でき、より広い範囲の論理記述を

可能とする。

また主体についてだが、人狼は相手の思考を読むゲームであり、各エージェントの記述の中に他のエージェントの信念・願望・意図を記述する必要がある。したがって、BDI の各演算子において、誰がその行為者であるかという主体を記述することがとても重要である。これらの主体は演算子の上に記述される。新出らの拡張は従来の BDI 論理を包含しており[4]、本研究の記述法は新出らの記述を拡張したものである。本研究での拡張も既存の BDI 論理で表す範囲を全て表すことが可能である。

表 1 BDI 論理のための基本オペレータ
Table 1 Basic operators for BDI logic

	演算子	説明
様相	BEL ^a BEL ^a [num %]	a が信じる
	DESIRE ^a DESIRE ^a [num %]	a が望む (可能不可能問わず)
	INTEND ^a INTEND ^a [num %]	a が意図する (可能なもののみ)
時相	A	全ての可能世界で
	E	ある可能世界で
	X	次の時点で
	G	現在時点を含み永遠に
	F	現在を含む時点のいつか
	U	条件が成立する時点まで
	B	現在を含まない前の時点で

3.2 ゲーム状態の記述

通常の BDI 論理では、各心的状態の中の記述はしばしば自然言語で書かれる。心的状態の中の記述を自然言語で書く場合、自然言語の中に BDI オペレータが含まれてしまい、外部から操作がしづらくなる可能性がある。また、人狼で行われる会話は限定的であるため、記述の簡略化が可能である[2]。筆者らの以前の研究を元に、人狼の状態を記述する演算子を表 2 の 2 文に絞り込んだ

表 2 情景描写と行為のためのオペレータ
Table 1 Operators for descriptions and actions

	演算子	説明
Is 文	Is(character, role)	character が role である
	Is(character/role, compositant)	character/role が compositant である
Do 文	Do(character/role, verb, character/role)	character/role が character/role を verb する
	Do(character/role, verb, act)	character/role が act を verb する
	Do(character/role, verb, IS())	character/role が IS 文を verb する

3.3 人狼ゲーム記述のために使われる記述語

人狼における使用単語を表 3 に記述する。

表 3 基礎語
Table 3 basic words

	名詞	説明
character	"NAME"	人物の名前が入る
	anychar	∀character
	who	前に述べられている人物
role	villager	村人 (市民)
	seer	占い師 (予言者)
	medium	霊能者 (霊媒師)
	hunter	狩人 (守護者)
	freemason	共有者
	wolf	狼
	lunatic	狂人
	HUMAN	villager ∪ seer ∪ medium ∪ hunter ∪ freemason ∪ lunatic
	VILLAGESIDE	villager ∪ seer ∪ medium ∪ hunter ∪ freemason
	WOLFSIDE	wolfside ∪ lunatic
	GIFTED	seer ∪ medium ∪ hunter ∪ freemason
ANYROLE	seer ∪ medium ∪ hunter ∪ freemason ∪ lunatic	
act	inspection	占い行為
	inquestion	霊能行為
	protection	護衛
	execution	処刑
	voting	投票
	attack	襲撃
	comingout	役職宣言(CO)
	speech NUM	発話番号
verb (能動形)	divine	対象を占う
	guard	守る
	execute	処刑する
	vote	投票する
	attack	襲撃する
	tell	IS 文を伝える
	know	知る
	discuss	議論する
	decide	決定する
	avoid	避ける
	comingout	役職公開(CO)させる

	estimate	推測する
composant (受動形)	divined	占われる
	inquested	霊能結果を告げられる
	guarded	守られる
	executed	処刑される
	voted	投票される
	attacked	襲撃される
	comingouted	役職公開(CO)した
	speak	話す
	exist	存在する
	win	勝つ
	lose	負ける
	suddendead	突然死する
	KILLED	attacked ∪ executed

4. 人狼 BBS のログを用いたモデル記述の検証

記述例として、人狼 BBS 村 1「疑心暗鬼の村」のログを用いた[2]。人狼 BBS の初めの村は、村人と人狼の読み合いが面白い名勝負として記録されており、他の村のプレイに影響を与えたログとして知られている[5]。従って、この村のプレイヤーの流れを記述できることは重要であると考えられる。

4.1 人狼 BBS 村 1 の展開の記述

人狼 BBS 村 1 の村の流れを、3 章の記述を用いて表す。これは将棋で言う棋譜にあたる。ただしスペースの都合上、各自の投票先など、全ての行動は記述しない。

村 1 は初日犠牲者のノンプレイヤーキャラクターが 1 名、村人が 7 名(Liesa, Nicholas, Thomas, Dieter, Katharina, Pamela, Jacob)、人狼が 2 名(Simson, Otto)、占い師が 1 名(Regina)、霊能者が 1 名(Joachim)、狂人が 1 名(Albin)、狩人が 1 名(Molitz)、共有者が 2 名(Walter, Peter)の構成となる。人狼 BBS 初期のログであるため、人狼の構成は 3 名ではなく 2 名である。

Is(Liesa, villager), Is(Nicholas, villager), Is(Thomas, villager)
Is(Dieter, villager), Is(Katharina, villager), Is(Pamela, villager)
Is(Jacob, villager), Is(Simson, wolf), Is(Otto, wolf)
Is(Regina, seer), Is(Joachim, medium), Is(Albin, lunatic)
Is(Molitz, hunter), Is(Walter, freemason), Is(Peter, freemason)

- 1 日目に共有者の Walter と Peter が共有者を CO (役職宣言) する。

Do(Walter, tell, IS(Walter, freemason))

Do(Walter, tell, IS(Peter, freemason))

Do(Peter, tell, IS(Peter, freemason))

- 2 日目に狂人 Albin が占い師を CO し、村人 Liesa が人間であることを宣言する。次に占い師 Regina が占い師を CO し、狼 Simson の占い結果が人狼であることを宣言する。最後に霊能者 Joachim が霊能者を CO

する。狼 Simson が処刑される。狼 Otto が狂人 Albin を襲撃するが、狩人 Molitz が Albin を護衛する。無発言により共有者ペーターが突然死する。

Do(Albin, tell, IS(Albin, seer)), Do(Albin, tell, IS(Liesa, HUMAN))
 Do(Regina, tell, IS(Regina, seer)), Do(Regina, tell, IS(Simson, wolf))
 Do(Joachim, tell, IS(Joachim, medium))
 Is(Simson, executed)
 Do(Molitz, protect, Albin), Do(Otto, attack, Albin)
 Is(Albin, guarded), Is(Peter, suddendeath)

- 3 日目に霊能者 Joachim が Simson 狼を宣言し、占い師 Regina が狼 Otto の占い結果が人狼であることを宣言する。狂人 Albin が Joachim の占い結果が狼であることを宣言する。霊能者 Joachim が処刑され、共有者 Walter が襲撃される

Do(Joachim, tell, IS(Simson, wolf))
 Do(Regina, tell, IS(Otto, wolf))
 Do(Albin, tell, IS(Joachim, wolf))
 Is(Joachim, executed), Is(Walter, attacked)

- 4 日目に狂人 Albin が狼 Otto の占い結果が人間であることを宣言する。村人 Pamela が処刑され狩人 Molitz が襲撃される

Do(Albin, tell, IS(Otto, HUMAN))
 Is(Pamela, executed), Is(Molitz, attacked)

- 5 日目に狂人 Albin が狩人・死亡者 Molitz の占い結果が人間であることを宣言する。占い師 Regina が処刑され村人 Liesa が襲撃される

Do(Albin, tell, IS(Molitz, HUMAN))
 Is(Regina, executed), Is(Liesa, attacked)

- 6 日目に狂人 Albin が村人 Nicholas の占い結果が人間であることを宣言する。村人 Jacob が処刑され村人 Nicholas が襲撃される

Do(Albin, tell, IS(Molitz, HUMAN))
 Is(Jacob, executed), Is(Nicholas, attacked)

- 7 日目に狂人 Albin が村人 Katharina の占い結果が人間であることを宣言する。狼 Otto が処刑され、村側勝利。

Do(Albin, tell, IS(Katharina, HUMAN))
 Is(Otto, executed), Is(VILLAGESIDE, win)

4.2 人狼 BBS 村 1 での推理の記述

BDI 論理記述を用いることで、各々のエージェントの発言を BDI 論理で記述することが可能となる。例えば、図 1 は Joachim の 2 日目の発言で、占い師を宣言した Albin の後に、他の村人に向かって対抗の占い師が居ないかどうか確認する発言である。こういった、CO を要請する発言は、人狼ゲームにおいてよく行われる。この Joachim の発言は下記のとおり記述できる。

$BEL^{Joachim} \left(AX \left(\neg Do(\forall \text{people except Albin, tell, Is(who, seer)}) \right) \rightarrow Is(Albin, seer) \right)$

これを他プレイヤーに対する願望の発露と考えれば、以下のように簡潔に記述可能である。

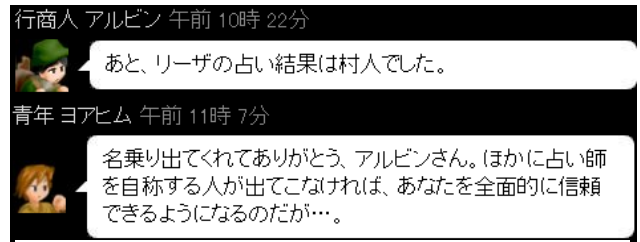


図 1 二日目の Albin のリーザへの占い結果と、Joachim の占い師に対する条件付きの印象評価

Fig. 1 Albin's information about Liesa and Joachim's recommendation about seer's comingout for others in 2nd day.

$DESIRE^{Joachim}(\forall \text{people except Albin, tell, Is(who, seer)})$

また、「自分が他のプレイヤーの占い師 CO を意図し、そのために行動したい」という意味であれば、下記のように記述できる。

$INTEND^{Joachim}(\forall \text{people except Albin, tell, Is(who, seer)})$

上記 3 つの発言は、発言の与える意味合いが異なるが、このような微妙な違いは BDI 論理でも再現可能である。

4.2 複雑な推論の記述

本村での特徴的な点としてまず、狼 Otto が 3 日目の襲撃失敗から、狩人 Molitz の正体を見抜いた点が挙げられる(図 2)。これは本試合ログのまとめでも特徴的な出来事として挙げられている。この図の推論を 3 章の BDI 論理で表すと以下のとおりである。

$BEL^{Otto}(BEL^{Molitz}(Is(Otto, VILLAGESIDE)) \rightarrow Is(Molitz, hunter)) \quad (1)$

これは人狼のルールから導かれる論理規則と、プレイヤー間で共通する規則と、狩人の思考に対する Otto 個人の BDI 論理規則から導かれる。人狼のルールから導かれる論理は以下の 2 点である。これらの規則はプレイヤーの立場に関わらず、真である。

- ・ x が狩人で、かつ y が護衛されているなら、y が襲撃されたとき x は知っている。また、y が襲撃されたとき x が知っているということは、x は狩人であり、かつ y が護衛されている。

$IS(x, hunter) \cap B(IS(y, guarded)) \leftrightarrow$

$Do(x, know, IS(y, attacked)) \quad (2)$

- ・ x が襲撃されたにも関わらず、全ての次の可能世界で x が生きている場合、x が護衛された

$Is(x, attacked) \cap AX(Is(x, live)) \rightarrow Is(x, guarded) \quad (3)$

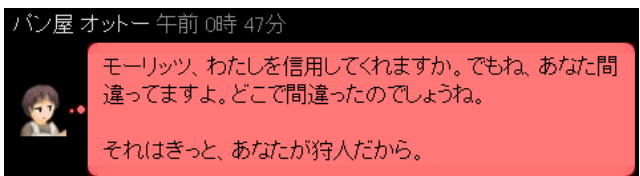


図 2 三日目の Otto の独り言 (エージェント内部の推論)

Fig. 2 Otto's monologue on a 3rd day's log

狩人に対して想定できる BDI 論理規則は、以下の 1 点である。これはプレイヤー個人の思考法に依存せず、合理的なプレイヤーであれば推論可能なルールである。

- someone が x を村人側だと信じる場合、some は x を人間だと信じている

$$BEL^{someone}(Is(x, VILLAGESIDE)) \rightarrow$$

$$BEL^{someone}(Is(x, HUMAN)) \quad (4)$$

Otto 個人が持つ合理的な BDI 論理規則は、以下の 4 点である。これはプレイヤー個人の思考方法に依存する。

- x が占い師で、y を占って人間だと言っている場合、x は人間だ、とプレイヤーは信じる（占い師が嘘をつく可能性を考えると、これは必ずしも真ではないが、合理的な思考である）

$$BEL^{player}(Is(x, seer) \cap Do(x, tell, Is(y, HUMAN))) \rightarrow$$

$$IS(x, HUMAN) \quad (5)$$

- y が人間の時、占い師宣言した x が y を人間と宣言し、占い師宣言した x 以外の z が y を狼と宣言した場合、x が占い師である（占い師が嘘をつく可能性を考えると、これは必ずしも真ではないが、合理的な思考である）

$$BEL^{player}(Is(y, HUMAN), Do(x, tell, Is(x, seer))) \cap$$

$$Do(x, tell, Is(y, HUMAN)) \cap Do(\forall z \text{ except } x, tell, Is(z, seer)) \cap$$

$$Do(z, tell, Is(y, wolf)) \rightarrow IS(x, seer) \quad (6)$$

- x が護衛されていれば、狩人は x を村側の人間だと信じる、とプレイヤーは信じる

$$BEL^{player}(Is(x, attacked)) \rightarrow BEL^{hunter}(x, VILLAGESIDE)$$

$$(7)$$

- x が村側だと y が信じていることは、y が x が襲撃されていると知っていることを示唆する、とプレイヤーは信じる

$$BEL^{player}(BEL^y(Is(x, VILLAGESIDE)) \rightarrow$$

$$DO(y, know, Is(x, attacked))) \quad (8)$$

式(2)~(8)より、式(1)が導出可能となる。

Otto は Molitz の発言より、自身が信頼されていることを理解する。図 1 中の以下の発言は式(9)を支持する。

「モーリッツ、私を信用してくれますか」

$$BEL^{Otto}(BEL^{Molitz}(Is(Otto, VILLAGESIDE))) \quad (9)$$

本式と式(4)より $BEL^{Otto}(BEL^{Molitz}(Is(Otto, HUMAN)))$ が導ける。式(3)より Albin が護衛されたことが狼 Otto にはわかる。よって式(7)より、hunter が Albin を村側と判断していることが推察できる。また式(9),(5),(6)より、Albin が占い師であると Molitz が信じていることが導ける。

$$BEL^{Molitz}(Is(Albin, seer)) \quad (10)$$

定義より、 $Is(Albin, seer) \rightarrow Is(Albin, VILLAGESIDE)$ これと式(8)より、

$$BEL^{Otto}(Do(Molitz, know, Is(Albin, attacked))) \quad (9)$$

式(2)及び Albin が襲われた事実より、より Molitz が hunter であることが導かれる。よって式(1)が導かれた。こうした複雑な推論を扱うことも、本研究で表した拡張型 BDI 論理で可能となる。

5. 考察とゲーム研究における貢献

本研究で行ったように、BDI 論理を用いていくつかの記述からプレイヤー同士の複雑な推論、入れ子構造になった推

論を再現できることが明らかになった。BDI 論理で記述された事前知識、推論エンジンを各エージェントが保持することで、人狼をプレイする人工知能が作成できることがわかる。実際に作成される人工知能エージェントが、これと同様の推論モデルを持つ必要は無いが、人間と同様にオンライン人狼をプレイするためには、このモデルで記述されるのと同様の課題を処理できる必要があることが推測できる。

人狼において特徴的なのは、相手の信念に対する自分の信念、という入れ子構造の推論が頻繁に登場することである。こうした入れ子型の信念を処理する課題は、認知科学の心の理論などにおいては頻繁に登場する[6]が、完全情報ゲームなどでは必ずしも必要とされない情報処理である。相手の心的モデルを解くことがゲームの勝利に結びつく（他の要件が極めて少ない）という点で、人狼を解く人工知能エージェントを作ることは、ゲーム研究分野に新しい課題を提示することに繋がる。また、相手の思考を読む、という意味で、対話・説得エージェントなど、応用範囲の広い研究を生み出すことが想定される。

6. 結論

本研究では人狼におけるエージェントの推論モデルとして BDI 論理を用い、実際の人狼ゲームで行われている複雑な推論が BDI 論理を用いて記述できることを検証した。

謝辞 本研究は JSPS 科研費 26118006A の助成を受けたものです。

参考文献

- [1] J. L. Austin, "How To Do Things With Words. The William James Lectures delivered at Harvard University in 1955," *J. Symb. Log.*, vol. 36, p. 513, 1962.
- [2] ninjin, "人狼 BBS," 2004. [Online]. Available: <http://ninjinix.x0.com/wolf/>.
- [3] A. S. Rao and M. P. Georgeff, "Modeling Rational Agents within a BDI-Architecture," in *2nd International Conference on Principles of Knowledge Representation and Reasoning*, 1991, pp. 473–484.
- [4] 新出尚之, 高田司郎, and 藤田恵, "拡張 BDI 論理 TOMATO を用いた確率的状態遷移のモデル化とその応用," *情報処理学会研究報告. BIO, バイオ情報学*, vol. 2010, no. 23, pp. 1–9, Dec. 2010.
- [5] "人狼 BBS まとめサイト まとめサイト - 1 村." [Online]. Available: <http://wolfbbs.jp/1%C2%BC.html>.
- [6] D. Premack and G. Woodruff, "Does the chimpanzee have a theory of mind?," *Behav. Brain Sci.*, vol. 1, no. 04, p. 515, Feb. 2010.