

## 人間同士の人狼ゲームで用いられる戦術を反映させた人狼知能の研究

伊藤 幹太<sup>1,a)</sup> Reijer Grimbergen<sup>2,b)</sup>

概要：囲碁のトップ棋士に人工知能が勝利し、完全情報ゲームにおける強い人工知能の研究はおおむね決着がついた。次の題材として不完全情報ゲームである人狼ゲームが注目されており、人狼知能の研究が行われている。人狼ゲームは一般的に有効とされる通説が存在するが、人間が用いる戦術を分析した研究はあまり行われていない。そこで本研究では人間同士で行われた人狼ゲームの対戦ログの分析を行い、2つの戦術を発見した。1つ目は自身に処刑の投票を誘導する戦術である。2つ目は村人騙りといった戦術である。また、既存の人狼知能ではこれらの戦術を用いていないことを示し、人狼知能に導入する際の方法について提案手法を述べた。その結果、自身に処刑の投票を誘導する戦術は既存の人狼と比較して勝率が5%増加することを示したが、村人騙り戦術は勝率が変化せず有効性を見出すに至らなかった。しかし、勝率が変化しなかったことから有効ではないと断言することはできないが、戦術についてより考察を行うことで有効的に使用できるのではないかと考えた。

## Research into a Werewolf Game AI Based on Strategies Used by Human Players

Kanta Ito<sup>1,a)</sup> Reijer Grimbergen<sup>2,b)</sup>

Abstract: In complete information games like Go, AI programs have become stronger than the top human players, so it can be argued that there is not much room for improvements in this area. Research into game AI seems to be shifting to incomplete information games like the Werewolf game which has become a focus of attention in recent years. In Werewolf, a strategy based on “common belief” is considered to be effective, but there has been little research into strategies used by human players. In this research, we have analyzed the logs of games played by humans and found two interesting strategies that have not been used in Werewolf AI. The first is the self-sacrifice, asking to be the next victim. The second strategy is villagers pretending to have a more important role in order to confuse the werewolf pack. We have confirmed that both of these strategies are being used only rarely by current Werewolf AI. We have added the strategies to the strongest Werewolf AI and compared the percentages the villagers won the game. The results of the experiments were that the self-sacrifice strategy increased the winning percentage of the villagers by 5%. However, the confusing strategy did not show any increase in performance. Further research will be needed to conclude whether the confusing strategy is ineffective or if its application rules need to be refined.

---

<sup>1</sup> 東京工科大学大学院 バイオ・情報メディア研究科 コンピュータサイエンス専攻

<sup>2</sup> 東京工科大学 コンピュータサイエンス学部

<sup>a</sup> Tokyo University of Technology Graduate School  
Bionics, Computer and Media Science, Entrepreneurship Program

<sup>b</sup> Tokyo University of Technology  
School of Computer Science

## 1. はじめに

近年、囲碁のトップ棋士に人工知能が勝利をした。このことから完全情報ゲームにおける強い人工知能の研究はおおむね決着がついたと言える。そこで、人工知能研究の次の題材として不完全情報ゲーム、「人狼ゲーム」が注目されており、「人狼知能プロジェクト」[1]といったイベントが立ち上がっている。また、ゲームにおける人工知能以外に医療の分野においても、人工知能が人間とコミュニケーションを取ることが期待されている。人狼ゲームはプレイヤー同士で会話を行い、誰が人狼であるか推理を行う。その中で説得力のある発言や円滑に議論を進めるためにコミュニケーション能力が重要である。人狼ゲームをコミュニケーション能力や論理的な思考、リーダーシップなどを測るため、就職試験に採用する企業も存在する。従って、人狼ゲームが上手な人はコミュニケーション能力が高いと言える。このことから、人工知能の題材として適していると考えられる。

人狼知能プロジェクトでは毎年人狼ゲームエージェント(以下、人狼知能と呼称する)同士で対戦大会が行われ、有効的な戦術の研究が行われている。人狼ゲームは一般的に有効とされる戦術が通説として存在する。通説の有効性については検証が行われている[2]。しかし、人間が用いる戦術、考え方を考慮した人狼知能の研究はあまり行われていない。人間が用いる戦術は通説以外に様々なものが想定されるが、現在の人狼知能では通説に沿ったプレイを前提として作成されている。したがって、人間と人狼ゲームをプレイするためには不十分であると言える。また、人狼知能同士で対戦する研究に対し、人狼知能と人間が対戦する研究はあまり行われていない。

ゲームにおける人工知能の最終目標の1つに、人間と対戦を行うことが挙げられる。その際、人狼知能では人間らしい行動を取ることが求められる。ここで、人間らしい行動の1つに例として「ゲームが終了するまで何も話さない」といった行動が考えられる。人狼知能がこのような行動を取るとは考えにくく、人間らしい行動であると考えられる。しかし、このような行動はゲーム性を損なう行動でもあり、人狼知能に導入した場合、相手に不快感を与えてしまう可能性がある。そこで、本研究では既存の人狼知能では用いられておらず、人間が用いる戦術を人間らしい行動として研究対象とする。

本研究では初めに人間同士の人狼ゲームで用いられていた戦術を分析し、人狼知能に導入する。その後、既存の人狼知能と対戦を行い、戦術の有効性について分析を行う。また、人間同士の戦術を導入した人狼知能と人間を対戦させ、人狼知能が人間と違和感なく対戦が行えるか従来の人狼知能と比較する。

## 2. 本研究で用いるルール

本研究では、人狼知能大会のレギュレーションおよび、先行研究で用いられた選択回答式の人狼ゲーム[3]を元にルールを定めた。本研究で用いる役職を2.1、通常の人狼ゲームと大きく異なる点を2.2および2.3に示す。

### 2.1 役職

使用する役職は人狼、占い師、霊媒師、狩人、狂人、村人の6種類とする。プレイヤーは15人とし、内訳は人狼3人、占い師1人、霊媒師1人、狩人1人、狂人1人、村人8人とする[1]。

### 2.2 発言テンプレートの導入

人狼知能が人狼ゲームを行う上で、日本語特有の解釈の差が問題となる。例として「あなたは人狼じゃない?」という発言が行われたとする。この発言に対し、発言の対象が人狼であることを肯定する解釈と否定する解釈の2通りが考えられる。これは発言内容の「じゃない」の部分に対し解釈が2通り存在することが原因である。前者の場合、「じゃない」を推定の意味として捉えており、後者の場合、否定を意味する形容詞「ではない」の転化として捉えられている。このような解釈の差を避けるため、本研究ではあらかじめ決められた発言内容を発言テンプレートとし、その一部分を変更して発言を行う。

昼時間の発言は、全て発言テンプレートを用いて行う。発言テンプレートはあらかじめ発言内容が決められており、「私は【ROLE】です」といったテンプレートの一部分に該当する語句を入れ発言を行う。先行研究ではテンプレートの数は9種類であったが、本研究では人間同士の人狼ゲームで頻繁に使われる会話を元に最大22種類のテンプレートを用いる。使用例を図2.2、代表的なテンプレートを表2.2に示す。

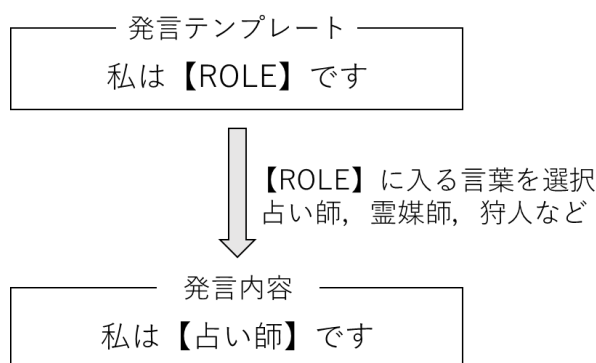


図2.2 発言テンプレートの使用例

表2.2 代表的な発言テンプレート

【WHO】に投票します
【WHO】は【ROLE】だと思えます
【WHO】に賛成します
【WHO】に反対します
私は【ROLE】です
占いの結果は【WHO】は【Werewolf/Villager】です
霊媒の結果は【WHO】は【Werewolf/Villager】です

図2.2では人狼ゲームで頻繁に使われる発言、カミングアウト(以下, CO)を元に作成された発言テンプレートである。COとは、自身の役職を発言することである。プレイヤーはCOを行う際、「私は【ROLE】です」という発言テンプレートを選択する。その際、【ROLE】の部分に役職を入れることでCOを現す発言を行うことができる。その他に例として表2.2にある「占いの結果は【WHO】は【Werewolf/Villager】です」という発言テンプレートが存在する。これは占い師の結果、誰がどのような結果になったかを発言する際に用いる。【WHO】の部分にAさんのようにプレイヤーの名前を入れ、【Werewolf/Villager】に人狼、又は村人を入れることで「占いの結果はAさんは人狼です」といった発言を行うことができる。

### 2.3 ターン制の導入

通常の人狼は昼時間が5分、10分など、時間で区切られている。しかし、発言テンプレートを用いる場合、テンプレートを選択する速度などによって個人差が発生し、プレイヤーの意図しないタイミングで発言が行われる可能性がある。そのため、本研究では時間ではなく人狼知能大会で用いられたターン制を導入する。人狼知能大会で用いられたルールと同じく昼時間は全20ターンとなっており、各プレイヤーはターンごとに1回発言を行うことができる。また、発言を行わないことも可能である。各プレイヤーの発言は同時に行われるため、ターン内の発言順番に意味はない。発言は1日に10回まで行うことができる[1]。

## 3. 人間同士の人狼ゲームの分析結果

人間同士の人狼ゲームを分析したところ、以下ののような戦術が見られた。

### 3.1 自身に処刑の投票を誘導する

「自分を処刑したい」といった旨の発言が人間同士の人狼ゲームで見られた。この戦術の意図は村人陣営の特殊能力を持った役職を誤って処刑しないために用いられる。占い師や霊媒師といった役職は村人陣営にとって有益な情報をもたらす。そのような役職

を誤って処刑しないために、あえて自身が犠牲となって処刑され、能力を持った役職を守るといった戦術である。人間同士の人狼ゲームでは、主に村人が用いており、ゲームの序盤に使用されていた。

2017年に行われた人狼知能大会では10種類の人狼知能で10万回対戦が行われた。対戦ログを自作のツールを用いて分析を行ったところ、この戦術が用いられたケースは170戦で使用され、合計で239回発言が行われていた。この結果は全体の約0.17%であった。

人間同士の人狼ゲームでは50戦中29戦、58%の割合で自身に処刑の投票を誘導する戦術が行われており、「自分は処刑しても問題ない」といった消極的な誘導が22戦、全体の44%の割合で用いられていた。また、「自分を処刑して欲しい」といった積極的な誘導が7戦、全体の14%の割合で用いられた。また、この戦術を用いる役職は主に村人であり、ゲーム開始の初日に多く用いられていた。

このことから、人間同士の人狼ゲームと比較して人狼知能がこの戦術を用いた割合約0.17%は非常に低いと考えられる。また、この戦術を用いたAIは1種類であり、役職が狩人の時のみ用いる、ゲーム終盤にもかかわらず使用するという行動を行っていることから、人間が使用する状況と異なっていた。

人狼知能が用いた状況が適切か分析を行うため、自作のシミュレーションツールを作成し、確率的に有効であるか検証を行った。検証時に以下の条件を加え、村人陣営の勝率を比較した。

- 人狼は役職騙りをしない。
- 狂人は占い師COを行う。
- 昼時間の会話はCO及び役職の能力結果のみ発言する。
- 投票は原則としてCOしていない人からランダムに処刑する。ただし、ゲームの終盤で本物の占い師が判明していない場合、占い師COをしている人から人狼判定を受けている人を優先的に処刑する。
- 人狼はCOしていない人からランダムに襲撃する。
- 占い師はCOしていない人からランダムに占う。
- 霊媒師は人狼判定が発生した場合、COを行い情報を開示する。
- 狩人は占い師を護衛する。
- 狩人は指定した日にちに自身に処刑の投票を誘導する。このとき、他のプレイヤーは誘導に従ったこととし、狩人を処刑する。

以上の条件で、ゲームの何日目に狩人が投票を誘導した場合の勝率について比較を行い、その結果を図3.1に示す。

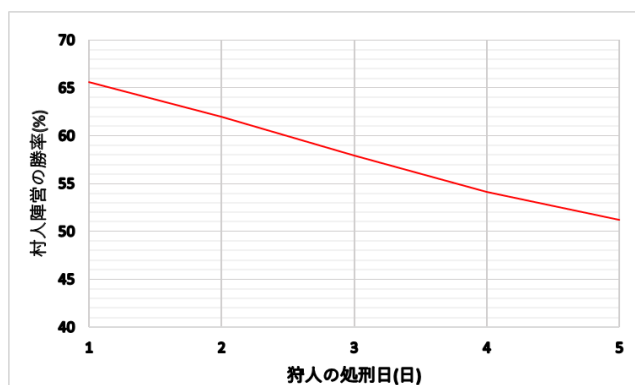


図3.1 戦術を使用した日にち毎の村人陣営の勝率

図3.1より、狩人が自身に処刑の投票を誘導する戦術は、ゲーム終盤に用いた場合、勝率が減少することが分かる。このことから、人狼知能が有効的に使用しているとは考えにくく、現在の人狼知能ではこの戦術を使用していないと言える。

### 3.2 村人騙り

人間同士の人狼ゲームで見られた戦術として、「村人騙り」というものが存在する。通常、人狼ゲームでは嘘を言う人は人狼、または狂人のどちらかであると考えられる。しかし、嘘を言わないとされる村人陣営の役職があえて嘘を言う戦術が村人騙りである。

村人騙りは人狼陣営を騙すことによって、人狼陣営の誤解を招く行動であるが、同時に村人陣営も騙してしまうため、村人陣営を混乱させてしまう可能性もある。

既存の人狼知能では、嘘を言うのは人狼陣営である前提でアルゴリズムが考えられており、村人騙りのケースを想定していない。

## 4. 提案手法

本研究では、第2回人狼知能大会で優勝し、文献[4]で解説されている人狼知能「餛飩」を土台とし作成する。また、人間同士の人狼ゲームで行われる戦術は幾つか存在するが、より顕著に表れる自身に処刑の投票を誘導する戦術について優先的に研究を行う。

### 4.1 自身に投票するよう誘導する

人間同士の人狼ゲームを分析した結果、自身に処刑の投票を行うよう誘導する戦術を用いる際の条件は3つ存在する。優先度の高い順に以下に示す。

1つ目はゲームの終盤ではないことである。自身が処刑された場合、占い師などを誤って処刑しない利点がある反面、自分の視点では確実に人狼以外を処刑してしまう欠点がある。ゲーム終盤では、人狼以外を処刑する猶予が無いため、用いることができない。

2つ目は情報が無い人からランダムに処刑する「グレーランダム」と呼ばれる戦術で処刑の対象を決める場合である。情報が無い人はCOを行っていない人

が対象となる。COを行っていない人の中に特殊能力を持った役職が存在する場合、誤って処刑してしまう可能性があるため、この戦術を用いる場合がある。

3つ目は自身の役職が村人である点である。自身に処刑の投票を行うよう誘導する際の利点は、特殊能力を持った役職を誤って処刑しないためである。従って、この戦術を用いる前提として自身は特殊能力を持っていない村人であることが挙げられる。

以上のことから、以下の3点がこの戦術を用いる状況となる。

1. ゲームの序盤
2. グレーランダムで処刑の対象を決める
3. 自身の役職が村人

また、自身の役職が人狼であった場合、人狼が処刑されてしまうことは人狼陣営にとって大きな不利益であるため、通常は用いないと考えられる。しかし、敢えて用いることで自身が村人陣営であることの主張、仲間の人狼が役職騙る際の補助になる場合がある。そのため、自身の役職が人狼である場合はこの戦術を用いる可能性があるものとする。

### 4.2 村人騙り

村人騙りは人狼陣営のみならず、村人陣営を騙す行動であるため、リスクの高い戦術である。しかし、限定的に村人騙りを使用することでリスクを軽減できると考える。そこで、人狼の役職騙り阻止の目的で使用される。例として、プレイヤー2名が占い師COし、そのうち1名が村人騙りであった場合、村人視点、人狼視点それぞれの内訳を表4.2に示す。

表4.2 2名が占い師CO場合の内訳(順不同)

視点	占い師 A	占い師 B
村人	占い師	人狼 or 狂人
人狼	占い師	狂人

このとき、占い師COしていた村人がCOを撤回した場合、村人視点では占い師が確定する。また、人狼視点では、狂人が占い師COしていたと勘違いさせることができ、人狼の占い師COを阻止することができる。このように、短期間において村人騙りを用いることは有効であると考えられる。

## 5. 実験方法及び結果

本研究では、人狼知能「餛飩」に提案手法を導入し、実験を行った。

### 5.1 自身に投票するよう誘導する戦術の有効性

自身に投票するよう誘導する戦術を餛飩に追加し、有効性について検証を行う。餛飩同士で対戦を行い、



戦術を用いた場合と用いていない場合で勝率の比較を行う。このとき、比較対象となる餽鈍は常に村人とする。また、この戦術はゲーム開始の初日であり、占い師が人狼を特定していない場合に用いるものとする。

通常の戦術を用いて常に村人となる餽鈍を固定餽鈍、投票を誘導する戦術を用いて常に村人となる餽鈍を誘導餽鈍とし、それぞれの場合で1000回対戦を行った。固定餽鈍、誘導餽鈍ともに自身が村人陣営になった際の勝率を表5.1に示す。

表5.1 固定餽鈍と誘導餽鈍の勝率

人狼知能	村人陣営勝率
固定餽鈍	79.1%
誘導餽鈍	84.1%

表5.1より、誘導餽鈍は固定餽鈍より5%勝率が増加しており、自身に処刑の投票を誘導する戦術は有効的であると考えられる。しかし、占い師 CO をする人が多く、特定の役職 CO を行っている人を全員処刑する役職ローラーと呼ばれる戦術が有効的である場合にこの戦術を使用するといった、効果的に用いていないケースが見られた。従って、「ゲームの序盤」「グレーランダムで指名されたとき」「自身の役職が村人」の条件では効果的に使用できていないと考えられる。特にグレーランダムが有効的である状況に効果を発揮しやすい戦術であるため、どのような状況でグレーランダムが有効であるか検証が必要である。

## 5.2 村人騙り戦術の有効性

村人騙り戦術を餽鈍に追加し、有効性について検証を行う。餽鈍同士で対戦を行い、戦術を用いた場合と用いていない場合で勝率の比較を行う。このとき、比較対象となる餽鈍は常に村人とする。また、検証時に村人騙りの前に他のプレイヤーが CO を行うことを防ぐため、村人騙りを行う餽鈍以外は初日の最初のターンに限り、CO ができないものとする。

通常の戦術を用いて常に村人となる餽鈍を5.1と同じ固定餽鈍、投票を誘導する戦術を用いて常に村人となる餽鈍を村人騙り餽鈍とし、村人騙り餽鈍を使用して500回対戦を行った。固定餽鈍、村人騙り餽鈍ともに自身が村人陣営になった際の勝率を表5.2に示す。

表5.2 固定餽鈍と村人騙り餽鈍の勝率

人狼知能	村人陣営勝率
固定餽鈍	79.1%
村人騙り餽鈍	79.6%

表5.2より、村人騙りを用いた戦術は村人陣営の勝率に大きく影響しないことが分かる。実際に餽鈍同士の対戦を見ると、村人騙りを用いたことにより人狼の騙りを制限することに成功した。しかし、人狼の騙りを制限してしまったため、人狼が誰であるか情報が減少してしまったことが考えられる。仮に占い師 CO が3人いた場合、その内訳は一般的に占い師1人、人狼1人、狂人1人となり、占い師 CO を行っている人の中に人狼が紛れ込んでいることが分かる。すなわち、占い師 CO をしている人からランダムに処刑した場合、約33%の確率で人狼を処刑することができる。しかしながら、村人騙りを用いた場合、占い師 CO をしていない人の中から人狼が3人いるといった状況になり、プレイヤー人数が15人であった場合、占い師を除いた14人中3人が人狼となる。従って、この中からランダムに処刑した場合、約21%の確率で人狼を処刑することになり、村人騙りを用いていない場合の方が高い確率で人狼を当てることができる。

以上のことから、村人騙り戦術は有効であると断言することはできない。しかし、勝率が大きく変化していないため無効であると断言することはできない。従って、村人騙りを用いる方法や状況についてより詳細に考察することで、有効な戦術となることが期待できる。

## 5. おわりに

本研究では、実際に人間同士で行われた人狼ゲームの分析を行った。その結果、自身に処刑の投票を誘導する、村人騙りといった2つの戦術が見られ、この戦術は既存の人狼知能では用いられていないことを示した。また、2つの戦術を有効的に用いる方法について考察を行い、戦術の有効性について検証を行った。実験の結果、自身に処刑の投票を誘導する戦術は5%勝率が増加し有効であると考えられるが、村人騙り戦術は勝率が変化せず有効性を見出すに至らなかった。しかし、村人騙り戦術を用いた場合でも勝率が変化しないことから、今後の改善次第で有効的に使用することが可能ではないかと考えられる。

今後の課題として、村人騙り戦術の有効的な使用方法について考察を行うとともに、これらの戦術を用いた人狼知能が人間と対戦を行った場合にどのような挙動となるか検証が必要であると考えられる。また、本研究では餽鈍を用いて実験を行ったが、餽鈍以外を用いた場合や異なる人狼知能を対戦させた場合についてどのような挙動となるか検証が必要である。

## 参考文献

- [1] 人狼知能プロジェクト,  
 <<http://aiwolf.org/>>(2017年7月12日参照)  
 [2] 神田直樹,伊藤毅志(2015)「人狼サーバによる自動対戦を用いた通説の検証~人狼は占い師を騙るべきか~」ゲームプログラミングワークショップ2015論文

集,pp.20-24

[3] 高田和磨,杉原太郎,五福明夫(2015)「人狼ゲームにおける人間らしいエージェントの要素の分析:騙りと同調行動の影響」人工知能学会全国大会論文集,20,pp.529-532

[4] 狩野芳伸,大槻恭士,園田亜斗夢,中田洋平,箕輪峻,鳥海不二夫(2017)「人狼知能で学ぶ AI プログラミング」株式会社マイナビ出版