

# 人狼知能エージェント同士の対戦における 自プレイヤーの役職別の役職推定

石川達也<sup>†</sup> 渡邊裕司<sup>†</sup>

名古屋市立大学理学研究科<sup>‡</sup>

## 1. はじめに

人狼知能プロジェクト[1]とは、不完全情報ゲームの一種である人狼ゲームをプレイする人工知能（人狼知能）について研究するプロジェクトである。本プロジェクトは、「人間と自然なコミュニケーションをとりながら人狼をプレイできるエージェントの構築」を目標とする。その前段階として、人狼知能エージェント同士の対戦におけるモデル化された行動の研究が行われている。人狼ゲームで勝利するためには「役職推定」と「戦略」が重要である。人狼知能エージェント同士の対戦を対象として役職推定を行った大川らの研究[2]では、多層パーセプトロンを用いて71.9%の正答率を実現した。

本研究では、先行研究の特徴に対して、新たな特徴として、全役職共通の特徴を二つ、村人以外の役職にはその役職でしか得られない「役職別特徴」をそれぞれ一つ追加する。そして、自プレイヤーの役職別に分けて推定を行い、正答率が向上することを示す。

## 2. 人狼ゲーム

人狼知能プロジェクトでは、エージェントの性能評価の場として人狼知能国際大会が開かれ、5人と15人の部門が存在する。本研究では5人人狼の部門を対象とする。5人人狼ゲームのルールを以下に説明する。

ゲーム開始時、3人の村人陣営と2人の人狼陣営に分かれ、各プレイヤーには表1に示す役職が与えられる。占い師と人狼は能力を有する。

その後は昼のターンと夜のターンを繰り返す。昼のターンは、生存プレイヤー全員で会話を行うフェーズと、生存者全員の投票によって人狼容疑者を決定し処刑するフェーズで構成される。一方、夜のターンは、能力を持つプレイヤーが能力を使用するフェーズと、各陣営の勝利条件が満たされているかどうかを判定するフェーズで構成される。勝敗判定のフェーズでどちらかの陣営の勝利条件が満たされていれば、ゲームが終了する。村人陣営の勝利条件は人狼を追放することであり、人狼陣営の勝利条件は生存している人狼の数が村全体の生存者の過半数になることである。

昼のターンでの会話には、あらかじめ設定された定型文が用いられる。この定型文を発話プロトコルという。例えば、VOTE、AGREEなどがある。

Estimation of role by the player's role in battles between AI werewolf agents

<sup>†</sup> Ishikawa Tatsuya <sup>†</sup> Watanabe Yuji

<sup>‡</sup> Graduate School of Science, Nagoya City University

表1 5人人狼ゲームにおける役職一覧

役職	陣営	能力	数
村人	村人	なし	2
占い師	村人	生存者から指定した1人を占い、その役職が人狼かどうかを知ることができる。	1
人狼	人狼	指定した1人を襲撃することができる。	1
狂人	人狼	なし	1

## 3. 先行研究

大川らの先行研究[2]では、表2に示す全役職共通の特徴を使用して、多層パーセプトロンを用いて役職推定を行った。学習とテストに人狼知能プレ大会@GAT2017の5人人狼の対戦ログを使用して71.9%の正答率を実現した。

表2 先行研究のプレイヤーXの特徴

特徴	詳細
日にち	現在何日目か
占い師の数	占い師であるとカミングアウトしたプレイヤーの数
被占い結果	Xが人間判定された数と人狼判定された数
何番目の占い師	Xが何番目に占い師と名乗り出たか
占い結果	Xが報告した人間判定の数と人狼判定の数
投票変更数	XがVOTEの発言の対象にしたプレイヤーと、Xが行った投票の対象が異なった回数
生死	Xが生きているか処刑されたか襲撃されたか
肯定的意見の数	Xが別のプレイヤーYに対して村人陣営であると推定した、またはYの会話にAGREEの発言をした数
否定的意見の数	Xが別のプレイヤーYに対して人狼陣営であると推定した、またはYの会話にDISAGREEの発言をした数

## 4. 提案手法

本研究では、先行研究の特徴に対して、表3に示す全役職共通の特徴を二つ、村人以外の役職にはその役職でしか得ることができない特徴をそれぞれ一つ追加する。このような特徴を本研究では「役職別特徴」と呼ぶ。以下では各特徴について詳しく説明する。

「発言回数」は、議論への参加度に役職毎の傾

向が表れると考えて追加する。

「生存日数」は、人狼が処刑されるとゲームが終了することから、生存日数が日数よりも小さい時にプレイヤーXが人狼ではないことがわかるため追加する。

「占い師からの被占い結果」は、占い師視点の推定でのみ使用する特徴である。自分は占い師であると嘘をついている他プレイヤーからの被占い結果を取り除くことができ、より正確な被占い結果として使用できると考えて追加する。

「人狼への肯定的意見、否定的意見の数」は、人狼視点の推定でのみ使用する特徴である。人狼への意見の傾向から、人狼であるプレイヤーを疑っているのかどうかと、疑っている場合陣営はどちらなのかを推定できると考えて追加する。

「狂人からの被占い結果」は、狂人視点の推定でのみ使用する特徴である。自分は占い師であると嘘をつく戦略をとることが多い狂人からの被占い結果を取り除くことができ、より正確な被占い結果として使用できると考えて追加する。

学習とテストには、2020年度第2回人狼知能国際大会の5人人狼の対戦ログ1万試合分から生成したデータを用いる。学習モデルには機械学習のXGBoostを使用し、5分割の交差検証を行う。本研究では新たに「役職別特徴」を追加するため、自プレイヤーの役職別に分けて役職推定を行い、その正答率を算出する。

表3 提案するプレイヤーXの特徴

特徴	詳細
発言回数	XがSKIPとOVER以外の発言をした回数
生存日数	Xが生存していた日数
占い師からの被占い結果 (占い師)	Xが占い師から人間判定された数、人狼判定された数
人狼への肯定的意見、否定的意見の数 (人狼)	Xが人狼であるプレイヤーYに対して村人陣営であると推定した、またはYの会話にAGREEの発言をした数と、人狼陣営であると推定した、またはYの会話にDISAGREEの発言をした数
狂人以外からの被占い結果 (狂人)	Xが狂人以外から人間判定された数と人狼判定された数

5. 結果

表4に自プレイヤーの役職別の推定の正答率を示す。使用したデータと学習モデルが異なるため、先行研究の71.9%の正答率と単純に比較できないものの、全ての役職において80%以上の高い正答率を実現した。特に占い師では90%近い正答率を達成した。

また、追加した特徴の有効性を検証するために、各特徴を加えた場合と加えない場合について、同様に推定を行った。その推定結果を表5に示す。同表で○はその特徴を追加したことを表し、○が全てついていない場合は先行研究の特徴だけを用いたことになる。全ての役職において、先行研

究の特徴に対して、提案する特徴全てを追加すると正答率が著しく向上することが分かった。

表4 自プレイヤーの役職別正答率

自プレイヤーの役職	正答率
村人	0.809020
占い師	0.895875
人狼	0.833725
狂人	0.886175

表5 各特徴の有効性検証

役職	発言回数	生存日数	役職別特徴	正答率
村人			なし	0.67756
	○			0.74660
		○		0.75062
	○	○		<b>0.80902</b>
占い師				0.78675
	○			0.81408
		○		0.87025
			○	0.80625
	○	○		0.88628
	○		○	0.83053
	○	○	○	<b>0.88170</b>
人狼				0.77288
	○			0.83118
		○		0.77295
			○	0.77443
	○	○		0.83318
	○		○	0.83258
狂人				0.77490
	○	○	○	<b>0.83373</b>
				0.76595
	○			0.80358
		○		0.85185
			○	0.77730
	○	○		0.87828
	○		○	0.81648
○	○	○	0.86007	
	○	○	<b>0.88618</b>	

6. まとめ

本研究では、人狼知能エージェント同士の5人人狼の対戦における役職推定に対して、新たに「役職別特徴」などを追加することで高い正答率を実現した。今後の展望としては、本研究の役職推定をもとにした戦略を考え、エージェントを作成する。そして、作成したエージェントが勝率に与える影響を調べるために、人狼知能国際大会に参加したエージェントと実際に対戦させることで勝率を比較する必要がある。

参考文献

[1] 人狼知能プロジェクト, <http://aiwolf.org/>  
 [2] 大川貴聖, 吉仲亮, 篠原歩, 「深層学習を用いた役職推定を行う人狼知能エージェントの開発」、ゲームプログラミングワークショップ 2017, pp. 50-55, 2017