

# チームcashドキュメント

2023/08/19 Kei Harada

## About me

aiwolfpyの作者のcashです。プロジェクトメンバーの1人です。久々に参加します。  
この4月から電気通信大学で教員をしています。[harada@uec.ac.jp](mailto:harada@uec.ac.jp)

## エージェント概要

基本的な思想は以前のもの(GAT2016ごろ)と同様。5460パターンのロジスティック回帰。(Udonメソッドの確率化)

パラメータ等の細かい部分はソースコード(<https://github.com/k-harada/AI-WolfPy>)を参照。

## 今回工夫している点

- パラメータを実際のデータを元に再計算した
  - 特に、村陣営の投票先が初日からかなり正しいということが興味深い結果
- Agentごとにパラメータを推定し、100ゲーム中にAgentを推定する
  - Agentの発言を分類して、その割合でクラスタリングする
    - かなり適当に実施したがこれでも結構精度が良い
  - 推定結果を加味してロジスティック回帰の係数を更新
- 人狼陣営のときに勝手なことをする
  - 人狼陣営の勝率が全体的に悪いので、少し好き勝手する
  - 狂人の場合に初日から黒出しをする(対象は人狼陣営の勝率が最も悪いエージェント)
  - 人狼の場合に、一定程度敵対した相手に対して”IDENTIFIED Agent[XX] WEREWOLF”と宣言する(人狼勝率の高いチームprontoさんがそうしているのを見つけたので真似ているだけ)
  - PPする(主な想定は5人村)
  - 残り3人で狂人COがいる場合に人狼COする(村陣営でも)

## 相変わらず工夫していない点

- CO
  - 占いと霊媒は開幕CO
  - 狂人の場合も開幕で占いCO
  - 騎士と人狼の場合は沈黙(PPは例外)
    - 騎士COは本来すべきケースがあるのは理解しつつ、難しいので諦め
- 難しいことを発言しない
  - COMINGOUT, VOTE, 結果の報告くらいしかしない
- 難しい発言を理解しない
  - ANDとかが入るとこのエージェントは理解できない

## エージェントの詳細

### Udonメソッドの確率化について

時々聞かれるので再掲。

人狼と狂人の配置を考えると、15人村では5460通り、5人村では20通りある。このパターンのどれに該当するかを確率モデル化する。

一様分布から出発して、ゲームの進行にともなって情報が得られるごとに、確率分布を更新する。例えばAがBに投票したらAとBがともに人狼であるcaseの確率が下がり、Cが襲撃されたらCが人狼であるcaseの可能性が0になる。

情報はAgent1人で完結する情報と、Agentの組(順序あり)に対して付与される情報とに分けて管理する。

例えばAがBを占って人狼と出た場合、AとBの関係性が「人間-人間、人間-狂人、人間-人狼、狂人-人間、狂人-人狼、人狼-人間、人狼-狂人、人狼-人狼」の8通りのどれかを、この情報のみから推定する。この推定のロジットを全ペアに対して足し合わせる形で推定値(のロジット)を差分更新する。

ちなみに、毎ターン入る情報は最悪30程度なので、 $15 * 15 = 225$ を全て更新するよりも高速であり、差分更新にすることで地味であるが以前のものよりも高速化している。

この設計の意図は対称性の問題を回避して機械学習をすることにある。普通に深層学習等で役職推定すると、agentのidに関する対称性の問題があり、idによるリークを避けるための工夫が必要になる。5人村の場合は $5! = 60$ 倍にデータ拡張することが考えられるが、15人村では到底現実的ではない。

この設計によって見えなくなる問題がなくはない(例えば3者の関係で起きる矛盾)。だが、役職推定精度が著しく劣るわけではないと考え、気にしないことにした。

### エージェントごとのパラメータ推定について

上記のA-Bの関係の推定を決勝進出エージェントごとに実施した。実際のゲーム進行上は、エージェントの行動傾向をクラスタリングして、そのクラスタの係数をあてることに相当する。(構築時にはエージェントごとの集計をしているが)

試合での適用においては、エージェントが誰であるかの確率分布を試合ごとに更新している。今回は時間の都合でかなり適当に実施したが、それでも高精度であった。これにより、「霊媒COした場合は人狼」などのエージェントごとの癖も含めた推定が可能と考える。

## 今後の課題について

### エージェントの特定について

これは回避することは無理だと思われる。そもそも同じセットで100ゲーム実施する以上は、禁止することもできない。参加者は自分のエージェントが特定されたとしても役職がバレないように行動することを心がける必要がある。

複雑な発話は特定をより容易にするため、さらに難しい。

### 「だまされる」のか？

以前よりも役職推定の精度が高まっていると感じる。これは裏を返せば他のエージェントの行動を見ているということで、騙されやすい可能性がある。エージェントごとの騙し方などの、研究の余地があると思われる。