

Title	花札の「こいこい」ゲームの強化学習によるコンピュータプレイヤ
Author(s)	佐藤, 直之; 上原, 隆平; 池田, 心
Citation	情報処理学会研究報告. GI, 研究報告ゲーム情報学, 2017-GI-38(6): 1-7
Issue Date	2017-07-08
Type	Journal Article
Text version	publisher
URL	http://hdl.handle.net/10119/15102
Rights	<p>社団法人 情報処理学会, 佐藤 直之, 上原 隆平, 池田 心, 情報処理学会研究報告. GI, 研究報告ゲーム情報学, 2017-GI-38(6), 2017, 1-7. ここに掲載した著作物の利用に関する注意: 本著作物の著作権は(社)情報処理学会に帰属します。本著作物は著作権者である情報処理学会の許可のもとに掲載するものです。ご利用に当たっては「著作権法」ならびに「情報処理学会倫理綱領」に従うことをお願いいたします。</p> <p>Notice for the use of this material: The copyright of this material is retained by the Information Processing Society of Japan (IPSJ). This material is published on this web site with the agreement of the author (s) and the IPSJ. Please be complied with Copyright Law of Japan and the Code of Ethics of the IPSJ if any users wish to reproduce, make derivative work, distribute or make available to the public any part or whole thereof. All Rights Reserved, Copyright (C) Information Processing Society of Japan.</p>
Description	

花札の「こいこい」ゲームの強化学習による コンピュータプレイヤ

佐藤 直之^{1,a)} 上原 隆平^{b)} 池田 心^{1,c)}

概要: 花札の「こいこい」ゲームは交互2人零和不完全情報ゲームの一種で、様々な媒体で多くの人に遊ばれているが研究例が少なく、人間の上級者に匹敵する人工プレイヤが開発されたという話も聞かない。そのため我々は強化学習の方策勾配法を用いて強い「こいこい」プレイヤの実装を試みた。まずはゲーム知識に基づいた高級な特徴量を人間が設計し、その重み付き線形和モデルで状態行動の価値を推測して学習を行った。その結果、ランダム行動プレイヤとルールベースプレイヤを上回る強さを獲得した。さらに我々は、状態行動の価値により複雑なモデルを適用すれば更に高い性能が引き出せると考えて、その準備のための実験のみを本稿で行った。ゲームに関する低級な特徴量を設計して、それがANNの学習を通じて適切にゲームの最終スコア予測のために利用できそうな事を確かめた。

SATO NAOYUKI^{1,a)} RYUHEI UEHARA^{b)} IKEDA KOKOLO^{1,c)}

1. はじめに

花札は日本で古来から親しまれてきたカードゲームの1つである。簡単なルールと手ごろなゲームサイズを持ち、スマートフォンのアプリとして手軽に遊ばれたり、ビデオゲームの商業タイトルの中でのミニゲームとして登場したりする。しかし一方で花札を対象とした研究は例が非常に少なく、人間の上級者より十分に強い人工プレイヤが作られた例も我々の知る限りでは無い。

そこで我々は強化学習により強い花札の人工プレイヤ作成を目指す。花札を使った遊び方のうち我々は特にルールが簡明な「こいこい」ゲームに着目する。これは交互2人ゼロ和不完全情報ゲームで、同様の不完全情報ゲームでは麻雀に形式が似ている。既存の麻雀プレイヤ研究では上級者棋譜の教師あり学習がまず基礎の部分に適用された [1] が、花札ではそうした上級者の棋譜が大量に用意しづらい。よって我々は強化学習を用い、適度な強さの単純なルール

ベースプレイヤを相手に少しずつ訓練する事で強いプレイヤの獲得を目指す。

本研究では強化学習の方策勾配法を用いる。この手法ではパラメタライズされた方策を持つエージェントの受け取る報酬を観察し、その獲得報酬の期待値が上昇するようにパラメータを調整していく。

この方策中の“目的関数（状態行動の良さを評価する関数）”として我々は2種類のセッティングを考えている。まずは実装が簡便な方法として、高級な少数の特徴量を重みづけた線形和関数を試みる。次に低級な多数の特徴量を入力信号にとった人工ニューラルネットワーク（以下ANNと呼ぶ）を使う事を想定する。ただし後者の複雑なセッティングについてはまだ人工プレイヤの作成まで行わず、そのための準備実験として特徴量がゲーム結果を正しく反映できるか確かめるだけに本稿ではとどめる。

2. 花札の「こいこい」ゲーム

対象ゲームについて説明する。花札はトランプのように様々な種類の遊び方を持つが、中でも最も有名な遊び方の「こいこい」を我々は扱う。図1はその局面の例である。このゲームは2人のプレイヤが交互に自分の手番に手札から

¹ 北陸先端科学技術大学院大学
JAIST, Nomi, Ishikawa 923-1211, Japan

a) satonao@jaist.ac.jp

b) uehara@jaist.ac.jp

c) kokolo@jaist.ac.jp

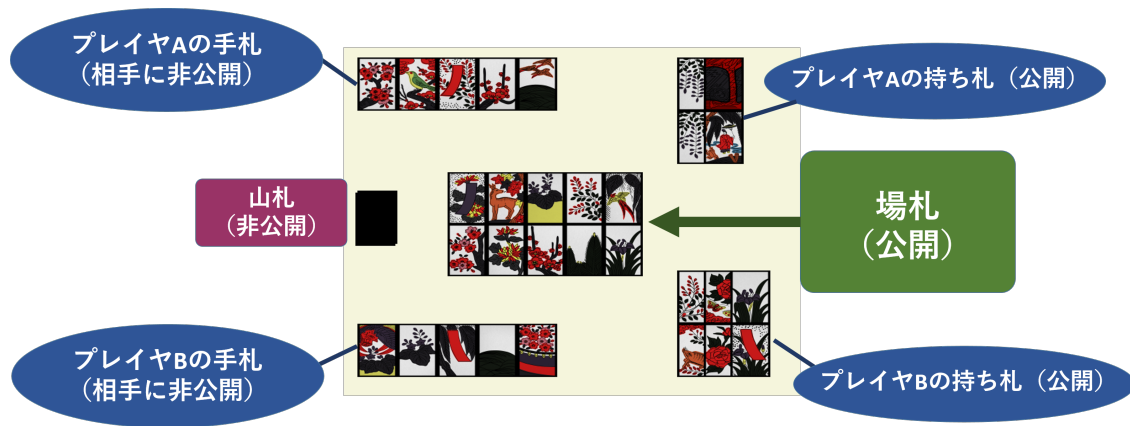


図 1 at.

カードを場に出して、「持ち札(手札とは異なるフィールドである)」を増やす。持ち札のカードの中で「役」と呼ばれる一定のパターンが完成するとそのプレイヤーは「あがり」によって得点をもらう事ができる。以下に詳細を示す。

- (1) 初期状態：各プレイヤーは8枚の手札を持ち、場には場札が8枚表向きに公開されている。また互いに持ち札は0枚ずつである。残りのカードは山札として伏せられた状態で場に積まれる。カードは全48枚で、12種類の花が描かれた札4枚ずつから構成される。
- (2) 先手プレイヤー行動-手札の提出：まず先手のプレイヤーが手札から好きなカードを場に出す。場札に、それと同じ「花」が描かれたカードがある場合には自分の持ち札に加える。その加え方のルールはやや複雑で、場に同じ花のカードが1枚または3枚だけあった場合は、自分の出した札と併せてそれら全てを持ち札に加える。しかし同じ花のカードが2枚だけあった場合のみ、それらのうち好きな方1枚と自分の出した札を自分の持ち札に加える。そして場に同じ花のカードが1枚も無かった場合は自分は何も持ち札に加えず、自分が出したカードも新たな場札として追加される。
- (3) 先手プレイヤー行動-山札めくり：続けて先手プレイヤーは山札の一番上にある札を表向きにめくる。その札についても先ほどと同様の処理に基づき自分の持ち札を増やす。つまりめくったカードと同じ花が描かれたカードが場に1枚か3枚あるときは、めくったカードもあわせてそれら全てを持ち札に加える。2枚だけあるときはその片方とめくったカードを持ち札に加え、1枚もなければ単に場札に加える。
- (4) 後手プレイヤー行動：次に後手プレイヤーが同様に手札の提出と山札のカードをめくる手続きを行う。以上の手続きは両プレイヤー交互に、どちらかのプレイヤーの持ち札に「役」が完成するか提出できる手札が無くなるまで繰り返される。
- (5) 「あがり」と「こいこい」の選択：片方のプレイヤーが

持ち札に「役」を完成させた時、そのプレイヤーはただちに「あがる」事によって役に応じた得点をもらうか「こいこい」を宣言する事によってゲームを継続する事を選ぶ事ができる。「こいこい」は通常、自分が更に高い点数の役を狙えそうでなおかつ相手が役の完成から遠そうな場合に選ばれるオプションである。

- (6) ゲームの終わり：片方のプレイヤーがあがった場合に、そのプレイヤーに得点が与えられてゲームが終わる。あるいはどちらのプレイヤーもあがらないまま8枚の手札を使いきった場合にゲームが終わり、この場合はどちらのプレイヤーにも得点が与えられない。このゲームは確率的な要素も大きく、通常は何回もゲームを繰り返してその合計スコアを競い合う。

以上がこいこいゲームのルールであるが、役の種類や特定条件下での手札の配り直し等のローカルルールが加えられる事もよくある。特に近代的なオンラインゲームとして提供される場合には、既存のものとの差別化のためか、かなり大がかりな独自の特殊ルールが導入されている事もある。

そのように様々なルールがある中で本稿で採用するルールはかなりシンプルで、認められる役は五光(10点)四光(8点)雨四光(7点)三光(5点)赤短・青短・猪鹿蝶(各5点)タネ・短冊・カス(各1点から1枚増加で1点追加)のみである。しばしば用いられる「手四」と「くつつき」は無しとする。^{*1}また得点の倍増に関するルールも取り入れない。

3. 既存研究

不完全情報のカードゲームに関する研究はポーカーを対象にしたものが盛んである。MartinらはCounter Factual Regretと呼ばれる指標の最小化を強化学習で行う事でHeads-up limit Texas holdemルールで ϵ ナッシュ均衡の導出に成功した[2]。またTammelinはその発展形とし

^{*1} 初期状態で手四の形ができていたらゲームはやり直しとし、くつつきはそのままゲームを続ける。

ての指標である CFR+ を提案し学習収束速度の向上を示した [3]。さらにゲーム規模を 2 人用に限定していない多人数ポーカーにて Counter Factual Regret の適用を試みた研究 [4] やナッシュ均衡から敵モデルを構築して搾取する研究がある [5]。

また花札のこいこいと同様に、役作りとあがり、あがりを延期し更なる高得点を狙うオプションを備えた不完全情報ゲームには麻雀があるが、麻雀では上級者棋譜が使えるため教師あり学習がしばしば用いられる。水上らは教師あり学習からの麻雀プレイヤーの作成 [1] とその発展形としてモンテカルロシミュレーションを加えた手法を提案した [6]。さらに麻雀では人工プレイヤーの技術発展を目的とした競技用のサーバーが用意されている [7]。

4. 適用手法

我々はアプローチとして強化学習の方策勾配法 [8] を選んだ。ポーカーは各状態での行動がコール・レイズ・フォールドの 3 つだけだが花札は全カードについて 48 種の行動が想定されるため行動政策の定式化が難しい。また麻雀のように利用可能な上級者棋譜が見当たらないため教師あり学習も困難である。そこで強化学習による動的な強さ向上を目指した。また TD 学習ではなく方策勾配法を選んだのは一般に方策勾配法のほうが学習の難度が下がると言われている（状態価値や状態行動価値を推定しなくても、報酬を最大化する行動のみ求められれば良いため）ためである [10]。

方策勾配法は「パラメタライズされた方策で行動するエージェント」の得る期待報酬を、報酬に対する勾配方向に各パラメータを動かす事で増大させようとする。一口に方策勾配法といっても様々な流儀のものがあるが、我々が用いる方法は五十嵐らが提案した手法を参考にしている [9]。

まずエピソード開始から t 回目の行動を行う状態 s_t で行動 a_t をエージェントが選択する確率（つまり方策）を

$$\pi_a(a_t|s_t; \vec{\omega}) = \frac{\exp(E_a(a_t, s_t; \vec{\omega})/T_a)}{Z_s} \quad (1)$$

と定めているものとする。ここで $E_a(a_t, s_t; \vec{\omega})$ は、 s_t での a_t の選ばれやすさを表す「目的関数」と呼ばれる指標である。 T_a は温度パラメータで、方策で選ばれる行動のバラつきに影響する。 Z_s は s_t でのエージェントの全可能行動についての選択確率値を 1 以下、総和 1 にし正規化するための項であり、式 (1) の右辺の分子を全可能行動に対し足し合わせる事で求められる。

そしてある 1 回のエピソードに対応した報酬を r とし、 i 回目のエピソードの報酬を r_i 、エピソードを重ねていった場合の報酬合計 R の期待値を $E[R; \vec{\omega}]$ と表すとき、この期待値を重みベクトルの調整によって極大化しようとするためのパラメータ更新式を考える。方策が式 (1) のとき $\nabla_{\vec{\omega}} E(R; \vec{\omega})$ は一般に、

$$\begin{aligned} \nabla_{\vec{\omega}} E(R; \vec{\omega}) &= \frac{1}{T} \sum_i r_i p(x_i; \vec{\omega}) \sum_{t=1}^{L_i} \{ \nabla_{\vec{\omega}} E(a_{i,t}, s_{i,t}; \vec{\omega}) \\ &\quad - \sum_{a' \in A_{s_{i,t}}} \pi(a'|s_{i,t}; \vec{\omega}) \nabla_{\vec{\omega}} E(a', s_{i,t}; \vec{\omega}) \} \end{aligned}$$

という形で計算される。ここで $p(x_i; \vec{\omega})$ はそのエピソードの生成確率で、 $\prod_{t=1}^{L_i} \pi(a_{i,t}|s_{i,t}; \vec{\omega})$ に等しい（ただし L_i はエピソードの行動ステップ数）。これを用いて

$$\vec{\omega} \leftarrow \vec{\omega} + \eta \nabla_{\vec{\omega}} E(R; \vec{\omega}) \quad (2)$$

と示されるパラメータ更新を行えばよい。ただし η は小さな正の定数である。

このように式 (2) のような更新式を、エピソードを繰り返しながら重みパラメータに適用し続けて、なんらかの終了条件を満たしたときに繰り返しを打ち切るのが本稿で利用する方策勾配法の概要である。

5. 実験 1：高級な少数の特徴量による線形和

まず我々は、ゲーム知識に依存した少数の特徴量を用いた方策勾配法の人工プレイヤーの性能を評価するため対戦実験を試みた。

5.1 使用特徴量

ゲームの特徴量を以下に示す。まず花札のゲーム状態を 4 種に分類し、行動（手札から 1 枚選んで場に出す）に 8 つの特徴を設ける。これらによって後の実験で花札の状態行動を $4 \times 8 (= 32)$ 種の特徴で表す。まず状態の分類を以下のように行った。

- s_x : 自分も相手もあと 1 行動でアがりうる
- s_y : 自分のみあと 1 行動でアがりうる
- s_z : 相手のみあと 1 行動でアがりうる
- s_w : 自分も相手もアガりにあと 2 行動以上必要とする
次に行動の特徴は、場の札を取れる（持ち札に加える）行動と場の札を取れない行動について設けた。場の札を取る行動が持ちうる特徴は以下 4 つである。
- f_{a-g1} : 自分にとって“高得点貢献度”が最高の札を取る手である
- f_{a-g2} : 自分にとって“早上り貢献度”が最高の札を取る手である
- f_{a-g3} : 相手にとって“高得点貢献度”が最高の札を取る手である
- f_{a-g4} : 相手にとって“早上り貢献度”が最高の札を取る手である

この“高得点貢献度”と“早上り貢献度”の詳細な定義は付録に譲るが、つまり「高い得点がもらえる役の完成」と「(点が安くても) なるべく早い役の完成」に貢献する度合いとして我々が適当に定めた指標である。

そして場の札を取れない行動についての特徴は以下 5 つ

である。

- f_{a_ng1} :ゲームに2枚残って(どちらの持ち札にもなっていない)いて自分が(手札に)2枚持っている花の札を出す
- f_{a_ng2} :ゲームに4枚残っていて自分が3枚持っている花の札を出す
- f_{a_ng3} :ゲームに2枚残っていて自分が1枚持っている花の札を出す
- f_{a_ng4} :ゲームに4枚残っていて自分が2枚持っている花の札を出す
- f_{a_ng5} :ゲームに4枚残っていて自分が1枚持っている花の札を出す

これらの分類と特徴を使うと「こいこい」ゲームの全ての状態行動は s_x, s_y, s_z, s_w のいずれか1つの状態において行う, $f_{a_g1}, \dots, f_{a_g4}$ のうち0個以上4個以下の特徴を備えた行動であるか, $f_{a_ng1}, \dots, f_{a_ng5}$ のうち1個の特徴を備えた行動である. よって状態行動は s_x からの f_{a_g1} , s_x からの f_{a_g2}, \dots, s_w からの f_{a_ng5} という32種のバイナリな特徴量で表現できて, その重みづけ線形和を方策の目的関数とし, その重みをパラメータとした.

5.2 実験条件

方策勾配法による人工プレイヤーが用いる特徴量は前項に示した通りである. 1回の対戦を1エピソードとして, 各エピソードは不連続なものとして独立に扱われる. エピソードに割り当てられた報酬は, アガリを達成したときの役の点数を用い, 自分があがった場合はプラス, 相手があがった場合はマイナスの符号をつける. 引き分けは報酬0である.

対戦相手も人工プレイヤーで, ベースラインとしてのランダム行動プレイヤーと, ゲーム知識に基づくIf-thenルールで処理を書き下したルールベースプレイヤーである. このルールベースプレイヤーはランダムプレイヤーと1,000戦して平均獲得点数2.52点をおさめる程度の強さがあった. 方策勾配法プレイヤーはこれらのプレイヤーとそれぞれ30,000エピソード(GPU無しのマシンで高速化処理なしで処理時間約10分程度)にわたって対戦した. ただし方策勾配法プレイヤーは常に先手番でゲームを始める. 実験環境は自作プラットフォームを用いた.

方策勾配法プレイヤーは学習率 η を $0.5 \times 100 / (100 + Epi)$ とした. ただし Epi は現在までに経験したエピソード数を示すため, η は0.5から最終的に約0.001まで下降する. 重みパラメータの初期値は0.0以上1.0以下の一様ランダムで定めた.

なお, このゲームの醍醐味である「こいこい」をするかしないかの判断はどのプレイヤーも原始モンテカルロ手法のシミュレータにより行う. すなわちこの実験に登場する人工プレイヤーはどれも, 役を完成させた時に「こいこい」を

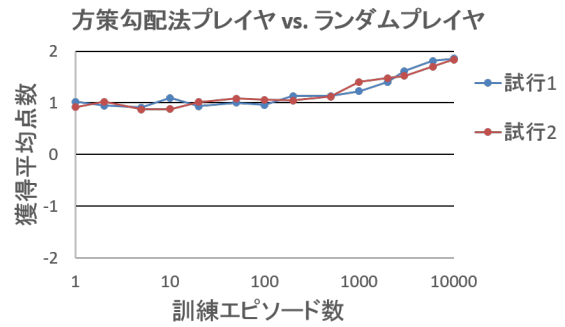


図2 at.

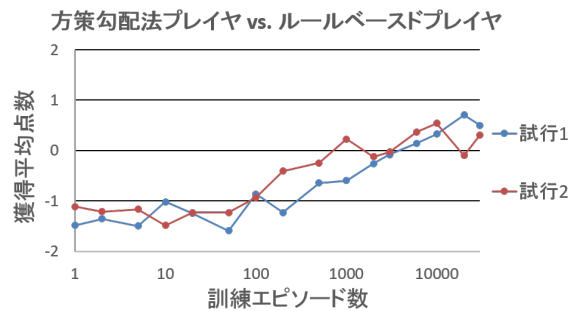


図3 at.

宣言した後のゲーム展開を300回シミュレーションする. そのシミュレーションの中でゲームを進めるのは先手後手ともにランダム行動プレイヤーだがそのプレイヤーはもはや「こいこい」はせず, あがれるチャンスには必ずあがる. このシミュレーションの平均獲得スコアが, ただちに「アガリ」を選んだ場合より高い場合は実際のゲームで「こいこい」を宣言する.

5.3 結果

対戦の結果, 図2と図3のような結果が得られた. 方策勾配法プレイヤーはランダムプレイヤーに対して搾取できる点数を上昇させていき, ルールベースプレイヤーに対しても互角以上に戦うようになった. この対戦で方策勾配法プレイヤーは常に先手をとり, 花札は先手番が少し有利なゲームであるが, ルールベース型同士が対戦すると1000戦したときの先手番の平均報酬値は0.051だった. そのため方策勾配法プレイヤーは0.5点分程度ルールベースより性能が上回ったことになる.

6. 実験2: 低級で多数の特徴量による人工ニューラルネットワーク

項5の方策勾配法プレイヤーは我々が用意したルールベース型を性能で上回ったが, しかし特徴量は36個と少なく, 高級とはいってもそれぞれゲーム状態の数値の四則演算から得られる程度のものである. こうした特徴量の線形和関数がゲームのあらゆる局面での最善行動を精度よく, 例えば人間の上級者に匹敵するほどの精度で表現できるよ

うな表現力を方策に与えられるとは考え難い。

そこで我々は低級な特徴量を多く用いた入力信号を生成して ANN に渡し、その入出力で目的関数を近似することで更なる性能の向上可能性を考えた。本稿ではそのための準備実験として、ANN がゲームの盤面特徴量から結果得点（の平均値）を高い精度で推測できるかを試験した。この実験によって、我々が設計した低級な特徴量とその後のゲームの結果得点を正しく近似できる程の表現力をニューラルネットに与えるかを確かめる事を狙う。

6.1 使用特徴量

以下に述べる 252 個のバイナリな状態特徴量を用いてゲーム中の状態行動を表現し、「ある状態行動の特徴量ベクトル」を入力にして「その後のゲーム終了時に（先手プレイヤーが）得る得点」を出力とする機械学習問題を ANN で試みた。

状態行動の特徴量づけには、いわゆる事後状態の考え方をを用いた。つまり状態行動 s, a の特徴量を、 s に a を適用した直後の状態（より厳密には不確定要素や相手行動の影響を受ける直前までゲームを進行させた状態とした）の状態特徴量とした。この状態特徴量は 240 個の「ゲームの現状状態」に関する特徴と 12 個の「ゲームの行動履歴」に関する特徴から成る。

まずは事後状態の持つ、ゲームの現状状態に関する特徴量について述べる。花札に使用される 48 のカードそれぞれに ID をつけて、行動判断の主体となるエージェント（自分）から見てそれぞれがどの場所にあるかで状態を表現する。その場所とは以下の 5 種である。

- 自分の手札の中
- 自分の持ち札の中
- 相手の持ち札の中（※持ち札は場に公開される）
- 場札の中
- それ以外の場所、つまり自分からは見えないどこか（※具体的には相手の手札の中か山札の中）

これをそれぞれのカードごとに設けて「ID1 のカードは自分の手札にあるか?」、「ID1 のカードは自分の持ち札に含まれるか?」、..., 「ID48 のカードは自分から見えない位置にあるか?」という $5 \times 48 (= 240)$ 種のバイナリ特徴が事後状態の特徴量に含まれる。

そして事後状態の到達に至った、ゲームの行動履歴に関する特徴量も我々は用意した。これは敵の手札内容へのヒントに結び付くものであり、「場に“松”の花の札がある状況で相手はどの札も取らない行動をした事がある」、「場に“梅”の花の札があるのに...」というような全 12 種の花に関する特徴である。花札の「こいこい」ゲームでは普通、場札を取れる行動があるときに他の行動を選ぶ事に利点がない。つまり相手が場の“松”の札を見逃して、そして場のどの札も取らないで何か適当な札を場に提供した場合は

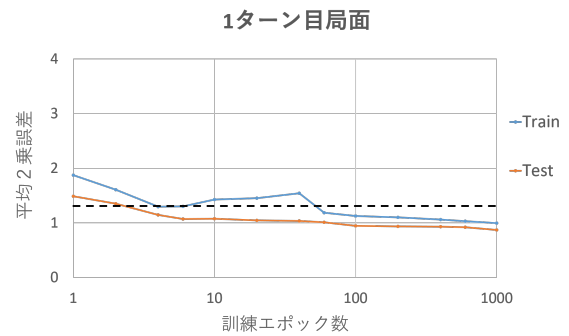


図 4 at.

相手は“松”の札を 1 枚も手札に持っていない可能性が濃厚になる。

この行動履歴に関する特徴量は ANN に利用する低級な特徴量としてはややゲーム特有の知識に頼ったものかもしれない。しかしこういう情報は花札においてプレイヤーがよく利用し、また判断への影響力も大きいものなので現段階ではひとまず導入の上での実験を行った。

6.2 実験条件

ANN の入力に使う状態特徴量は前節に示した通りである。出力となる結果得点であるが、ゲームのランダム性を考慮して複数回シミュレーションした結果の平均値をとった。具体的には、入力の局面に含まれる不完全情報の部分をランダムにシャッフルしながら 300 回の対戦をルールベース型（5 章で使ったものと同じ）同士の間で行って、入力局面で手番だったプレイヤーの獲得点数にプラスをつけてその相手の獲得点数にマイナスをつけたものの平均を ANN 出力用の教師信号にした。

ANN は入力層と隠れ層と出力層 1 個ずつの 3 層でそれぞれニューロンの数は 252, 200, 1 である。隠れ層の発火関数はシグモイド関数で出力層の発火関数は線形関数にした。学習率は 0.05 で L2 正規化項を用いてその係数は 0.1、データは 1 ターン目、7 ターン目、11 ターン目の局面のみを使用した。訓練データを 10,000 個、テストデータを 300 個ずつ用意した。

6.3 結果

図 x と x に 2 乗誤差の学習曲線を示す。

点線はベースラインで、出力をある 1 つに固定（この場合全ての出力値の平均値）した場合の 2 乗平均誤差である。よって我々の用意した特徴量は入出力の係数に意味のある情報をとらえられていると考えられる。

7. まとめ

我々は花札の「こいこい」ゲームを対象に、方策勾配法による人工プレイヤー実装を試みた。その結果、初歩的なルールベース型に獲得点数期待値で勝ち越す結果が観察され

7ターン目局面

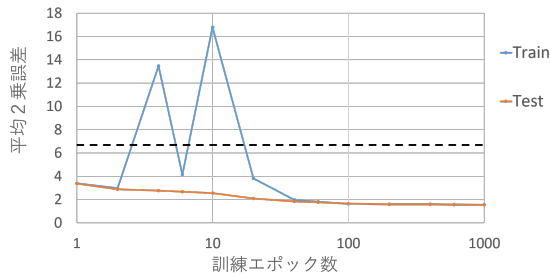


図 5 at.

11ターン目局面

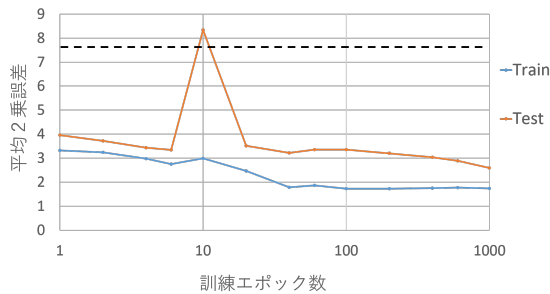


図 6 at.

Turn1 出力-教師信号散布図

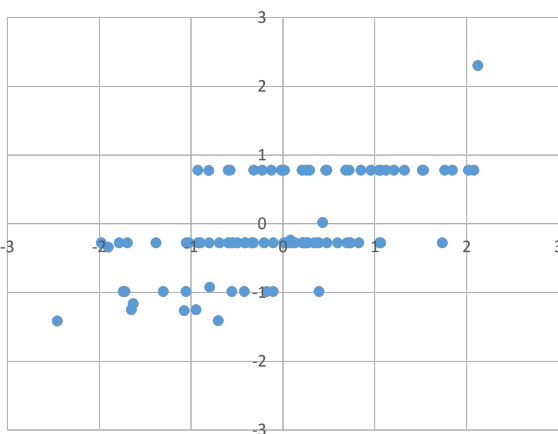


図 7 at.

た。また我々は、より高い表現力を持つ方策による性能向上を見据えて、局面の低級な特徴量と結果得点による ANN の学習を試した。これも、健全に学習の精度が高まっていく様子を観察できた。

これらの結果から、次に我々は ANN による方策勾配法プレイヤーの実装、そして市販タイトルなどに含まれる強い花札人工プレイヤーとの対戦実験と評価を行う予定である。

謝辞 謝辞

参考文献

[1] 水上直紀, 中張遼太郎, 浦晃, 三輪誠, 鶴岡慶雅, 近山隆. 降りるべき局面の認識による 1 人麻雀プレイヤーの 4 人

Turn7 出力-教師信号散布図

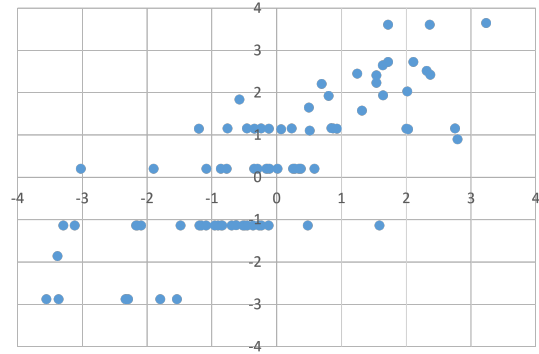


図 8 at.

Turn11 出力-教師信号散布図

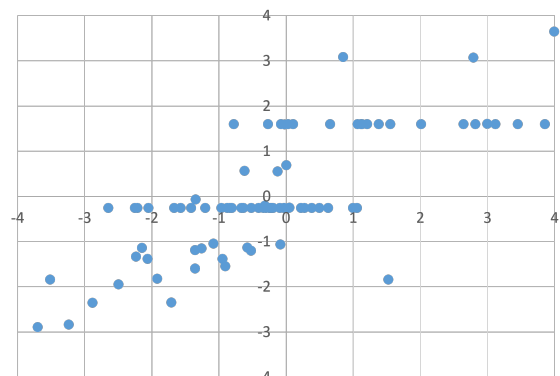


図 9 at.

麻雀への適用. The 18th Game Programming Workshop 2013, pp.1-7 (2013).

[2] Martin Zinkevich, Michael Bowling, Michael Johanson, and Carmelo Piccione . Regret Minimization in Games with Incomplete Information. Advances in neural information processing systems 2007, pp.1729-1736 (2007).

[3] Tammelin Oskari. Solving large imperfect information games using CFR+. arXiv preprint arXiv:1407.5042 (2014).

[4] Risk Nick Abou, Szfron Duane. Using counterfactual regret minimization to create competitive multiplayer poker agents. Proceedings of the 9th International Conference on Autonomous Agents and Multiagent Systems: volume 1. International Foundation for Autonomous Agents and Multiagent Systems 2010. p. 159-166 (2010).

[5] GANZFRIED Sam, SANDHOLM Tuomas. Game theory-based opponent modeling in large imperfect-information games. The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 2. International Foundation for Autonomous Agents and Multiagent Systems, pp.533-540 (2011).

[6] 水上 直紀, 鶴岡 慶雅. 牌譜を用いた対戦相手のモデル化とモンテカルロ法によるコンピュータ麻雀プレイヤーの構築, The 19th Game Programming Workshop 2014, pp.48-55 (2014).

[7] 「麻雀サーバーの紹介」, http://www.logos.ic.i.u-tokyo.ac.jp/~mizukami/slide/majong_server.pdf (accessed 2017-06-22).

[8] Williams Ronald J. Simple statistical gradient-following

algorithms for connectionist reinforcement learning. Machine learning, 8(3-4) pp.229-256 (1992).

- [9] 五十嵐治一, 石原聖司, 木村昌臣. 非マルコフ決定過程における強化学習—特徴的適正度の統計的性質—. 電子情報通信学会論文誌 D 90.9 pp.2271-2280 (2007).
- [10] Reinforcement Learning: An Introduction Second edition, in progress ****Draft****. <http://ufal.mff.cuni.cz/straka/courses/npfl114/2016/sutton-bookdraft2016sep.pdf> (accessed 2017-06-22).

付 録

A.1 付録