

繰り返し囚人のジレンマゲームにおける プレイヤーの信頼度を用いた行動決定のシミュレーション

大阪大学大学院 工学研究科 電子情報エネルギー工学専攻

森下 紗枝 (Sae Morishita), 巽 啓司 (Keiji Tatsumi), 谷野 哲三 (Tetsuzo Tanino)

Department of Electronics and Information Systems,

Graduate School of Engineering, Osaka University

1 はじめに

人間が社会生活を営む上で「信頼」は必要不可欠である。山岸らは繰り返し囚人のジレンマゲームを用いた被験者実験により、信頼関係についての研究を行った [6]。囚人のジレンマゲームにおいて、むやみに相手を信頼するプレイヤーは逆に搾取される可能性が高く、協調行動は成り立たないと直感的には思われる。しかし山岸らの実験において高得点をあげたのは他者を信頼し、囚人のジレンマゲームでは最初は協調行動を選択した人であった。これらの他者への信頼が高い人々は、同時に相手の信頼性についての情報に敏感であり、相手が信頼できないと判断した場合は協調行動を控える傾向があることがわかった。

本研究では、信頼研究を参考に、構築した相手を評価し行動を決定するプレイヤーのモデルを提案する。このプレイヤーは、過去の相手の行動によって相手を評価し、その評価と自分内部の評価基準に基づいて対戦での行動を確率的に決定する。このモデルを用いて、計算機上のプレイヤーによる進化的なシミュレーションを実行しその結果を解析する。またこのモデルをより現実の問題に適用しやすくするために、 n 人ゲームへと拡張する。同時にモデルの特徴を明らかにするため、過去の履歴を参照し行動を決定する Lindgren モデル [2] を使用したシミュレーションも行い、挙動を比較する。

2 Lindgren モデル

繰り返し囚人のジレンマゲームを用いた進化的なシミュレーションでよく使用される Lindgren モデルについて述べる。Lindgren モデルでは、 n 人のプレイヤーが2人ゲームの対戦を行う。プレイヤーは戦略と呼ばれる行動指針に従って対戦での行動を決定する。この戦略はプレイヤーの遺伝子として表され、同じ戦略(遺伝子)をもつプレイヤーが複数存在する。各プレイヤーは他のすべてのプレイヤーと1対1の繰り返し対戦を行い利得を得る。全プレイヤーが対戦を終えたら、ある戦略 i をもつプレイヤーが獲得した利得 g_i と平均利得 g_{ave} を比べる。平均利得よりも大きい利得を獲得したプレイヤーは多くの子孫を残せるので、次の世代では獲得利得が高い戦略をもつプレイヤーが増加する。また、一定の確率でプレイヤーの遺伝子に突然変異が起き、新たな戦略が生成される。

2.1 対戦と利得

Lindgren モデルでは、プレイヤーの対戦で得られる利得は表1のような利得表を用いる。表1のような

表 1: 囚人のジレンマゲームの利得表

囚人 A \ 囚人 B	黙秘 (協調)	自白 (裏切り)
黙秘 (協調)	3 \ 3	0 \ 5
自白 (裏切り)	5 \ 0	1 \ 1

利得の大小関係では、利己的なプレイヤー同士では(自白, 自白)がゲームの解となるが、(黙秘, 黙秘)の行動の組み合わせの方が両者にとって望ましい。しかし自分だけが黙秘に行動を変えても相手が行動を変えなければ搾取されるだけである。このような状態を囚人のジレンマという。囚人のジレンマは現実世界の様々な局面で見られるものである。本研究では Lindgren モデルと信頼度モデルの両方でこの囚人のジレンマゲームをプレイヤー間の相互作用として用いる。

2.2 プレイヤーの戦略と履歴

Lindgren モデルでは、プレイヤーが持つ戦略テーブルと、相手と自分の行動の履歴を参照して今回の行動を決定する。記憶長(履歴の長さ)が2の場合は、前回の相手の行動と前回の自分の行動の2つを記憶しておき、その行動の組合せを戦略テーブルのアドレスとして今回の行動を決定する。記憶長が2の場合のすべての履歴と代表的な戦略テーブルを表2に示す。

表 2: 記憶長 2 のときの履歴と戦略テーブル一覧

履歴		戦略テーブル			
前回の自分の行動	前回の相手の行動	AllD	TFT	ATFT	AllC
D	D	D	D	C	C
D	C	D	C	D	C
C	D	D	D	C	C
C	C	D	C	D	C

2.3 個体数の変化

個体数とは、ある戦略をとるプレイヤーの数である。ある戦略の次の世代の個体数は、以下のように決定する。ある戦略 i が戦略 j と対戦したときに得られる利得を g_{ij} 、戦略 i の個体数が全プレイヤーに占める割合を x_i とすると、戦略 i が他の戦略と対戦して獲得する総利得 g_i は $g_i = \sum_j g_{ij}x_j$ となる。また全戦略の平均利得を g_{ave} は $g_{ave} = \sum_i g_i x_i$ となる。世代 t から世代 $t+1$ へ進むとき、戦略 i の全プレイヤーに占める個体数は次の式に従って変化する。

$$x_i(t+1) - x_i(t) = d_{agent}(g_i - g_{ave})x_i(t) \quad (1)$$

ただし d_{agent} は個体増加率である。

2.4 戦略の突然変異

プレイヤーの戦略テーブルは、世代交代時に一定の確率でランダムに1箇所が書き換えられる。これを突然変異という。このとき、もともとのテーブル内容がCだった場合はDに、Dだった場合はCに変化する

2.5 n 人ゲームへの拡張

2.5.1 n 人ゲームにおける利得決定

n 人ゲームでの利得は、Lindgren の n 人ゲームの研究 [3] を参考に以下の式を用いた。

$$V(C|n_C) = \frac{Rn_C}{n-1} + \frac{S(n-n_C-1)}{n-1} \quad (2)$$

$$V(D|n_C) = \frac{Tn_C}{n-1} + \frac{P(n-n_C-1)}{n-1} \quad (3)$$

ただし $V(C|i)$ は自分を除いて i 人が協調行動をとり、自分が協調行動をとったときの利得、 $V(D|i)$ は自分を除いて i 人が協調行動をとり、自分が裏切り行動をとったときの利得、 n_C は自分を除いた協調 C を選択したプレイヤーの数とする。

2.6 n 人ゲーム用 Lindgren モデル

本研究では、Lindgren モデルを n 人ゲームへ拡張するにあたり次のような行動決定方法を考案した。 n 人で行う対戦においてグループ内で過半数を占めた行動をそのグループの選択した行動とする。各プレイ

ヤーはグループの行動を2人ゲームにおける相手の行動とみなして、戦略を決定する。またこのグループの行動と自分の行動によって自分が獲得する利得を決定する。

3 信頼度モデル

実社会におけるゲーム的状况での意思決定では、一般的に次のような2段階の手順を踏むと考えられる。まず、相手について事前に知っている情報、相手との過去の対戦結果、相手の第一印象などから相手の次の行動を予測する。次に、その予測した相手の行動に対して、自分がどう行動するかを決定する。

本研究では以上の意思決定方法をモデル化した信頼度モデルを提案する。第1段階の相手の行動の予測として、相手との過去の対戦結果に応じて、相手に対する信頼度を設定する。この信頼度とは、相手が協調行動をとるだろうという自分の相手に対する期待の高さを表す数値である。

第2段階における相手の評価に基づく行動決定として、相手に対する信頼度と自分内部の基準である行動決定関数によって対戦で選択する行動を決定する。行動決定関数とは信頼度によって自分が協調行動をとる確率を決める関数である。この関数は各プレイヤー個別のもので、相手の信頼度に対する各プレイヤーの特徴を表すものである。

3.1 信頼度を用いた行動決定

プレイヤー*i*はプレイヤー*j*に対し信頼度 $t_{ij}(i, j = 1, \dots, n, i \neq j)$ という値を保持する。対戦時、プレイヤー*i*がプレイヤー*j*に対し協調行動*C*をとる確率 p_{iC} は次のように決定する。

$$p_{iC} = f_i(t_{ij}) \quad (4)$$

f_i をプレイヤー*i*の行動決定関数と呼ぶ。行動決定関数は、プレイヤー*i*が持つグラフ形状 G_i 、軸 K_i の2つの要素によって決まる。今回の実験では、グラフ形状としては右上がりの高信頼度-協調型、右下がりの高信頼度-裏切り型、山型の限定-協調型、谷型の限定-裏切り型の4種類の形状を使用した。4つの行動決定関数の軸 $K = 0$ 、傾き $a = \frac{2}{7}$ のときのグラフを図1に示す。

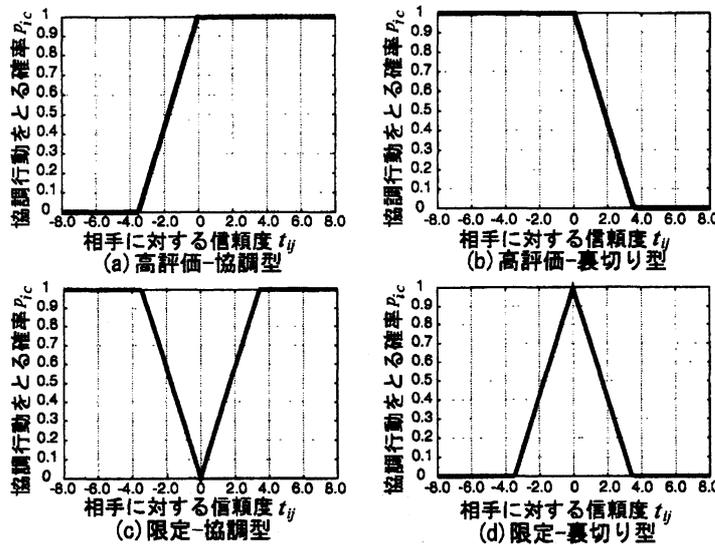


図 1: 行動決定関数一覧

1対戦ごとに行動の組み合わせによって利得とは別に信頼度が増減する。信頼度の増減は表3のようにする。対戦の結果が共に協調行動であった場合、互いの信頼度は T_{CC} だけ変化する。また共に裏切り行動であった場合は、信頼度は T_{DD} だけ変化する。信頼度の定義より、 $T_{CC} > 0$ 、 $T_{DD} < 0$ とする。

表 3: プレイヤーの信頼度増減表

自分 \ 相手	C	D
C	T_{CC}	T_{CD}
D	T_{DC}	T_{DD}

3.2 n 人ゲームへの拡張

自分の保持する他の各プレイヤーに対する信頼度に基づき、信頼度モデルを用いた2人ゲームと同様にして行動を決定し、 $n-1$ の行動のうち数が多い方をその対戦での行動として採用する。グループ内で過半数を占めた行動をグループの行動とし、グループの行動と自分の行動から利得を決定する。その対戦での各プレイヤーの行動と自分の行動に基づき、相手に対する信頼度を更新する。

4 シミュレーション実験

今回は、繰り返し囚人のジレンマゲームを用いたシミュレーションとして、Lindgren モデルと信頼度モデルのそれぞれについて、対戦を行うグループの大きさが2人の場合と5人の場合の2種類を行った。

以下に n 人ゲームを行う場合のシミュレーションの流れを大まかに示す。

1. その世代で生存している(その戦略をとるプレイヤーが存在している)戦略からゲームを行う n 個の戦略を重複を許して選び、対戦を行う。
2. 1を生存しているすべての戦略について行い、戦略のすべての組合せについて、利得を計算する。
3. 対戦で獲得した利得から適応値を計算し、個体数の増減を行う。
4. 一定の確率で突然変異を行う。
5. 以上を1世代とし、繰り返す。

4.1 実験条件

4.1.1 共通条件

シミュレーションで使用したパラメータの一覧を表4に、利得表を表5に示す。

表 4: 使用したパラメータ一覧

最低対戦回数	100
未来係数 ω	0.98
人口増減係数 d_{agent}	0.1
新戦略への人口分割比 d_{div}	0.1
突然変異発生確率 p_{mute}	0.001
世代数	50000

表 5: 実験で用いる囚人のジレンマゲームの利得表

$i \setminus j$	C(協調)	D(裏切り)
C(協調)	1.0 \ 1.0	0 \ 1.5
D(裏切り)	1.5 \ 0	0.2 \ 0.2

4.1.2 Lindgren モデルの実験条件

Lindgren モデルでは実験条件として以下を用いた。

- 初期履歴について

繰り返し対戦の初回の対戦時に各プレイヤーがもっている履歴は、(前回の自分の行動, 前回の相手の行動)の組み合わせとして(C, C), (C, D), (D, C), (D, D)の4種類からランダムに1つ選んだものとした。

- 初期戦略について
シミュレーション開始時に存在する戦略は1種類とし、AllCの場合、AllDの場合、TFTの場合を行った。
- 戦略の表記
Lindgrenモデルの戦略は戦略テーブルというビット列で表記する。本実験では前回の対戦における自分と相手の行動の組合せに基づいて次の対戦での行動を決定するので、前回の対戦でとり得る4種類の行動の組合せにそれぞれ対応した4つの行動の並びが戦略となる。ビット列は $[a_3a_2a_1a_0]$ のように表し、 a_3 は前回の対戦が(C, C)であったときに次の対戦でとる行動、 a_2 は前回が(C, D)のとき、 a_1 は前回が(D, C)のとき、 a_0 は前回が(D, D)のときに次の対戦でとる行動とする。ここでCを1, Dを0として、1と0の並びで1つの戦略を表す。例えば、対戦では常に裏切り行動Dをとる戦略であるAllDは[0000]と表される。前回の相手の行動を繰り返すTFTは[1010]となる

4.1.3 信頼度モデルの実験条件

信頼度モデルにおいて、今回のシミュレーション実験の条件では遺伝子的に可能な戦略は20種類である。各戦略は戦略番号によって識別できる。戦略と戦略番号の一覧を表6に示す。

表 6: 戦略と戦略番号一覧

		グラフ形状			
		高評価-協調型	高評価-裏切り型	限定-協調型	限定-裏切り型
軸	-4	0	5	10	15
	-2	1	6	11	16
	0	2	7	12	17
	2	3	8	13	18
	4	4	9	14	19

信頼度モデルでは実験条件として以下を用いた。

- 初期戦略について
シミュレーション開始時に存在する戦略は4種類のグラフ形状についてそれぞれ軸 $K=0$ の場合の4通りを行った。
- 信頼度の変化について
対戦後の信頼度増減としては、表3で $T_{CC}=2, T_{DD}=-2, T_{CD}=T_{DC}=-1$ とした。

4.2 実験結果

4.2.1 2人Lindgrenモデル

最終的に平均利得が0.99になった場合と0.20となった場合の2種類に分かれた。ただし最終的な平均利得とは、30000世代から50000世代までの全戦略の利得の平均値である。初期戦略が協調的なものほど最終的に協調社会を築いている。

平均利得が0.99となった場合では戦略[1001]がプレイヤーのほとんどを占めている。この戦略同士が対戦した場合、どのような履歴の組合せで対戦が始まっても遅くとも3対戦目には相互の協調行動(C, C)を確立し、その後ずっと協調行動が続く。よって戦略[1001]が多数を占める場合は全体の平均利得はほぼ1.0となる。ただし、AllD([0000])に対しては一方向的に搾取されるので、AllDのような裏切りやすい戦略がプレイヤーの多数を占めている場合は、戦略[1001]は増加しない。

平均利得が0.20となった場合では先に AllD が増加してしまったために戦略 [1001] は増加できなくなった。代わりに戦略 [1000] が最終的には AllD とほぼ同数を占めている。しかしほとんどの場合戦略 [1000] は AllD と同じ行動をとるため (D, D) となる対戦が大多数となり、全体の利得もほぼ0.2となった。

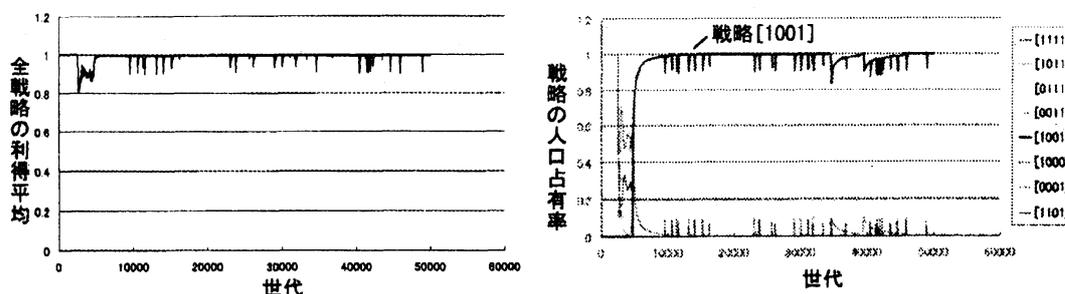


図 2: 2人 Lindgren モデル 利得の推移および戦略分布

4.2.2 2人信頼度モデル

初期戦略を 2(高評価-協調型, 軸 $K=0$), 7(高評価-裏切り型, 軸 $K=0$), 12(限定-協調型, 軸 $K=0$), 17(限定-裏切り型, 軸 $K=0$) の場合で行ったが, 派生する戦略の順番や時期が異なるだけで, 全体的にはシステムの挙動の違いは見られなかった。代表的な全戦略の利得の平均値の推移を図 3 に示す。

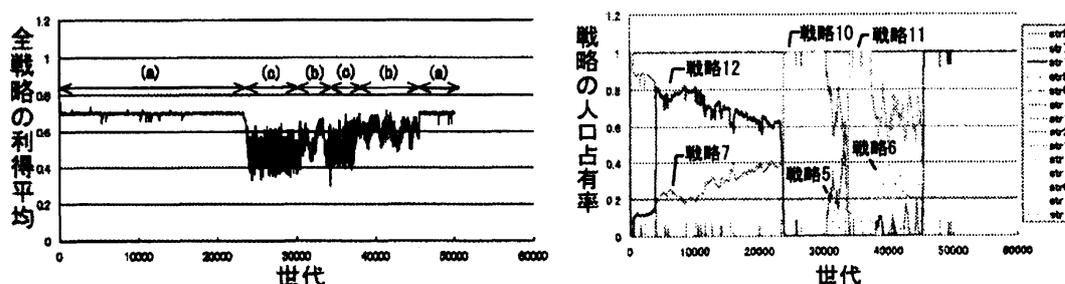


図 3: 2人信頼度モデル 利得の推移および戦略分布

全戦略の利得平均の推移をみると, 初期遷移時の状態以外に次の 3 種類の状態があった。

- 平均利得 0.7 (図 3 (a))
戦略 12 と 7 が集団を占めているとき, 利得はほぼ 0.7 となる。この平均利得 0.7 の状態は他の状態に比べると長期間安定しやすく, 50000 世代のシミュレーション中, 初期遷移にかかった期間 1000 世代を除いた 49000 世代でこの状態が続くこともある。
- 平均利得 0.5-0.7 (図 3 (b))
戦略 10 と 5, または戦略 11 と 6 の組み合わせがプレイヤーを占めている場合, 平均利得は 0.5 から 0.7 の間を振動する。1 世代ごとの利得の振動幅は 0.1 程度だが, それとは別により長い世代にわたる利得の変化による波のようなものが見られる。
- 平均利得 0.4-0.6 (図 3 (c))
戦略 10 または 11 が全プレイヤーを占めている場合, 平均利得は 0.4 から 0.6 の間を振動する。上記の 0.5 から 0.7 の振動と異なり, この振動には周期性は見られない。

どの状態も安定ではなく, 初期戦略に関わらずどの状態にも遷移する可能性がある。基本的には初期状態 (初期戦略に依存する) から 5000 世代以内に図 3 の (a), (b), (c) いずれかの状態に遷移し, その後は (a), (b), (c) を規則性なしに遷移を続ける。

4.2.3 5人 Lindgren モデル

初期戦略が AIIC, AIID, TFT の場合についてシミュレーションを行ったが、派生する戦略の順番が異なるだけで、全体的にはシステムの挙動に違いは見られなかった。

発生する戦略のうち生きのびる戦略は, AIID, [0010], [1000], [1010] の4種類のみであった。他の戦略は突然変異によって発生してもすぐに絶滅している。AIID 以外の戦略は, 前回の対戦時にグループの過半数が協調的でなければ協調行動をとらないという用心深い戦略といえる。初期戦略によって初期に多数を占める戦略は異なるが, 最終的にはこの4つの戦略がそれぞれ全プレイヤーの0.25ずつを占め安定する。全戦略の平均利得はシミュレーション開始時からほぼ0.2となった。

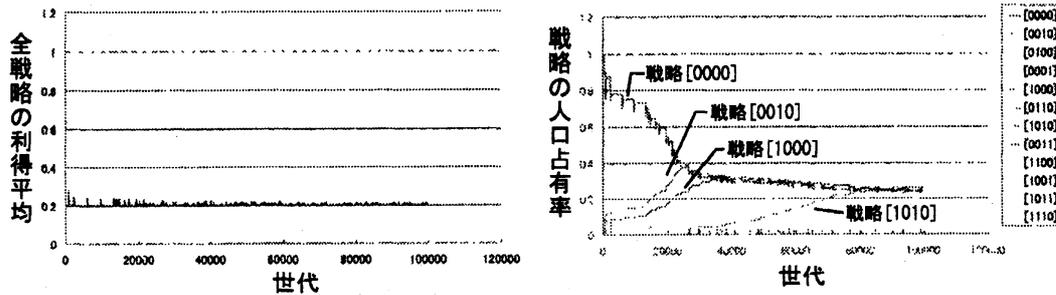


図 4: 5人 Lindgren モデル 利得の推移および戦略分布

4.2.4 5人信頼度モデル

最終的な平均利得が0.64となる場合と0.24となる場合の2種類の状態があった。どちらになるかは, 初期戦略に依存していると思われる。

最終的な平均利得が0.64になる場合は, 戦略の分布としては, 戦略2が約0.5, 戦略12が約0.35を占め, 残りを戦略4, 1, 0がほぼ同数ずつとなった。全体的に高評価-協調型が多く, 5人 Lindgren モデルよりも協調的な社会を築いているといえる。

最終的な平均利得が0.24となる場合は, 戦略12がプレイヤーのほとんどを占めてしまい, 他の戦略は人口を伸ばせなかった。行動が食い違くと信頼度が下がることから, いずれすべてのプレイヤーの互いの信頼度は0以下になり, すべてのプレイヤーが常にDをとる状態になる。しかし軸がストッパーの役割を果たしているため, 平均利得は0.2ではなく0.24になっていると考えられる。

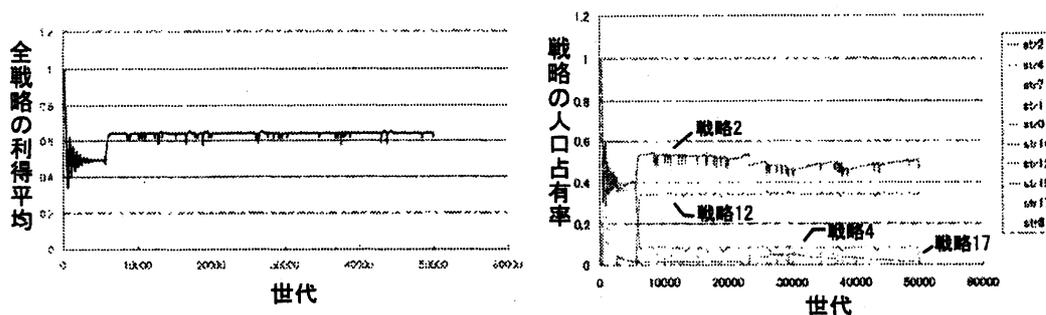


図 5: 5人信頼度モデル 利得の推移および戦略分布

4.3 考察

各実験の試行回数と, 全試行の最終的な利得の平均の一覧を表7に示す。また, 各モデルの2人ゲーム, 5人ゲームでの代表的な利得の推移を図6に示す。Lindgren モデルを用いた2人ゲームでは, 初期戦略が

表 7: 実験条件別利得一覧 (四人のジレンマゲーム)

使用モデル	ゲーム人数	初期戦略	試行回数	平均利得	備考
Lindgren	2人	AllC	10	0.75	0.99が8回, 0.20が2回
Lindgren	2人	TFT	5	0.67	0.99が3回, 0.20が2回
Lindgren	2人	AllD	5	0.20	
信頼度	2人	2	5	0.66	
信頼度	2人	7	5	0.63	
信頼度	2人	12	5	0.68	やや0.7の期間が長い
信頼度	2人	17	5	0.68	1.0になる期間がある
Lindgren	5人	AllC	10	0.20	
Lindgren	5人	TFT	5	0.20	
Lindgren	5人	AllD	5	0.20	
信頼度	5人	2	10	0.56	0.64が8回, 0.24が2回
信頼度	5人	7	5	0.24	
信頼度	5人	12	5	0.24	
信頼度	5人	17	5	0.32	0.24が4回, 0.64が1回

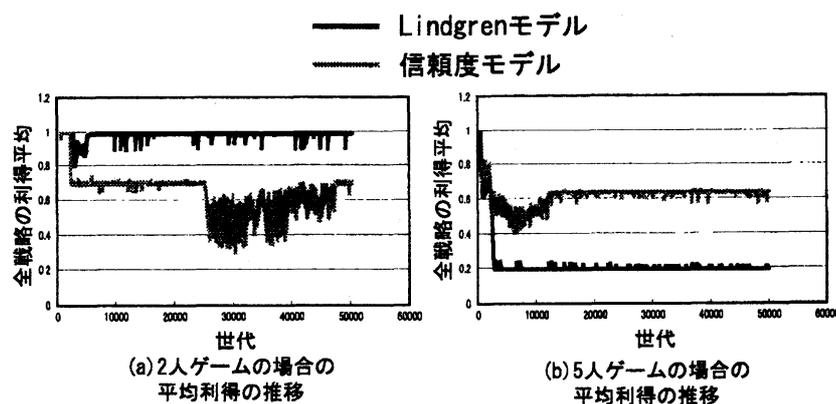


図 6: 2人ゲームと5人ゲームにおける利得の推移

協調的なものほど平均利得が高くなっている。これは初期戦略が協調的なものほど協調的な社会を達成した回数が多いためである。Lindgren モデルでは利得は0.99か0.20のどちらかとなり、協調的か非協調的かがはっきりしている。0.99を達成している場合、戦略分布としては[1001]が全プレイヤーの9割以上を占めている。

一方信頼度モデルを用いた2人ゲームでは、初期戦略に関わらず平均利得は0.6台となっており、初期戦略の影響が小さい。利得の推移を見ると、0.7となる期間と振動する期間が入り混じり、安定することはなかった。これはLindgrenモデルの戦略[1001]のように一人勝ちする戦略が存在せず、常に対抗する戦略が勢力を伸ばす可能性があることを示す。

Lindgrenモデルを用いた5人ゲームでは、初期戦略に関わらず平均利得は常に0.2となった。これは、今回用いたLindgrenモデルの n 人ゲームへの拡張方法として多数決を用いたことによると考えられる。この方法では5人中3人以上が協調行動をとらなければその対戦でのグループの決定は協調行動とはみなされず、各プレイヤーからは対戦相手は協調行動をとりにくいプレイヤーのように見える。非協力的なプレイヤーの割合が半数を超えた時点で、各プレイヤーはAllDのプレイヤーと対戦していることと同じ状態となるため、協調行動をとる動機が2人ゲームの場合より低くなると考えられる。

信頼度モデルを用いた5人ゲームでは、最低でも利得は0.24、初期戦略が協調的な戦略であれば最終的な利得が0.64となることもあり、Lindgrenモデルを用いた5人ゲームより高い利得を達成している。最も高い利得0.64を達成したときの戦略分布では、戦略2が全体の約半分、戦略12が約3割となっていた。戦略2は相手が信頼できれば自分も協調行動をとるという素直な戦略であり、戦略12はある程度以上に相手が協調的であればその裏をかき、裏切って利得を得ようとする戦略である。これらの戦略がバランスをとって一定の利得に落ち着くという結果は非常に興味深い。

5 まとめ

本研究では、繰り返し n 人ゲームを行うモデルとして信頼度モデルを提案し、Lindgrenモデルとの比較を行った。このとき、多数決の考え方をもとにLindgrenモデルを n 人ゲームへ拡張した。信頼度モデルとは、相手の行動を評価し、その評価に基づき行動を決定するプレイヤーを表現したモデルである。それぞれのモデルについて、2人ゲーム、5人ゲームの両方について利得表として囚人のジレンマゲームを用いたシミュレーション実験を行い、結果を比較した。

信頼度モデルを構築するにあたり参考にしたのは、山岸らによる人間同士の信頼関係の研究結果であった。被験者実験からは、特に情報がない場合の他者一般への信頼の度合いが高い人ほど他者の信頼性についての情報に敏感で、搾取されにくい性質を示すという結果が得られている。

本研究で提案した信頼度モデルでは、高評価-協調型の関数を持つプレイヤーは、相手に対する信頼度が増加すれば協調しやすくなり、信頼度が下がれば協調しにくくなるという最も単純な行動パターンを持つ。言い換えればこれは他者の信頼性に対して敏感なプレイヤーである。信頼度モデルによるシミュレーション実験において平均利得が1.0程度となる協調的な社会では、高評価-協調型やそれに似た行動をするプレイヤーが集団を占めていた。この結果と山岸らの研究を合わせて考えると、信頼度モデルは現実の意思決定状況の一面をある程度の妥当性を持って表現していると言える。

今後の課題としては、まずLindgrenモデルと信頼度モデルの n 人ゲームでの対戦方法の改良が考えられる。特に今回のシミュレーションでは2つのモデルの対戦方法に違いがあったため、 n 人ゲームの実験結果の差が顕著であった。両方のモデルに適用できる汎用的な n 人ゲームの対戦方法を考案したい。信頼度モデルで用いた行動決定関数についても、現実問題との整合性を考えて検討していきたい。また信頼度モデルにおけるパラメータのより詳細な解析が必要である。

参考文献

- [1] C. Fang, S. O. Kimbrough, A. Valluri, Z. Zheng and S. Pace "On adaptive emergence of trust behavior in the game of stag hunt," *Group Decision and Negotiation*. vol. 11, pp. 449-467, (2002).
- [2] K. Lindgren, "Evolutionary Phenomena in Simple Dynamics," *Artificial Life II*, pp. 295-312, (1991).
- [3] K. Lindgren and J. Johansson, "Coevolution of strategies in n-person prisoner's dilemma," in J. Crutchfield and P. Schuster, *Evolutionary Dynamics - Exploring the Interplay of Selection, Neutrality, Accident, and Function*, (Addison-Wesley, 2001).
- [4] 岡田章, "ゲーム理論," 有斐閣 (1996).
- [5] R. Suzuki and T. Arita, "Evolutionary analysis on spatial locality in the N-person iterated prisoner's dilemma," *International Journal of Computational Intelligence and Applications*, vol. 3, No. 2, pp. 177-188, (2003).
- [6] 山岸俊男, "信頼の構造-こころと社会の進化ゲーム," 東京大学出版会 (1998).