

私的観測下の繰り返し囚人のジレンマにおける協力のダイナミクス

西野上 和真 *
Kazuma Nishinoue五十嵐 瞭平 *
Ryohei Igarashi岩崎 敦 *
Atsushi Iwasaki

概要

本論文は私的観測下の繰り返し囚人のジレンマにおける協力のダイナミクスを分析した。私的観測は、プレイヤーが相手の行動についてノイズを含むシグナルを観測し、そのシグナルを他のプレイヤーは観測できないという特徴をもつ。ここで、どんな戦略の組が均衡になるかはゲーム理論の有名な未解決問題の一つであり、本研究では戦略空間を状態数2以下の有限状態機械に限定したレプリケータダイナミクスの帰結から、どのような戦略が生き残るかを吟味した。その結果、利得構造に応じて、4つの社会（非協力、不寛容、相互協力、周期協力）が現れることがわかった。とくに周期協力社会で、非専門家の間では有効と信じられてきたしっぺ返し戦略が最大多数になりうるが、他の戦略と共存しなければならないという不安定さをもつ。一方他の社会では、ある特定の戦略が人口のほぼ全てを占めるようになる。さらにノイズと突然変異率に関する感度分析から十分広いパラメータにおいて同じ傾向を保つことがわかった。

1 はじめに

無限回繰り返しゲームは、長期的関係にあるプレイヤー間の（暗黙の）協調を説明するためのモデルである [1]。主に経済学分野で企業間の談合といった協調行動を分析するために発展してきた [2]。暗黙の協調を実現するには、プレイヤーが相手の行動をある程度観測できることが前提となる。これまで、相手の行動が完全に観測できる完全観測 (perfect monitoring) のケースについては多く論じられている [3, 4, 5]。しかし、現実には相手の行動が完全に観測できない不完全観測 (imperfect monitoring) のケース、つまり、プレイヤーが相手の行動についてノイズを含むシグナルを観測し、そのシグナルを他のプレイヤーは観測できない場合がある。これはとくに、不完全私的観測 (imperfect private monitoring) のケースと呼ばれる [6, 7, 8]。不完全私的観測付き無限回繰り返しゲーム (infinite repeated games with imperfect private monitoring) の特徴は、プレイヤーが相手の行動に関してノイズを含む観測 (シグナル) を私的に受け取ると仮定する点にある。いいかえると、あるプレイヤーが相手の行動について観測したシグナルと異なるシグナルを他のプレイヤーが観測しているかもしれない。不完全私的観測付き無限回繰り返しゲームにおいてどのような振る舞い (戦略)

が最適なのかについては、ゲーム理論における代表的なゲームである囚人のジレンマの例でさえ十分にわかっていない。例えば、部分観測可能マルコフ決定過程 (Partially Observable Markov Decision Process, POMDP) を用いて均衡を計算する手法 [1] が知られているが、その計算量は一般には決定不能と知られている。

ここで本論文では、均衡の代わりに突然変異付きのレプリケータダイナミクス [9, 10] の帰結を用いて、私的観測下の繰り返し囚人のジレンマでどんな戦略が生き残るかを分析する。レプリケータダイナミクスは、進化ゲーム理論でよく用いられるダイナミクスの1つであり、頻度依存淘汰モデルを用いて最適な戦略を探るため、その帰結が比較的計算しやすい。無限に大きな集団を仮定し、各戦略をとるプレイヤーの頻度の時間的変化を計算する。無限集団を仮定することでモデルの持つ確率性を無視できるので、そのダイナミクスは決定論的となり、微分方程式で記述できる [11]。利得が高くなる戦略をとるプレイヤーの人口は増加し、低くなる戦略をとる人口はより良い戦略へ取って代わられてやがて絶滅するといった具合に自然淘汰の過程を表現する。厳密には、均衡とダイナミクスの帰結の2つに包含関係はない。ある戦略の組が均衡になったとしても、それがダイナミクスの帰結で最大多数を占めるとは限らない。その逆も必ずしも言えない。また、均衡は複数存在しうるので、ダイナミクスがどの均衡に収束するを事前に予測できないし、そもそも均衡を構成する戦略が存在しない場合もある。そのような場合でも、自然淘汰のダイナミクスはどんな戦略が生き残るかを示せる。

理論生物学や進化ゲームの文脈では、自分の行動を出し間違える振動 (trembling-hand) は盛んに研究されてきた [5]。しかし、その重要性にも関わらず、相手の行動を見間違える私的観測を進化ゲームの文脈で網羅的に分析することは非常に難しいと考えられてきた。その理由の1つとして、振動と私的観測は戦略と情報の構造が異なるため、従来の成果が適用できないことに挙げられる。また、一般には複雑な行動計画となる繰り返しゲームの戦略を有限状態機械 (Finite State Automaton, FSA) で記述するとき、私的観測をどのようにモデル化し期待利得を計算するかよくわかっていなかった。

ここで本論文では、まずプレイヤーが取りうる戦略を状態数2以下のFSAに限定する。つまり、プレイヤーの今日とった行動と観測したシグナルから明日の行動への写像を考える。振動の場合はこれに加えて、自分が行動を取り間違えた後の振る舞いを考慮しなければならなくなる。戦略の状態数が同じ

* 電気通信大学大学院情報理工学研究所

であるとき、私的観測より摂動の方が戦略空間を制限することになる。戦略を FSA に限定したときの期待利得はマルコフ決定過程に基づいて計算し、その利得表をもとに突然変異付きレプリケータダイナミクス [12] を計算する。

その結果、利得構造に応じて、4 つの社会（非協力、不寛容、相互協力、周期協力）が現れることがわかった。非協力社会では、常に裏切る戦略 (ALLD) のみが生き残る、つまり単独の戦略で人口のほぼ全てを占める、次に不寛容社会ではトリガー (Grim trigger, GRIM) 戦略、はじめに協力し、相手が一度でも裏切ったら二度と許さない戦略が生き残る。さらに相互協力社会では、1 期相互処罰 (1-period Mutual Punishment, 1MP) という新しい戦略が生き残る [1]。最後に、非自明な均衡戦略がないときに発生する周期協力社会では、非専門家の間では有効と信じられてきたしっぺ返し戦略 (Tit-For-Tat, TFT) が最大多数になりうるが、他の戦略とサイクルを構成し共存しなければならない。単体の戦略として安定しないことがわかった。さらにノイズと突然変異率に関する感度分析で十分広いパラメータでも傾向が変わらないことを示した。

2 モデル

本章では文献 [1] に基づいて、2 人私的観測付き無限回繰返しゲームをモデル化する。ここでプレイヤー $i \in \{1, 2\}$ は成分ゲームを無限期間 $t = 0, 1, 2, \dots$ に渡って繰り返す。各期においてプレイヤー i は有限集合 A から行動 a_i を選択し、その行動の組を $\mathbf{a} = (a_1, a_2) \in A^2$ とする。次に、プレイヤー i は \mathbf{a} に関する私的なシグナル $\omega_i \in \Omega$ を観測する。 \mathbf{w} をシグナルの組 $(\omega_1, \omega_2) \in \Omega^2$ とする。また、プレイヤーが \mathbf{a} を選択したとき \mathbf{w} が生起する同時確率を $o(\mathbf{w} | \mathbf{a})$ とし、この同時確率を与える分布のことをシグナル分布と呼ぶ。成分ゲームは無限回繰返し行われるので、プレイヤー i の割引利得和は割引因子 $\delta \in (0, 1)$ により $\sum_{t=1}^{\infty} \delta^t g_i(\mathbf{a}^t)$ となる。ただし、 $g_i(\cdot)$ の値は利得表によって定められた値に従う。

本論文では利得表として表 1 に示す囚人のジレンマを用いる。表中の C は協力行為を、 D は裏切り行為を表す。囚人のジレンマの利得構造は $g > 0, l > 0$ であり、このとき D は厳密な支配戦略となる。また、囚人にジレンマでは $|g - l| < 1$ が要求される。もしこの条件が成り立たないとすると、繰返し囚人のジレンマにおいて協力和裏切りを交互に出すほうが、純粋な協力よりも利得が高くなってしまい、純粋な協力が維持できなくなる。

次にプレイヤー 2 の行動に関するプレイヤー 1 のノイズを含む観測をプレイヤー 1 の私的シグナルとし、 $\omega \in \{g, b\}$ (good, bad) とする。正しい観測ではプレイヤー 2 が C を選択した際のプレイヤー 1 の私的シグナルは g 、 D を選択した際の私的シグナルは b となる。プレイヤー 2 についても同様である。よく使われる不完全私的観測のシグナル分布にほぼ完全観測がある。ここでは、両プレイヤーが正しいシグナルを観測する確率は p 、片方のプレイヤーが間違っただけのシグナルを観測する確率はそれぞれ q とする。また、 $1 - p - 2q$ の確率で両方のプレイヤーが間違っただけ

表 1: 囚人のジレンマ ($g > 0, l > 0$, および $|g - l| < 1$)

| | | |
|-----------|-------------|-------------|
| | $a_2 = C$ | $a_2 = D$ |
| $a_1 = C$ | 1, 1 | $-l, 1 + g$ |
| $a_1 = D$ | $1 + g, -l$ | 0, 0 |

表 2: (C, C) のときのシグナル分布

| | | |
|-----------|-----------|--------------|
| | $w_2 = g$ | $w_2 = b$ |
| $w_1 = g$ | p | q |
| $w_1 = b$ | q | $1 - p - 2q$ |

シグナルを観測する。例として、 (C, C) が実現した場合のシグナル分布を表 2 に示す。ただし、両プレイヤーが正しいシグナルを観測する確率 p が最も高くなるように設定する。

プレイヤーの戦略は、そのプレイヤーの過去の行動と受け取ったシグナルから現在の行動への写像で表現される。FSA は繰返しゲームの戦略を簡略に表記する方法であり、本研究では、状態数 2 以下の非同相な 26 個の FSA を用いる。

このような数ある戦略の中から有効な戦略を発見する方法の 1 つとして、レプリケータダイナミクスがある。ゲームを行うプレイヤーの集団を考え、プレイヤーはいくつかの戦略の中からランダムに戦略を選択し、他のプレイヤーとゲームを行い利得を得る。その後、戦略の集団に対する利得と集団全体の平均利得との差に応じて戦略の人口比を増減させる [5]。本論文では突然変異の概念を導入したレプリケータダイナミクスを用いる。ここで、戦略の集団 \vec{x} の中で戦略 j が占める割合を x_j とし、 \vec{x} に対して戦略 j が得る利得を $f_j(\vec{x})$ とする。また、 $\sum_{j=1}^n q_{ij} = 1$ を満たすような q_{ij} を戦略 i の子孫が戦略 j となる確率とおく。このとき、突然変異付きのレプリケータ方程式は以下のように表される。

$$\dot{x}_i = \sum_{j=1}^n x_j f_j q_{ji} - x_i \phi, \quad i = 1, \dots, n$$

$\phi(\cdot)$ を全ての戦略の利得の平均 $\sum_j x_j f_j(\vec{x})$ 、 $f_j(\cdot)$ を $\sum_m x_j a_{jm}$ とする。ただし、 a_{jm} は戦略 j をとるプレイヤーが戦略 m を取るプレイヤーと無限回プレイしたときの割引利得和である。

数値実験では、割引利得 ($\delta = 0.9$) を固定した上で、 g, l を $[0.05, 3.00]$ の範囲で 0.05 刻みで変化させた。戦略として状態数 2 以下の FSA 26 個を用いる。また、初期時点において、各戦略の人口は一律に分布、つまり、各戦略の存在比率は全て等しいものとする。さらに、突然変異を起こす確率 $\sum_{i \neq j} q_{ij} = \mu$ を 0.01 とした。このとき、戦略 j から戦略 $i \neq j$ に突然変異する確率 q_{ji} は $\mu / (26 - 1)$ の等確率で起こるとする。この微分方程式を解く際は期数の刻み幅 Δt を可変とした Dormand-Prince 法 [13, 14] を採用し、その実装には scipy を用いた [15]。さらに Δt の最大値は 0.5 とした。ダイナミクスは全ての戦略

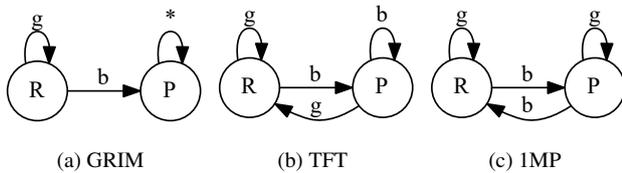


図1: 主要な FSA

の1期あたりの人口の変化量 $|x|$ が 10^{-5} 以下となった時点で収束と判定した。また、50000期までに収束と判定されなかった場合は計算を終了する。

3 主要な戦略とナッシュ均衡

本章では、繰り返し囚人のジレンマにおいて重要な戦略とそのナッシュ均衡を概説する。繰り返しゲームの戦略は過去の行動と観測の履歴から現在の行動への写像で定義される。一般には複雑になる戦略でも FSA を用いて簡略に表記できる。FSA の状態は、 R (reward, 報酬) と P (punishment, 処罰) の2つに区別され、プレイヤー i は状態 R で行動 $a_i = C$ を選び、状態 P で行動 $a_i = D$ を選ぶ。状態数1の戦略には ALLC と ALLD の2つが存在し、ALLC は状態 R のみを持ち毎期必ず協力する戦略、ALLD は状態 P のみを持ち毎期必ず裏切る戦略である。一方で状態 R と P をもつ状態数2の著名な戦略としては、まず最初に無限期罰則のトリガー戦略 (grim trigger, GRIM) が挙げられる (図1a)。GRIM は最初に協力し、相手の裏切りを観測するとそれ以降裏切り続ける戦略であり、多くの場合 GRIM は完全観測、不完全観測の両方の下で均衡を構成できる。別の戦略としては“しっぺ返し” (tit-for-tat, TFT) がある (図1b)。TFT は、状態 R からスタートし、相手の協力を観測した次の期では協力を、裏切りを観測した次の期には裏切りを行う戦略である。完全観測下では協力関係を維持できる一方で、不完全観測下では、いったん相手が裏切ったというシグナルを観測すると再び協力状態に戻るの難しくなる。他にも重要な戦略として、“1期相互処罰” (1 mutual punishment, IMP) が存在する (図1c)。IMP [16] は、従来“Pavlov”または“win-stay, lose-shift” [17] として知られている。IMP は状態 R からスタートし、相手の協力を観測したときは同じ状態に留まり、裏切りを観測するともう一つの状態へと遷移する。裏切りを観測してから協力に戻るのは一見不自然に見えるが、お互いを処罰してから協力に戻ることで、見間違えのある環境で TFT より協力状態を維持しやすくなっている。最後に、“一回処罰” (one-shot punishment, OSP) もしばしばダイナミクスに含まれる戦略である。OSP は状態 R からスタートし、相手の裏切りを観測した次の期のみ裏切る (状態 P に遷移する) が、その後は何を観測しても協力に戻る (状態 R に遷移する) 戦略である。

次に各プレイヤーの戦略空間を26個のFSAに限定した2人ゲームのナッシュ均衡を考える。相手があるFSAにしたがってプレイするとき、自分の割引利得和を最大化するFSAを最

適反応 FSA と呼ぶ。ある FSA の組がナッシュ均衡になるとは、その組がお互いに最適反応となる FSA になっていることを言う [18]。完全観測の場合、割引因子 δ が十分に大きければ、ALLD や GRIM, IMP, TFT を含む多くの戦略が均衡を構成する。不完全観測の場合、厳密な均衡条件を解析的に求めるのは難しいが、均衡を構成するのは ALLD, GRIM, IMP および状態 P から始める IMP の4種類のみである。とくに TFT が不完全観測で均衡を構成することはない [1]。例えば、 $p = 0.95, q = 0.01, \delta = 0.9$ のとき、ALLD は常に均衡を構成し、GRIM は l がおよそ 0.15 より大きければ均衡を構成する。状態 R もしくは P から始める IMP は g が 0.75 より小さければ均衡を構成する。先に述べたように均衡とダイナミクスの帰結に包含関係はないが、ダイナミクスの帰結でどんな均衡戦略 (もしくは均衡でない戦略) が生き残るかを知ることが協力の仕組みを理解する上で重要である。

4 2 状態戦略間のダイナミクス

4.1 完全観測下におけるダイナミクス

図2に完全観測のダイナミクスの帰結を示す。ここで、シグナル分布パラメータを $p = 1$ および $q = 0$ とする。それぞれの図の横軸は自分の裏切りによる利得の増分 g 、縦軸は相手の裏切りによる損失 l に対応し、0.05刻みで $[0.05, 3.00]$ をプロットした。図2aに収束時に最も多くの人口を獲得した戦略、すなわち最大多数戦略を、図2bは協力率を示している。これは収束時の戦略人口比に対して無限回繰り返しゲームを行うとして実現する (C, C) の頻度から計算した値である。残りの図2c-2hは主要6戦略の収束時の人口比を示している。

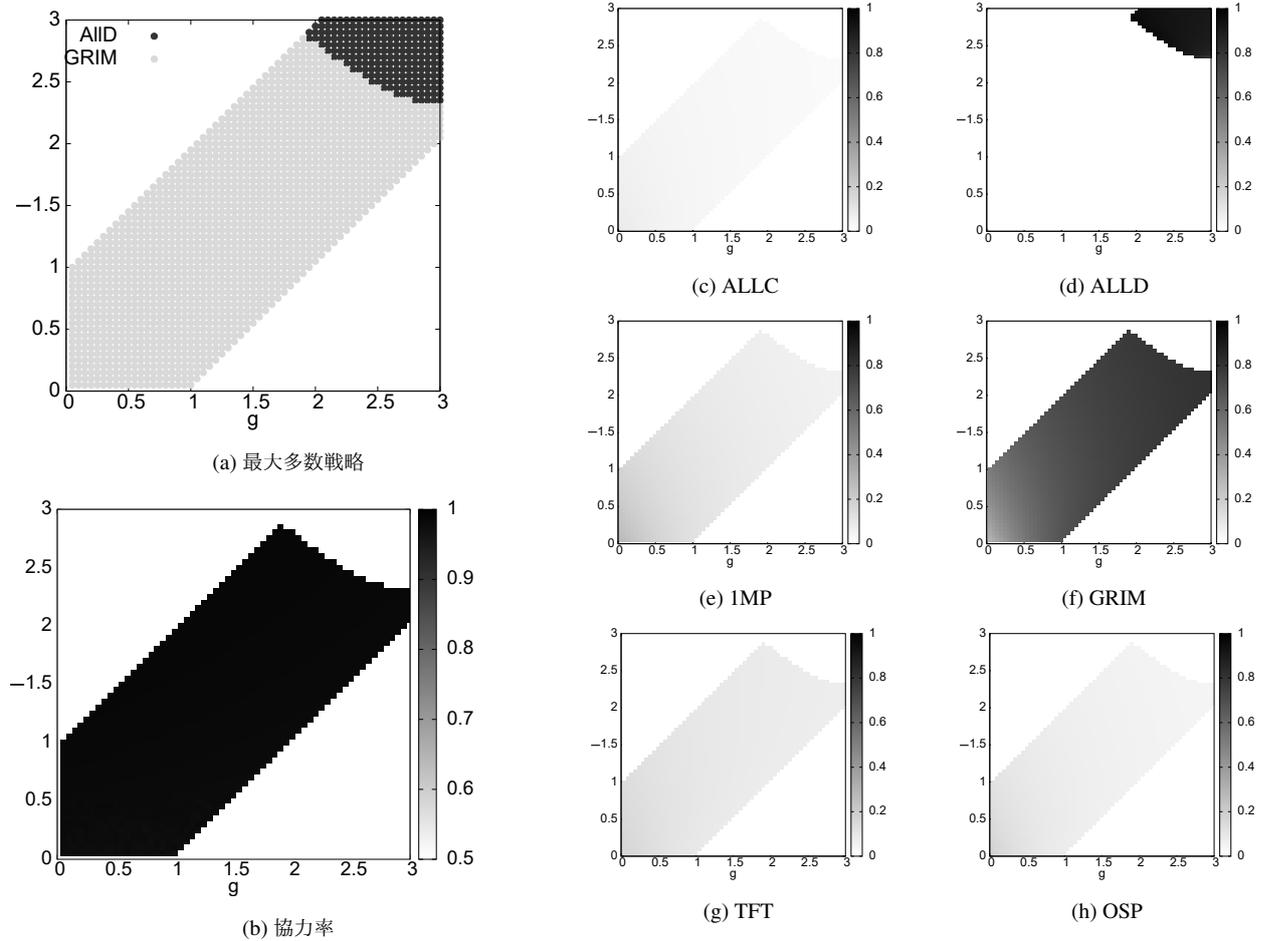
図2aでは、 g と l が十分大きい領域では ALLD が、それ以外の領域では GRIM が最大多数戦略となる。ALLD の人口比は約9割に到達する。一方で、GRIM が最大多数となるとき、他の4つの戦略 (ALLC, IMP, TFT, OSP) とそれなりの割合で共存する。どれくらいの割合で共存するかは g, l の値に依存し、 g, l が大きくなるにつれて GRIM の占める割合が増加する。図2bでは、ALLD が最大多数となるときの協力率はほぼ0である一方、GRIM が最大多数となるときは0.97を上回る。これは図2c-2hにあるように GRIM を含む5つの戦略はお互いに恒久的な協力関係を実現するためである。

4.2 ほぼ完全観測と4種の社会

図3にほぼ完全観測のダイナミクスの帰結を示す。ここで、シグナル分布パラメータを $p = 0.95$ および $q = 0.01$ とする。完全観測のときと同様に、 g および l を変化させながら図3aおよび3bのそれぞれに最大多数戦略と協力率を、図3c-3hに主要戦略の人口比率を示した。

図3aが示すように、どんな戦略が生き残るかは利得構造に依存し、おおまかに4つの領域に分けることができる。本論文ではこの4つの領域を以下の4つの社会に分類する。

- 非協力社会： ALLD が最大多数となる領域
- 不寛容社会： GRIM が最大多数となる領域

図 2: 完全観測のダイナミクス: $p = 1.00, q = 0.00$

相互協力社会: IMP が最大多数となる領域

周期協力社会: 複数の戦略が共存もしくは周期を構成する領域

非協力社会は g および l が十分大きいときに発生する。このとき、裏切る誘引や裏切られることによる損失が大きいため、他のどの戦略も協力を維持するに十分な将来利得を獲得できない。不寛容社会は g および l がそこそこの大きさでかつ、 g が l より小さいときに発生する。ここでは GRIM が最大多数を占めるため、最初はお互いに協力するが、一度でも裏切りが発生すると、永遠に裏切り続けることになり、相手を許すことはなくなる。実際 l が十分に大きいときは、裏切られることによる損失が大きくなるため、寛容な戦略で協力を回復する誘引を提供しにくくなる。GRIM は見間違えが起こるまでは協力状態を維持できるため、図 3b に示すようにその協力率は 0.680 程度となる。全体としては TFT も若干生き残るがその人口比率はわずか 0.01 程度にしかない。

次に相互協力社会は g および l が十分小さいときに発生する。IMP 同士の対戦では、どちらか一方のプレイヤーが *bad* を

観測して協力状態が途切れた後も、互いに裏切り合う相互処罰を経て、協力状態に簡単に戻ることができる。互いに罰を与えることで相互協力に戻るの、一見直感に反するが、相互処罰がうまく協力に戻るタイミングを明確にしている。そのため、IMP の協力率は見間違えのある状況でも約 0.928 と非常に高く、急速に人口を獲得する。また、図 3h にあるように OSP がわずかに存在するが、その比率は 0.05 以下にしかない。

最後に、周期協力社会は g および l がそこそこの大きさでかつ、 g が l より大きいときに発生する。ここで最大多数となる戦略は GRIM もしくは TFT のいずれかになるが、いずれも他の社会ほど大きな人口比率を獲得することない。代わりにこれらの戦略が共存もしくは循環する一方で、その協力率は図 3b にあるように 0.680 に安定する。

周期協力社会以外の社会では、ダイナミクスが収束するにつれて、ある戦略が単独で最大多数を占めるようになる。これに対して周期協力社会では、いくつかの戦略の人口比率が振動し、ある一定の比率に収束する、もしくは最大多数戦略が周期的に入れ替わるサイクルが続くようになる。この戦略

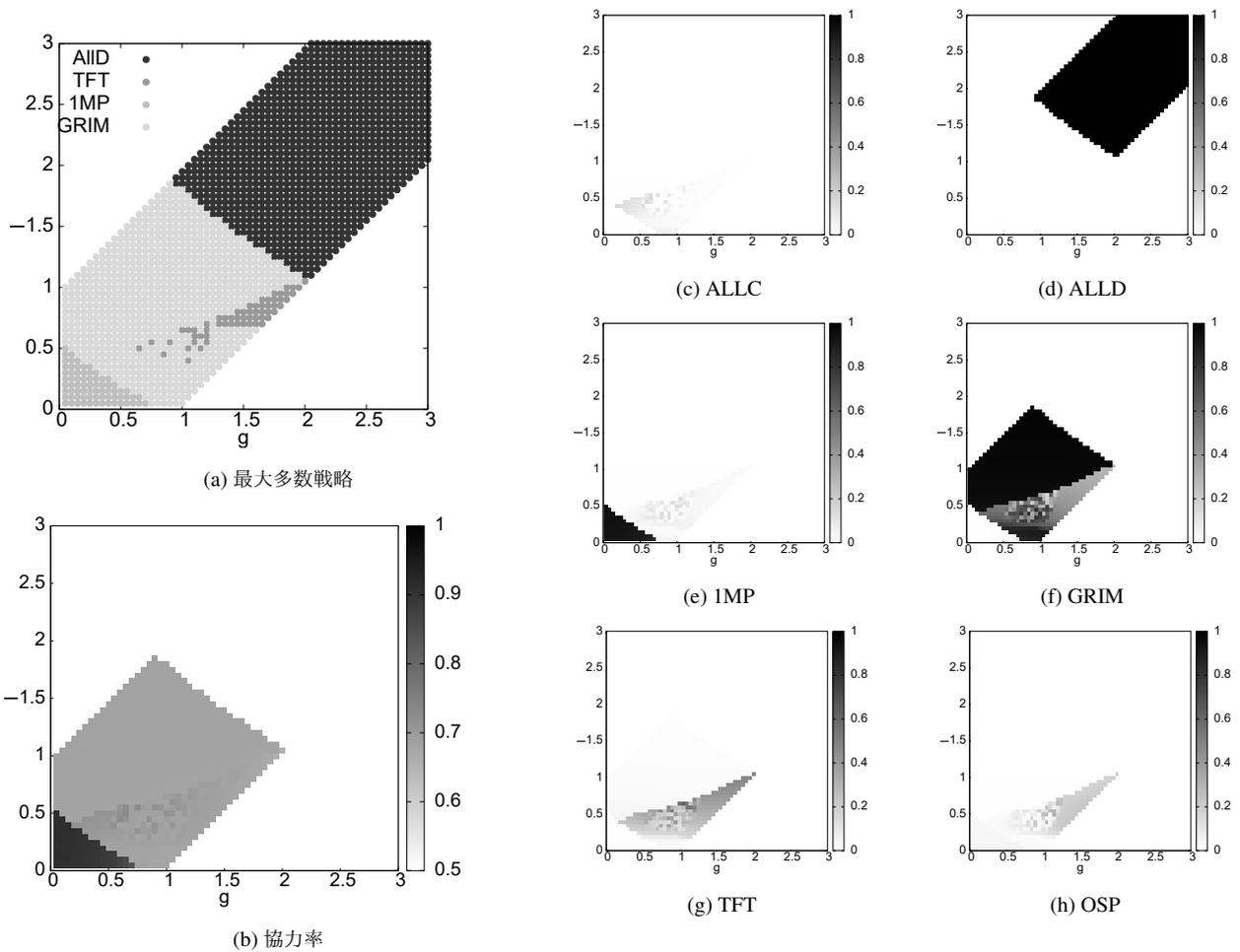


図3: ほぼ完全観測のダイナミクス: $p = 0.95, q = 0.01$

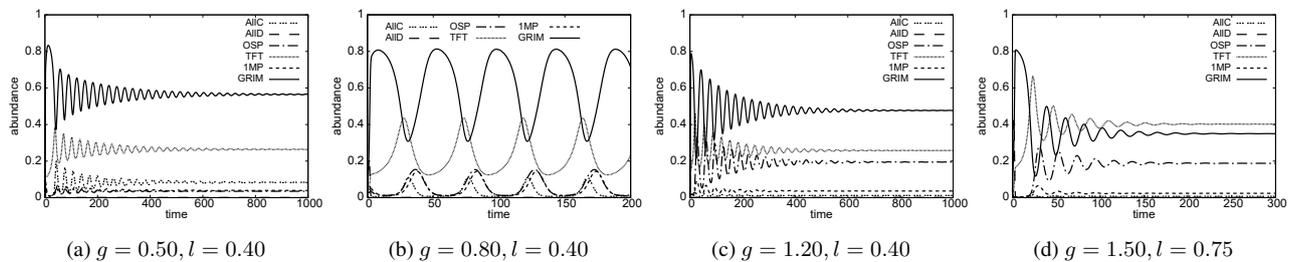


図4: 複数戦略が共存する点における人口割合の時間変化

の人口比率における協力率はほぼ一定で, ALLC, GRIM, TFT, IMP, OSP の5つの戦略を含む。ただし, g が小さいときには ALLC が, 大きいときには OSP の比率が増加する傾向がある。実際, 裏切りによる利得の増分が増えると ALLC のような相手を処罰しない戦略は簡単に搾取されてしまう。こうした協力の移り変わりを理解するために, そのダイナミクスをいくつか吟味する。

図4aに $g = 0.50$ および $l = 0.40$ における戦略の人口比率の時間変化を示す。それぞれの戦略はまず振動を繰り返すが,

徐々に振幅は小さくなり, 1000期以降はほとんど変化しない。このとき, 先に述べた5つの戦略が生き残っており, 比率の大きい順に GRIM (0.565), TFT (0.264), ALLC (0.083), IMP (0.039), OSP (0.035) となっている。括弧内にその戦略の人口比率を示す。また, l をわずかに大きくしても同様のダイナミクスを観察するが, それ以上大きくすると不寛容社会への移行し, GRIM がすぐに全ての人口を獲得するようになる。

図4bでは, l を0.4に固定したまま, g を0.8に増加させた。このとき, 5つの戦略の人口比率は循環し, 一定の比率に安定

しない。その最大多数戦略はほとんどが GRIM だが、定期的に TFT の比率が GRIM より高くなる。いつ計算を打ち切るかによってどの戦略が最大多数となるかが変わるため、図 3a で TFT が最大多数となる領域が飛び地を形成する。ただし、図 4d の $g = 1.5, l = 0.75$ のダイナミクスが示すように、TFT が最大多数となる人口比率に収束する場合も存在する。

さらに l を 0.4 に固定したまま、 g を 1.2 に増加させたのが図 4c である。ここでは、図 4b で観察した循環はなくなり、図 4a と同じように戦略の人口比率がほぼ一定に収束する。ただし、その比率は大きい順に GRIM (0.477), TFT (0.258), OSP (0.197), IMP (0.036), ALLC (0.016) となる。図 4a での ALLC の代わりに OSP が GRIM, TFT に続く人口比率を獲得する。

ここまで見間違いのあるほぼ完全観測でのダイナミクスの帰結が、囚人のジレンマの利得構造に応じて 4 つの社会に分類できることを示した。一方で見間違いのない完全観測では、非協力と不寛容の 2 種類の社会しか発生しない。つまり IMP や TFT といった非自明な協力的戦略が生き残らない。

完全観測では、状態 R から始まる戦略の組はナッシュ均衡になる。TFT や IMP のような協力を回復させようとする戦略は、状態 P から始まる戦略 (26 個中 13 個) や状態 R から始まるが、相手の協力を観測すると状態 P に遷移するような“ずい”戦略 (26 個中 8 個) に搾取される。例えば、状態 R から始まる TFT と状態 P から始まる IMP の組を考える。その積状態は RP から PP, PR, RP と遷移していく。こうしたお互いに報復する過程が TFT の期待利得を減少させ、 g と l が十分大きいとき、その将来利得は、相互協力を継続して得られる利得 1 より小さくなる。このため、不寛容な戦略 (GRIM) が相対的に高い利得を実現する。実際、状態 R から始まる GRIM と状態 P から始まる IMP の積状態は RP から PP, PR, PP, ... と遷移していく。このとき、IMP は 2 回に 1 回は裏切られて損をするため、GRIM や GRIM と似た戦略に容易に搾取される。

一方、ほぼ完全観測では、見間違いが起きるため、不寛容な戦略では相手を処罰しすぎてしまい、利得や協力率を下げてしまう。このため、 g および l が十分小さいときは IMP が、 g が l に比べて大きいときは TFT が最大多数となりえる。つまり IMP や TFT がもつ見間違いの後に協力を回復させる仕組みが機能するようになる。したがって、わずかでもお互いに相手の行動を見間違えるというノイズが、人が協力をどのように維持するかに多様性を与えているといえる。

5 議論

5.1 周期協力社会における三すくみ

周期協力社会では、GRIM もしくは TFT が最大多数戦略となり、OSP, IMP, ALLC を加えた 5 つの戦略が共存する。ただ、GRIM と TFT 以外の戦略の人口比率はかなり小さい。実際、IMP の比率が 1% を超えることはめったにない。また、 g が小さいときは、GRIM と TFT に次いで ALLC の人口比率が大きくなり、 g が大きいときは、ALLC の代わりに OSP の

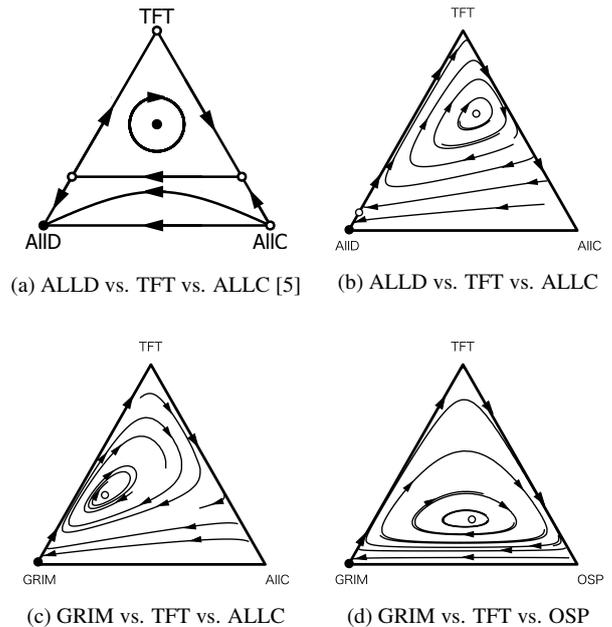


図 5: 3 戦略間のレプリケータダイナミクス: 図 5a は摂動におけるダイナミクスを文献 [5] より引用

人口比率が大きくなる。そこで本章では GRIM vs. ALLC vs. TFT と GRIM vs. OSP vs. TFT の 2 つの 3 戦略の組の関係を分析する。

図 5a に Sigmund [5] でよく知られている ALLD vs. ALLC vs. TFT のダイナミクスを示す。図中の矢印はダイナミクスが進む方向を、黒点は安定な不動点を、白点是不安定な不動点を表す。ここではプレイヤーが行動を取り間違える (trembling-hand) とき、これらの 3 戦略が三すくみを形成する。TFT の人口が少ないときは ALLD へ収束する一方で、TFT の人口が一定数を超えると 3 戦略による周期が発生する。

ほぼ完全観測 ($p = 0.95$ および $q = 0.01$) において $g = 1.20$, $l = 0.40$ とする。図 5b に、ALLD, ALLC, TFT の 3 戦略間のダイナミクスを示す。ここで不安定な不動点 (ALLD, ALLC, TFT) = (0.150, 0.266, 0.585) を中心にサイクルが形成される。図 5a と同じように TFT の人口比率が一定数を下回ると ALLD に収束するようになる。次に、ALLD を GRIM に入れ替えた結果を図 5c に示す。サイクルの中心が不安定な不動点 (GRIM, ALLC, TFT) = (0.527, 0.125, 0.347) に移った以外は、図 5b と同じような結果となった。さらに、ALLC を OSP に入れ替えた結果を図 5d に示す。ここでサイクルの中心は (GRIM, OSP, TFT) = (0.348, 0.427, 0.225) に移り、図 5a とほぼ同じダイナミクスを異なる戦略の組で形成する。

すでに見てきたように周期協力社会では GRIM もしくは TFT が最大多数となる。見間違いのあるとき、TFT は ALLD や GRIM を支配する。ALLC もしくは OSP は TFT を支配するが、どちらの戦略も GRIM に支配される。26 戦略間のダイ

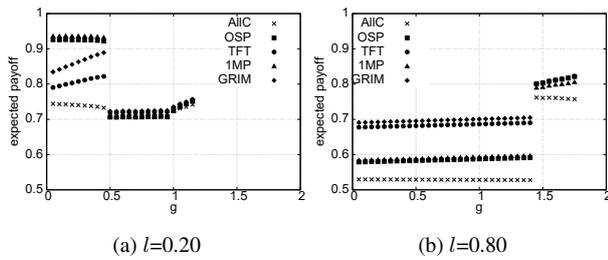


図6: 収束時の人口比における期待利得の変化

ナミクスの帰結で ALLD は生き残らないが、代わりに GRIM が (行動を取り間違えるときの) ALLD の役割を果たす。また g が大きいと裏切りへの誘因が強くなるので、ALLC が生き残りにくくなるとともに、相手を1回だけ処罰する OSP が ALLC にとって代わる。 g が大きすぎもせず、小さすぎもしないときに図 4b のような複雑な周期を観測し、戦略の安定的な共存が成立しなくなる。

摂動とほぼ完全観測では戦略の定義が異なるため、単純に結果を比較することはできないが、摂動における ALLD vs. ALLC vs. TFT に相当する GRIM vs. OSP vs. TFT という三すくみ状態を、私的観測において世界で初めて発見した。

5.2 感度分析

本節では、利得、シグナル分布、そして突然変異率といったパラメータに関する感度分析を行う。まず、利得パラメータ g, l が期待利得に与える影響を図 6 に示す。ただし、ここでの期待利得は収束時の人口比における各戦略の期待利得の平均と定義する。図 6a に $l = 0.2$ に固定して g を変化させたときの平均利得を示す。 g が約 0.5 よりも小さい、つまり 1MP が最大多数となると、1MP の利得は最も高く、その値は 0.9 を超える。 g が 0.5 よりも大きくなると、GRIM が最大多数を占める不寛容社会が実現する。このとき 5 戦略すべての利得が 0.7 程度になり、その差がほとんどなくなる。次に、図 6b に $l = 0.8$ に固定して g を変化させたときの平均利得を示す。 g が 1.4 よりも小さいとき、GRIM と TFT の利得はともに 0.7 程度であるが、GRIM が最大多数となる。 g が 1.4 を超えると、GRIM, TFT, OSP の期待利得は 0.8 程度となり、1MP が 0.75 程度と若干低くなるこのときは、これら 3 つの戦略が多数を占める共存が発生する。

次に、シグナル分布パラメータを変化させ、ダイナミクスの帰結を観察したところ、最大多数となる戦略それぞれの領域の大きさや、GRIM がどんな戦略と共存するかが変化した。一方で、有効な戦略の傾向に大きな変化は見られなかった。これを確認するために、図 7 に q を 0.01 に固定したまま p を変化させたときの、図 8 に p を 0.90 に固定したまま q を変化させたときの収束時 GRIM の人口割合を示した。ただし、GRIM が人口を獲得していない g および l が小さい領域では 1MP が、 g および l 大きい領域では ALLD がほとんどすべての人口を獲得し最大多数となっている。

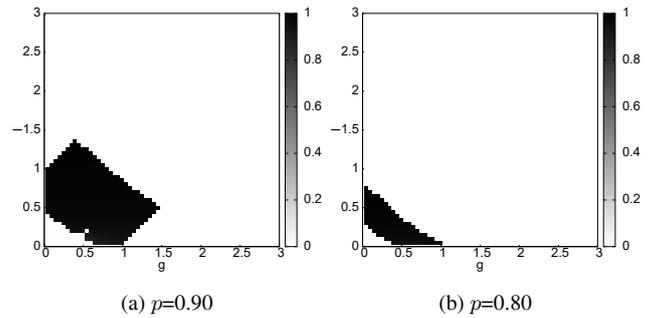


図7: $q=0.01$ で固定したときの GRIM の人口割合の変化

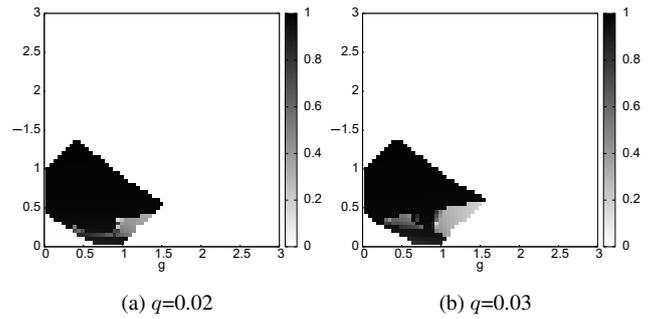


図8: $p=0.90$ で固定した GRIM の人口割合の変化

図 7 では p が小さくなるにつれて、ALLD が最大多数となる g および l の領域が増加していることがわかる。同じように GRIM も g および l の値が小さいとき最大多数となるが、ALLD と比べるとその領域は小さくなっている。一方で、1MP は g および l がさらに小さくないと最大多数にならなくなる。この傾向は $p = 0.95, q = 0.01$ とした図 3f でも確認できる。そして GRIM の人口割合に注目すると、 $p = 0.95$ では $g > l$ の範囲において GRIM が他戦略と共存する傾向が見られたが、 $p = 0.90, 0.80$ においてはその傾向はみられなかった。

一方で、図 8 では各戦略が最大多数となる領域はほとんど変化していないが、 q の値が大きくなるにつれて $g > l$ における GRIM と他戦略の共存領域が増加している。この傾向は $p = 0.90, q = 0.01$ とした図 7a でも確認できる。

ここで再び、GRIM が他の戦略と共存しうる 3 つのパラメータ (図 3f, 8a, 8b) に注目すると、両方のプレイヤーが間違っただけのシグナルを受け取る確率 $r (r = 1 - p - 2q)$ が比較的小さいという共通点が見つかる。具体的には $r \leq 0.06$ では GRIM が他戦略と共存するが、 $r \geq 0.08$ では他戦略とほとんど共存しない。したがって、 r の大きさが GRIM と他戦略の共存度合に影響を及ぼすと考えられる。

次に p が領域の変化に与える影響を見るために、 $g = 0.25, l = 0.25, q = 0.01, \mu = 0.01$ に対して p を変化させたときの協力率と期待利得の推移を図 9 に示す。 p が大きくなる、つまりシグナルの正確さが増加するにつれて最大多数戦略が ALLD, GRIM, 1MP と変化し、協力率と期待利得が段階的に上昇する。

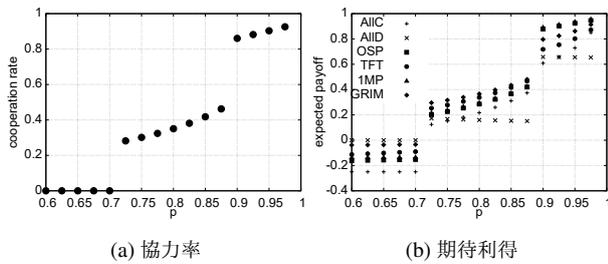


図9: $g=0.25, l=0.25$ における協力率および利得の変化

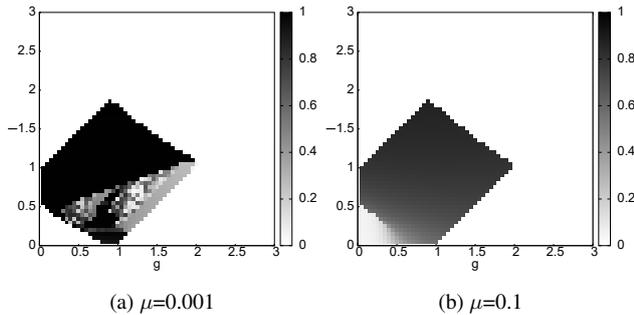


図10: μ を変化させたときの GRIM の人口割合の変化

最後に突然変異確率の影響を吟味するため $p = 0.95, q = 0.01$ のほぼ完全観測で GRIM が最大多数となる領域を図10に示す. 突然変異確率 $\mu = 0.01$ の図3f に対して, 突然変異確率を0.001に小さくすると, 図10aにあるように $g > l$ において GRIM が他戦略と共存する領域が増加する. 逆に, 突然変異確率を0.1に大きくすると, その人口比率は安定的に大きくなる. その結果, 周期協力社会で見たような最大多数戦略の周期的な変化が起こりにくくなる. 言い換えると突然変異率が高いとき, 多数を占めていない戦略が発生する確率が高くなる. その結果, GRIM のような不寛容な戦略の方が安定した利得を実現しやすくなるためである.

6 おわりに

本論文は私的観測付き繰り返し囚人のジレンマを突然変異付きレプリケータダイナミクスで分析した. 完全観測では, ALLD もしくは GRIM のいずれかの戦略しか生き残らなかったのに対して, 私的観測では, 見間違えの後でも協力に戻りやすい IMP という比較的新しい戦略が生き残るようになった. さらに非自明な均衡戦略を持たないような利得構造では, TFT が GRIM, OSP とともに三すくみになり, 周期協力社会を構成することを世界で始めて明らかにした. 今後の課題として状態数3以下のFSAや確率的FSAを含む戦略空間におけるダイナミクスを分析することなどが挙げられる.

参考文献

[1] ジョヨンジュン, 岩崎敦, 神取道宏, 小原一郎, 横尾真. 部分観測可能マルコフ決定過程を用いた私的観測付き繰り返し

返しゲームにおける均衡分析プログラム. 情報処理学会論文誌, pp. 1234–1246, 2012.

[2] 岡田章. ゲーム理論 新版. 有斐閣, 2011.

[3] Robert Axelrod. *Genetic Algorithms and Simulated Annealing*, chapter The Evolution of Strategies in the Iterated Prisoner's Dilemma, pp. 32–41. Morgan Kaufman, Los Altos, CA, 1987.

[4] Martin Nowak. *Evolutionary Dynamics: Exploring the Equations of Life*. Harvard University Press, 2006.

[5] Karl Sigmund. *The Calculus of Selfishness*. Princeton University Press, 2010.

[6] Michihiro Kandori. Repeated games. In Steven N. Durlauf and Lawrence E. Blume, editors, *Game theory*, pp. 286–299. Palgrave Macmillan, 2010.

[7] 関口格. 経済セミナー増刊:ゲーム理論プラス, 「協調達成のための正しいお仕置きの方」. 日本評論社, 2007.

[8] 松島齊. ゲーム理論の新展開. 勁草書房, 2002. 第4章: 「繰り返しゲームの新展開:私的モニタリングによる暗黙の協調」, pp.89-114.

[9] Peter D. Taylor and Leo B. Jonker. Evolutionarily stable strategies and game dynamics. *Mathematical Biosciences*, pp. 145–156, 1978.

[10] Josef Hofbauer and Karl Sigmund. *Evolutionary Games and Population Dynamics*. Cambridge University Press, Cambridge, 1998.

[11] 大槻久. 有限集団における進化ゲーム理論の発展. 特集「多様性と進化の統計解析」, 第60-2章, pp. 251–262. 統計数理研究所, 2012.

[12] Drew Fudenberg Loren A. Imhof and Martin A. Nowak. Evolutionary cycles of cooperation and defection. in *Proceedings of the National Academy of Sciences*, Vol. 102, No. 31, pp. 10797–10800, 2005.

[13] P. J. Prince R. Dormand. A family of embedded runge-kutta formulae. *Journal of Computational and Applied Mathematics*, Vol. 6, No. 1, pp. 19–26, 1980.

[14] L. W. Shampine. Some practical runge-kutta formulas. *Mathematics of Computation*, Vol. 46, No. 2, pp. 135–150, 1986.

[15] Eric Jones, Travis Oliphant, Pearu Peterson, et al. SciPy: Open source scientific tools for Python, 2001–. [accessed 1 February 2019].

[16] David Kraines and Vivian Kraines. Pavlov and the prisoner's dilemma. *Theory and Decision*, Vol. 26, pp. 47–79, 1989.

[17] Martin Nowak and Karl Sigmund. A strategy of win-stay, lose-shift that outperforms tit for tat in prisoner's dilemma. *Nature*, Vol. 364, pp. 56–58, 1993.

[18] 神取道宏. ミクロ経済学の力. 日本評論社, 2014.