

# 道徳の起源に関する理論的研究\*

——動学的相互依存性理論と制度分析アプローチによる考察——

清 水 裕 士\*\*

藤 原 武 弘\*\*\*

## 1 節 道徳の起源についての研究

本論文は、道徳が発生する起源について、文化や制度の進化の枠組みから整理することを目的とした理論論文である。その目的のために、Kelley & Thibaut (1978) の相互依存性理論を用いて、文化と制度がいかんして進化しうるかについての試論を行い、道徳の成立について論じる。

道徳判断についての研究は、道徳心理学や社会心理学、教育心理学、発達心理学など、多くの分野で行われている。多くの道徳判断の研究は、人が道徳的な判断を行う心理プロセスや、人がいかんして道徳的判断を獲得するかといった発達プロセスが扱われていることが多い。前者については、トロッコ問題などを用いた神経科学的なアプローチ(芋坂, 2012)が、後者については Kohlberg (1976) などによる道徳性の発達のな獲得についての研究がある。

本論文では、これらの研究とは異なり、現代社会において人々がなぜわれわれが行っているような道徳的判断を行うのか、その起源について答えることに関心がある。そのため、まずは道徳の起源を扱った道徳心理学研究とその批判点について簡単にレビューし、その後本論文の理路の方向性を示そう。

### 道徳基板理論

道徳心理学者の Haidt (2012) は、道徳基板理

論において、道徳的判断が直観的で、進化的に獲得された心の特性であることを主張している。Haidt (2001) では道徳判断の基準となるのは理性ではなく直観であるとし、またその判断の基盤が通文化的なものであることを主張した。具体的には、危害・擁護、公正、内集団への忠誠、権威への敬意、純潔・神聖の5つを挙げている。Haidt (2012) では6つ目として自由も加えている。Haidt は特に前者の二つ、危害・擁護と公正(場合によっては自由も)はどの文化でも存在するかなり普遍的な道徳であると考えており、後者の集団主義的な3つの道徳はどちらかという文化的な影響があると考えられている。

Haidt (2012) は6つの基準の起源に多少の違いを認めながらも、道徳を進化心理学的な基盤によって構成されていると考えている。しかし、その論点についてはいくつか批判点もある。たとえば、Barash (2007) は集団への忠誠や権威への服従といった集団主義的な道徳の進化的基盤について、集団淘汰理論、あるいはマルチレベル淘汰理論を援用しているが、その議論が成立する根拠が不十分であると指摘している。具体的には、集団淘汰あるいはマルチレベル淘汰が遺伝子に与える影響は非常に限定的、あるいは限られた状況(病原菌などによる大量の死者が出る状況や、戦争などの深刻な集団間葛藤状況)であるとし、Haidt の議論はマルチレベル淘汰の可能性を過大に評価しすぎているとしている。それ以外にも、集団主義的な道徳を必要とするほどの大きな集団が形成

\*キーワード：道徳の起源、制度、動学的相互依存性理論

\*\*広島大学大学院総合科学研究科助教

\*\*\*関西学院大学社会学部教授

されて、およそ1万1千年 (Diamond, 1997) ほどしか経っていないにもかかわらず、遺伝的な要素のみで道徳を説明するのは無理があるだろう。このように Hadit の議論は進化心理学的な説明を行っているにもかかわらず、その成立根拠についての議論が不十分である。

本研究では、これらの批判点に基づいて、道徳の基盤として遺伝的要因以外に、制度的な要因を考察する。

### 道徳の起源における遺伝的要因と社会文化的要因

Haidt (2012) の道徳基板理論では、道徳の起源において個体レベルの淘汰と集団レベルの淘汰の両方を仮定しているとはいえ、どちらも遺伝子レベルの進化プロセスを前提としている。しかし、すでに指摘したように、集団への忠誠心や権威への服従といった集団主義的な道徳は、人類の歴史において非常に最近 (1 万年程度) であり、病原菌などによる強い淘汰がかからない限りは遺伝子の分布にはそれほど大きな影響があるとは考えにくい。

一方、遺伝子による進化ではなく、文化そのものが進化することを想定した文化進化理論がある (たとえば Aunger, 2000)。その中でも二重継承理論は、文化が遺伝的なプロセスと社会的な学習プロセスの両方を重視するという立場である。たとえば Richerson & Boyd (2005) や Tomasello (1999) では、文化は人間が遺伝的に獲得した学習能力によって、他の種ではできない高精度の観察学習によって伝達されるという。

二重継承理論では、文化や道徳といった社会集団における規範的性質を、遺伝子からの継承で考えるわけではなく、人々の相互作用 (観察学習) をも継承の単位として考えている点において Haidt (2012) の道徳基板理論よりは妥当な説明が可能である。しかし、二重継承理論においても、社会レベルの文化の継承を、観察学習を中心としたメカニズムを想定している点において、道徳の起源を考察する上で限界がある。なぜなら、道徳はある種の行動が斉一的に実践されていることのみを指すわけではない。もし、行動の斉一性や集団内の戦略の均衡としてのみ道徳を位置づけるならば、突然変異以外の違反が生じえないこと

になり、そもそも規範や道徳についての信念が不要になるだろう。よって、道徳は行動的な規則性のみではなく、行動を実践するときに事前に主体が持っているだろう当為的判断、具体的には「～すべきである」といった内的な判断が前提とされる必要がある。これは、チャルディーニがいうところの「記述的規範」と「命令的規範」の区別に対応する。社会がある行動を斉一的に実践している状態が記述的規範であり、「～すべきである」という命令的規範を含む規範的な性質として道徳が位置づけられると考えられるのである。

そこで本研究では、当為的な判断を伴う道徳が社会に共有されるメカニズムを理論的に考察する。その手がかりとして、制度分析 (Aoki, 2000) に注目する。制度分析アプローチは、社会集団が実践している戦略の集合が、ゲーム理論でいうところの「均衡」であると認める点については進化ゲーム理論と同じである。だが一方で、社会集団における相互作用の利得構造そのものの変化の可能性を考慮に入れている点が異なっており、またその特徴となっている。つまり制度分析アプローチは、進化ゲームと同様により利得を得た個体の戦略が増え、均衡に至るという考えを持ちながら、均衡そのものが変化し、進化していくという立場にたっているのである。この戦略の均衡を変化させる装置こそが本論文の想定する制度である。

ここでいう制度とは、近代的な社会制度だけを意味するのではなく、「人々が社会集団において相互作用を円滑にするために定式化された社会装置」、という広い意味で用いている。また制度は、ゲーム理論的な意味での戦略のように個体を持つ特性に還元できるものの集合ではない。むしろ、個体のレベルを超えた、新たな戦略を個体に可能とするような社会的な仕組みのことを指している。本論文では、道徳は個人の進化生物学的な観点を超えた、社会的な要因に注目しているのである。このような制度と道徳・規範の捉え方は、Aoki (2000) 以外にも、たとえば Heath (2008) や Gintis (2009)、そして河合 (2013) などにも見られる。

ただし本論文では、道徳が制度そのものであると考えているわけではない。なぜなら制度が成立

することと、当為的な道徳が成立することは直接的な関係性はないからである。道徳について議論するためには、当為的な信念が成立するプロセスについて明らかにしなければならない。そこで本論文では、道徳は人々が作り上げてきた制度に適應するために共有されてきた判断基準である、という立場をとっている。つまり、制度をスムーズに運営するために作り上げたのが道徳である、ということである。この結論に至るために、本論文では以下の様に議論を展開する。

### 本論文の目的と理路

本論文の目的は、制度分析アプローチの観点から、道徳起源の制度的な要因について考察することである。

この制度分析アプローチをゲーム理論的に考察する理論的道具立てとして、本論文では Kelley & Thibaut (1978) の相互依存性理論に焦点を当てる。相互依存性理論では、個人間の利得についての相互依存性のパターンを分析し、人々が主観的に利得構造をどのように変換しているかを考察している理論である。2節では相互依存性理論を動学モデルに発展させて、3節の制度分析に接続させる準備を行う。

そして3節では制度分析アプローチが想定する制度の変化や進化を、相互依存性理論で議論されているところの「相互依存性パターンの変換」という観点から考察することで、制度の進化について理論的に展開することが可能になることを示す。

続いて4節では、道徳が進化した制度のものでいかにして発生するかについて、言語による伝達を用いた教育の可能性について議論する。ここでは、道徳が社会秩序の原因なのではなく、社会秩序に人々が適應するために獲得した心の特性であることを示す。

## 2節 動学的相互依存性理論

### 相互依存性理論

ゲーム理論による対人的相互作用の研究は数多いが、なかでも Kelley & Thibaut (1978) の相互依存性理論は、社会心理学的に対人行動を分析す

る上で重要な含蓄を有している。それは、以下に述べるように対人的相互作用における相互依存性のパターンを分析することができること、そしてその相互依存性のパターンを変換するというアイデアが盛り込まれている点にある。

相互依存性理論は、二者の相互作用を  $2 \times 2$  の利得行列によって表現し、その相互依存性を分析したものである。なかでも、相互依存性理論が他のゲーム理論と異なる点は、ゲーム構造を三つの要素に分解し、その要素の組み合わせによって対人的相互作用を分析しようとするところにある。

例として、囚人のジレンマゲーム (Prisoner's Dilemma Game: PDG) を用いて解説しよう。要素の分解は、図1のようになる。すなわち、プレイヤー A と B の利得行列をそれぞれの利得で分割して表示し、それを平均 (Grand Means: GM) と三つの分散要素に分解する。分散要素は、自己の効果 (Self Control: SC)、相手の効果 (Partner's Control: PC)、そして共同の効果 (Joint Control: JC) と呼ばれ (Kelley, et al. 2003)、二要因分散分析でいうところの各要因 (自分と相手) の主効果と交互作用に対応する。SC は自分の選択によって自らの利得が変化する効果、PC は相手の選択によって利得が変化する効果、JC は自分と相手の選択の組み合わせによって利得が変化する効果である。

なお、各効果は以下のように計算される。図1の A の利得 (上の方) を例にするならば、まず SC は A 自身の選択による分散を表していることから、行ごとの平均の差によって算出される。具体的には、協力は  $(8+0)/2=4$ 、裏切りは  $(12+4)/2=8$ 、となり、協力と裏切りには期待値として4ポイントの差があることになる。ここで、各効果の平均を0とすると、SC は、協力の期待値の  $-2$  と裏切りの期待値の  $+2$  として表現することができる。さらに、PC は列ごとの平均、JC は対角項と非対角項の平均によって同様に計算することが可能である。また、本論文では PD ゲームの SC を  $2$ 、PC を  $-4$  としているが、わかりやすさのために SC を正の方向に基準をあわせているだけであり、これを SC が  $-2$ 、PC が  $+4$  としても、本質的には利得構造に違いはない。

このように、 $2 \times 2$  のゲーム構造は GM、SC、

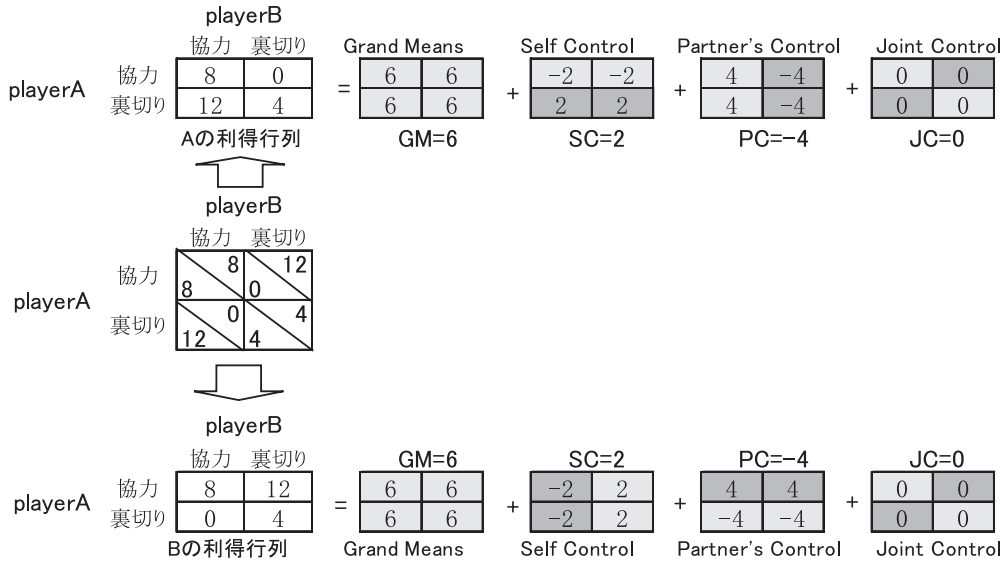


図1 相互依存性理論による2×2の利得構造の要素への分解

PC、JCの4つの要素によって規定されるが、実質的に相互依存性のパターンを決定している要素はGMを除く3つの要素である。なぜなら、GMは全体的な利得の大きさを決めていて、戦略を変えることによる利得の変化は生じないからである。事実、相互依存性理論においてもSCとPCとJC野組み合わせによる相互依存性のパターンの考察が行われている。

動学的な相互依存性のパターン

ただし、Kelleyらの相互依存性理論は一時点における相互作用の相互依存性のパターンが考察されており、進化、文化進化ゲーム理論が想定するような、何世代も繰り返して相互作用することを前提とする枠組みとは異なっている。本論文では、相互依存性理論が対象としてこなかった、動学的な相互依存性のパターンについて考察しよう。

進化ゲーム理論では、相互依存性のパターンは進化的安定戦略 (Evolutionary Stable Strategy ; ESS) と呼ばれる、均衡概念と対で考察されてきた。ESSとは、ある戦略のみを行う個体同士が繰り返して、何世代にわたって相互作用した場合に、進化的に生き残る戦略の組のことである。たとえば囚人のジレンマゲームでは非協力的な戦略がESSとなり、タカハトゲームではタカ戦略と

ハト戦略が一定の頻度割合で共存する状態がESSとなる。ESSは経済ゲーム理論におけるナッシュ均衡に非常に近い概念であるが、不安定なナッシュ均衡はESSとはならないことが知られており、ナッシュ均衡の部分的な戦略集合である。

しかし、本論文では社会集団が成立したあとの文化や道徳の進化に関心があるため、進化ゲーム理論では考察のスケールが一致しない。よって、以降では学習ゲーム理論に基づく均衡概念を対象に相互依存性のパターンを考察しよう。学習ゲーム理論は、自然淘汰による戦略分布の変化を考えるのではなく、個体が得る利益や罰に基づいて行動を変化させる、学習心理学的な観点で戦略分布の変化を考えるゲーム理論である。試行錯誤ダイナミクスとも呼ばれる。試行錯誤ダイナミクスは、ゲーム理論に対して動学的なアプローチを採用し、微分方程式や差分方程式を用いて戦略の均衡がどのような様相を示すかを明らかにするための理論的ツールである。

相互依存性理論では、相互依存性のパターンを、1時点のゲームにおける戦略のあり方を踏まえて考察されていた。それに対して本論文では、試行錯誤ダイナミクスを用いて、集団内のメンバーが同じゲームを繰り返したときの均衡が、相互依存性パターンによってどのように異なるかを検

討しよう。それによって、集団内における相互作用の根源的なパターンを見出すことができるだろう。

### 差分集団型試行錯誤ダイナミクス

対称  $2 \times 2$  の集団型試行錯誤ダイナミクスを考えよう。対称ゲームとは、すべての個人間の相互依存性が等しい、という意味である。あるサイズの集団内の二人がランダムに組になり、一度ゲームを行う。集団内の個体は変わらず、ゲームで得られた利益に基づいて行動傾向のみを変化させる。行動傾向とは、2 戦略のどちらの行動を選択するか、その確率で表現される。A と not A 戦略があれば、A をとる確率が高くなったり低くなったりする、ということである。その試行を何度も繰り返して各個体の行動傾向および集団全体の戦略分布がどのように変化するかを検討することを考えよう。

ここで Roth & Erev、(1995) の差分方程式モデルを、大浦 (2008) を参考に、対称  $2 \times 2$  の集団型試行錯誤ダイナミクスを相互依存性理論に当てはめてみよう。対称  $2 \times 2$  のゲームは、成果が 4 パターンあり、それぞれ相互依存性理論に基づいて SC と PC と JC によって表現することができる。ここで、符号の不定性を避けるため SC は 0 より大きい、つまり常に正であることを仮定しよう。この条件に基づいて試行錯誤ダイナミクスの均衡点を考える。数理的な展開は Appendix に載せているので、以下に結論のみを提示する。各個体の行動選択確率変化の期待値は、

$$E[\Delta x_i] = x_i(1 - x_i)(-JC + SC + 2JCx)/k$$

で表される。なお、 $x_i$  は個体  $i$  が、SC が正の方向の戦略を選択する確率を表している。また  $x$  は全成員の  $x_i$  の期待値を、 $k$  は定数 (Appendix 参照) を意味する。どの個体も同じ利得行列をもつ対称ゲームであるので、これを解くと、集団の平均的な行動傾向が変化しなくなる均衡点は  $x = 0$  か 1、あるいは  $(JC - SC)/2JC$  となる。ここで、 $(JC - SC)/2JC$  が 0 から 1 の範囲に収まり、SC が 0 以上であることを前提に展開すると、 $|JC| > |SC|$  という結果が得られる。この結果が意味するところは、JC よりも SC の絶対値が大きい場合に、内的均衡点に、それ以外の時に

は  $x$  が 0 か 1 の状態が均衡点となることを意味している。ただし、内的均衡点のパターンには 2 つある。JC > SC の場合には不安定な均衡点となり、すこしのゆらぎで戦略分布が  $x = 0$  か 1 のどちらかに移行する。もう一つは  $-JC > SC$  の場合に安定な均衡点となり、 $(JC - SC)/2JC$  に収束する。また  $|JC| < |SC|$  の場合は、 $SC \geq 0$  を前提にすると、均衡点は常に  $x = 1$  となる。

これらの結果から得られる知見として、均衡点のパターンは不安定な均衡を除くと、 $x = 1$  になる場合 ( $|JC| < |SC|$ )、 $x = 0$  か 1 になる場合 (JC > SC)、そして内的均衡点である  $x = (JC - SC)/2JC$  となる場合 ( $-JC > SC$ ) の 3 パターンが存在することがわかる。また、均衡点に影響する相互依存性の要素は SC と JC だけであり、PC は全く関与しないこともわかる。これは、PC が自己の行動を変化させることによって利得が影響を受けない要素であるからである。

一方、均衡点における平均利得  $u$  は次のように表される。

$$u = (GM + (2x - 1)(PC + SC + JC(2x - 1)))/2$$

このように、均衡点の平均利得には PC や GM なども影響しうることがわかる。

### ジレンマと不確実性による 3 種類の不調和のパターン

動学的試行錯誤ダイナミクスの分析から、相互依存性のパターンには 3 つのタイプがあることが明らかとなった。つまり、SC が JC よりも影響力が強い「SC 優位パターン ( $|JC| < |SC|$ )」、SC と JC の符号が一致しており JC のほうが大きい「正の JC 優位パターン (JC > SC)」、そして SC と JC の符号が一致しておらず JC のほうが大きい「負の JC 優位パターン ( $-JC > SC$ )」の 3 つである。一方、均衡点の利得は相互依存性の要素すべてが関与することから、均衡点は常にパレート効率的ではないことがわかる。具体的には、SC 優位パターンの場合、 $-PC > SC$  の場合は、均衡点は常に  $x=0$  の戦略の組に対してパレート不効率となる。これは典型的な囚人のジレンマの状況である。同様のことは、他のパターンにも当てはまり、JC 優位パターンであっても、PC が JC よりも十分大きい場合、均衡点よりパレート効率的な

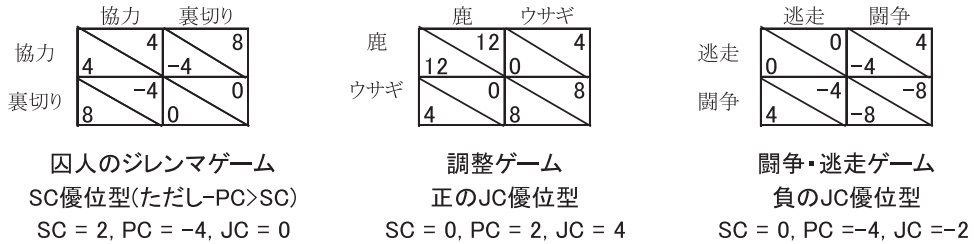


図2 3つの不調和な相互依存性のパターンの例

戦略が少なくとも1つ存在することになる。  
 正のJC優位型のゲームにおいても、均衡点がパレート不効率になる場合がある。それは正のJC優位型のゲームには不確実性が存在することによる。正のJC優位型ゲームでは、 $x=0$ か1が均衡点であるがどちらの均衡点になるかは、初期値に大きく依存する。よって、 $x=1$ がパレート効率的な戦略の組(両方とも鹿を狩る)であったとしても、 $x=0$ に収束してしまう(結果的に二人ともウサギを狩る)ことがある。このように、パレート効率的な均衡点に到達できない可能性があるという意味で、正のJC優位型ゲームは不確実性がある。

この不確実性に類似した問題は、負のJC優位型ゲームにおいても存在する。図のように負のJC優位型ゲームは、チキンゲームやタカハトゲームと呼ばれる闘争・逃走(fight or flight)型のゲームである。具体的には、片方が闘争を、もう片方が逃走を選択した場合、闘争を選んだほうが有利だが、両方が闘争を選ぶと大きな損害がある、というゲームである。この場合、闘争による損害をさけるためには相手がどちら戦略(闘争か逃走)をとる傾向を持っているかを事前に知っている必要がある。相手が闘争型なら逃げたほうがよいし、逃走型なら戦ったほうが良いからである。均衡点は闘争型と逃走型の両方が混在している状態であるので、特定の割合で戦いが行われていることになる。進化ゲーム理論では、このような闘争型のゲームにおいて、所有権方略がESSになりうるということが指摘されている。これはある資源に先に辿り着いた場合に闘争を、あとの場合は逃走を選ぶ、ということである。この戦略は、先にたどり着くという偶然的な手がかりを手がかりに闘争を避けるという方略である。このような戦略による均衡を相関均衡(Aumann, 1987)と呼

ぶ。負のJC優位型ゲームの均衡点は、相関均衡点に比べてパレート非効率になる。

このように、これら3つの相互依存性パターンについては、それぞれ異なるかたちで、均衡点よりもパレート優位な戦略の組があり得ることが明らかとなった。均衡点よりもパレート優位な戦略の組があるゲームの事を、相互依存性理論にならない「不調和な」ゲームと呼ぼう。またこれら3種類の不調和のパターンは、代表的にはそれぞれジレンマゲーム、調整ゲーム、チキンゲーム(闘争ゲーム)と呼ばれるものを一般化したものであるといえる。Kelley (1979)の表現を使えば、SC優位型で $-PC > SC$ のパターンは「交換型」、JC優位型は「調整型」である。Kelley (1979)ではJCの正と負の区別をつけていないため、ここでは正のJC優位型を調整型、負のJC優位型のゲームを闘争型と呼ぼう。なお相関均衡を範囲に入れ場合、調和的な(均衡点とパレート効率的な戦略の組が一致する)ゲームは、SC優位型で $-PC < SC$ の場合のみである。

以降では、交換、調整、闘争型それぞれの不調和パターンについて、均衡点をパレート効率的な均衡点に変換する仕組みについて考察する。

### 利得行列の変換

相互依存性理論の特徴は、ゲーム構造の要素の分解によるパターン分けだけではなく、不調和な利得構造を調和的な利得構造に変換することについての考察が行われている点にもある。利得構造の変換は、まず二人が与えられた利得構造(これを所与行列と呼ぶ)を主観的に変換し、変換された利得構造(これを実効行列と呼ぶ)に基づいて行動を行う、とされる。変換方略には、大きく分けて3つあげられており、成果変換と転置変換、そして系列型変換である。

成果変換は、得られる利得をそのまま変換する変換方略である。例えば相手に協力することはいいことだと捉え直すことによって、囚人のジレンマゲームは調和的な安心ゲームに変換される、といったものである。成果変換は、交換型ゲームを調和的にするために使われる変換であると考えられている。次に転置変換は、同期ゲーム（同時に選択するゲーム）を非同期ゲームに変えるものである。例えば調整ゲームにおいて相手がどちらの手をとったかわかっているならば、容易にパレート優位な戦略の組を選択することができるだろう。このように、転置変換は調整ゲームを調和的にするための変換方略である。最後に系列変換は、1時点の闘争ゲームを繰り返しの闘争ゲームとして考え、片方が闘争をとっても片方が逃走を選んだら、次のゲームでは逆の手をとる、というやりとりを繰り返すようにする、というものである。つまり闘争型の不調和ゲームに時系列を導入することで、平等な資源分配が行えるように変換する変換方略であるといえる。

このように、相互依存性理論は基本的には1時点の二者間の相互作用を対象としており、利得構造の変換も「主観的に」変換されると考えられている。しかし、進化・学習的な観点に立った場合、主観的に利得構造を変換したとしてもダイナミクスに影響すると考えることは難しい。頭のなかだけで利得が増えたと考えたとしても、実際に得をしているわけではないからである。そこで、主観的ではなく、相互作用の仕方そのものを変換することで不調和なゲームを調和的にするための変換方略を動学的に考察する必要がある。

### 動学的な利得行列の変換

相互依存性理論で提案されている交換型の変換方略である成果変換は、あくまで主観的なものであるため、動学的なゲーム理論においては有効ではない。一方、進化ゲーム理論においてはAxelrod (1984) などが提案している互恵性戦略 (Tit For Tat: TFT 戦略) がある。互恵性とは、協力してくれた相手に対して協力をするという戦略であるが、協力的な行動を選択することが、その時は損をするが次以降のゲームで得をするという変換を行っているといえる。これはつまり、ある時

点のゲームにおける選択が後のゲームにおける利得と時系列的に連結されているような変換であるといえる。これを、「連結変換」と呼ぼう。TFT戦略は集団内における二者関係の直接的な互恵性を達成するための変換であるが、間接的互恵性 (Nowak, 1998) のように、二者関係から拡張された利得構造も同様に連結変換の結果と考えることができる。

複数の均衡点があることによって不確実性が存在する調整型の不調和ゲームにおいては、相互依存性理論では転置変換によって調和なゲームに変換されていた。しかし、転置変換を用いる場合であっても、どちらが先に手を決めるかについて同様の不確実性が存在する。よって、不確実性そのものを無くすための変換が必要である。そのためには、相手との意思疎通をはかったり、どちらの均衡点がより価値が高いかについての共有知識が手かかりとして用いたりすることが有効だろう。具体的には、コミュニケーションによってどちらの均衡がパレート効率的であるかを確認する、あるいはどちらの行動を選択すべきかについてあらかじめ集団内で合意を得る (規則を作る) といった事が考えられる。つまり、転置変換を確実に機能させるための仕組みがさらに必要となる、ということである。

第三の闘争型ゲームにおける変換は、系列変換があげられていた。しかし、集団型試行錯誤ダイナミクスは同じ二者が繰り返し相互作用するわけではなく、ランダムマッチングを仮定しているため、同様の変換を用いることができない。しかし、系列変換の本質が資源の平等分配にあるのであれば、分配可能な資源を平等に分割する、といった変換が可能であろう。また、もし資源の分割が不可能な場合、数土 (1998) が示したように儀礼的闘争による闘争の回避も機能すると考えられる。儀礼的闘争とはコストの比較的かからない競争 (木を揺らすなど) によっても相手の力の強さが正直なシグナルとして伝わることを利用した相互作用である。これをコストリーシグナリングとも呼ぶ。儀礼的闘争により、相手より自分が弱いと判断したら闘争を、強いと判断したら逃走を選択するという戦略をとることが可能になる。つまり、コストリーシグナリングによってどちらが逃

走すべきかの手がかりを共有することができるのである。

本節では、動学的な相互依存性理論に基づいて、不調和なゲームのパターンの分類と、それぞれの調和性を解決するための利得構造の変換について議論した。以降の節では、これらの変換を制度分析アプローチから捉え直す作業を行う。

### 3節 相互依存パターンの制度的変換

#### 制度分析アプローチ

本節では、動学的相互依存性理論の考え方を制度分析アプローチに拡張することを目標とする。そこで、まず制度分析アプローチを支える比較制度分析について簡単にレビューしよう。

比較制度分析は Aoki (2000) が提唱した、比較的新しい経済学理論である。制度は、比較制度分析においては「均衡の要約表現と予測」というように定式化されている。つまり、ある均衡がどのような状態であるかについての、暗黙あるいはシンボリックな共有された信念体系を指す。すなわち、制度分析アプローチでは制度を個人の戦略の集合として考えているのではなく、その均衡状態を人々がどのように捉えているか、という側面を強調している。また、制度分析では、制度の変化についてその認知的な側面や規範などの社会的な側面など多岐にわたって考察がされている。このように、制度分析では進化あるいは経済ゲーム理論が均衡のみを考察の対象にするのに対して、均衡についての認知的な側面やその変化をも分析対象としている点について画期的である。

しかし、本論文では制度を、均衡そのものや、均衡についての要約表現として捉えているわけではない。本論文では相互依存性理論が主張するように、所与行列と実効行列を区別する点が制度分析と異なっている。対比的に述べれば、制度分析における制度が、変換された後の実効行列に基づく均衡を指しているのに対し、本論文では、所与行列から実効行列に変換させる、あるいは所与行列による学習ダイナミクスの均衡から実効行列による学習ダイナミクスの均衡にシフトさせる社会的な仕組みのことを指している。

では、どの状態を所与行列と考えればい

うか。次項では所与行列と実効行列の関係について考察し、本論文における制度の中身について詳述する。

#### 資源獲得における原初的な相互依存性のパターンと制度

Jacobs (1992) は、人々が資源を獲得する方法は、交換することと採取することの2パターンしかないことを指摘し、道徳や倫理がこの二つの相互依存性に基づいて形成されていることを指摘した。この交換と採取は、第2節でわれわれが見出した3つのパターンに対応させることができる。具体的には、交換は交換型の不調和パターン、採取は調整型と闘争型の不調和パターンに分類することができる。後者については、資源の採取は人々が共同で採取する場合と、資源を奪い合って採取するという二つの状況が考えられるからである。資源の採取が一人では不可能な場合は協同で採取し、一人で採取可能な場合は取り合いになる、ということである。Jacobs の指摘を踏まえれば、この3つのパターンは、社会集団において資源を共有するときに生じる本質的な相互依存性パターンを意味していると考えられる。これらを原初的相互依存性パターンと呼ぶことにしよう。

これらを踏まえ、本論文では以下の想定を行う。すなわち、人類が集団を形成し始めた当初は、このような原初的な不調和パターンに遭遇していた。それに対して、社会集団を形成した人々はより効率的な相互作用の仕組みを生み出していくことで、よりパレート効率的な均衡に至ることができるようになっていったと考えるのである。このように考えると、原初的な不調和パターンを所与行列、そしてパレート効率的な均衡をもたらす実効行列に変換するための装置を制度と呼ぶことができる。ここで、変換後の実効行列のことを制度的相互依存性パターンと呼ぶことにしよう。まとめると、人々は、原初的な相互依存性パターンから制度的相互依存性パターンに利得行列を変換することによって、よりパレート効率的な均衡に至るように制度を発展させてきた、ということである。

制度による変換を具体的に説明すると、交換型のジレンマ構造が、連結的変換によって間接的互



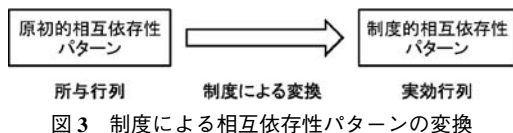


図3 制度による相互依存性パターンの変換

互惠性が達成され、相互協力的な社会状態が達成される、などが挙げられるだろう。まず資源を交換するという原初的な相互依存性パターンがあり、そのまま相互非協力が非協力となる。それに対して連結変換、つまり協力を選択すると、次のゲームで相手が協力する確率が高くなるといった変換（これは具体的には評判制度によって達成される）によって、間接的互惠性が成立し、相互協力状態が均衡になりえる（Nowak, 1998；Takahashi & Mashima, 2006）。このとき、制度的な変換を支えるのは協力する人を識別する個人の判断傾向だけではなく、評判をうまく共有するための社会的な仕組みも同時に必要となるだろう。たとえば Yamagishi & Kiyonari (2000) では、間接的互惠性を支える器として集団が機能した、という論を展開している。この議論を踏まえると、集団あるいは集団主義的な制度が、間接的互惠性（連結変換）を維持するための仕組みであるという考え方も可能であろう。具体的には、協力的な他者を「仲間」、非協力的な他者（これは非協力者に協力的な人も含む）を「裏切り者」や「よその」として排除するような、村八分的な集団主義的な制度が、間接的互惠性を維持している、ということである。このように、制度は、「個人のもつ行動傾向」というよりは、「人々が共有している相互作用のやり方（河合、2013）」という考え方のほうが妥当であると考えられる。

制度が個人の戦略に還元できない場合を考察する例として、さらに複雑な社会制度を考えることもできる。たとえば、間接的互惠性をより広範囲に、一般的に、効率的なものにするための制度として貨幣制度が挙げられる。貨幣は、他者に利他的に振る舞った人（商品を提供した人）が所有できる、いわば評判を数量的に把握するための装置である。つまり貨幣を使えば他者からの協力行動を引き出すことができ、逆に他者に対して利他行動をすれば貨幣を手に入れることができる。このように貨幣は、そのものに価値がなくとも、一度

流通してしまえば間接的互惠性の評判システムをより効率的に支えるメディアとして機能するのである。

ここで貨幣制度を、個人の遺伝的に継承された行動傾向によってのみ支えられていると考えるのは難しいだろう。なぜなら、貨幣そのものは社会によって外在的に用意された道具だからである。このように、制度分析アプローチは個人の戦略の集合としてではなく、外在的な社会装置をも分析の対象にすることができるのが利点である。

### 制度の進化

文化進化理論は、遺伝子以外の複製単位として他者の行動についての情報を挙げ、それが個体の学習能力という遺伝的な傾向性に支えられて進化すると考えている。つまり、個体学習に加えて模倣学習（社会学習）の二重の継承メカニズムによって適応的な戦略が文化として累積的に継承されていくことを仮定する。

それに対して本論文では、文化や規範の進化を考えるにあたって、複製の単位が制度であると考えられる。制度を複製子と考えるということは、制度そのものが複製され、社会集団の中に増えていく、ということである。このことを理解するためには、制度が個人の戦略に還元できない社会装置である、という本論文の前提を振り返ろう。

制度は個人の戦略に還元できない以上、個人を媒体として複製することはできない。少なくともふたり以上の個人が相互作用し、利得構造を原初的なパターンから制度的パターンに変換するというプロセスが必要である。よって、制度を複製子と考えるとき、媒体は社会的相互作用であると考えられる必要がある。より多くの人々が次々にその制度によって相互作用するとき、制度が複製されていると考えることができるのである。

ここで、原初的な相互依存性をよりパレート効率的に変換できる方法が複数あると考えよう。例えば、貨幣制度 A と貨幣制度 B があり、どちらも間接的互惠性を維持することができる均衡点を可能にする制度であると仮定する。しかしこのとき、貨幣制度 A において使われている貨幣がとても重く使いづらく、貨幣制度 B で用いられている貨幣が比較的軽くて使いやすいとしよう。す

ると、相互作用するときのコストは貨幣制度 B のほうが小さくなり、貨幣制度 A の均衡点は貨幣制度 B の均衡点よりもパレート不効率になる。このとき、貨幣制度に参加している個人がどちらの貨幣制度をより採用するようになるだろうか。おそらく、貨幣制度 B のほうの貨幣を用いて相互作用するようになると考えられる。

この例をもとに、制度の進化を理論化しよう。複数の制度が同じ社会集団で併存しているとき、制度の複製は、その制度を採用している個人の多さ（正確には制度を採用した相互作用の多さ）によって決定される。そして、制度の利用は制度による変換後の実効行列の均衡点の利得に依存する。すなわち、よりパレート効率的な均衡点を可能にする制度がより複製される、つまり進化するということである。よって、貨幣はより使いやすく、偽札が作りにくく、製造コストがかからないようなものに進化していくはずである。

ここで注意が必要なのは、制度が所与行列の均衡点から別の均衡点に変換するための装置であるため、集団淘汰理論で批判される淘汰スピードの問題（集団レベルの淘汰よりも個体レベルの淘汰のほうがスピードは早いため、個人レベルで成立しない均衡は集団レベルでも成立しない問題）は生じない。しかし、人々が制度を選択するというプロセスをまだ厳密にはダイナミクスに取り入れられてはいない。この議論を精緻化するためには、個人だけではなく、制度をプレイヤーとした利得行列に基づく高次のゲームダイナミクスを構築する必要があるだろう。この点について本論文ではまだ十分な検討はできていない。

### 不調和パターンごとの制度進化

この制度進化のロジックは、交換型の不調和パターンだけに当てはまるものではなく、調整型や闘争型であっても同様である。変換後の均衡点の期待利得のパレート効率性によってその制度が支配的になるからである。

しかし、不調和パターンによって、進化する制度は当然異なったものになるだろう。それぞれの不調和パターンをより効率的に変換する制度がどのようなものであるかについては、実証的な問いとなるため、本論文では十分に答えをだすことは

できない。ただ、いくつかの歴史的な制度の変遷を手がかりに、不十分ながらもデモンストレーションすることは可能だろう。

以後は、人類学で指摘されている、集団の規模の推移に基づいて制度の形態の変容を考察しよう。人類学的に、社会集団は狩猟採集社会、部族社会、首長社会、国家、という四段階で発展していったと考えられている（Diamond, 1997）。もちろんこれらの分類は便宜的なもので、もちうる性質には相互に重なりがある。狩猟採集社会は、50人程度の規模の、主に親族集団によって形成されていた小集団の社会である。そして部族社会は複数の親族集団がまとまった数百人規模の集団で、およそ1万年前に発生したと考えられている。首長社会はさらに大きく、数千人規模の社会である。国家はおよそ5000年前に発生した、数万人規模の社会である。ここでは、狩猟採集社会のような小規模な集団、そして部族社会や首長社会のような大規模な集団、そして国家のような更に大きな社会の3つに分けて制度の進化を考察していこう。

交換型はすでに述べたように、連結型の変換によって直接あるいは間接的な互惠性を築くことで相互協力の状態を作り上げてきたと考えることができる。小規模な集団においては二者関係における TFT 戦略のような直接的互惠性を形成し相互協力状態を築いてきたと考えられる。しかし、集団の規模が大きくなるに連れて、直接的な互惠性では相手が協力的な人間かをどうかを判断することが難しくなってくるだろう。そこで評判制度を用いた間接的互惠性によって集団単位の相互協力が達成されると考えられる。すでに述べたように、協力的な人であるかどうか、つまり仲間であるかどうかを互いに監視しあって、もし裏切り行動を見せたら即座に悪い評判を集団内に伝搬させる。そのような仕組みによって集団内の相互協力状態を維持していたのかもしれない。そして現代社会のように更に大きな社会では、評判が行き渡らない広範囲な集団や国家において貨幣制度を用いて間接的互惠性を成立させることに成功したと考えられる。

調整型は、二者関係においてはジェスチャーや言語などを用いたコミュニケーションによって転

置変換が行われていたと考えられる。コミュニケーションによって、自分があるいは相手が何に価値を持っているか、何がしたいのかを伝えることができるからである。しかし、大きな集団で調整を行う場合、いちいちコミュニケーションで合意を形成するのは難しい。集団の歴史があるならばこれまでの経験から慣習を作り上げて、規則や宗教などの価値体系を作ることが有用かもしれない。あるいは、のちに述べる権威のある人物による意思決定なども機能したかもしれない。そして現代では、学問によって事前にどの均衡点に至るのが妥当であるかについて知ることができているといえる。Luhmann (1984) は科学や学問が真理をシンボリックなメディアとするコミュニケーションシステムであると指摘した。科学や学問システムが、(科学的な) 事実を共有するための制度であると考えられることができるだろう。

最後に闘争型は、二者関係においては系列変換や、平等的な変換、あるいはコストリーシグナリングによる儀礼的闘争によって深刻な闘争を避けていたと考えられる。たとえば寺嶋 (2011) や Diamond (1997) は、人類が狩猟採集社会において非常に強い平等規範を持っていたことを指摘している。また Kameda, Takezawa, & Hastie (2003) も平等分配制度が闘争を回避することができる戦略の均衡の結果であるという考察を行っている。続いて、より大きな集団になった場合の闘争回避の制度として、たとえば Diamond (1997) は権威や階層制度が用いられていたと述べている。権威は、個体の腕力ではなく、財力や影響力についてのコストリーシグナリングによる儀礼的闘争であると考えられることができるかもしれない。つまり、多くの強い影響力を持つものほどより大きな人を動かすことができるため、その相手と闘争して勝つことが困難になるからである。そして現代では、闘争回避システムとして裁判制度、つまり法制度が用いられている。裁判制度では法によって争点となる資源をどちらに所有権がある、あるいはどれほど分配するのが妥当であるかを判断するシステムである。Luhmann (1975) は法が物理的暴力を抑える効力を持つ一方、暴力が国家に集中する (軍や警察など) というかたちで権力を持ちうることを指摘している。つまり、裁判や警

	二者関係・小集団	大規模集団	国家
交換型	直接的互惠性	評判制度・村八分	貨幣制度
調整型	言語・コミュニケーション	慣習・宗教	学問・科学
闘争型	平等・儀礼的闘争	権威・階層制度	法・裁判

図 4 制度の進化

察、軍といった国家権力は、社会集団の闘争を管理するための制度であるという見方が可能である。

このように、それぞれの不調和パターンに従って、制度は図のように集団の規模や時代によって進化してきたと考えることができるだろう。

#### 4 節 道徳と制度の共進化

##### 制度維持の潤滑油としての道徳

3 節では制度分析アプローチから制度の進化までを考察したが、本節では本論文の目的である、道徳の起源について考察しよう。

Smith (1970) では、経済制度が機能する背景として、人々の利他的な道徳感情があることを指摘している。Smith のこの主張を一般化すれば、道徳感情を持つ個人がゲームを行うことで初めて達成される均衡点がある、ということの意味している。しかし、進化ゲーム理論の立場にたてば、そもそも利他的な個人が進化する基盤そのものを説明しなければ、この議論は成立しない。Haidt (2012) の道徳基盤理論では、道徳感情が進化的に獲得されたものであるという立場で議論を行っているが、すでに指摘したように、集団淘汰を不十分な形で適用しているに過ぎず、十分な説得力を持つものではない。

本論文の立場は、Smith のように道徳感情を原因としていまの社会制度や道徳規範が維持されているという立場には立たない。また同時に、道徳感情が遺伝的な意味での進化の産物であり、それによって道徳規範が社会に共有されているという立場にも立たない。本論文では、まず制度があり、それに適応するかたちで人々は道徳規範を成立させた、という立場である。つまり、道徳は制度の原因ではなく、制度に人々がうまく適応するための潤滑油としての機能を持っている、という

立場である。

### 制度に適応するための学習方略としての道徳

制度は所与行列による均衡を、実効行列の均衡にシフトさせる社会装置であった。よって、均衡点に至っている個人は道徳や規範意識などがなくとも所与行列の不調和は解決されているはずである。つまり、道徳がなくとも制度が機能していれば、相互協力や闘争回避などは実現されていることになる。では、なぜ道徳が必要なのだろうか。

制度は長いスパンでみれば変化するが、個人レベルで見れば長期にわたって維持されているものであると考えられる。よって、制度に参加する個人は世代交代が行われる一方、制度そのものは大きくは変わらない状態が続く。このような環境においては、個人は自らが制度に適応するだけでなく、新規成員に対しても制度に適応するように教育することが動機づけられるはずである。なぜなら、ある制度において相互作用をするときに、相手が制度を理解していなければ、不調和を解決することができなくなるからである。もし制度を知らない個人が少数でもいた場合に不調和が解決されないような利得構造であるなら、N人チキンゲームと同じ利得構造になり、少数の個人が教育コストをかけて全体の利益を向上させるような頻度依存型の均衡状態になる。つまり、少数の教育者が生まれることになる。また、包括適応度を考慮に入れた場合、自分子どもなどの血縁他者が制度に参加してうまく適応できるようになるためには、幼少期からその制度に適応するように教育する必要があるだろう。

教育者が学習する側に対して制度を学習させる場合、罰や賞を与えることで行動傾向を是正する方略が第一に考えられる。しかし、もし学習者が教育者の罰を予期することができて、かつ、教育者が罰を実際に与えずに正直なシグナルによって伝達することができれば、実際に罰（あるいは賞）が行使されずとも学習が行われる可能性がある（清水・嘉志摩、2011）。具体的には、教育者が怒りを表出させながら「～してはいけない」と伝達したとしよう。このとき、学習者が教育者の怒りを正直なシグナルであると判断できる能力があり、それによって罰を予期できるなら、罰を受

ける前に行動傾向を変化させることが適応的できるということである。つまり、教育者は罰のコストをかけずとも、感情をこめて「～してはいけない」という伝達を教育者にするだけで、学習者は実際に罰を受けずに当為的な道徳を学習することができる。

もしこのような言語伝達による制度の教育が機能するなら、人々は幼少期から制度に適応するための行動傾向そのものを身につけるだけではなく、どの行動をすると罰を受けるのか（叱られるのか）、という判断基準を内在化するようになるだろう。行動に対する逐次的な罰は、実際に行った行動と罰が連合されることによる学習であるので、道徳のような抽象的な行動のカテゴリ化や判断基準の内在化には至らないだろう。しかし、言語伝達による教育は、言語表現による行動の抽象化を媒介することで、「していい行動」と「してはいけない行動」にカテゴリ化する機能があると考えられる。しかし、これらの考察はまだ経験的なデータがない予想であるので、今後実証的な検討が必要である。

上記の議論が仮に正しいとすると、道徳が特定の制度を学習・教育する仮定で生まれる個人内の判断基準である、ということが出来る。ただし制度はパレート効率的な均衡点を与えるので、教育される道徳は同じ制度のもとでは共有されることになるだろう。

また、ある制度のもとで成員に道徳が共有されると、それによって新しい制度へのシフトへの道が開かれる可能性も考えられる。つまり、制度が道徳を作り、道徳を持った個体が多く集まることによってさらに効率的な制度の発明を促進する、という共進化の関係がある可能性がある。この観点は、本論文においてほとんど触れられていない「制度の発明」に対して、重要な視点を与えてくれるかもしれない。

### 制度と道徳の対応関係

図に記した各制度と、その制度下において発生するだろう道徳の対応関係を図に記そう。これらは、Haidt (2001) の道徳基盤理論で指摘されている道徳（公正や集団主義、権威への敬意など）と重複している部分もあるが、貨幣や科学につい

	二者関係・小集団	大規模集団	国家
交換型	直接的互惠性 →利他性	評判制度・村八分 →集団主義	貨幣制度 →儉約
調整型	言語・コミュニケーション →信頼	慣習・宗教 →慣習や宗教の遵守	学問・科学 →知の受容
闘争型	平等・儀礼的闘争 →平等・公正	権威・階層制度 →権威への敬意	法・裁判 →正義

図5 制度と道徳の対応関係

での道徳などについては新たに付け加えられている。

図を見ると、制度の進化にともなって、その道徳がどのように人々に習得されているかに違いが伺える。小集団において成立した制度に伴う道徳は、Haidt (2012) が指摘するような、生得的な道徳感情によって支えられている部分もあるかもしれない。しかし一方、大規模集団における制度については、生得的というよりは、幼少期の躾や教育によって形成される道徳（仲間を大切にしない、目上の人のことを聞きなさい）と対応しているように見える。そして国家における社会制度は、初等教育以降において学習する道徳に対応すると見ることができるとも思われる。つまり制度が複雑になるにつれて、生得的な道徳から、発達の道徳、そしてより高次の道徳へと変遷している、ということである。これは Kohlberg (1976) が前慣習的道德、慣習的道德、後慣習的道德と表現した道徳発達の理論とも整合的である点も興味深い。

しかし、これらの道徳と制度の対応についてはあくまで予測であって、実証的な証拠や精緻に論理的な証明ができるわけではない。よって、ここでは各制度と道徳の具体的な対応関係について詳述はしない。しかし、この予測がどれほど妥当であるかについては経験的な検証を必要とするだろう。

## 5 節 結論

本論文では、動学的相互依存性理論と制度分析アプローチのアイデアをもとに、道徳の起源について考察してきた。

相互依存性理論を動学的に拡張したモデルによ

って、社会集団における相互依存性には、3つの不調和な相互依存性が存在することが明らかとなった。また、それぞれの不調和パターンに対するパレート効率性を高める変換方略が分析された。

続いて、制度分析アプローチから、利得行列の変換こそを制度として定義し、その進化の可能性について議論を行った。その結果、制度は個体を媒体とするのではなく、社会的相互作用を媒体とする複製子として進化しうることが議論した。そして、狩猟採集社会から現代社会にかけての制度の進化について不十分ながら、その推移の描写を試みた。

最後に、道徳は人々が制度により効率よく適応するための教育によって発生した心のメカニズムであり、同じ制度内で相互作用する人々に共有されうることを示した。最終的に、道徳と制度の対応関係と共進化についての予測を述べた。

以上の議論から、本研究では、道徳は制度の進化と、制度に適応する教育方略の両方によって発生した人々の適応的な行動判断システムである、という結論にいたることができた。本論文が主張する結論は、道徳と制度の対応関係の結果だけを見ればそれほど目新しいものではないかもしれない。しかし、制度に適応するための言語伝達の教育によって道徳が発生するという主張は、これまでの道徳研究にはない結論である。

ただし、まだまだ途中の論理展開も飛躍し、検証すべき経験的な問いも多く残されている。ただ、本論文の結論はこれだけによって成立する成果というよりは、道徳研究に新しい見方を提供し、さらなる研究の積み重ねを狙った理論的研究である。よって、今後は、検証可能な仮説について実証的なデータを収集し、検証していく必要があるだろう。

**Appendix**

大浦 (2008) を参考に対称 2×2 の集団型試行錯誤ダイナミクスの均衡点について考える。

ここで、SC が正の方向の戦略を A、反対の戦略を B とする。このとき、各戦略の組の利得は以下のようになる。

$$\begin{aligned} AA &= (SC + PC + JC + GM)/2 \\ AB &= (SC - PC - JC + GM)/2 \\ BA &= (-SC + PC - JC + GM)/2 \\ BB &= (-SC - PC + JC + GM)/2 \end{aligned} \quad \text{式 1}$$

ここで Roth & Erev (1995) の試行錯誤ダイナミクスに従い、以下の仮定を置く。

まずプレーヤー  $i$  が戦略 A をとる行動傾向を  $P_i a$ 、戦略 B をとる行動傾向を  $P_i b$  とする。このときプレーヤー  $i$  が時間  $t$  に置いて戦略 A を取る確率を  $x_i$  は

$$x_i(t) = P_i a(t)/(P_i a(t) + P_i b(t)) \quad \text{式 2}$$

で表されるとする。また、時間  $t$  におけるプレーヤー  $i$  の戦略 A を選択する傾向性  $P_i(t)$  は、

$$P_i(t+1) = (1 - \phi)P_i(t) + R_i(t) \quad \text{式 3}$$

で表されるとする。なお、 $\phi$  は忘却率、 $R_i(t)$  は時間  $t$  におけるプレーヤー  $i$  に与えられた強化値である。ただし、今回の結論に忘却率は影響しないので、以降  $\phi = 1$  として省略する。

ここでプレーヤー  $i$  が戦略 A をとる確率の時点  $t+1$  と  $t$  の変化量である  $\Delta x_i$  を考えよう。また、 $P_i(t) = (P_i a(t) + P_i b(t))$  とおくと、

$$\begin{aligned} \Delta x_i &= x_i(t+1) - x_i(t) \\ &= Pa(t+1)/P(t+1) - x_i(t) \\ &= ((1-x(t))Pa(t+1) - x_i(t)P_i b(t+1))/P_i(t+1) \end{aligned} \quad \text{式 4}$$

ここで式 2 から  $P_i a(t) = x_i(t)P_i(t)$ 、 $P_i b(t) = (1 - x_i(t))P_i(t)$  であることから、これお式 4 の分子に代入して展開すると、

$$\begin{aligned} &((1-x(t))Pa(t+1) - x_i(t)P_i b(t+1)) \\ &= R_i a - x(t)(R_i a + R_i b) \end{aligned} \quad \text{式 5}$$

となる。以上より、

$$\Delta x_i = (R_i a - x(t)(R_i a + R_i b))/P_i(t+1) \quad \text{式 6}$$

となる。

ここで、 $x$  を集団中の  $x_i$  の平均とすると、集団中からランダムに一人選んだ相手が戦略 A をとる確率は  $x$  となる。このとき、プレーヤー  $i$  が A や B をとろうとする傾向性である  $P_i a(t)$ 、 $P_i b(t)$  に与えられる強化値  $R_i a$ 、 $R_i b$  はそれぞれ、

$$\begin{aligned} \text{確率 } x_i x \text{ で } R_i a &= AA, R_i b = 0 \\ \text{確率 } x_i(1-x) \text{ で } R_i a &= AB, R_i b = 0 \\ \text{確率 } (1-x_i)x \text{ で } R_i a &= 0, R_i b = BA \\ \text{確率 } (1-x_i)(1-x) \text{ で } R_i a &= 0, R_i b = BB \end{aligned} \quad \text{式 7}$$

となる。式 6 に式 7 を代入すれば、 $\Delta x_i$  は

$$\begin{aligned} \text{確率 } x_i x \text{ で } &(1-x_i) AA/(P_i a(t) + AA) \\ \text{確率 } x_i(1-x) \text{ で } &(1-x_i) AB/(P_i a(t) + AB) \\ \text{確率 } (1-x_i)x \text{ で } &-x_i BA/(P_i b(t) + BA) \\ \text{確率 } (1-x_i)(1-x) \text{ で } &-x_i BB/(P_i b(t) + BB) \end{aligned} \quad \text{式 8}$$

となる。ここで、行動傾向性  $P_i a$  や  $P_i b$  に対して一度で得られる利得の大きさは相対的に小さいと考えられるため、分母を概ね同じと考え (大浦、2008)、定数  $k$  とする。そして分子だけに注目して、 $\Delta x_i$  の期待値  $E[\Delta x_i]$  の分子を求めると、

$$\begin{aligned} E[\Delta x_i] \text{ の分子} &= x_i x(1-x_i) AA + x_i(1-x)(1-x_i) AB \\ &\quad - (1-x_i)x x_i BA - (1-x_i)(1-x)x_i BB \\ &= x_i(1-x_i)(x AA + (1-x) AB - x BA \\ &\quad - (1-x) BB) \end{aligned} \quad \text{式 9}$$

式 9 を展開して、式 1 のように相互依存性理論の要素 SC、PC、JC、GM で表現し直すと、

$$E[\Delta x_i] = x_i(1-x_i)(-JC + SC + 2JCx)/k \quad \text{式 10}$$

となる。この差分方程式を解くと、 $x_i$  の均衡点と一致するため、

$$x = 0, 1, (JC - SC)/2JC \text{ の 3 つとなる。}$$

対称ゲームであるので、 $x_i$  の均衡点は、 $x$  の均衡点と一致する。

**引用文献**

Aoki, M. (2001). *Toward a Comparative Institutional Analysis*. MIR Press.

Aumann, R. J. (1987). Correlated equilibrium as an expression of bayesian rationality. *Econometrica*, 55, 1-18.

Aunger, R. ed (2000). *Darwinizing culture : the status of memetics as a science*. Oxford University Press.

Axelrod, R. (1984). *The Evolution of Cooperation*, Basic Books.

Barash, D. P. (2007). *Evolution and Group Selection*. Comment on The new synthesis in moral psychology. *Science*, 317, 596-597.

Diamond, J. (1997). *Guns, Germs, and Steel : The Fates of Human Societies*. W. W. Norton & Co.

Gintis, H. (2009). *The Bounds of Reason : Game Theory and the Unification of the Behavioral Sciences*. Princeton University Press. (ギンタス, H. 成田悠輔・小川一仁, ・川越敏司・佐々木俊一郎 訳 (2011). *ゲーム理論による社会科学の統合* NTT 出版)

Haidt, J. (2001). *The Emotional Dog and Its Rational Tail :*

- A Social Intuitionist Approach to Moral Judgement. *Psychological Review* 108, 814–834.
- Haidt, J. (2012). *The Righteous Mind: Why Good People Are Divided by Politics and Religion*. Pantheon.
- Heath, J. (2008). *Following the Rules: Practical Reasoning and Deontic Constraint* Oxford University Press. (ヒース, J. 瀧澤弘和 (2013). ルールに従う NTT 出版)
- 芋坂直行編 道徳の神経哲学 神経倫理から見た社会意識の形成 新曜社
- Jacobs, J. (1992). *Systems of survival: a dialogue on the moral foundations of commerce and politics*, Random House. (ジェイコブス, J. 香西泰訳 (2003). 市場の倫理 統治の倫理 日本ビジネス人文庫)
- Kameda, T., Takezawa, M., & Hastie, R. (2003). The logic of social sharing: An evolutionary game analysis of adaptive norm development. *Personality and Social Psychology Review*, 7, 2–19.
- 河合香吏 (2013). 制度 人類社会の進化 京都大学学術出版会
- Kelley, H. H. (1979). *Personal relationships: The structure and processes*. Lawrence Erlbaum Association, Inc. (ケリー, H. H. 黒川正流・藤原武弘訳 親密な二人についての社会心理学 パーソナル・リレーションシップ ナカニシヤ出版)
- Kelley, H. H., Holmes, J. G., Kerr, N. L., Reis, H. T., Rusbult, C. E., & Van Lange, P. A. (2003). *An atlas of interpersonal situations*. New York: Cambridge.
- Kelley, H. H. & Thibaut, J. H. (1978). *Interpersonal relations: A theory of interdependence*. Wiley-Interscience. (ケリー, H. H.・テイボー, J. H. 黒川正流 監訳 (1995). 対人関係論 誠信書房)
- Kohlberg, L. (1976). Moral stages and moralization: The cognitive-developmental approach. Lickona, T ed. *Moral Development and Behavior: Theory, Research and Social Issues*. Holt, NY: Rinehart and Winston.
- Luhmann, N. (1975). *Macht*. Ferdinand Enke Verlag. (ルーマン, N. 長岡克行訳 (1986). 権力 勁草書房)
- Luhmann, N. (1984). *Soziale Systeme. Grundriss einer allgemeinen Theorie*. Suhrkamp Verlag, Frankfurt am Main. (佐藤勉監訳 (1993). 社会システム論 上 恒星社厚生閣)
- Nowak, M. A. & Sigmund, K. (1998). Evolution of indirect reciprocity by image scoring, *Nature*, 393, 573–577.
- 大浦宏邦 (2008). 社会科学者のための進化ゲーム理論 基礎から応用まで 勁草書房.
- Richerson, P. J. & Boyd, R. (2005). *Not by Genes Alone: How Culture Transformed Human Evolution*. Univ. of Chicago Press, Chicago.
- Roth, A. E. & Erev, I. (1995). Learning in extensive-form games: Experimental data and simple dynamic models in the intermediate term, *Games and Economic Behavior*, 8, 164–212.
- Smith, A. (1759). *The Theory of Moral Sentiments*. (スミス, A. 高哲男訳 (2013). 道徳感情論 講談社学術文庫)
- 数土直紀 (1998). 権力構造の発生モデル 理論と方法, 12, 167–179.
- Takahashi, N. & Mashima, R. (2006). The importance of subjectivity in perceptual errors on the emergence of indirect reciprocity. *Journal of Theoretical Biology*, 243, 418–436.
- 寺嶋秀明 平等論 霊長類と人における社会と平等性の進化 ナカニシヤ出版
- Tomasello, M. (1999). *The cultural origins of human cognition*. Cambridge, MA: Harvard University Press. (トマセロ, M. 大堀壽夫・中澤恒子・西村義樹・本多 啓 (訳) (2006). 心とことばの起源を探る: 文化と認知 勁草書房)
- Yamagishi, T. & Kiyonari, T. (2000). The group as the container of generalized reciprocity. *Social Psychology Quarterly*, 116–132.

## A Theoretical Approach to the Origin of Morality: The Dynamic Interdependence Theory and Institutional Analysis

### ABSTRACT

This article discusses the origin of morality on the basis of dynamic interdependence theory and institutional analysis. The moral foundation theory explains the origin of morality through a multilevel selection theory, a genetic evolutionary approach taking both individual and group selection into consideration. However, this theory, similar to other cultural evolutionary theories, has certain limitations in explaining how injunctive morality arose in social groups. Therefore, we offer an institutional evolutionary explanation of the origin of morality, which is based on dynamic interdependence theory and comparative institutional analysis. This approach gives us certain important advantages in discussing the origin of morality. First, three fundamental patterns of interdependence are identified on the basis of the dynamic interdependence theory. Then, transformation methods of outcome matrices corresponding to the patterns of interdependence are developed. Next, institutional analysis allows us to describe how social institutions operate with these transformation methods. These theoretical components indicate that social institutions (authority, indirect reciprocity, and traditional religion, among others) constructed an environment producing selection pressure on morality. Finally, we discuss how injunctive morality as an emotion functions in social groups.

**Key Words:** the origin of morality, institutions, dynamic interdependence theory