

# 裏切りの効果による協調の進化的安定

## Evolution of cooperation in the meta-norms game by a social vaccine

山本仁志\*<sup>1</sup> 岡田勇\*<sup>2</sup>

Hitoshi Yamamoto and Isamu Okada

\*<sup>1</sup> 立正大学 Rissho University

\*<sup>2</sup> 創価大学 Soka University

**要旨:** 集団における規範維持のモデルとして良く知られたメタ規範ゲームは、 $n$  人囚人のジレンマの拡張モデルとして、国際問題における協調問題など中央集権的でない集団においていかに規範を維持するかを検討する上で優れたモデルである。進化論的な分析によって、規範ゲームでは協調は維持されないが、メタ規範を導入することで協調が維持されることが知られている。しかし近年、メタ規範がシミュレーションの世代数に対して脆弱であるとの指摘がなされている。我々は、様々なシミュレーション条件におけるメタ規範の成立条件を精査し、脆弱性のメカニズムを探る。更に我々は、ジレンマ状態にある集団における協調の頑健にするための「社会的ワクチン」を提案し、その効果を検討した。社会的ワクチンとは、集団の中にごく少数の常に裏切り行為をとるエージェントが存在することで、集団全体の規範を高く維持することができる効果をいう。社会的ワクチンを導入することで、協調は広範なパラメータ空間で頑健に維持されることがわかった。

**キーワード:** 社会的ワクチン、メタ規範、協調の進化、Agent-Based Simulation

**Abstract:** Norms game and metanorms game are well known models for maintaining order in a group. As an extension of the  $n$ -person prisoner's dilemma game, the norms game introduces the behavioral principle of non-cooperation in group participants. We argue that Axelrod's findings have limits, we go one step further by attempting to extract sufficient conditions for making cooperation stable. We have discovered that cooperation can be robustly maintained by introducing into the group a small number of agents who are always behaving in a non-cooperative manner. We call this the "social vaccine" effect.

**Keywords:** Social Vaccine, Metanorms, Evolution of Cooperation, Agent-Based Simulation

### 1. はじめに

Axelrod[1]の規範ゲーム・メタ規範ゲームは、集団における秩序の維持に対するよく知られたモデル化である。規範ゲームは、 $n$  人囚人のジレンマを拡張したものであり「非協調者を罰する」という行動原理を集団の参加者に導入している。しかし、この行動原理を導入するだけでは規範は達成されず、非協調が支配的な戦略になることが示されている。そこで Axelrod は、「非協調者を罰しなかったものを罰する」というメタ規範を導入した。この導入によって、集団内で協調が維持されることがシミュレーションによって示された。Deguchi[2] はレプリケータダイナミクスを用いてメタ規範を分析して協調の安定性を支持している。また Heckathorn[5] や Horne and Cutlip[6] は、心理学的な実験を行いメタ規範が存在することを示している。

しかし、Axelrod の枠組みにはいくつかの根強い批判が存在する。Yamashita et. al. [10] や Galan et. al. [4] では、メタ規範として集団全員の相互監視を敷いているモデル化は、認知限界から集団数の上限をもたらし、また相互監視というシステムは非現実的な強い制約であることを主張している。そのため、メタ規範ゲームを部分集団に拡張する研究(Prietula and Conway [9])や、スモールワールドネットワーク上での相互監視に限定した研究での研究(Newth, [7])が提案されている。

また、Axelrod の知見はごく限られたパラメータ空間でしか成立していないことも指摘されている。織田[8]では、初期の懲罰確率によってはメタ規範においても協調が成立しないと述べている。さらに、Galan and Izquierdo[3] は Axelrod [1] をコンピュータシミュレーションと数理解析で精査した結果、メタ規範が協調を安定させるパラメータ空間は限定的であることを指摘した。

彼らによると、協調の維持には大きく3つの制約が存在している。一つは世代制約である。高々100世代では協調が維持されるが、100万世代までシミュレーションを実施するとメタ規範でも協調が崩壊する。メタ規範崩壊のメカニズムとしては後述のように、二重の壁による協調の不安定性と、崩壊社会の方が安定的であることが原因である。これは、そもそもメタ規範が規範社会より協調安定的とはいえ、何らかの契機によって崩壊すると不可逆性を持つという危険性があることを指摘するもので、メタ規範の非頑健性を示す重要な知見である。

二つ目に、メタ規範における懲罰損失と懲罰コストを0.1倍にすることでも崩壊することを指摘している。メタ懲罰損失が減ることはメタ規範的な裏切りにメリットを上昇させるので、直観的にも支持されるものであろう。

最後に突然変異確率を0.1倍にすることでも協調が崩壊されることを指摘している。これはメタ規範の維持には、ある程度の突然変異が必要であることを予想させる。しかし、彼らは協調が崩壊することを示しているが、では協調の安定を実現するためにはどうすればよいのかに対して答えていない。我々は、Axelrod の限界を指摘する意味では彼らと同じ主張をしているが、もう一步踏み込んで、協調を安定させる十分条件を抽出しようとするものである。

我々は新たに、集団に少数の常に非協調行動をとるエージェントを導入することで頑健に協調が維持されることを発見した。我々はこの効果を社会的ワクチン効果と呼ぶ。

## 2. モデル

ここで、Axelrod の規範ゲーム・メタ規範ゲームを整理し、社会的ワクチンの導入をおこなう。

規範ゲームは  $n$  人囚人のジレンマゲームの拡張としてとらえることができる。 $N$  人のエージェントで構成される集団を考える。エージェント  $i$  は裏切るか協調するか二つの行為を選択することができる。裏切る確率を  $B_i$  (大胆さ) で表現する。 $i$  が裏切ると、 $i$  は  $T(=3)$  の利得を得ることができる。残りの  $(N-1)$  人エージェントは  $H(=-1)$  の利得を得る。 $i$  が協調すれば、すべてのエージェントの利得は  $0$  である。

ここまでは  $n$  人囚人のジレンマゲームであるが、規範ゲームでは、このあと  $(N-1)$  人のエージェントに懲罰のチャンスがある。エージェント  $j$  は確率  $s$  で  $i$  の裏切りを発見する。発見しなかった場合、なにも起こらず  $i, j$  いずれの利得も変化しない。 $j$  が  $i$  の裏切りを発見した場合、 $j$  は自身の持つ復讐度  $V_j$  の確率によって  $i$  を罰する。 $j$  が  $i$  を罰した場合、 $i$  は  $P(=-9)$  の利得を  $j$  は  $E(=-2)$  の利得を得る。罰しなかった場合、 $i, j$  の利得に変化はない。

ここまでは規範ゲームである。メタ規範とは、エージェント  $j$  が  $i$  の裏切りを発見し、更に  $j$  が  $i$  を罰しなかったことをエージェント  $k$  が発見したときに  $k$  が  $j$  を罰するという構造を導入したものであるこのとき、 $k$  が  $j$  を罰すれば、 $j$  は  $P(=-9)$  の利得、 $k$  は  $E(=-2)$  の利得を得る。

社会的ワクチンは、常に  $(B, V)=(1, 0)$  の戦略をとるエージェントをさす。集団内で少数の社会的ワクチンエージェントが存在することが、規範の維持にどのような効果をもたらすかを検討することが本稿の目的である。また、社会的ワクチンエージェントの戦略を  $(B, V)=(0, 1)$  とすることで、純然たる規範維持のための監視者が存在することの効果を観察できるなど、メタ規範に社会的ワクチンを適用することの応用範囲は広い。

## 3. メタ規範の脆弱性

ここでは、集団の規模  $N$  を 20 から 100 まで変化させ、更に世代数も 100 から 100,000 まで変化させた実験を行った。実験は 50 回の試行を行い、最終世代の大胆さ  $B$  の平均値をプロットしている (図 1)。復讐度  $V$  は、 $B$  と強い負の相関があり  $B$  の値を観察することで  $V$  の挙動もわかるため、本論文では大胆さ  $(B)$  のみを観察する。

規範ゲームにおいて、世代数が増えるとほぼ裏切り支配になることがわかる。これはもともと規範ゲームが懲罰に対するフリーライドを容易にしている構造のため、長期的には裏切りが優位になるためである。

また、集団の規模が大きくなると協調が維持されやすくなっている。これは、集団の規模が大きくなることで、裏切りが発見される回数も増え、裏切ることで得られる利得より裏切りを発見されて集中的に罰せられることで裏切りが不利になるためである。ただしこれは大規模な集団での完全な相互監視を意味しており、現実的には非常に厳しい制約であると考えられる。

メタ規範ゲームにおいては、規範ゲームと比較すると協調が支配的となるが、集団規模が小さい範囲において世代数を長くすることで規範が崩壊していることが分かる。集団規模が少し大きくなると協調が安定していることは、先に述べたように完全な相互監視がメタ規範のレベルで徹底しているために、非常に厳しい監視社会となり協調が維持されている。

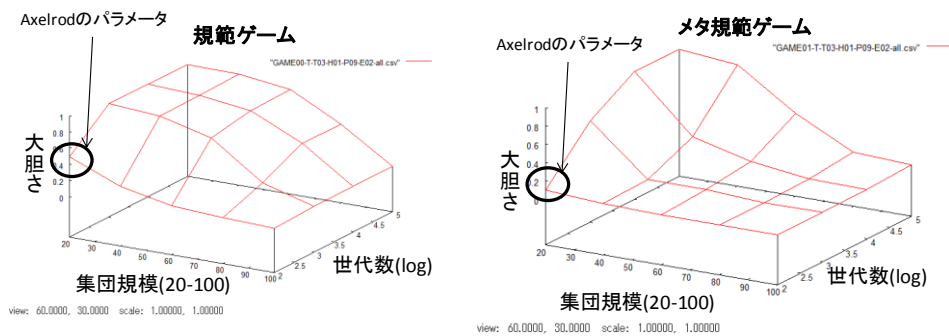


図 1：集団規模と世代を変化させた GA モデル

メタ規範において崩壊が早期に発生していることが観察される。これは、Axelrod モデルでは新たな形質は突然変異によるものしかなかったが、GA を採用したことによって、交叉によって新たな形質が生じる可能性が高くなり、崩壊の起こるタイミングが頻発するためと考えられる。また、いずれの場合も集団規模が大きくなることは協調に優位に働くことが分かる。

続いて、本研究では集団規模は基本である  $N=20$  に固定し、突然変異率を変化させた実験を行う。図 2 は、規範ゲーム・メタ規範において、突然変異率、世代数を変化させたものである。突然変異率が 0% と 5% 以上のときに協調が成立していることが観察される。

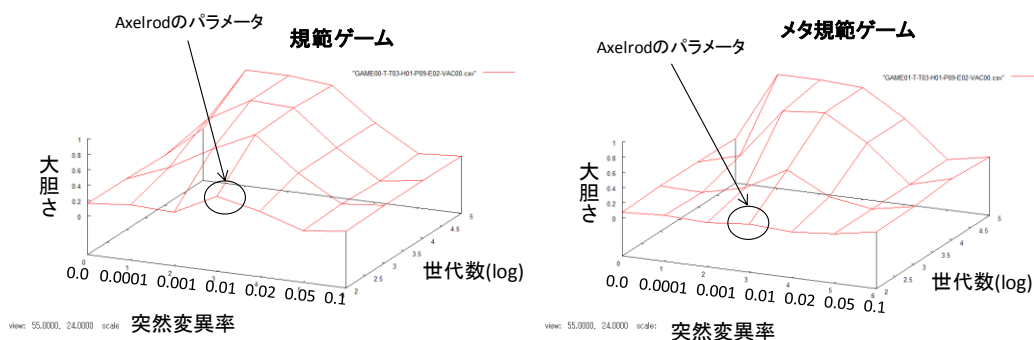


図 2：突然変異率と世代を変化させた GA モデル

突然変異率が 5% のときに協調が達成されているように観察されるが、時系列の推移を観察してみると非常にランダムな世界であることが分かる。突然変異率が 5% のときの時系列変化を観察する。図 3 はそれぞれ規範ゲーム、メタ規範ゲームの突然変異率 5% における時系列推移である。メタ規範において大胆さ(B)は 0.15 から 0.5 までの間を不規則に推移する。

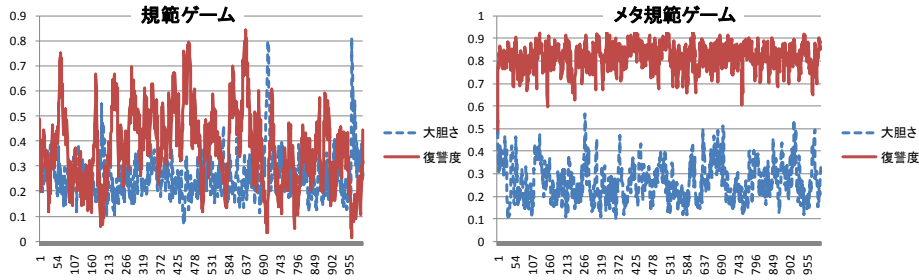


図 3：突然変異率が5%のときのGAモデルの挙動

突然変異率が0%のときは、いったん戦略が安定すると変化する要因がないため最終的な結果は安定的となる。規範ゲームにおいても、初期に高まった復讐度の高さが協調の安定を促し、早い世代で戦略が一樣になるため裏切りの侵入がなく安定的な結果となる。ただし、一樣になる直前に交叉で裏切りが発生して裏切りが支配的になることも低い頻度である。その結果大胆さの平均値は0.2前後で安定している。

#### 4. 社会的ワクチンの導入

前節までの結果をまとめると以下のように整理できる。規範ゲームでは集団規模を増やさない限り裏切りに収束する。メタ規範ゲーム (Agent=20) も、超長期では裏切りになる。突然変異率が0%や5%では平均的な裏切率は抑制されたが、0%は進化ゲームとしては不自然であるし、5%の状態はランダム性が大きい。

我々は、頑健に協調を維持するための方策として「社会的ワクチン」の導入を提案する。ワクチンとは一般的に弱毒化した病原体を接種することで抗体をつくり病原体への感染を予防することをいう。社会的ワクチンとは、集団の中にごく少数の常に裏切り行為をとるエージェントが存在することで、集団全体の規範を高く維持することができる効果をいう。

図4は集団に集団規模の5%のワクチンエージェント (常に裏切るエージェント) を導入した際の大胆さの平均値である。集団規模と世代数を変化させた。5%の理由は、最小規模の20に導入できる最低数1が5%であることによる。

規範ゲームにおいては、裏切りが支配的となっているが、メタ規範ゲームにおいては世代数を変化させても協調が安定的に維持されることが分かった。

メタ規範が崩壊する理由は、協調達成時に復讐度の低いエージェントが侵入してきても、裏切り行為がないため、それを発見できないため復讐度の低いエージェントが広まってしまうためであった。しかし、ワクチンエージェントがいることで、復讐度の低いエージェントは発見されやすくなり、集団全体の復讐度が下がることを防ぐことができる。

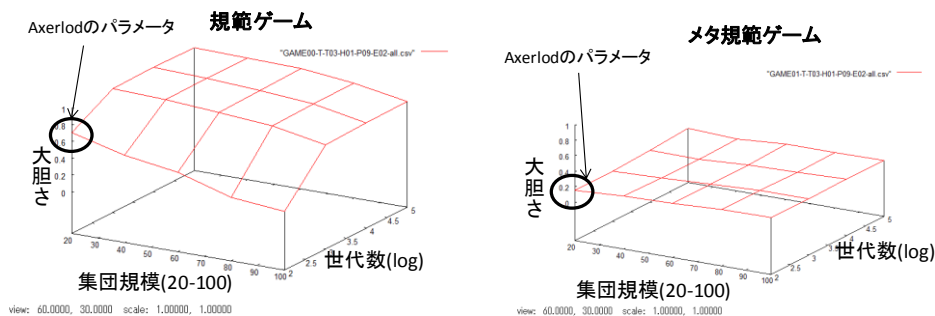


図 4：社会的ワクチンを導入したことによるメタ規範の安定 (集団規模の変化)

続いて、突然変異率と世代を変化させた実験の結果を図5に示す。メタ規範ゲームにおいて、突然変異率に対しても頑健に協調が維持されていることが分かった。

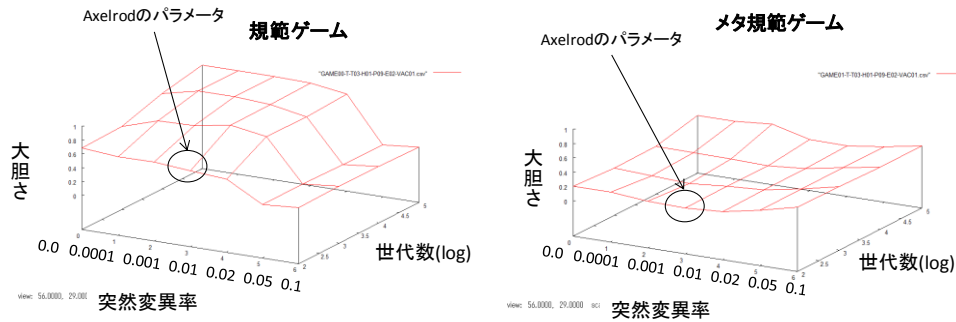


図 5 : 社会的ワクチンを導入したことによるメタ規範の安定 (突然変異率)

## 5. 社会的ワクチンの効果

前節では、社会的ワクチンを集団規模の 5% として実験をおこなった。本節では、社会的ワクチンの適切な濃度が存在するのかどうかをさまざまなパラメータを用いて観察する。

### 5.1. 社会的ワクチンの濃度

ここではワクチンエージェントの数を変化させることで社会的ワクチンが効果的に機能する条件を精査する。図 6 左は、Agent=20、突然変異=0.01 に固定した結果であり、図 6 右は、Agent=100、突然変異=0.01 に固定した結果である。ワクチンエージェントは 0-20 で実行した。図 6 左では、社会的ワクチンなしでは長期効果による裏切り出現している。これはメタ規範の基本モデルの挙動と同一である。また、ワクチン数最少の 5% (1 エージェント) で最も高いレベルの協調を達成している。ワクチンはごく少数で効果的に協調に貢献することを示唆している。一方で図 6 右ではワクチンエージェントが存在しないときに最も協調が達成されているが、これはそもそも N=100 では非常に強い相互監視がしかかれており、ワクチンが存在しなくても協調の達成が可能となるからである。

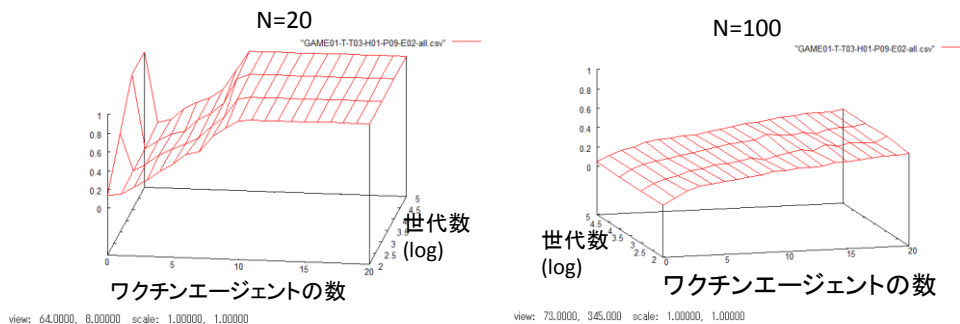


図 6 : ワクチンエージェントの数による協調の安定

### 5.2. 小集団における社会的ワクチンの効果

本節では、ここまで基本的な集団規模としてきた N=20 より小さい集団でメタ規範ゲームがどのような挙動を示すのか、集団として安定した挙動をもたらすのは最低どの程度の規模が必要なのかを分析する。また社会的ワクチンの効果はどの程度の集団規模で出現するののかも併せて明らかにする。

図 7 は、左列に社会的ワクチンの存在しないメタ規範の挙動を時系列で観察したものを示し、右列に社会的ワクチンを集団中に 1 エージェントのみ投入した場合の時系列の挙動を示している。

N=5 では、ワクチンの有無に関わらずランダムな挙動となり集団的に安定した挙動は観察できない。N=10 で両者ともランダム性の強いふるまいであるが、社会的ワクチンがない場合は、協調的な社会と非協調的な社会が周期的にあらわれる。ワクチンがある場合では大胆さが低く復讐度の高い状況が安定的に観察できる。N=15 では振る舞いは安定的となる。ワクチンがない場合、基本的なメタ規範の振る舞い

と同様、協調的な社会が達成されたのちに協調が崩壊し、裏切りが支配的となり安定する。ワクチンがある場合には、協調が安定的に維持された。

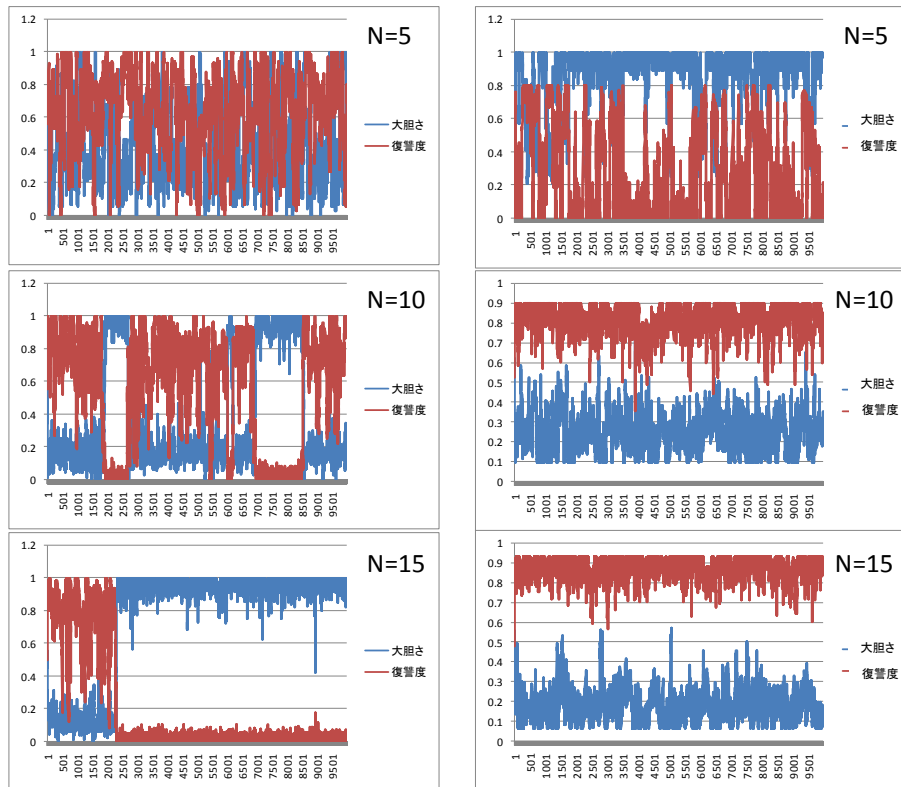


図7：小集団における社会的ワクチンの効果

## 6. まとめ

我々は、メタ規範が協調を安定させる条件を探るためにシミュレーション実験をおこなった。その結果多くのパラメータ環境において、協調が崩壊することを示した。また、我々は従来協調が崩壊するといわれているパラメータ空間においても協調が頑健に維持されるための方策として「社会的ワクチン」の導入を提案した。社会的ワクチンを導入することでメタ規範における超長期および様々な突然変率における安定達成を可能とした。

## 文 献

- 1) Axelrod, R.M., An Evolutionary Approach to Norms, *American Political Science Review*, 80(4), 1095-1111, 1986.
- 2) Deguchi, H., Norm Game and Indirect Regulation of Multi Agent Society, *Proc. of Computational Social and Organizational Science Conference*, 2000.
- 3) Galan, J.M. and L.R. Izquierdo, Appearances Can Be Deceiving: Lessons Learned Re-Implementing Axelrod's Evolutionary Approach to Norms', *Journal of Artificial Societies and Social Simulation* 8(3), <http://jasss.soc.surrey.ac.uk/8/3/2.html>, 2005.
- 4) Galan, J.M., M. Latek, M. Tsvetovat, and S. Rizi, Axelrod's Metanorm Games on Complex Networks, *Proc. of Agent 2007 Conference*, 271-280, 2007.
- 5) Heckathorn, D.D., Collective Sanctions and Compliance Norms: A Formal Theory of Group-Mediated Social Control, *American Sociological Review*, 55(3), 366-384, 1990.
- 6) Home, C., and A. Cutlip, Sanctioning Costs and Norm Enforcement: An Experimental Test, *Rationality and Society* 14(285), DOI: 10.1177/1043463102014003002, 2002
- 7) Newth, D., Altruistic Punishment, Social Structure and the Enforcement of Social Norms, in R. Khosla et al. (Eds.): *KES 2005, LNAI 3683*, 806-812, 2005.
- 8) 織田輝哉, 秩序問題への進化論的アプローチ-メタ規範ゲームの展開, *理論と方法*, 5(1), 81-99, 1990.
- 9) Prietula, M.J. and D. Conway, The evolution of metanorms: quis custodiet ipsos custodes?, *Computational Mathematical Organization Theory*, DOI 10.1007/s10588-009-9056-4, 2009.
- 10) Yamashita, T., H. Kawamura, M. Yamamoto, and A. Ohuchi, Effects of Proportion of Metanorm Players on Establishment of Norm, *Fourth International Conference on Computational Intelligence and Multimedia Applications (ICCIMA01)*, 2001.