

# 間接互惠状況における評判ダイナミクスの分析：理論・実験・シミュレーションの統合

## An analysis of reputation dynamics on indirect reciprocity: theory, experiment and simulation

山本仁志<sup>1\*</sup> 岡田勇<sup>2</sup> 鈴木貴久<sup>3</sup> 梅谷凌平<sup>1</sup>

Hitoshi Yamamoto<sup>1</sup>, Isamu Okada<sup>2</sup>, Takahisa Suzuki<sup>3</sup> and Ryohei Umetani<sup>1</sup>

<sup>1</sup> 立正大学<sup>1</sup> Rissho University

<sup>2</sup> 創価大学<sup>2</sup> Soka University

<sup>3</sup> 津田塾大学<sup>3</sup> Tsuda University

**Abstract:** Indirect reciprocity is one of the keystones to explain the evolution of human cooperation. For sustainable cooperation by indirect reciprocity, a mechanism needs to be in place to ensure that benefits are returned to the individuals who cooperate. Indirect reciprocity relies on social norms that distinguish the good from the bad, and the results of the assessment are called reputation. Extensive studies on indirect reciprocity in the framework of evolutionary game theory have discovered various types of norms or reputation dynamics that promote the evolution of cooperation. However, few is known about norms that human adopts in their real life. Most previous studies are based on the assumptions that people use binary assessment rule, which is limited to good or bad. Here we reveal reputation dynamics that people adopt in real life situations and analyze robustness and evolutionary stability using mathematical analysis and agent based simulation. Our results show that people adopt ternary values when they assess others. Notably, justified defection, defined as defection against a bad reputation, is justified only when a person of a good reputation did it.

## 1 はじめに

協力とは行為者がコストを支払って他者の利得を増加させることである。協力は社会的に望ましい事だが、協力することにはコストがかかるため適応的な観点からは無条件に協力し続ける個体はやがて淘汰されてしまうはずである。しかし人々は他者の行動を評価し評判として共有することで、協力行動が非協力者によって搾取されずに安定的に存続できる社会を築いてきた。評判を通じた協力は間接互惠性として理解することができる [Sugden 86, Alexander 87]。間接互惠性とは、誰かを助けた行為者にたいして、協力した相手からではなく別の他者から協力されることをさす。情けは人のためにならずという言葉が示すように社会で広く受け入れられている。間接互惠性が機能するためには協力した個人へと利益が還元される仕組みが存在しなければならない。人間が善悪を判別する機能（規範）は、そのような仕組みの一つとして働いていると考えられ、

どのような規範が進化し定着するのかについて数学、物理学、経済学、心理学など幅広い分野で研究がなされてきた。

間接的互惠性の最も単純な規範はイメージスコアリング (IS) [Nowak 98a, Nowak 98b] であり、これは一次情報のみを参照する。しかし、イメージスコアリングは、エラーやフリーライダーの侵入に弱く進化的に安定していない [Sigmund 10]。そのため、間接的な互惠性の中で協力を持続させるためには、二次情報が必要となる。2次情報を利用する規範としては、主にスターンジャッジング (SJ) [Kandori 92, Pacheco 06, Santos 18]、シンプルスタンディング (SS) [Sugden 86, Milinski 01, Panchanathan 03]、シャニング (SH) [Takahashi 2006] などがある。2次情報については、レシピエントの評判が悪い場合の評価ルールを考えることが重要である。なぜなら、評判の良いレシピエントへの行動の評価は容易だからである。当然のことながら、良いレシピエントへの協力と良いレシピエントへの非協力は、それぞれ善と悪と評価されるべきである。

間接的互惠性に関する実験は数多く行われている

\*連絡先：立正大学  
東京都 品川区 大崎 4-2-16  
hitoshi@ris.ac.jp

が [Wedekind 00, Milinski 01, Milinski 02, Bolton 05, Ule 09, Yoeli 13]、人間が実際にどのように二次情報を利用してはいるのかは不明である。一方で、Milinski ら [Milinski 01] の研究では、人は複雑な情報を処理することを嫌うため、一次情報のみを利用する傾向があると報告されている。一方、他の研究では [Swakman 16, Okada 18]、情報を得ることにコストがかかっても、人は 1 次情報と 2 次情報の両方で意思決定を行うことが示されている。さらに、悪意のある相手への行動を無視する評価ルールが、進化的に安定した協力の規範として提案されている。また、実験的には悪い人の評価は良い人の評価よりも不確実性が高いことが示されている [Siegel 18]。

上述したように、評判の悪いレシピエントに対する行動の評価ルールについては、理論からの予測と実験結果との間に乖離がある。特に、現実社会における悪いレシピエントに対する行動の評価ルールの分布は明らかになっていない。理論的には協調性を維持するために必要とされる二次情報が、人間が他者を評価する際にどのように利用されているかを分析することで、人間社会における罰の省略と正当化された逸脱の機能が明らかになる。

本研究では、間接的互惠関係における人の評価ルールの分布を明らかにするために実験をデザインする。実験では、シナリオに二人の人物が登場する。1 人は、レシピエントに協力するか協力しないかを問うドナーであった。もう一人は、ドナーに協力を依頼するレシピエントであった。参加者には、ドナーがレシピエントに協力するか協力しないかを選択したときの様子を観察して、ドナーを Good か Bad かの評価をしてもらった。レシピエントには Good か Bad のどちらかのラベルが貼られている。参加者は、ドナーの行動とレシピエントのラベルを組み合わせた 4 つのシーンを評価した。本研究の目的は、現実社会で人々がどのような規範を採用しているかを調べることであり、経済的な実験ではなくシナリオベースの実験を採用した。

これまでの理論研究では、ある集団には一つの規範が浸透しており、結果としてすべての人が同じ規範を使うという仮定が採用されてきたが、この仮定は明らかに現実を単純化しすぎている。近年では、この前提に基づき、母集団の中に様々な規範を共存させたモデルを開発することで、この前提に取り組む研究もある [Yamamoto 17, Uchida 18]。しかし、理論研究で提案された規範が実際に人々に採用されているのであれば、オブザーバーによるドナーの評価は広く分布しているはずである。例えば、評判の悪いレシピエントに協力しない「正当化された裏切り」の評価に良い評価と悪い評価が混在している場合、正当化された裏切りの評価のばらつきは大きくなるはずである。本研究では、人々の規範がこれまでの理論研究で提示された規範の

範囲内にとどまっているのか、あるいは、人々が理論研究で想定された規範とは異なる規範を採用しているのかを探る。さらには、観察された規範が理論的に安定であるのかどうかを数理モデルとエージェント・ベースド・シミュレーションによって検証する。

## 2 実験 1

我々は間接互惠状況において人々が採用する評価ルールの分布を探るための 3 つの実験を設計した。

参加者はドナーの行動（協力/非協力）とレシピエントの評判（Good/Bad）の組合せた 4 通りの場面を評価した (図 1)。順序効果を検討するために、参加者は CtoG、DtoG の順序でシナリオを評価し、続いて CtoB と DtoB の順序を入れ替えた 2 つのケースを用意した。参加者はレストランで働いている状況を想定したシナリオを提示された。具体的には、同僚の Alice (レシピエント) が別の同僚 Bob (ドナー) にレストランの夜勤を交代して欲しいと依頼し、Bob が承諾する場面（協力）と拒否する場面（非協力）を想定する。更に Alice の評判をコントロールした。良い評判のケースは“ Alice は働きぶりも真面目で、他の人が夜勤に来れないときにはいつも快く代わりを引受けてくれます。そのため Alice はあなたを含め職場の人達からとても好かれています”である。悪い評判のケースは” Alice は働きぶりも不真面目で、何だかんだと理由をつけてはよく夜勤を休んでいます。プライベートな理由で休むことも多く、Alice あなたを含め職場の人達から良く思われていません。”である。(実際には Alice と Bob の名前は日本の一般的な苗字にランダムに変換した)。シナリオを読んだ後に参加者は Bob の行動を“ Bob は信頼できる人物か”、“ Bob に好感を持つか”、“ Bob に親しみを感じるか”の 3 項目を 5 件法で評価した。

また、特定のシナリオの依存性を排除して結果を一般化するために異なるシナリオによって実験をおこなった。追加したシナリオは以下の 2 つである。シナリオ 2 : シナリオをレストラン従業員から近所の住民に変更し、レシピエントがドナーに個人的な相談に乗ってほしいと要請する。シナリオ 3 : 日常生活における評価だけでなく、経済実験における行動も検討するために Giving ゲームの場面を設定した。実験は実際の金銭インセンティブを想定せずに、参加者は、Giving ゲームにおけるドナーの行動を（協力か非協力）シナリオ 1,2 と同様に評価した。レシピエントの評判は、彼が過去 5 ラウンドでドナーとしてプレイしたときに何回寄付を選んだかによってコントロールした (5 回 (Good)、0 回 (Bad))。匿名状況での経済ゲームであることを参加者に示すために、プレイヤー A・プレイヤー B という表記で説明した。

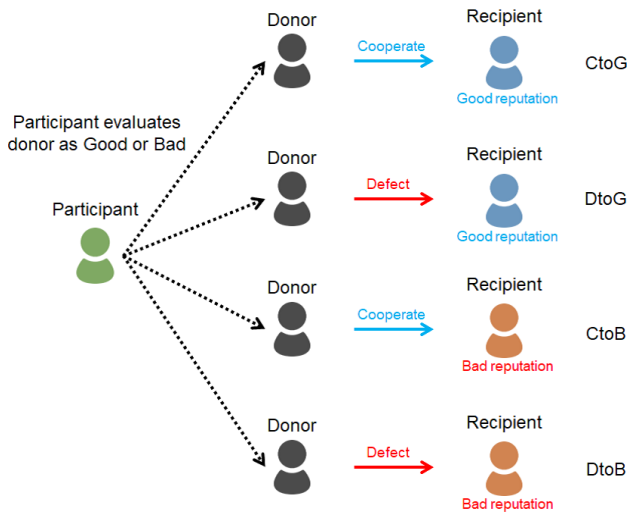


図 1: 実験の概要。参加者は 4 つのシーンを評価する。

### 3 結果 1

第一に我々は各場面におけるドナーの評価の分布を観察する (図 2)。パネル A は良いレシピエントへの協力、パネル B は良いレシピエントへの非協力、パネル C は悪いレシピエントへの協力、パネル D は悪いレシピエントへの非協力をそれぞれ表す。パネル A と B は自明な結果であり、CtoG は Good と評価され DtoG は Bad と評価されている。パネル C からわかるように CtoB も Good と評価されていることがわかる。これらのことから多くの参加者は IS か SS 規範のいずれかを採用していると考えられる。一方でパネル D は非常に興味深い結果を示している。評価の分布が中央に集中しており、参加者は DtoB を Good とも Bad とも評価していない。IS と SS が混在しているのであれば 2 つのピークが表れるはずなので参加者の多くは IS も SS も採用していないといえる。つまり、DtoB については参加者は判断を保留し評判をアップデートしないという規範を採用している [Yamamoto 20]。

### 4 実験 2

正当化される裏切りは正当化も不当化もされないことが実験 1 より明らかとなった。では、正当化される裏切りが正当化されることはないのでしょうか？ また正当化されるのであればどのような条件で正当化されるのでしょうか？ これまで間接互惠性の研究で扱われてきた情報は、主にドナーの行動とレシピエントの評判であった。前者は 1 次情報と呼ばれ、後者は 2 次情報と呼ばれてきた。他方で評価されるドナー自身の評判は 3 次情報と呼ばれいくつかの研究では扱われてい

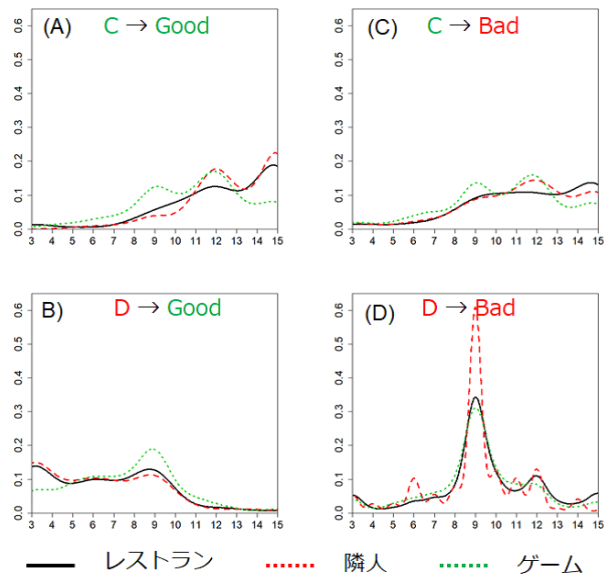


図 2: 各シナリオのドナーに対する評価の分布。 [Yamamoto 20] より改訂

るが [Ohtsuki 06]、特に実験研究では扱われることがほとんどなかった。我々はドナーの評判を導入することとする。また、実験 1 の結果のように Good/Bad の 2 値ではなく中立という評価も考える必要が生じるために中立 (Neutral) という評価も導入する。加えて、実験 1 ではドナーについては評判に関する情報を与えていない。この状況も併せて中立として、1 次、2 次、3 次情報に加えて、Good/Neutral/Bad の 3 値を想定したシナリオ実験をおこなう。シナリオは実験 1 のレストランシナリオを用いた。シーンはドナーの評判 3 種、レシピエントの評判 3 種、ドナーの行動 2 種の計 18 シーンであり実験はすべて被験者間計画によって実施した。

### 5 結果 2

実験の結果、18 シーンはいずれも図 2 のパネル A, B, D のように右にゆがんだ分布 (良いという評価)、左にゆがんだ分布 (悪いという評価)、中央にピークのある分布 (中立という評価) の 3 つに分類できることが分かった。表 1 にそれぞれのシーンの分類の結果を示す。結果からわかることは、ドナーの評判によらず協力行動によって評価は 1 段階良い方向に上がることが分かった。すでに良い評判を持っているドナーの協力行動は当然良いと評価されるが、今回の実験では 5 件法による評価のため、評価がより良くなったのかよいのままとどまっているのかは判別できない。評判を 3 値で表

現すべきなのか、より多段階へと拡張すべきなのかは今後の課題である。

他方で非協力行動に関しては、基本的には1段階悪い方向へさがることがわかった。しかし悪いレシピエントへの非協力だけは、現在の評判がアップデートされずそのままとどまっている。すなわち、中立なドナーの正当化される裏切りは中立のままであり、良い評判のドナーの正当化される裏切りは良いと評価された。つまり、正当化される裏切りは良い評判のドナーが行ったときにのみ正当化されている。

表 1: ドナーの評判あり 3 値実験の結果

Who	Act	To whom	Evaluation
G	C	G	G
		N	G
		B	G
N		G	G
		N	G
		B	G
B		G	N
		N	N
		B	N
G	D	G	N
		N	N
		B	G
N		G	B
		N	B
		B	N
B		G	B
		N	B
		B	B

## 6 シミュレーション

実験 2 で示された評判ダイナミクスは安定的に協力を維持することができるのであろうか？本節ではエージェントシミュレーションによってこの評判ダイナミクスの性質を分析する。評判ダイナミクスを分析するためには行動ルールを検討する必要がある。直観的に Good には協力、Bad には非協力という行動ルールを想定することが自然である。これまでの Good/Bad の 2 値を採用した理論モデルでは上記の行動ルールを採用し、様々な評価ルール（規範）の無条件協力、無条件非協力に対する安定性を検証してきた。しかし、3 値モデルを検討する差には、中立的な評判を持つ相手に対して協力するのか否かの決定をしなくてはならない。そのため行動ルールをあらかじめ与えることは不

自然であり、行動ルールの真価についても検討する必要がある。考える行動ルールは以下の 8 通りで表現できる (表 2)。本モデルの行動ルールを採用することで、無条件協力 (ALL-C)、無条件非協力 (ALL-D) もモデルに含めることができるので網羅的な検証が可能となる。注目すべきは、G-DISC, S-DISC で表現される 2 つのディスクリミネータの安定性である。これらの行動ルールは Bad には非協力、Good には協力という向社会的ルールを保持しているが、Neutral に対する行動ルールが異なる。これらの行動ルールが支配的となり協力行動が安定するかどうかを検証する。

表 2: 3 値モデルにおける行動ルール

Reputation	Bad	Neutral	Good
AllD	D	D	D
S-DISC	D	D	C
	D	C	D
G-DISC	D	C	C
	C	D	D
AllC	C	D	C
	C	C	D
	C	C	C

### 6.1 モデル

集団は  $N = 1000$  個体のエージェントからなる。エージェントは上述した 8 つの行動戦略のいずれかを初期値としてランダムに割り振られる。また各エージェントの評判もランダムに割り振られる。集団からランダムに選ばれた 2 個体のエージェントがドナーとレシピエントに割り振られギビングゲームをおこなう。ドナーの行動は協力か非協力であり、協力を選択するとコスト  $c$  を失うがレシピエントは  $b (b > c > 0)$  の利得を得る。非協力の場合は両者の利得は変化しない。他の  $N - 2$  個体のエージェントはこのゲームを観察し表 1 の評判ダイナミクスを用いてドナーの評判をアップデートする。

1 世代の中で全てのエージェントが 100 回ドナーとなる。つまり 1 世代の間に  $1,000 \times 100 = 100,000$  回のギビングゲームがおこなわれる。1 世代の最後にエージェントは行動戦略を進化させる。エージェント  $i$  は集団からランダムに選んだエージェント  $j$  とそれぞれの利得  $U_i, U_j$  を比較し確率  $P_{ij}$  でエージェント  $j$  の戦略を模倣する。模倣しない場合には自身の戦略はアップデートされない。シミュレーションは 10,000 世代実行し、乱数種の異なる 50 回の試行をおこなった。

$$P_{ij} = \frac{1}{1 + \exp\left(\frac{U_i - U_j}{S}\right)} \quad (1)$$

## 6.2 シミュレーション結果

図3に示すように実験2から得られた評判ダイナミクスの下では協力が安定的に維持できることがわかった。また、行動ルールとしてはG-DISCとS-DISCが混在した状況で安定することがわかった。実社会で観察される規範はシミュレーションにおいても協力を進化させる。

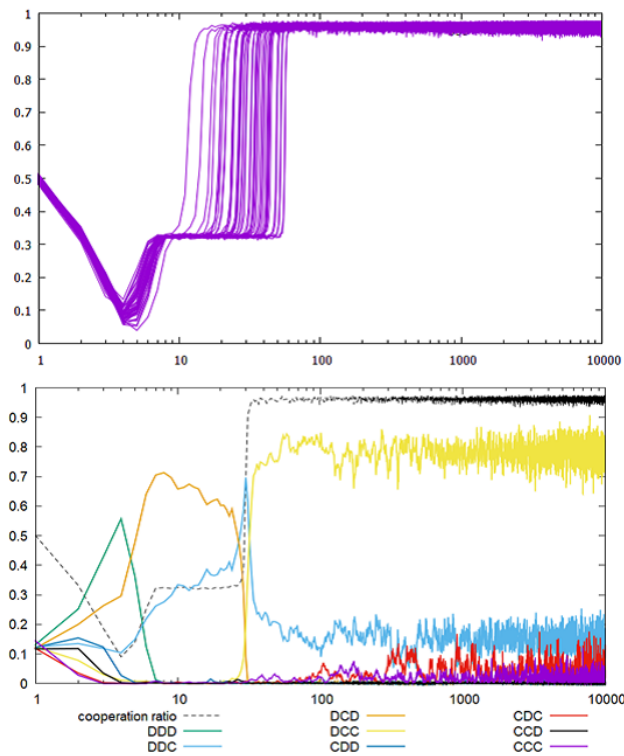


図3: シミュレーション結果。上パネル: 50 試行の協力率の推移。下パネル: 1 試行の中の各行動戦略の推移。

## 7 まとめ

本研究では、人間の協力の進化を説明する鍵となる間接的互惠性において人々が採用する評判ダイナミクスを明らかにした。進化ゲーム理論の枠組みでの間接互惠性に関する広範な研究により、協力の進化を促進する様々なタイプの規範や評判ダイナミクスが発見されてきた。しかし、人間が実際の生活の中で採用している規範についてはほとんど知られていない。これまでの研究の多くは、人間が良いか悪いかという二項対立的な評価ルールを用いていることを前提にしていた。

本研究では、人間が実生活の中で採用している評判のダイナミクスを明らかにし、被験者実験とエージェントベースのシミュレーションを用いて、頑健性と進化的安定性を解析した。その結果、人は他人を評価する際に3値を用いたな価値観を採用していることがわかった。特に、正当化された裏切りは良い評判の人がおこなった場合にのみ正当化されることがわかった。

## 謝辞

本研究の一部は科研費17H02044, 18H03498, 19H02376, 19K21570の助成を受けている。

## 参考文献

- [Alexander 87] Alexander, R.: *The Biology of Moral Systems*, Aldine de Gruyter, New York (1987)
- [Bolton 05] Bolton, G. E., Katok, E., and Ockenfels, A.: Cooperation among strangers with limited information about reputation, *J. Public Econ.*, Vol. 89, No. 8, pp. 1457–1468 (2005)
- [Kandori 92] Kandori, M.: Social norms and community enforcement, *Rev. Econ. Stud.*, Vol. 59, No. 1, pp. 63–80 (1992)
- [Milinski 01] Milinski, M., Semmann, D., Bakker, T. C. M., and Krambeck, H.-J.: Cooperation through indirect reciprocity: image scoring or standing strategy?, *Proc. R. Soc. B Biol. Sci.*, Vol. 268, No. 1484, pp. 2495–2501 (2001)
- [Milinski 02] Milinski, M., Semmann, D., and Krambeck, H.-J.: Reputation helps solve the 'tragedy of the commons', *Nature*, Vol. 415, No. 6870, pp. 424–426 (2002)
- [Nowak 98a] Nowak, M. A. and Sigmund, K.: Evolution of indirect reciprocity by image scoring, *Nature*, Vol. 393, No. June, pp. 573–577 (1998)
- [Nowak 98b] Nowak, M. A. and Sigmund, K.: The dynamics of indirect reciprocity., *J. Theor. Biol.*, Vol. 194, No. 4, pp. 561–574 (1998)
- [Ohtsuki 06] Ohtsuki, H. and Iwasa, Y.: The leading eight: social norms that can maintain cooperation by indirect reciprocity., *Journal of theoretical biology*, Vol. 239, No. 4, pp. 435–44 (2006)

- [Okada 18] Okada, I., Yamamoto, H., Sato, Y., Uchida, S., and Sasaki, T.: Experimental evidence of selective inattention in reputation-based cooperation, *Scientific Reports*, Vol. 8, No. 1, p. 14813 (2018)
- [Pacheco 06] Pacheco, J. M., Santos, F. C., and Chalub, F. A. C.: Stern-judging: A simple, successful norm which promotes cooperation under indirect reciprocity, *PLoS Comput. Biol.*, Vol. 2, No. 12, p. e178 (2006)
- [Panchanathan 03] Panchanathan, K.: A tale of two defectors: the importance of standing for evolution of indirect reciprocity, *J. Theor. Biol.*, Vol. 224, No. 1, pp. 115–126 (2003)
- [Santos 18] Santos, F. P., Santos, F. C., and Pacheco, J. M.: Social norm complexity and past reputations in the evolution of cooperation, *Nature*, Vol. 555, No. 7695, pp. 242–245 (2018)
- [Siegel 18] Siegel, J. Z., Mathys, C., Rutledge, R. B., and Crockett, M. J.: Beliefs about bad people are volatile, *Nature Human Behaviour*, Vol. 2, No. 10, pp. 750–756 (2018)
- [Sigmund 10] Sigmund, K.: *The calculus of selfishness*, Vol. 6, Princeton University Press (2010)
- [Sugden 86] Sugden, R.: *The Economics of Rights, Cooperation and Welfare*, Oxford: Basil Blackwell (1986)
- [Swakman 16] Swakman, V., Molleman, L., Ule, A., and Egas, M.: Reputation-based cooperation: Empirical evidence for behavioral strategies, *Evol. Hum. Behav.*, Vol. 37, No. 3, pp. 230–235 (2016)
- [Uchida 18] Uchida, S., Yamamoto, H., Okada, I., and Sasaki, T.: A Theoretical Approach to Norm Ecosystems : Two Adaptive Architectures of Indirect Reciprocity Show Different Paths to the Evolution of Cooperation, *Front. Phys.*, Vol. 6, No. February, p. 14 (2018)
- [Ule 09] Ule, A., Schram, A., Riedl, A., and Cason, T. N.: Indirect punishment and generosity toward strangers, *Science*, Vol. 326, No. 5960, pp. 1701–1704 (2009)
- [Wedekind 00] Wedekind, C. and Milinski, M.: Cooperation through image scoring in humans, *Science*, Vol. 288, No. 5467, pp. 850–852 (2000)
- [Yamamoto 17] Yamamoto, H., Okada, I., Uchida, S., and Sasaki, T.: A norm knockout method on indirect reciprocity to reveal indispensable norms, *Sci. Rep.*, Vol. 7, No. March, p. 44146 (2017)
- [Yamamoto 20] Yamamoto, H., Suzuki, T., and Umetani, R.: Justified defection is neither justified nor unjustified in indirect reciprocity, *PLOS ONE*, Vol. 15, No. 6, p. e0235137 (2020)
- [Yoeli 13] Yoeli, E., Hoffman, M., Rand, D. G., and Nowak, M. A.: Powering up with indirect reciprocity in a large-scale field experiment, *Proceedings of the National Academy of Sciences*, Vol. 110, No. Supplement 2, pp. 10424–10429 (2013)