

# Self-interest, trust, and cooperation

Kirbi Joe<sup>1</sup>, Ronan Arthur<sup>2</sup>

1 University of California, Irvine; 2 Stanford University

June 21, 2018

## 1 Introduction

In the age of digital technology, information is more easily accessible to the general public than ever before. Given the wealth of information that is available to the average person, it is important that people evaluate the quality of the information they receive and the source from which they receive it. In an environment where information can simultaneously be shared and consumed, entities participating in such an exchange may benefit from acting cooperatively with one another—that is, outputting “good” information in order to incentivize others to reciprocate and also output “good” information back to them.

However, cooperative or altruistic human behavior intuitively flies in the face of one of the primary drivers of evolution: selfishness. Yet, cooperation in groups to complete large tasks or support one another is often essential to the fitness of the constituent members of that group [1]. Key to motivating cooperation in human communities is trust, the belief that one’s partners in an endeavor will fulfill their part of the bargain. Consequently, individuals are more likely to exchange information with each other if they trust one another to not take advantage of them or exploit their vulnerabilities [2]. Key to trust is the track record of an individual’s decisions when faced with the option to expend effort to support the group or selfishly refuse. An important question, then, is what happens if the track record of an individual is not completely known or delayed? Can trusting and cooperative communities still emerge?

To answer these questions, we use a game theoretic approach to model a system of agents engaging in a cooperative endeavor with delayed information about the trustworthiness of the other members of the group.

## 2 Model specifications

Agents are initialized on a Bernoulli random graph  $\sim G(N, p)$  (e.g. Fig. 1), where  $N$  represents the number of agents and  $p$  represents the probability of a tie being formed between two agents (i.e. network density).

We assume that each of these agents are endowed with some kind of information and they want to share this information with each of their neighbors. Agents have the choice to either send high (H) quality information or low (L) quality information. In this framework, it is favorable to receive high quality information, so those who send high quality information reliably are deemed more

“trustworthy.” At each time step, agents simultaneously and privately choose to give either H or L quality information to each of their neighbors. Choices are not revealed until  $\Delta$  rounds later.

Agent  $i$ 's *trust* of agent  $j$  is determined based on their history of past play. Agents will be more trusting of others if those others frequently choose to give them H quality information. Trust is defined to be the probability with which agent  $i$  believes that agent  $j$  will give him H quality information in the current round; that is, trust is equal to the proportion of previous rounds played (up until  $\Delta$  rounds ago) in which  $j$  gave  $i$  H quality information. Therefore, agent  $i$ 's trust,  $T$ , of agent  $j$  at time  $t$  is given by:

$$T_{i,j}(t) = \frac{H(t - \Delta)}{H(t - \Delta) + L(t - \Delta)} \quad (2.1)$$

where  $t$  is the round number,  $t - \Delta$  is the round number counting the time delay (the agent's most recently revealed round),  $H : \{t_0, t_1, \dots, t_n\} \rightarrow \mathbb{N}$  is the function that returns the number of rounds up to  $t_k$  in which agent  $i$  received H quality information from agent  $j$ , and  $L : \{t_0, t_1, \dots, t_n\} \rightarrow \mathbb{N}$  is a similar function which counts the number of Ls received by  $i$  from  $j$ .

Table 1: Payoff structure for agents

		Alter	
		H	L
Ego	H	a, a	d, b
	L	b, d	c, c

The incentives for agents to choose which type of information to send is modeled using the payoff structure detailed in Table 1. This payoff structure takes after the popular game the Stag Hunt, which imposes the restriction  $a > b \geq c > d$  on the individual payoffs. The Stag Hunt was theoretically a good fit for our model because it captures the notions of social cooperation as well as aspects of trust that causes agents to choose to engage in a riskier or more costly activity in order to gain a higher reward.

Because of the time delay which prevents players from seeing the choices of their partners until  $\Delta$  rounds later, agents will make a choice to play H randomly according to the probability  $P_H$  for the first  $\Delta$  many rounds of the simulation. After the first  $\Delta$  rounds have passed, players will choose to send H or L to all their partners by electing the higher of the expected values for each scenario (as in equations 2.2 and 2.3), based on the summed expected payoffs of each interaction and the proportional trustworthiness of a partner to play H.

$$E_i[H] = \sum_{j \in N(i)} [T_{i,j} * U(H, H) + (1 - T_{i,j}) * U(H, L)] \quad (2.2)$$

$$E_i[L] = \sum_{j \in N(i)} [T_{i,j} * U(L, H) + (1 - T_{i,j}) * U(L, L)] \quad (2.3)$$

where  $T_{i,j}$  is as defined in equation 2.1,  $U(x, y)$  is the utility function of the ego given the choice he made,  $x$ , and the choice his alter made,  $y$ . For example,  $U(H, H) = 3$ .

### 3 Results

A series of parameters were systematically altered in various simulation runs in order to evaluate the variation in convergence behavior relative to these alterations. The parameters of interest are:

- $N$ , the number of agents
- $p$ , the network density
- $\Delta$ , the time delay
- $a, b, c, d$ , the individual payoffs as shown in Table 1
- $P_H$ , the probability that an agent chooses H in the first  $\Delta$  rounds

Each simulation was run for a total of 25 rounds. All simulations were found to converge to an equilibrium in which all agents were playing the same strategy (either all playing H or all playing L), provided they were not isolated in the network. For every set of parameters tested, 500 repetitions of each simulation with that particular parameter set was run.

In order to systematically test the effects of parameter changes on model convergence, a set of default parameters were chosen to remain fixed while one of the variables was altered. The default parameters were:

$$N = 25, p = 0.5, \Delta = 1, (a, b, c, d) = (3, 2, 2, 0), P_H = 0.65$$

A network with example parameters  $N$  and  $p$  is represented in Figure 1.

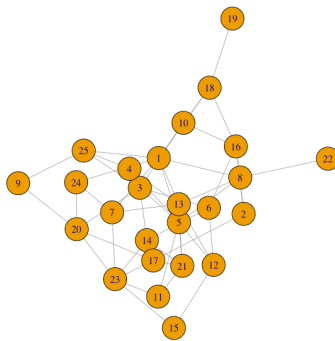


Figure 1: A Bernoulli random graph generated with  $N=15$  and  $p=0.25$

For 5000 repetitions of simulations using the default parameters, the proportion of runs that ended with an equilibrium of all players playing H equaled 0.20. Figure 2 shows the number of agents who play the strategy H over time for 10 different simulation runs (of which 2 ended in the H equilibrium and 8 ended in the L equilibrium, consistent with the results from the high volume simulations with the same parameters). Varying the default parameters can drive the system from a convergence to all playing L to all playing H. An example of this is illustrated in Figure 2, where  $\Delta$  is 4.

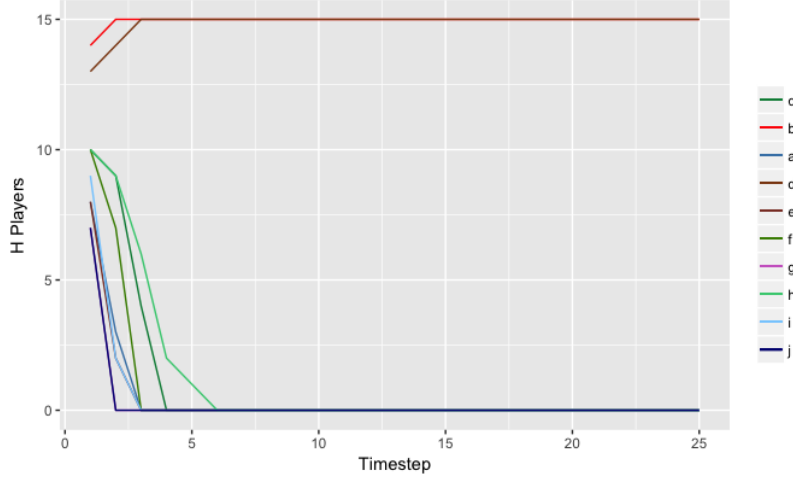


Figure 2: Set of 10 simulations using the default parameter conditions

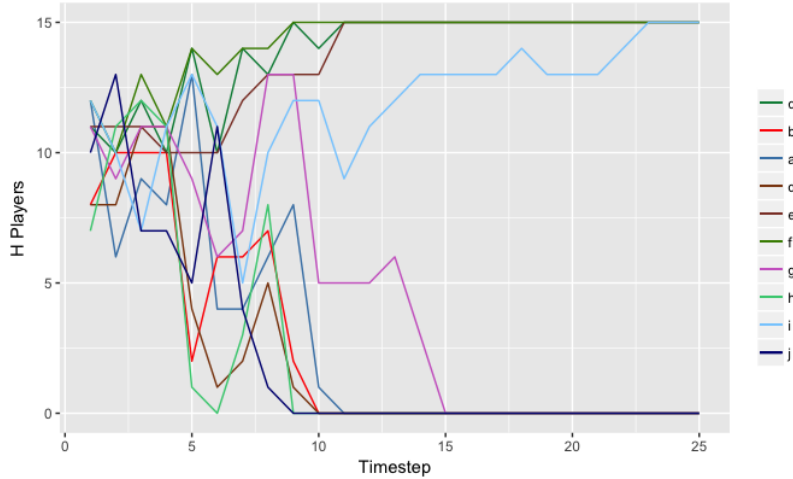


Figure 3: Set of 10 simulations using the default parameter conditions with  $\Delta = 4$

The following tables detail the parameter sets used to run various simulations as well as the proportion of simulations using those parameters which resulted in an equilibrium of H, represented by the final column in each table titled “Runs w/ eq=H”. Key findings in these simulations include a divergence of systems, ones that converge to networks of agents that all play H and networks of agents that all play L. Several of the variables exhibit threshold dynamics, capable of pushing the system from likely to end up all H to likely to end up all L. We chose default parameters near this threshold to demonstrate how the manipulation of some of these parameters had non-linear effects on the percentage of all H simulations.

First, consider Table 2, which examines the effects of varying the number of agents in the network. This appears to have little to no impact on the proportion of runs that end with all H. There is a small increase for an intermediate number of players, potentially due to the structure of such networks that may allow for its quicker alignment and trust formation. Table 3 looks at the variation of the probability of edge formation. This, again, affects network structure, increasing the density

Table 2: Results for Simulations with Varied  $N$

$N$	$p$	$\Delta$	$a, b, c, d$	$P_H$	Runs w/ eq=H
<b>10</b>	0.65	1	3,2,2,0	0.65	0.199
<b>15</b>	0.65	1	3,2,2,0	0.65	0.200
<b>20</b>	0.65	1	3,2,2,0	0.65	0.238
<b>25</b>	0.65	1	3,2,2,0	0.65	0.233
<b>100</b>	0.65	1	3,2,2,0	0.65	0.180

as  $p$  increases. This increase is not linear, rather its growth is greater in the middle of this scale and slower near 0.25 and 1.

Table 3: Results for Simulations with Varied  $p$

$N$	$p$	$\Delta$	$a, b, c, d$	$P_H$	Runs w/ eq=H
15	<b>0.25</b>	1	3,2,2,0	0.65	0.091
15	<b>0.35</b>	1	3,2,2,0	0.65	0.184
15	<b>0.45</b>	1	3,2,2,0	0.65	0.200
15	<b>0.55</b>	1	3,2,2,0	0.65	0.220
15	<b>0.65</b>	1	3,2,2,0	0.65	0.26
15	<b>0.75</b>	1	3,2,2,0	0.65	0.315
15	<b>0.85</b>	1	3,2,2,0	0.65	0.335
15	<b>0.95</b>	1	3,2,2,0	0.65	0.345
15	<b>1</b>	1	3,2,2,0	0.65	0.380

Table 4: Results for Simulations with Varied  $\Delta$

$N$	$p$	$\Delta$	$a, b, c, d$	$P_H$	Runs w/ eq=H
15	0.5	<b>1</b>	3,2,2,0	0.65	0.220
15	0.5	<b>2</b>	3,2,2,0	0.65	0.193
15	0.5	<b>3</b>	3,2,2,0	0.65	0.140
15	0.5	<b>4</b>	3,2,2,0	0.65	0.127
15	0.5	<b>5</b>	3,2,2,0	0.65	0.127
15	0.5	<b>6</b>	3,2,2,0	0.65	0.107
15	0.5	<b>7</b>	3,2,2,0	0.65	0.140

Table 4 looks at variations of time-delay delta. With an increase in delta, the number of runs that go to H decrease but increase again at 7. This initial decrease is likely due to an increase in noise during the first delta rounds of randomized assignments of L and H. As agents sum their expectations, this may then lead them to a higher likelihood to conclude they must play L.

Table 5 shows the results for simulations with varying  $b$  and  $c$ , which we held as equal, the payoffs associated with playing L. This trend was very clearly delineated: as  $b$  and  $c$  are increased the run proportions of all H decrease, and the rate of this decrease increases until it approaches zero. This can be qualitatively explained somewhat simply - as payoffs for playing L increase, the expected value of playing L across all transactions also increases to the point where, at 2.5, it is very rare that a run can find a trusting equilibrium where all players play H.

Varying  $P_H$  has equally clear consequences as shown in Table 6.  $P_H$  is the probability that agents will play H in each of their first delta time-steps, and since the process of establishing trust is very dependent on these initial delta time-steps, the system will lock into L or H from there. Therefore, it was expected that as we increased  $P_H$  we saw a significant increase in the H proportion.

Table 5: Results for Simulations with Varied  $b, c$  ( $b = c$ )

$N$	$p$	$\Delta$	$a, b, c, d$	$P_H$	Runs w/ eq=H
15	0.5	1	<b>3,1.5,1.5,0</b>	0.65	0.816
15	0.5	1	<b>3,1.6,1.6,0</b>	0.65	0.793
15	0.5	1	<b>3,1.65,1.65,0</b>	0.65	0.718
15	0.5	1	<b>3,1.7,1.7,0</b>	0.65	0.656
15	0.5	1	<b>3,1.75,1.75,0</b>	0.65	0.622
15	0.5	1	<b>3,1.8,1.8,0</b>	0.65	0.456
15	0.5	1	<b>3,2.2,2.2,0</b>	0.65	0.0.12
15	0.5	1	<b>3,2.5,2.5,0</b>	0.65	0.0059

Table 6: Results for Simulations with Varied  $P_H$

$N$	$p$	$\Delta$	$a, b, c, d$	$P_H$	Runs w/ eq=H
15	0.5	1	3,2,2,0	<b>0.55</b>	0.055
15	0.5	1	3,2,2,0	<b>0.6</b>	0.158
15	0.5	1	3,2,2,0	<b>0.7</b>	0.368
15	0.5	1	3,2,2,0	<b>0.75</b>	0.518
15	0.5	1	3,2,2,0	<b>0.8</b>	0.683
15	0.5	1	3,2,2,0	<b>0.85</b>	0.855
15	0.5	1	3,2,2,0	<b>0.9</b>	0.955
15	0.5	1	3,2,2,0	<b>0.95</b>	1

## 4 Discussion

These results suggest a number of qualitative takeaways. However, we discuss these with caution as the model and its results have limitations, chiefly that they depend on structural mathematical assumptions. We assume the community is connected in a random graph network, though we know that network structures in social communities are not random [3], which may affect trust and cooperation. We assume a plausible payoff structure, but one which may not always represent the way social interaction and exchange function. And finally, we assume perfect knowledge and rational agent behavior, that all agents know about the delayed history of all other agents and that they rationalize their subsequent actions by taking the higher of two expected value functions.

Despite these limitations, the model suggests promising results about how trust might be built in social exchange. First, we have seen that trust, or the establishment of reputation, can lead to cooperative environments, where all agents trust one another and act in good faith. Second, cooperation emerges from self-centered behavior, so generosity or self-sacrifice is not necessary to discourage agents from free-riding on one another. This is to be expected given that the agents are payoff-maximizing and (H,H) and (L,L) are the pure Nash equilibria of the game. They may

instead cooperate to the benefit of all. With the introduction of a relative point structure, where agents compare their payoffs and compete with one another, we would likely see the environment converge to all agents playing L.

Trust is instrumental to the future of the network community. In high trust scenarios, agents feel safe playing H, while in low trust scenarios, agents will defensively also play L. What then affects trust? In this model, trust is built on the history of what agents play. Since this is assigned randomly for the first delta time-steps, the history of play, and thus the future too, may be path-dependent and rely on how this randomization process works. Other variables and parameters had significant effect on the number of runs that converged to all H, but  $P_H$ , the probability that players would choose H in each of the first delta time-steps is the most influential (see Tables 2-6).

Future work will include incorporating structured networks such as small world or scale-free networks to see how network position might affect the agent's path to equilibrium. It has been suggested that the presence of a single leader or a small influential group could affect the dynamics of a group engaging in the formation of social contracts, as modeled by the Stag Hunt [4]. Some have argued that influential actors in a network occupy positions called "structural holes" [5], positions which bridge between two groups of actors, while others have argued that leaders emerge from central positions in the network [6]. Therefore, it may be interesting to consider various node-level indices, like degree centrality and betweenness centrality, to see if there is any correlation between such measures and the time in which they adopt the equilibrium strategy relative to their neighbors.

Another potential improvement that could be implemented to this model would be changing the behavior of agents in the first  $\Delta$  many rounds before they are informed of their neighbors' choices. Since the preliminary results show that changing the probability with which agents play H during these first rounds has significant effects on the likelihood of the system converging to H, altering this mechanism may provide more insight as to how initial behavior affects later actions. In all, the model and its results suggest further research into the mechanisms of establishing trust and its importance on social exchange choices is needed to better understand their dynamics, incentives, and patterns.

## References

- [1] E. Fehr and U. Fischbacher, "The nature of human altruism," *Nature*, vol. 425, no. 6960, p. 785, 2003.
- [2] J. K. Butler Jr, "Trust expectations, information sharing, climate of trust, and negotiation effectiveness and efficiency," *Group & Organization Management*, vol. 24, no. 2, pp. 217–238, 1999.
- [3] M. O. Jackson, *Social and economic networks*. Princeton university press, 2010.
- [4] B. Skyrms and R. Pemantle, "A dynamic model of social network formation," in *Adaptive networks*, pp. 231–251, Springer, 2009.
- [5] R. S. Burt, "Structural holes and good ideas," *American journal of sociology*, vol. 110, no. 2, pp. 349–399, 2004.
- [6] H. J. Leavitt, "Some effects of certain communication patterns on group performance.," *The Journal of Abnormal and Social Psychology*, vol. 46, no. 1, p. 38, 1951.