# The effects of imperfect information on social norms in networks

**Group 13**

João Carranca - 102600 || Javier Santa Maria - 99240

Masters in Computer Science and Engineering and Masters in Engineering and Data Science

**Network Science**

Instituto Superior Técnico

## Abstract

Indirect reciprocity is one of the most fundamental concepts in the study of population dynamics and cooperation dynamics. We can find the effects and marks of this mechanism all over the planet in a variety of species, including humans. Social norms are used to study rates of cooperation in populations, governing the nature of interactions between agents. In this work, we analyze the effect of imperfect information on the performance of discriminators using 4 different social norms relative to full cooperators, full defectors and inverse-discriminators, first in a well-mixed population and then in the context networks. We show that the choice of social norm yields considerably different results, with imperfect information generally significantly decreasing rates of cooperation. We also show that the the network topology has little effect on results.

## I. Introduction

*Reciprocity* can be used to describe several behavioural patterns found in nature. The idea of "You scratch my back, I scratch yours" has been present in several species for hundreds of thousands of years. Standard reciprocity as a societal model faces a logical problem [1]. A society as a whole does benefit more from being made up of cooperators than defectors, however, on an individual level, each defector can extract more benefits than each cooperator since helping someone incurs a cost. This thought process would lead to societies full of defectors and no cooperators. This doesn't, however, happen in real life. A variety of biological systems actively use cooperation in their organization [3].

Social norms allow us to craft more accurate models of interactions in populations, especially human ones, that address this collapse in cooperation [2].These social norms are derived from the idea of indirect reciprocity, which has proven to maintain mutual cooperation in populations [4] [5]. *Indirect reciprocity* makes use of the concept of reputations in individuals which adds a higher level of cognitive capacity to interactions [6]. An individual can, for example, choose to cooperate or not cooperate, based on the reputation of the individual he is interacting with. The way he might use this information on reputation will constitute a specific social norm. The degree of a norm depends on the amount of information used to decide on whether to cooperate or not cooperate, with each level of information constituting a "layer" [7].

Stern-Judging, Shunning, Simple-Standing and Image-Scoring are four of the most well studied social norms in the literature [18] [19]. Stern-Judging, Simple-Standing and Shunning are second degree social norms that, given an interaction between two individuals, one donor and one recipient, will take into account the reputations of both the donor and recipient as well as the donor´s action to calculate the donor´s new reputation [7] [18] [21] while Image-Scoring is a first degree norm [20]. Stern-Judging has proven to be the most successful social norm when it comes to maintaining high levels of cooperation across many population sizes [8], but when presented with things like assessment errors and imperfect information, it has shown low levels of resilience [22] [9] [24]. Are the other three norms more resilient? If we create a population mix with different types of agents like defectors cooperators, discriminators and inverse-discriminators for each norm, what will be the distribution of these agents at the state of equilibrium? What will the average fitness look like? How will cooperation evolve? Here we extend on previous work related to imperfect information in social norm models like [24], [**?**] and [11], testing first in a well-mixed population model and then on a variety of network topologies.

## II. Methods

### I Evolution Model

Our evolution model, which will be used to evaluate which agents dominate given certain simulation parameters, is grounded in *evolutionary game theory*. This evolution dynamic will be applied to the three types of agents we will be mixing in our populations: cooperators, defectors and discriminators, where discriminators will follow the Stern-Judging social norm. A cooperator can, for example, turn into a defector and vice-versa.

In evolutionary game theory, the success of a strategy is tied to its *fitness*—a measure of how well an individual with that strategy performs relative to others in the population. The higher the fitness, the more likely a strategy is to proliferate. The evolution of strategies in a population follows the *replicator equation*, which describes how the frequency of a given strategy changes over time based on its fitness.

The general form of the *replicator equation* is given by:

$$\dot{x}_i = x_i \left( f_i(\mathbf{x}) - \phi(\mathbf{x}) \right)$$

Where $\dot{x}_i$ is the rate of change in the frequency of strategy $i$, $x_i$ is the current frequency of strategy $i$, $f_i(\mathbf{x})$ is the fitness of individuals using strategy $i$, and $\phi(\mathbf{x})$ is the average fitness of the population. This equation implies that strategies whose fitness exceeds the population's average fitness will increase in frequency, while strategies with below-average fitness will decrease.

In addition to the natural selection process described by the replicator equation, our model also includes *mutation* based on fitness comparisons where each individual $A$ has a certain probability of imitating individual $B$ based on the difference in their fitness. That probability is given by the following equation:

$$p = \left[ 1 + e^{-\beta(f_B - f_A)} \right]^{-1}$$

Where $p$ is the probability, $f_B$ is the fitness of individual $B$ and $f_A$ is the fitness of individual $A$ and $\beta$ will have a value of 0.05.

## II  Well-Mixed population model, social norms and imperfect information

A large, well-mixed population of players is considered. From time to time, two players are chosen at random from the population and they engage in a one-shot prisoner´s dilemma [15] with the following payoff matrix:

|   | C | D |
|---|---|---|
| C | R,R | S,T |
| D | T,S | P,P |

|   | C | D |
|---|---|---|
| C | 0.4, 0.4 | −1, 0.6 |
| D | 0.6, −1 | 0.0, 0.0 |

The values are designed to reflect real-world scenarios where cooperation provides mutual benefits, but defection can be tempting. Each population has an assessment rule by which to judge the action of a donor. We assume a binary judgment: either the label 'good' or 'bad', represented by 1 and 0, is assigned to the donor. That will now be his reputation. The social norms mentioned are considered to label these actions 'good' or 'bad'. Discriminator agents make their decision to cooperate or defect based on the reputation value they have assigned to the agent they interact with, at the moment of the interaction, while cooperators and defectors always either cooperate or defect and inverse discriminators, as the name suggests, do the opposite of discriminators.

Stern-judging views those as good who, in their previous game, gave help to a good recipient or refused help to a bad recipient [7]. Image Scoring doesn´t care about the nature of the recipient and will only consider good those that that in the previous round cooperated. Shunning is more selective, only considering good those that cooperated with good recipients and Simple-Standing only considers bad those that defected with a good recipient [21]

When it comes to implementing the idea of imperfect information in this model, we have followed studies on the *Three Degrees of Influence* [10], where information has a certain probability of reaching individuals depending on their distance to the action. We have used the value of 61%, used in [10], as the experimental value of "influence" for first-degree neighbours. Only this value is needed since in a well-mixed population there are only first-degree neighbours. In the case of networks, the percentages for second and third-degree neighbours will be obtained from the same source. It should be noted that only discriminators make use of reputation to guide their decisions.

Mutations can occur in the population. After any given interaction, an individual has a chance $\mu$ of mutating into a another type of agent. In our simulation $\mu = 0.01\%$, based on analysis in [25] and our own tests.

## III  Reputation matrix

Let us the notation $r_{ij}$ to represent the reputation of player $j$ in the eyes of player $i$. The values $r_{ij} = 1$ and $r_{ij} = 0$ correspond to the situations where $i$ thinks that $j$ is good and bad, respectively. The matrix $(r_{ij})$ is called the *reputation matrix* .

The update rule is described as follows: let $u$, $v$, and $W$ denote two random individuals that interact, with $u$ deciding on an action towards $v$ , and the set of observers. The number of observers is $qN$, where $N$ is the total number of individuals in the population. For the analytical calculation, the population is assumed to be infinite.

After the game, the updated reputation matrix for discriminators is given by:

$$r'_{ij} = \begin{cases} r_{ij} & \text{if } i \notin W \text{ or } j \neq u \\ f(i)(a_u, r_{iv}) & \text{if } i \in W \text{ and } j = u \end{cases}$$

where $a_u$ is the action of $u$ toward $v$ in the game, and $r_{iv}$ is the current reputation of $v$ in the reputation matrix.

Using the abbreviations $C$ for "help" and $D$ for "refuse", the function $f(i)(a,b) \in \{0,1\}$ with $a \in \{C,D\}$ and $b \in \{0,1\}$ is the assessment, from the viewpoint of $i$, of the action of node $u$ towards the recipient $v$ . For *stern-judging*, the function $f(i)$ is defined as $f(i)(C,1) = f(i)(D,0) = 1$ , $f(i)(C,0) = f(i)(D,1) = 0$.

For simple-standing: $f(i)(C,1) = f(i)(C,0) = f(i)(D,0) = 1$ , $f(i)(D,1) = 0$.

For shunning: $f(i)(C,1) = 1$ , $f(i)(C,0) = f(i)(D,1) = f(i)(D,0) = 0$.

For image-scoring: $f(i)(C,1) = f(i)(C,0) = 1$, $f(i)(D,0) = f(i)(D,1) = 0$.

Since the action $a_u$ is probabilistically determined by $r_{uv}$, the new assessment $r'_{wu}$ in the eyes of observer $w \in W$ probabilistically depends on the assessments of recipient $v$ both in the eyes of $u$ ($r_{uv}$) and $w$ ($r_{wv}$) before the game. That is to say, the updated image matrix is probabilistically determined by the old image matrix.

## IV  Average Consensus Level for analysis of the reputation matrix

To evaluate the level of imperfect information and its propagation over time we use a metric called Average Consensus Level.

We compute the consensus level, based from the definition of the Rice index [26]. We assume that the label U does not count as consensus, and thus we define the Average Consensus Level parameter $\kappa$ as:

$$\kappa = R^{-1} \sum_{i=1}^{R} \kappa_i \quad \text{with} \quad \kappa_i = Z^{-1} \sum_{k=1}^{Z} \frac{G_k - B_k}{G_k + B_k}$$

Where $G_k$ is the number of individuals in the population that identify individual $k_i$ as Good, while $B_k$ is the number of individuals in the population that identify individual $k$ as Bad; in both cases, these labels are associated with the population status associated with the last time step of run $i$. $\kappa$ is maximal (100%) when all opinions coincide, and minimal (0%) for maximum population polarization, which happens when, for individual $i$, half the population that knows $i$ sees $i$ as Good while the other half sees $i$ as Bad.

## V  Simulation settings and step by step execution

We start by setting up our fixed parameters for the simulation. For the well-mixed model we use $N = 100$ for the nodes, $T = 250$ for the number of time steps or epochs, the number of simulation per input distribution which will be 30, the percentage associated with imperfect information, 61%, the number of interactions per epoch which will be the same as the number of nodes, $N$, the value of $\beta$ for the evolution function which is always 0.05, the mutation rate $\mu$ which is always 0.01%, and the parameters for the Prisoner´s Dilemma that have already been presented.

As input, our simulation function receives, 6 arguments: the initial distributions of all four types of agents, a boolean indicating if imperfect information is to be used or not and the social norm to be used.

We use arrays to keep track of values such as total fitness of the population per epoch, node distribution per epoch and percentage of cooperative actions per epoch. These will then be used to produce the graphs. The reputation matrix is initialized with all entries as "True". We also initialize a matrix with the name "standing" which will always have perfect information. We define a function "interaction" which

will simulate an interaction between two nodes in its entirety, deciding the payoff based on the decisions made by the two nodes, updating standing and updating reputation with the specifications already described. It will also run the nodes through mutation and evolution functions that work as described above. In each epoch (time step) $N$ interactions take place, so 100, which will guarantee that, on average, a node is having one interaction per epoch. We make sure a node can´t interact with itself.

At the end of the program we will have data representing the average of 30 simulations with each setting for the given input parameters.

## VI  Network Model

We test three types of networks with the purpose of comparing results to those achieved in a well-mixed setting. These are Scale-Free networks [13] [14] , Watts-Strogatz networks [12] and random networks.

We use graph generating functions from the python library *networkx* to setup these three type of graphs with the following parameters:

$$nx.gnp_random_graph(N, 4/N)$$

$$nx.barabasi_albert_graph(N, 2)$$

$$nx.watts_strogatz_graph(N, 4, 0.2)$$

To guarantee that all significant marks of all topologies, especially Scale-Free, are visible in the network we scale the number of nodes $N$ to 2000, keeping everything else the same.

The main changes come in the interaction model: Following the studies on the *Three Degrees of Influence* [10],once we randomly pick a first node for an interaction, the second node will be picked given the following probabilities: 61% of being a first degree neighbour, 29% of being a second degree neighbour and 10% of being a third degree neighbour.

Similarly, when updating reputations, for each of the interacting nodes, the probabilities of their neighbours receiving information about the interaction, follows the same distribution described above.

## VII  Evaluation Criteria

For the purposes of evaluating our simulation results, we shall be using four criteria: Average fitness of the population at the end of a simulation, the percentage of cooperative interactions per time-step in the simulation, the overall population mix in comparison to the starting mix and overall imperfection of information at the end of the simulation measured by the Average Consensus Level equation for the final reputation matrix.

# III. Results and Discussion

## I Well-mixed

We tested configurations with each agent starting out at 25% of the population, for all social norms and with both perfect and imperfect information. We also tested configurations with one group of agents (i.e discriminators) starting at 100%, to see how mutations affected settings where one group of agents dominates and is stable.

For perfect information, Stern-Judging (SJ) performs significantly better than every other norm, with the population at equilibrium reaching rates of 100% cooperation on average, while Simple-Standing (SS) doesn´t go beyond 40%, Image-Scoring (IS) stabilizes at less than 25% and Shunning(SH) quickly approaches 5%, where it remains for the remainder of the simulation. Fitness, at equilibrium, follows a trend of increase for all norms but IS. SS at first experiences a steady decrease, but once cooperation rates stabilize, the slope of the graph becomes positive and remains that way. For shunning, the increase in fitness can be justified by the fact that, at equilibrium, a stable amount of inverse discriminators remain, which mostly cooperate with each other and defect with discriminators, allowing for a slight increase in the overall fitness.
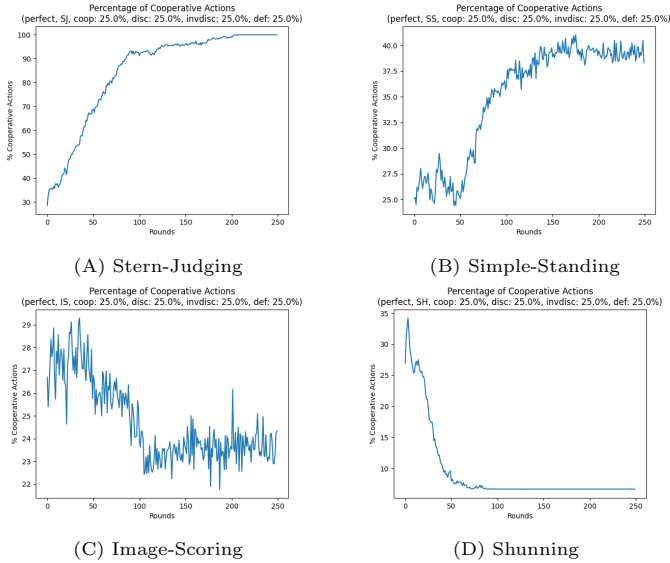


Figure 1: Well-mixed population: Evolution of cooperation rates for all social norms under perfect information with each agent constituting 25% of the population at the start

For the population distribution, we find that, in all cases, either discriminators or inverse discriminators come to dominate the population, with defectors and cooperators holding very minor stakes at equilibrium. This can be explained by the fact that discriminators and inverse discriminators are more complex in their decision making than cooperators and defectors. This theory is backed up by the fact that the population distribution for IS, the least sophisticated of the social norms analyzed, is the only one where all four types of agents are still present at equilibrium. As for some of the specific distributions at equilibrium, for SJ, for example, inverse discriminators, hold a majority stake of the population while for SH, discriminators crush everyone else very quickly. In the case of SJ, as the population evolves, inverse discrim-

inators may consistently find stable cooperation within their subset, even if it appears as a non-traditional form of cooperation. Their interactions may not target defectors in the same way discriminators do, allowing them to avoid direct conflict while fostering cooperation within their in-group.

For SH, discriminators not only punish defectors directly but also defect with individuals who interact with defectors, creating a cascading effect. This means that even inverse discriminators, who might engage with defectors, are heavily penalized by being shunned by discriminators.
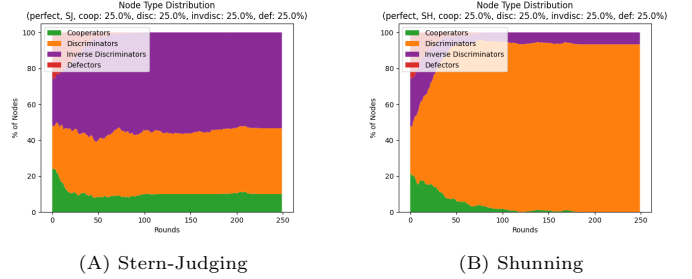


Figure 2: Well-mixed population: Population at equilibrium under perfect information with each agent constituting 25% of the population at the start for Stern-Judging and Shunning.

For imperfect information, the scenario changes significantly, with SJ going from the best performing norm on all metrics to the second worst performing one, only behind SH. Cooperation rates collapse quickly and don´t recover, despite defectors being rooted out and the population distribution at equilibrium being virtually the same (with the exception of there being no minor stake of cooperators). Cooperation in the SS population, on the other hand, does not suffer significantly with the addition of imperfect information while in SH, rates remain virtually unchanged and IS actually shows improvements.

IS is a relatively simple social norm: agents judge others based on their reputation or history of cooperation, but without considering the context of those actions. This lack of nuance means that under perfect information, agents might be more likely to judge others harshly for a single defection, even if that defection occurred in a context that might have been rational (e.g., refusing to help a defector). With imperfect information, agents no longer have access to complete and precise reputational data. As a result, they may make less extreme judgments and are less likely to severely penalize a single defection. This can lead to more lenient interactions, where agents are more willing to cooperate despite imperfect or uncertain reputational data. This can, paradoxically, stabilize cooperation in the population by reducing harsh punishments for occasional defection.

The vulnerabilities of SJ under imperfect information can also be simply explained. Despite the fact that we are not using assessment errors like in [9], with our model only stating that a part of the population does not receive information about any given interaction, SJ still heavily struggles. Let´s consider a situation where we have a majority of SJ discriminators and just a few defectors. With perfect information this won´t be a problem, all defectors will be found out after the first round and signaled as so to all discriminators. From there on out, discriminators will always defect with defectors and that will not be a problem because all

other discriminators are in agreement about the reputation of these "bad" individuals. Now let´s consider what happens when some discriminators don´t get updated on everything that happens. With the same setting as before, now some discriminators immediately signal the defectors, but others don´t. For these others, some defectors will still be "good guys" after the first round. A discriminator now interacts with a defector that they know to be "bad", and, according to SJ, they defect, but from the point of view of a discriminator that did not observe that first round interaction, an individual just defected on a "good guy" making that individual bad. The opposite can also be seen: individuals that have signaled defectors, when they see a discriminator cooperate with who they still think is a "good guy", will perceive an act of cooperation towards a "bad guy", making said discriminator also bad. This "misunderstanding" keeps on compiling, breaking down cooperation chains that don´t recover, even after those initial defectors are gone. These results can be better understood and verified by looking at the reputation matrices for all of these settings (see below) as well as the cooperation rates over time:



(A) Stern-Judging

(B) Simple-Standing

(C) Image-Scoring

(D) Shunning

(E) Stern-Judging

(F) Simple-Standing

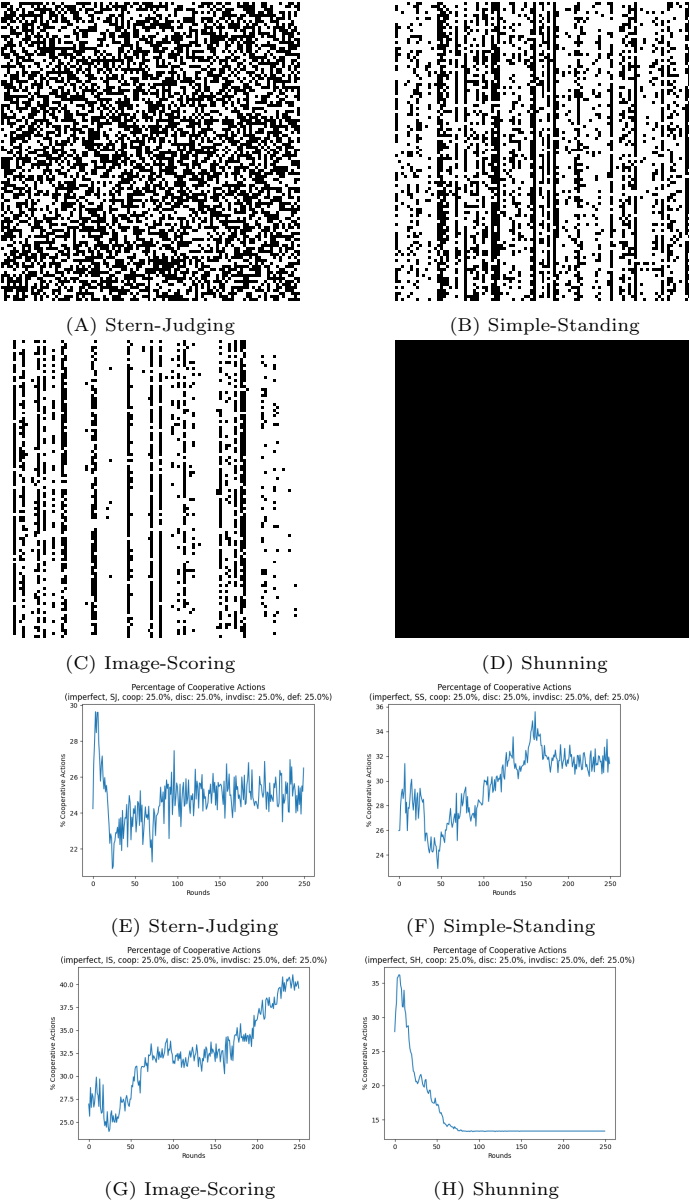(G) Image-Scoring

(H) Shunning

Figure 3: Well-mixed population: Reputation matrices [(A), (B), (C), (D)] and cooperation rates [(E), (F), (G), (H)] for all social norms under imperfect information with each agent constituting 25% of the population at the start

Using the Average Consensus Level (ACL) described in Chapter II on the matrices shown, we can verify the effect of imperfect information, with $\kappa = -0.02$ for SJ, $\kappa = 0.41$ for SS and $\kappa = 0.69$ for IS, with these values agreeing with the evolution of cooperation registered. For SH, calculating $\kappa$ would be of little interest, given the configuration of the matrix. The results shown and conclusions drawn from said results agree with those presented in [24]. Despite the differences in model, the same trends are exhibited.

Tests involving homogeneous mixes reinforce previous results and display just how vulnerable SJ is to imperfect information with defector infiltrations through mutation. In fact, only SS does well against the combination of slight defector infiltrations and imperfect information, being able to correct cooperation rates and easily root out defectors. IS and SH deal badly with small numbers of defectors, no matter the type of information, while SJ deals very well with them under perfect information and collapses completely under imperfect information. In all cases, defectors do not invade the population as they never gain any significant foothold and are quickly rooted out. Their temporary presence, however, can leave unrecoverable damage in cooperation.



(A) Stern-Judging

(B) Stern-Judging

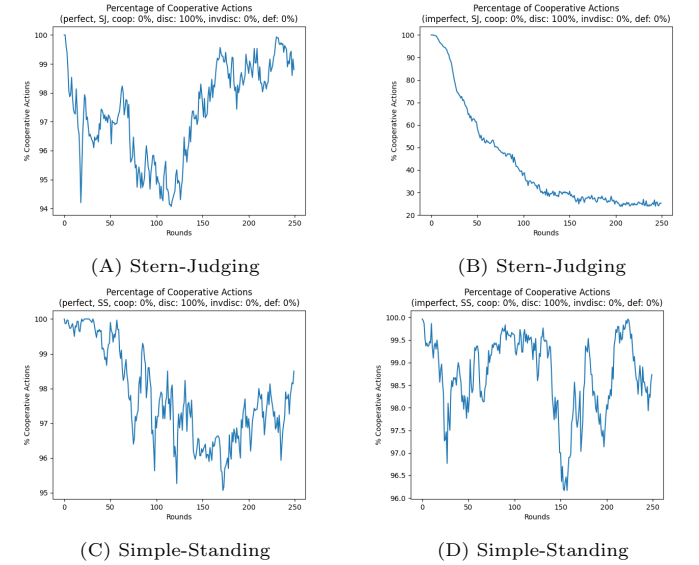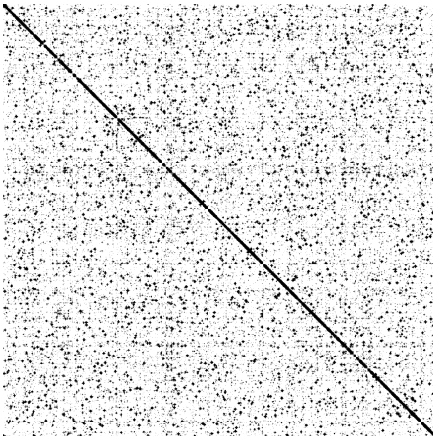(C) Simple-Standing

(D) Simple-Standing

Figure 4: Well-mixed population: Evolution of cooperation rates for Stern-Judging and Simple Standing, for a starting population of 100% discriminators with mutations, under perfect information (A) (C) and imperfect information (B) (D)
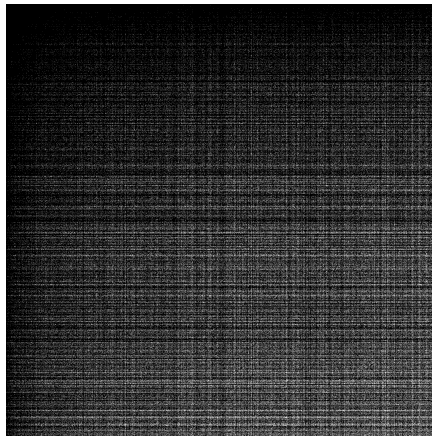
## II  Networks

For networks, we tested only configurations with each agent starting out at 25% of the population, for all social norms and all three types of networks. We find that the difference in topology, with our settings, doesn't yield any significant differences in fitness evolution, cooperation rates or node distribution. Only when we look at the reputation matrices, do we find the expected marks of the different topologies.

These patterns match the expectation for each type of network. For Watts-Strogatz, the sparse pattern shows that interactions or reputational updates are highly localized. This reflects the small-world property of the Watts-Strogatz graph, where most nodes are connected to close neighbors with occasional long-range connections due to rewiring. The
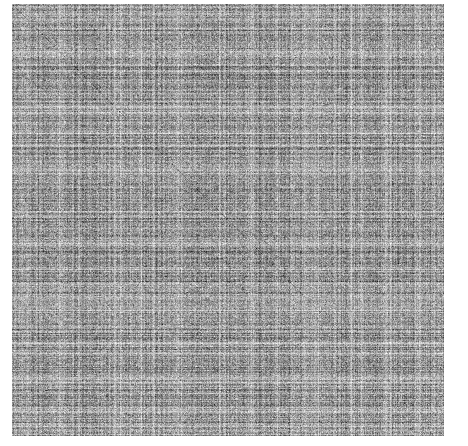
(A) Watts-Strogatz  (B) Barábasi-Albert  (C) Random Network

Figure 5: Reputation matrices for the three types of networks for the same setting: Shunning with perfect information and ach agent starting out at 25% of the population.

diagonal nature indicates that nodes mostly interact within their close-knit local groups, hence reputation updates remain sparse and clustered.

For the Barabási-Albert network, the darker, more saturated regions indicate that certain nodes (likely hubs) accumulate or interact with a disproportionately larger part of the population. This is characteristic of the scale-free nature of the Barabási-Albert graph, where highly connected hubs emerge. These hubs influence the reputational network heavily, and as a result, their reputation tends to propagate more widely, leading to the denser regions of interaction and reputation exchange.

For the random network, the matrix shows a more uniform distribution of interactions. The gray, evenly spread appearance reflects the randomness of connections in this network. Nodes interact more uniformly across the entire graph, leading to a less clustered and more homogeneous reputational structure. There is no clear clustering pattern, and the reputation updates seem to be distributed more evenly across the entire population.

As for why the remaining evaluation parameters seem to show no difference with the changes in topologies, that can be for a few reasons. All nodes are interacting frequently with their neighbors regardless of the topology and the decision to cooperate is mostly driven by reputation and behavior rather than network position. This can mean the network structure would have less impact.

On the other hand, the distribution of agents is done randomly throughout the network, when the initialization is carried out, meaning that no clusters of one agent type are likely to form from the beginning.

There are, however, differences between the results for networks and what was observed in the well-mixed population. The evolution of cooperation in SS differs in networks, although cooperation rates don´t change substantially. The shape of the cooperation evolution graphs is quite different with a lot less noise, despite the fact that the number of simulations per setting is the same for both networks and the well-mixed population. For other norms, differences are minor with both the shape and the values of the graphs being very similar.
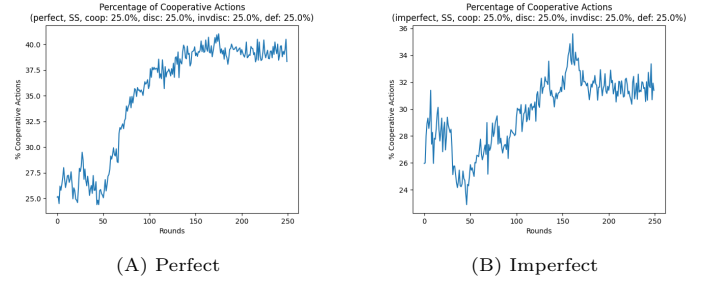


(A) Perfect  (B) Imperfect

Figure 6: Cooperation evolution for Simple-Standing for a well-mixed population
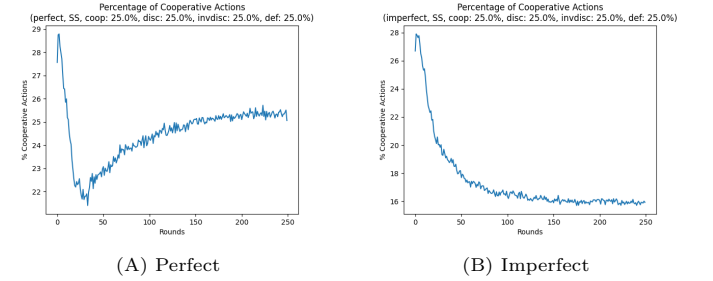


(A) Perfect  (B) Imperfect

Figure 7: Cooperation evolution for Simple-Standing for a Barabási-Albert network

Another trend that can be observed in networks is that for settings in the well-mixed simulations where one type of agent came to be very dominant at equilibrium, that process of domination occurs much more aggressively and completely, with no other agents remaining at equilibrium for the Barabási-Albert network, and only a minuscule percentage for the other two network types. In network settings, once one type of agent begins to dominate a region of the network, the spatial structure can cause a cascade where neighboring agents are "absorbed" by the dominant behavior. In a well-mixed population, even if one type dominates numerically, interactions are still more random, allowing minority types to persist at some level. In contrast, in a network, once dominance is established locally, it can become impossible for minority types to survive if they are isolated or cut off from other similar agents. This can justify the observed behaviour.
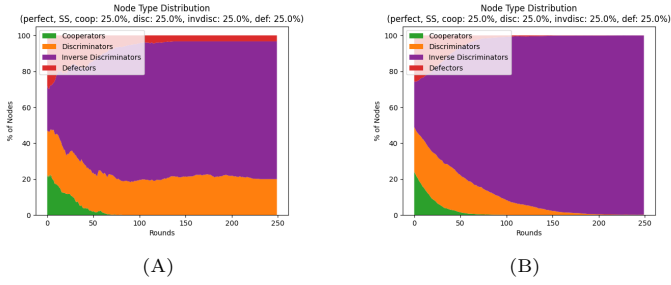
Figure 8: Population distribution in Simple-Standing with perfect information in a well-mixed setting (A) and for a Barabási-Albert graph (B)
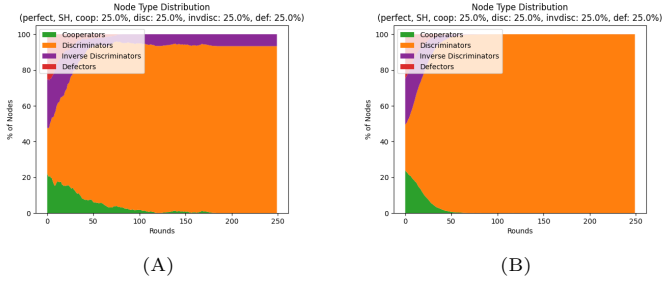


Figure 9: Population distribution in Shunning with perfect information in a well-mixed setting (A) and for a Barabási-Albert graph (B)

## IV. Conclusion

We managed to produce results that agree with those shown in previous works such as [24] and [11], despite the differences in model and the different approaches to the idea of imperfect information. Stern-Judging is shown to be very effective at fostering cooperation in populations when provided with perfect information but quickly collapses when faced with imperfect information. Simple-Standing and Shunning are barely affected by the change in information model while Image-Scoring actually improves. The extent of the frailties of Stern-Judging were further exemplified using mutations in an homogeneous discriminator population, with cooperation rates collapsing quickly and permanently.

Simple-Standing is shown to be the most consistent and best social norm analyzed, not wavering when subjected to imperfect information or infiltrations created by mutations on a stable discriminator population.

In networks, we show that topologies seem to have little effect on results but there are differences compared to the well-mixed model, with Simple-Standing cooperation rates shown to behave differently and stabilize quicker, and dominant agent types absorbing other agents quicker.

## V. Future work

In this work we built models to assess the robustness of social norms under the strain of imperfect information in a variety of different contexts. We tested 4 different social norms with imperfect information in a variety of different mixes. We also tested 4 different settings: well-mixed populations, a random network, a Scale-Free network and a Watts-Strogatz network. Other approaches could have been

taken, however. For example, we could have varied the parameters on each type of network. Dynamic networks could be studied in this context as well.

Other social norms could have be analyzed and compared, beyond the four chosen. Likewise, there are many possible approaches to imperfect information and many other models could be used. In [11], imperfect information is represented by a cost associated with acquiring information about interactions. Our model considers only the possibility of information not reaching certain individuals. Additionally, adding assessment errors can perhaps yield slightly different trends.

For networks, further degrees of neighbourhood could be considered and the probabilities of information propagation could be changed to investigate their impact. The values used in our model, associated with the Prisoner´s Dilemma, could be varied to resemble a Stag-Hunt game [17] or a Snowdrift game [16], for example.

We treated reputation as a binary variable, with individuals being perceived as either "good" (1) or "bad" (0). This was done to facilitate computation. However, reputation could be a discreet variable with a limited number of possible values ranging between 0 and 1 or even a continuous variable between 0 and 1. What trends could be derived from making such a change to the nature of reputation?

Imperfect information models, because of their versatility, can be studied in many different ways, using many different criteria with many different purposes and therefore, there is much that can still be explored.

## References

[1] Dawes, R. M. (1980). Social Dilemmas. Annual Review of Psychology, 31(1), 169–193

[2] Fehr, E., and Fischbacher, U. 2004. Social norms and human cooperation. Trends Cogn Sci 8(4):185–190.

[3] Smith JM, Szathmáry E (1995) The major transitions in evolution. Oxford: Freeman

[4] Trivers R. The evolution of reciprocal altruism. Quart Rev Biol 1971;46:35–57.

[5] Rand DG, Nowak MA. Human cooperation trends in cognitive sciences. Trends Cogn Sci 2013;17:413–25.

[6] Nowak MA, Sigmund K (2005) Evolution of indirect reciprocity. Nature 437: 1291–1298.

[7] Pacheco, J. M., Santos, F. C., Chalub, F. A. C. C. (2006). Stern-Judging: A Simple, Successful Norm Which Promotes Cooperation under Indirect Reciprocity. PLoS Computational Biology, 2(12), e178.

[8] Santos, F., Pacheco, J. Santos, F. Evolution of cooperation under indirect reciprocity and arbitrary exploration rates. Sci Rep 6, 37517 (2016).

[9] Uchida, S., Sasaki, T. (2013). Effect of assessment error and private information on stern-judging in indirect reciprocity. Chaos, Solitons Fractals, 56, 175–180.

[10] Christakis, N. A., Fowler, J. H. (2008). The Collective Dynamics of Smoking in a Large Social Network. New England Journal of Medicine, 358(21), 2249–2258.

[11] Fernando P. Santos, Jorge M. Pacheco, and Francisco C. Santos. 2018. Social norms of cooperation with costly reputation building. In Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence (AAAI'18/IAAI'18/EAAI'18). AAAI Press, Article 579, 4727–4734.

[12] Watts, D. J., Strogatz, S. H. (1998). Collective Dynamics of 'Small-World' Networks. Nature, 393(6684), 440–442.

[13] Barabási, A.-L., Albert, R. (1999). Emergence of Scaling in Random Networks. Science, 286(5439), 509–512.

[14] Albert, R., Barabási, A.-L. (2002). Statistical Mechanics of Complex Networks. Reviews of Modern Physics, 74(1), 47–97.

[15] Tucker, A. W. (1950). A Two-Person Dilemma.

[16] Nowak, M. A. (2006). Five Rules for the Evolution of Cooperation. Science, 314(5805), 1560–1563.

[17] Skyrms, B. (2004). The Stag Hunt and the Evolution of Social Structure. Cambridge University Press.

[18] Santos FP, Pacheco JM, Santos FC. 2021 The complexity of human cooperation under indirect reciprocity. Phil. Trans. R. Soc. B 376: 20200291.

[19] Yamamoto, H., Okada, I., Uchida, S., Sasaki, T. (2022). Exploring norms indispensable for both emergence and maintenance of cooperation in indirect reciprocity. , 10.

[20] Nax, H., Perc, M., Szolnoki, A., Helbing, D. (2015). Stability of cooperation under image scoring in group interactions. Scientific Reports, 5.

[21] Santos FP, Santos FC, Pacheco JM (2016) Social Norms of Cooperation in Small-Scale Societies. PLoS Comput Biol 12(1): e1004709.

[22] Hilbe, C., Schmid, L., Tkadlec, J., Chatterjee, K., Nowak, M. (2018). Indirect reciprocity with private, noisy, and incomplete information. Proceedings of the National Academy of Sciences, 115, 12241 - 12246.

[23] Olejarz, J., Ghang, W., Nowak, M. (2015). Indirect Reciprocity with Optional Interactions and Private Information. Games, 6, 438-457.

[24] Uchida, S. (2010). Effect of private information on indirect reciprocity.. Physical review. E, Statistical, nonlinear, and soft matter physics, 82 3 Pt 2, 036111 .

[25] Santos, F., Pacheco, J., Santos, F. (2016). Evolution of cooperation under indirect reciprocity and arbitrary exploration rates. Scientific Reports, 6.

[26] Rice, S. (1926). Some Applications of Statistical Method to Political Research. American Political Science Review, 20, 313 - 329.