

# Nash Equilibrium or Nash Bargaining? Choosing a Solution Concept for Multi-Agent Learning

Jeffrey L. Stimpson & Michael A. Goodrich\*

## ABSTRACT

Learning in many multi-agent settings is inherently *repeated play*. This calls into question the naive application of Nash equilibria in multi-agent learning and suggests, instead, the application of give-and-take principles of bargaining. We present an  $M$  action,  $N$  player social dilemma that encodes the key elements of the Prisoner's Dilemma and thereby serves to highlight the importance of cooperation in multi-agent systems. This game is instructive because it characterizes social dilemmas with more than two agents and more than two choices. We show how several different multi-agent learning algorithms behave in this social dilemma, including a satisficing algorithm based on [16] that is compatible with the bargaining perspective. This algorithm is a form of relaxation search that converges to a satisficing equilibrium without knowledge of other agents actions and payoffs. Finally, we present theoretical results that characterize the behavior of the algorithm.

## 1. INTRODUCTION

Many multi-agent learning problems can be viewed as social dilemmas. For example, in [11] we presented a multi-robot scenario that illustrated the difficulties in creating learning algorithms for environments where there are multiple learning agents and where games are non-zero sum. These difficulties arose because each robot needed to use a common resource, but if each robot tried to dominate the resource then every robot suffered. This is typical of prisoner's dilemma-like environments with ongoing interactions; the robots were required to act independently, but the solution concept of a single-play Nash equilibrium was inappropriate for these repeated interactions. Without the imposition of a central arbiter, however, cooperation is unlikely to emerge in repeated play using current algorithms and, even if it does emerge, cooperation may not be an attractor of the learning

\*Corresponding author: mike@cs.byu.edu. Computer Science Department, Brigham Young University, Provo, UT 84602

process.

In this paper, we introduce a multi-agent social dilemma game for which the single play Nash equilibrium solution may be undesirable. This game has the essential characteristics of the prisoner's dilemma, but is broad enough to represent social dilemmas with more than two actions (such as in extensive form games) and more than two agents (such as in multi-robot teams). We evaluate how several algorithms behave in this dilemma, and discuss whether they emphasize Nash equilibrium or Nash bargaining solutions. We then analyze a satisficing algorithm that is flexible enough to produce either solution, depending on the behavior of other agents. This algorithm and analysis contributes to an emerging body of work on learning to cooperate in repeated-play games [18, 19, 5].

## 2. RELATED LITERATURE

The literature in multi-agent choice is vast and space is limited, so we cite only a few. A more complete citation list can be found in [25]. Machine learning researchers have explored many approaches to learning in games. [13, 6] presented extensions to Q-learning for stochastic games that converge to Nash equilibrium solutions, and [27] has extended one of these algorithms to exploit the naive strategies of other agents. Complementing these papers is work from the economics literature [15, 9] that describes when and how model-based agents tend to converge to a Nash equilibrium. Others [22, 14], have explored algorithms based on reinforcement learning proceed when some of the assumptions required for convergence to a Nash equilibrium are violated.

Unfortunately, a lesson learned from the repeated play Prisoner's Dilemma game [1] is that strategies that tend to Nash equilibria are not always desirable when agents engage in repeated interactions. Attempts to generate cooperative solutions using algorithms with claims of bounded rationality have offered some insight into when cooperation is preferred to Nash equilibria [21]. From a machine learning perspective, augmenting state information with coordination-specific information can lead to cooperation [3].

In comparison to the standard prisoner's dilemma, literature related to multiple-player or multiple-action versions is smaller and far less unified. A formal discussion is provided in [12]. [20] and [2] briefly discuss a multiple-player, two-action prisoner's dilemma. In addition to multiple players, the prisoner's dilemma has been extended to continuous de-

grees of cooperation; much of this work is synthesized in [8], where a multiple-player, continuous prisoner’s dilemma is formulated.

### 3. A SOCIAL DILEMMA

In this section, we introduce the *multi-agent social dilemma* (MASD) which is a game with the same essential characteristics as the prisoner’s dilemma, but which allows for multiple players and actions. This game is useful for illustrating strengths and weaknesses of various multi-agent learning algorithms.

Consider a system consisting of  $N$  agents. At each iteration, every agent is faced with a decision of allocating  $M$  units of some discrete resource towards two possible goals  $S_i$  and  $G$ .  $S_i$  is some purely self-interested goal for agent  $i \in \{1, \dots, N\}$  and  $G$  is some group goal for all agents. Let  $u_i$  be the amount contributed by agent  $i$  towards the group goal  $G$  (and thus  $M - u_i$  is the amount contributed to the selfish goal  $S_i$ ). Let  $\mathbf{u} = [u_1, \dots, u_N]$  denote the vector of all actions taken by the agents. For each agent there are  $M + 1$  possible values for  $u_i \in \{0, 1, 2, \dots, M\}$ . Let each agent’s total utility be represented as a linear combination of the total amount contributed to the group goal  $G$  and the amount individually contributed to his or her own selfish goal  $S_i$ . The utility to agent  $i$  given the actions of all agents is

$$R_i(\mathbf{u}) = k_G \left[ \sum_{j=1}^N u_j \right] + k_{S_i} (M - u_i), \quad (1)$$

where  $k_{S_i}$  is agent  $i$ ’s weighting of his or her own selfish goal and  $k_G$  is agent  $i$ ’s weighting of the group goal.

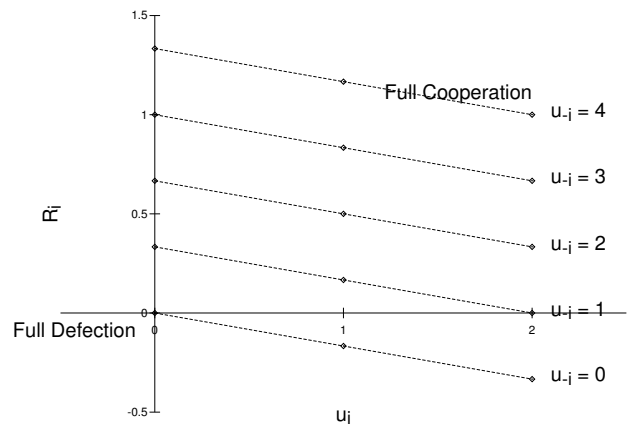
Suppose that all agents have the same  $k_S = k_{S_i}$  and  $k_G = k_G$ , the relative. Assuming that relative (not absolute) utilities are important, we can reduce the number of parameters by letting  $k_G = \frac{1}{NM}$  and  $k_S = \frac{k}{M}$  where  $k$  is a positive constant. When  $k < 1$  it means that each agent values a unit of contribution towards the selfish goal more than a unit of contribution to the group goal, and when  $k > \frac{1}{N}$  it means that there is a higher potential benefit from the group goal as long as enough agents contribute to the group. Thus, attention is restricted to the case where  $1 > k > \frac{1}{N}$ . Substituting this reparameterization into (1), dividing by  $M(1 - k)$ , and dropping a constant term from the end, gives

$$R_i(\mathbf{u}) = \frac{[\frac{1}{N} \sum_{j=1}^N u_j] - k u_i}{M(1 - k)}.$$

It will often be useful to examine the situation from the perspective of a single agent. For these circumstances, we define  $\mathbf{u}_{-i} \in \mathbf{U}_{-i}$  as the joint action of agent  $i$ ’s opponents. In the MASD,  $\mathbf{u}_{-i}$  can be reduced to a scalar integer because the reward function depends only on the sum of the actions of agent  $i$ ’s opponents whence  $u_{-i} = \sum_{j=1, j \neq i}^N u_j$ , whence

$$R_i(\mathbf{u}) = R_i(u_i, u_{-i}) = \frac{(1 - kN)u_i + u_{-i}}{NM(1 - k)}. \quad (2)$$

Figure 1 illustrates an example of the MASD reward structure when  $N = 3$ ,  $M = 2$ , and  $k = 0.5$ . The action of some agent  $i$  is shown along the x-axis. The reward for agent  $i$  is on the y-axis. For each of agent  $i$ ’s possible choices there are actually  $M(N - 1)$  possible rewards depending on how the other agents act.



**Figure 1: An example of payoff structure for social dilemma.**

From Figure 1 we can see a prisoner’s dilemma-like situation arising. The action  $u_i = 0$  corresponds to full “defection” by player  $i$  and  $u_i = M$  corresponds to full “cooperation” by player  $i$ . Clearly, no matter what  $u_i$  is chosen, agent  $i$  would receive the highest reward by choosing  $u_i = 0$ . However, if all the other agents also make this choice, then the bottom line is used and agent  $i$  (along with all other agents) receives a reward of 0. On the other hand, if the agent  $i$  chooses  $u_i = 2$  and the other agents fully cooperate, the top line is selected and agent  $i$  (along with the other agents) receives a reward of 1. However, by choosing  $u_i = M$ , agent  $i$  is exposed to possible exploitation; if the other agents all choose full defection ( $u_{-i} = 0$ ), then the bottom line is selected and agent  $i$  receives a reward of  $-\frac{1}{2}$ .

$\bar{R}(\mathbf{u}) = \frac{1}{NM} \sum_{i=1}^N u_i$ , is maximized at  $\bar{R} = 1$  by the joint action  $\mathbf{u}$  where  $\forall i u_i = M$ , and is minimized at  $\bar{R} = 0$  when  $\forall i u_i = 0$

3. **Nash Equilibrium** The joint action  $\mathbf{u}$  where  $\forall i u_i = 0$  is both strategically dominant and the unique Nash equilibrium.
4. **Nash Bargaining Solution** When the fallback position is defined as the strategically dominant solution, the joint action  $\mathbf{u}$  where  $\forall i u_i = M$  is the Nash Bargaining solution. It is therefore also Pareto optimal.

### 4. THE ALGORITHM

Herbert Simon introduced the term satisficing to mean “good enough” [23]. Although he discussed satisficing from several perspectives, a frequent perspective was one in which an agent searched through a set of possible decisions until a decision was found which had utility that exceeded an aspiration level. A formal treatment of this algorithm was

At each iteration  $t$

1. For each agent, compute

$$R_i(\mathbf{u}(t)) = \frac{[\frac{1}{N} \sum_{j=1}^N u_j(t)] - k u_i(t)}{M(1-k)}$$

2. Update the actions for satisficing agents

- If  $R_i(\mathbf{u}(t)) \geq \alpha_i(t)$  then  $u_i(t+1) = u_i(t)$  otherwise select  $u_i(t+1)$  from a uniform distribution over all actions.

3. Update the aspirations for satisficing agents

- $\alpha_i(t+1) = \lambda \alpha_i(t) + (1-\lambda)R_i(\mathbf{u}(t))$

**Figure 2: The satisficing algorithm for the MASD.**

analyzed in a prisoner’s dilemma context in [16] and further analyzed in [24] for deterministic updates. The conclusion of these papers is that a satisficing algorithm can lead to mutual cooperation in the prisoner’s dilemma under a broad variety of conditions.

#### 4.1 Extending Karandikar’s Algorithm to the MASD

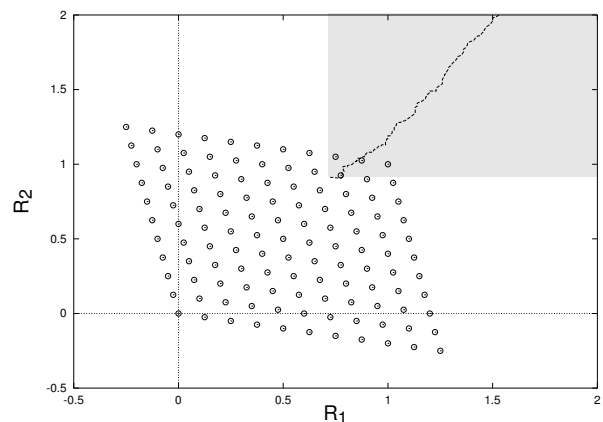
[16]’s algorithm works as follows: (a) when the aspiration level,  $\alpha$ , is not met, the agent switches actions, and (b) the aspiration level is updated as the convex combination of the old aspiration and the current reward via learning rate  $\lambda$ . In the prisoner’s dilemma, switching means simply switching to the other action. In the social dilemma, an agent must choose between an arbitrary number of actions. We adopt the simple method of selecting the next action randomly; more sophisticated techniques, such as policy hill climbing, are topics for future work. Figure 2 states the modified satisficing algorithm in the MASD context. For simplicity, we suppose that all agents use the same learning rate  $\lambda$ .

#### 4.2 An Example of Learning in the MASD

Figure 3 illustrates the satisficing learning process for two agents. The figure is shown for  $M = 10$  and initial aspirations are  $(\alpha_1, \alpha_2) = (1.5, 2.0)$ . In the figure, each open circle denotes a possible reward for some joint action, where the x-coordinate is Player 1’s reward and the y-coordinate is Player 2’s reward. At each time step, one of these rewards is determined from the joint action. The line that trails down from the upper right corner of the graph is a plot of the aspiration history for the agents. The gray area to the northeast of the aspiration level is termed the *satisficing region*, meaning that if a reward is selected that is in this region, both players will be satisfied and aspirations will converge to the chosen reward.

Initially, all actions produce rewards that are less than the aspiration levels of the agents. This causes aspirations to drop and, as a result, the agents are choosing randomly and thus the rewards are also randomly selected from any of the possibilities shown.

At the time shown in the figure, the satisficing region intersects the area of feasible rewards. It is now possible that



**Figure 3: An illustration of the satisficing learning process. For this example,  $M = 10$ ,  $k = 0.6$ , and  $\lambda = 0.99$  for both agents.**

a single agent may be satisfied with an action. However, because the aspirations are still quite high, most of those individually satisficing actions are likely to exploit the other agent. An agent is therefore unlikely to stay satisfied for more than a few iterations because the unsatisfied agent will constantly be changing. During this time, if both agents chose  $M = 10$  then they both receive a reward of (1,1) and they both continue to play this satisficing action ever after. Once this action is chosen, the aspiration vector approaches (1,1) until eventually the only action that is mutually satisfying is mutual cooperation.

This is the typical manner in which satisficing converges in the two-player, multiple-action MASD. Intuitively, we can see that mutual cooperation is the most probable outcome as long as aspirations start high and aspirations are updated slowly. In most cases,  $\mathbf{u} = (M, M)$  will be the first joint action that is satisficing to both agents.

## 5. SIMULATION RESULTS

In this section, we analyze expected average reward in self-play for (a) Bowling and Veloso’s WoLF algorithm [4], (b) two naive algorithms, (c) a general belief-based learner, (d) the Q-learning algorithm, and (e) the satisficing algorithm.

To compare the algorithms, it is useful to have a standard of comparison. In related work, Bowling and Veloso [27] state that a good learning algorithm should reach a Nash equilibrium in self-play and should find a best-response against inferior opponents. We flip these desiderata to identify two properties that are desirable from the bargaining perspective: a good learning algorithm should (I) reach a Pareto efficient solution in self-play and (II) should not be exploited by selfish agents. We will use average performance in self-play as a metric for measuring how often Pareto efficient solutions are obtained.

The selection of a Pareto efficient solution does not tell us which Pareto efficient solution should be selected. In the work presented herein, we adopt Nash’s definition of a fair bargain and select the Nash Bargaining solution [17, 20] as the most preferred Pareto efficient solution. We want to emphasize the importance of learning a Pareto efficient solution in general, and a Nash Bargaining solution in particular. When an algorithm is learning to play against itself, it seems unreasonable for a smart algorithm to learn to play the Nash equilibrium. Rather, a smart algorithm should learn that it is playing against another smart algorithm and therefore seek to find a solution that is beneficial to both algorithms. The repeated play nature of learning translates the problem from one of myopic optimizing to one of bargaining. Thus, notions of Pareto efficiency in self-play are important measurements of the applicability of the algorithm.

### 5.1 The WoLF algorithm

Bowling and Veloso’s Win or Learn Fast (WoLF) algorithm [4] seeks to learn to exploit inferior opponents, and to learn the Nash equilibrium against other opponents. Because the algorithm is based in Q-learning, the algorithm can successfully learn the best-response solution against stationary opponents; when these opponents play the Nash equilibrium solution, the WoLF algorithm converges to the Nash equilibrium solution. Thus, in self-play, the Nash equilibrium solution will be an attractor of the learning process.

In summary, the WoLF algorithm is unlikely to be exploited by selfish-agents, but it will fail to find a Pareto efficient solution in self-play.

### 5.2 Two Naive Algorithms

The simplest possible strategy would be to choose the same fixed action  $z$  at every iteration. In terms of average reward to a society of agents, the performance of this system when all agents use this strategy will be  $\frac{z}{M}$ . Thus, when  $z = M$ ,  $\bar{R}$  can be maximized at one. However, high values for  $z$  are open to exploitation by other agents, particularly agents that always play the Nash equilibrium. At the other end of the spectrum, if  $z = 0$ , the agents would never be exploited, but average reward would always be minimized. Intermediate choices for  $z$  would lead to less exposure to

exploitation, but also a lower average reward. Thus, fixed action strategies can never avoid exploitation while simultaneously guaranteeing that they learn a Pareto efficient solution in self-play.

Another possible consideration is to look at purely random strategies. For example, when all agents select their actions from a uniform distribution the expected average reward is  $E\{\bar{R}\} = 0.5$ . Again, however, agents playing the Nash equilibrium would exploit any individual playing a purely random strategy.

### 5.3 Belief-Based Learning

In this section, we present and discuss a general form of belief-based learning described in [7]. In this algorithm, player’s beliefs about an opponent’s play are characterized by a set of weights for each opponent action. At time  $t$  player  $i$  creates a probabilistic model,  $q_i(u_{-i}; t)$  of all other agents using standard techniques from fictitious play [9]. Given this opponent model, a player can compute the expected value,  $\hat{V}_i(u_i; t)$ , for each action  $u_i$  as  $\hat{V}_i(u_i; t) = \sum_{u_{-i} \in U_{-i}} R_i(u_i, u_{-i})q_i(u_{-i}; t)$ . A probability,  $p_i(u_i; t)$ , of choosing action  $u_i$  is then assigned as follows (thereby producing mixed strategies),

$$p_i(u_i; t) = \frac{\exp(\lambda_i \hat{V}_i(u_i; t))}{\sum_{u'_i \in U_i} \exp(\lambda_i \hat{V}_i(u'_i; t))},$$

where  $\lambda_i$  is the Boltzmann parameter that determines how optimally player  $i$  plays according to his beliefs. Note that this algorithm is a general case of many well-known belief-based learning algorithms including standard and cautious fictitious play [10].

Consider this learning model applied to the MASD. Substituting Equation (2) for  $R_i$  into the probability of choice and reducing leads to

$$p_i(u_i; t) = \frac{e^{-A\lambda_i u_i}}{\sum_{z=0}^M e^{-A\lambda_i z}}. \quad (3)$$

where  $A = \frac{1-kN}{NM(1-k)}$ . Note that these probabilities are completely independent of the opponent’s strategies or the player’s predictions about the probabilities of the opponents’ play. This means that learning models of this form are unable to adapt their behavior to their opponents in the MASD, and essentially reduce to a purely random strategy with the above distribution function. Furthermore, it can be shown that any dependence on state (whether from game history or player history) is eliminated in the final probability distribution.

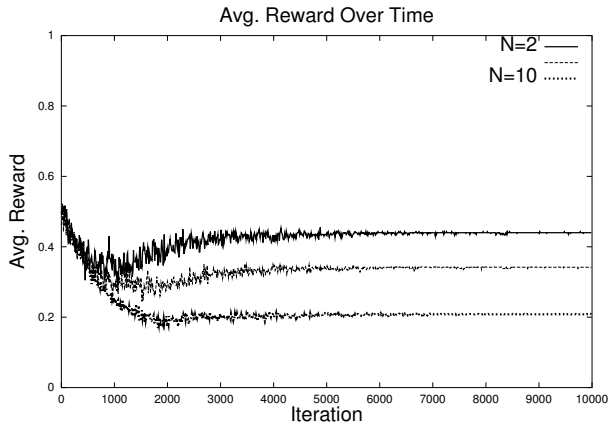
Consider the expected play for two extreme values of  $\lambda_i$ . When  $\lambda_i = 0$ , the expected play is  $\frac{M}{2}$ , and in the limit as  $\lambda_i \rightarrow \infty$ , then the expected play is 0. When all agents in a society use a learning strategy of this type,  $\bar{R}$  is bounded in  $[0, \frac{M}{2}]$  depending on the values for  $\lambda_i$ . Experiments were conducted for various parameter values where, in a given game, all agents used the same  $\lambda$ . When  $N = 5$ ,  $M = 3$ , and  $k = 0.6$ , the theoretical payoffs  $\bar{R}$  are as follows:  $\lambda = 0 \rightarrow \bar{R} = 0.5$ ,  $\lambda = 1 \rightarrow \bar{R} = 0.36$ , and  $\lambda = 10 \rightarrow \bar{R} = 0.012$ . Note that since these values are supported with the average empirical results, plots of the simulations are omitted.

In terms of the two desiderata, since the belief based agents act randomly, they will not learn mutual cooperation in self-play. They can, however, avoid exploitation by an appropriate choice of parameters.

## 5.4 Q-Learning

In strict terms, applying Q-learning to multiagent environments is not mathematically justified due to the fact that the transition function is not stationary when the other agents are able to learn and adapt their behavior. Such limitations are resolved in algorithms that adopt a stochastic games framework (such as WoLF [4]), but these algorithms emphasize convergence to Nash equilibria. Despite the theoretical difficulties, Q-learning has been shown to sometimes converge to Pareto efficient solution in some multiagent learning environments, and to best response solutions against selfish agents.

We designed several experiments to evaluate the performance of Q-learners in the MASD. The main results are presented in Figure 4, which displays the average rewards  $\bar{R}$  throughout the learning process for three different systems of Q-learning agents. As can be seen, in all cases, cooperation



**Figure 4:** The average reward of Q-learning agents over time in the MASD. The three lines represent three separate experiments with  $N=2$ ,  $N=3$ , and  $N=10$ . In all cases,  $M = 1$ . Each experiment consisted of averaging the rewards of all the agents over 200 trials. The game parameter  $k$  was chosen from a uniform random distribution over its legal range given  $N$ . The Q-learners used a fixed learning rate  $\alpha = 0.2$ , a discount factor  $\gamma = 0.9$  and Softmax exploration.

was relatively infrequent, although we found that the Q-learners always converged. In most cases, agents converged to the Nash equilibrium, but occasionally mutual cooperation emerged.

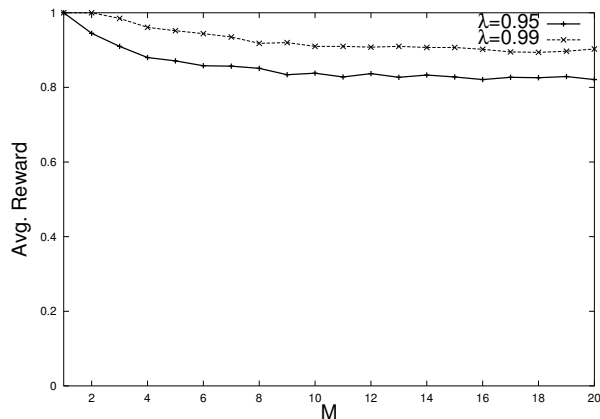
We varied both the parameters of the agents ( $\alpha$ ,  $\gamma$ , state representation) and the properties of the game ( $N$ ,  $M$ , and  $k$ ). Except for  $N$ , the performance of the Q-learning agents is not highly dependent on the agent parameters. For example, there is a wide range of values for both  $\alpha$ ,  $\gamma$ , and  $M$  that lead to similar results. We experimented with different state representations as well (account for the previous entire

joint action,  $\mathbf{u}$ , and account for the sum of the joint action,  $u_{-i}$ ) but found that it did not have a significant effect on the frequency of cooperation.

In terms of the two desiderata, the Q-learners tend to learn best responses to stationary strategies (as evidenced by the predominance of Nash equilibrium solutions), so they are unlikely to be exploited. However, they only rarely learn mutual cooperation.

## 5.5 Our Algorithm

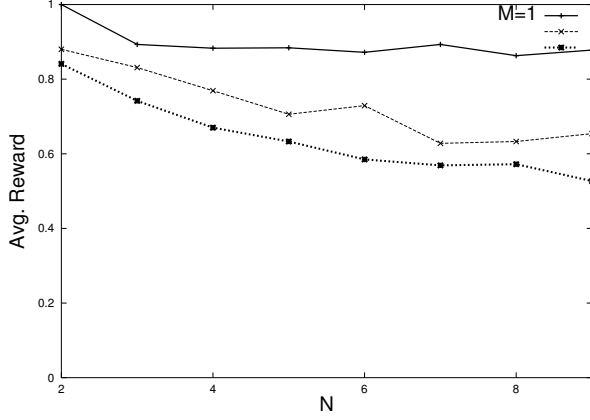
Figure 5 displays the average reward produced by the satisficing algorithm for two agents as a function of  $M$  in self-play. First, note that the performance is very high, meaning that in self-play the satisficing algorithm is likely to converge to a Pareto efficient solution. Note also that as  $M$  increases, the average reward decreases, but stays fairly high, even though the probability of guaranteed cooperation gets very small which means that the algorithm degrades gracefully as complexity increases. This can be accounted for by considering that mutual cooperation can still occur even when we cannot guarantee it. Also, the primary reason that  $\mathbf{u}^c$  becomes more difficult to obtain is not that bad solutions are found, but that fairly good solutions are found that are close to mutual cooperation.



**Figure 5:** The average reward for 500 games over the society of satisficing agents. In all games,  $k = 0.6$  and initial aspirations were randomly selected from the range  $[1.5, 2.0]$ .

Figure 6 compares the average reward over the society as  $N$  increases for three different values of  $M$ . We note that  $\bar{R}$  starts high, but falls off as  $N$  increases. By the time  $N = 10$ , for moderate values of  $M$ ,  $\bar{R}$  has significantly decreased and begins to approach 0.5.

In terms of the two desiderata, mutual cooperation is likely to emerge in self-play. Furthermore, we can prove that the algorithm is likely to converge to a Nash equilibrium when playing against a society of selfish agents, which means that the algorithm is not likely to be exploited. A corresponding theorem exists that states conditions that guarantee the algorithm will converge, with high probability, to mutual cooperation in self-play. The proof of the Nash equilibrium result is more straightforward, and contains the key elements



**Figure 6:** The average reward for 500 games of a society of satisficing agents as  $N$  increases.  $k$  was selected randomly from its legal range. Initial aspirations were chosen randomly between  $R_{\max}$  and  $2R_{\max}$ . All agents used a  $\lambda = 0.99$ .

of the second proof so, in the interest of space, we present only the first result in this paper.

One point of this paper is that naive application of the single play Nash equilibrium is sometimes inappropriate. The satisficing addresses the repeated play nature of learning by utilizing a history-dependent aspiration level. Thus, the aspiration level comes to encode and represent behavior that is considered acceptable in a repeated play context. Such acceptable behavior represents the bargain struck by the agents. Simulation results suggest that the Nash bargaining solution is the most likely solution selected from the set of Pareto efficient solutions provided that each agent begins with similar aspirations.

## 6. LEARNING A NASH EQUILIBRIUM

If a learning agent is facing agents that always attempt to exploit others, an effective learning algorithm should be able to learn the Nash equilibrium. In this section, we evaluate the ability of a satisficing agent to learn  $u_i = 0$  in such a society. Observe that  $u_{-i}$  will always be 0 for the satisficing agent. This means that the reward to the satisficing agent  $i$  for taking action  $u_i$  is

$$R_i(u_i(t), u_{-i}(t) = 0) = \frac{1 - kN}{NM(1 - k)}u_i(t). \quad (4)$$

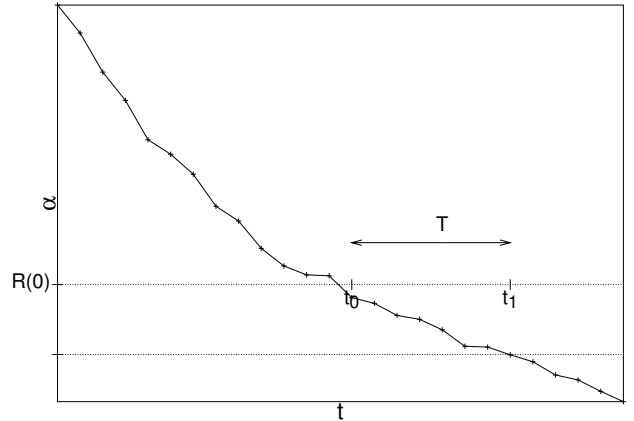
Note that  $R_i \leq 0$  which implies that  $\alpha(t)$  will always be decreasing until  $\alpha(t) \leq 0$ . At that point, whenever  $R_i \geq \alpha(t)$ , the agent will be satisfied indefinitely. As long as  $0 < \lambda < 1$ , the aspiration level  $\alpha(t)$  cannot fall below the minimum reward given and thus the algorithm will always converge to some action  $u^*$ .

### 6.1 Intuition Behind the Argument

The best response for a satisficing agent against  $u_{-i} = 0$  is  $u_i = 0$ ; this is the action that is most likely to be produced by the satisficing algorithm. This result requires that the initial aspiration  $\alpha(0) \geq R_i(0, 0) = 0$ . This means that the satisficing agent will be initially unsatisfied for several iterations while aspirations fall towards zero. Eventually, at

some  $t_0$ ,  $\alpha(t_0) < 0$ , after which if  $u_i(t) = 0$  is chosen for  $t > t_0$  then the agent will converge to the Nash equilibrium. However, it is possible at some time  $t_1$  for aspirations to fall below  $R_i(1, 0)$  before full defection is chosen. The trick is to make  $T = t_1 - t_0$  large enough that  $u_i = 0$  is chosen with a high probability.

Figure 7 illustrates this concept. The aspiration starts well above the reward for mutual defection denoted by  $R(0)$ . At each iteration, however,  $\alpha$  falls towards the received payoffs. At some point  $t_0$ ,  $\alpha(t)$  drops below the reward for playing the Nash equilibrium. At this point only the Nash equilibrium is satisfying. Eventually, at some time  $t_1 > t_0$ , if the agent does not play  $u_i = 0$ , the aspirations will fall below  $R_i(1, 0)$ , after which it is possible to converge to  $u = 1$ , and thus be indefinitely exploited.



**Figure 7:** An example of the aspirations of a single satisficing agent against a society of defecting agents over time.

### 6.2 Theorem

The critical factor in determining if the algorithm will converge to the Nash equilibrium is the length of the interval  $T = t_1 - t_0$ . We can place a lower limit on  $T$  by identifying the value of  $T$  that causes aspirations to fall most sharply between  $R_i(0, 0)$  and  $R_i(1, 0)$ . We refer the reader to [25] for a complete derivation. This bound is given by

$$T \geq \log_{\lambda} \left[ \frac{M-1}{M} \right] - 1,$$

which depends only on  $\lambda$  and  $M$ . Note that as  $\lambda$  approaches one  $T$  gets larger and approaches infinity, but as  $M$  goes to infinity,  $T$  goes to zero. We can now state and prove the following:

**LEMMA 6.1.** *Assuming  $\alpha(0) \geq R_i(0, 0)$ , then for any  $T' > 1$ , there exists a  $\lambda \in (0, 1)$  such that the shortest interval  $T$  in which  $\alpha(t_0 + T) > R_i(1, 0)$  satisfies  $T > T'$ .*

**PROOF.** Since  $T \geq \log_{\lambda} \left[ \frac{M-1}{M} \right] - 1$  it suffices to find a  $\lambda$  such that  $\log_{\lambda} \left[ \frac{M-1}{M} \right] - 1 \geq T'$  for any  $T' > 1$ . Such a  $\lambda$  must satisfy  $\frac{\ln \left( \frac{M-1}{M} \right)}{\ln \lambda} \geq T' + 1$  which is equivalent to  $\left( \frac{M-1}{M} \right)^{\frac{1}{T'+1}} \leq \lambda < 1$ . Thus, we can always choose a  $\lambda \in (0, 1)$  that will make  $T$  greater than any arbitrary  $T' > 1$ .  $\square$

We now know conditions on  $\lambda$  such that there is a time window of at least  $T$  iterations in which the Nash equilibrium action,  $u_i = 0$ , will be the only satisficing action. During this window, actions are selected from a uniform distribution where  $P[u = 0] = \frac{1}{M+1}$ . It follows, then, that the probability of the Nash equilibrium occurring in this window of length  $T$  is given by  $1 - (\frac{M}{M+1})^T$ . The Nash equilibrium could be reached in subsequent iterations (after  $\alpha(t) < R(1,0)$ ), but that will only increase the probability that  $u^* = 0$ . Thus, the probability that the agent learns the Nash equilibrium against a society of always-defecting agents can be bounded by

$$P[u_i^* = 0] \geq 1 - \left(\frac{M}{M+1}\right)^T. \quad (5)$$

**THEOREM 6.1.** *Consider a multiagent social dilemma specified by  $(N, M, k)$  played by a single satisficing agent  $i$  when  $u_{-i} = 0$ . Suppose that  $\alpha(0) \geq R_i(0, 0)$ . Then, for any  $\epsilon$  such that  $0 < \epsilon < 1$ , there exists a learning rate  $\lambda$  such that the probability of the single satisficing agent learning the Nash equilibrium is at least  $1 - \epsilon$ .*

**PROOF.** By Equation (5), we know that  $P[u^* = 0] \geq 1 - (\frac{M}{M+1})^T$ . Thus, if we can show that  $1 - (\frac{M}{M+1})^T \geq 1 - \epsilon$  then it follows that  $P[u^* = 0] \geq 1 - \epsilon$ . To satisfy this inequality,  $T$  must satisfy  $T \geq \log_{\frac{M}{M+1}}(\epsilon)$ . But by Lemma 6.1, we can always choose a  $T$  such that  $T \geq T' = \log_{\frac{M}{M+1}}(\epsilon)$ . Thus, for any  $\epsilon$  there exists a  $\lambda$  such that  $P[u_i^* = 0] \geq 1 - \epsilon$ .  $\square$

Empirical results confirm that  $P(u^* = 0)$  is indeed bounded by this limit. A similar proof can be used to show that a learning rate  $\lambda$  can be chosen such that a group of  $N$  satisficing agents will likely converge to mutual cooperation (the Nash bargaining solution) if they all begin with high and similar aspiration levels.

### 6.3 Graceful Degradation and Convergence Time

It is desirable for the algorithm to degrade gracefully in the presence of many possible actions ( $M$ ). Consider the system at any time  $t \geq t_1$ . At such times, the Nash equilibrium is at least as likely to be chosen as any other mutually satisficing action. Thus, the Nash equilibrium is not only possible earlier than higher values of  $u_i^*$ , but is always at least as likely as any other  $u_i^*$  as well. Furthermore, since  $R_i(u_i^*, 0)$  is proportional to the ratio  $\frac{u_i^*}{M}$ , as  $M$  increases the probability of missing the Nash equilibrium increases, but the cost of slightly missing  $u_i^* = 0$  decreases. Empirical results confirm that the average reward for a satisficing agent against  $u_{-i} = 0$  degrades gracefully. The trends are similar to those shown in Figure 5 so plots are omitted in the interest of space.

Time to converge is also a very important element of the performance of the satisficing algorithm. The order of convergence time is  $\left\lceil \frac{1}{\log \lambda} \right\rceil$ , which is obtained by taking the expected aspiration level at some time  $t$ . Thus, a high  $\lambda$  is required to make non-exploitation likely, but it also significantly increases convergence time.

## 7. DISCUSSION

We contend that when a learning agent interacts with other learning agents, the learning process is better treated as a bargaining problem than a search for the Nash equilibrium solution. This contention is based on the observation that learning is often inherently repeated play, so a give-and-take approach to adaptation is more appropriate than insisting on individual optimization.

We have presented an  $M$ -action,  $N$ -agent social dilemma that illustrates the importance of bargaining in environments with multiple learning agents. We evaluated the performance of several learning algorithms in this dilemma. Q-learning rarely converges to mutual cooperation in self play, belief-based learning generates random actions without regard to the game for this dilemma, naive solutions can either be exploited or cannot bargain efficiently, and WoLF learns to exploit others but cannot learn to cooperate with itself in self-play.

The satisficing algorithm, by contrast, usually converges to mutual cooperation in self-play, but usually avoids being exploited by selfish agents. Relaxing an agent's aspirations is one way to show respect for how an agent's choices affect other agents, and therefore acknowledges the repeated play nature of multi-agent learning. By relaxing aspirations, the satisficing algorithm assumes a bargaining perspective while avoiding being exploited by selfish agents.

One could argue that the bargaining perspective prevents a stable equilibrium from being learned, and therefore prevents convergence. Stirling has argued, however, that when all agents are "satisfied" then there is no incentive for any agent to change its choice [26]. Such agents are said to have reached a *satisficing equilibrium* which implies that the satisficing algorithm produces stable solutions in self-play. As demonstrated in this paper, such an equilibrium can represent what the agents perceive as an acceptable bargain based on their experiences with other agents.

## 8. REFERENCES

- [1] R. M. Axelrod. *The Evolution of Cooperation*. Basic Books, 1984.
- [2] J. Bendor and D. Mookherjee. Institutional structure and the logic of ongoing collective action. *The American Political Science Review*, 81(1):129–154, March 1987.
- [3] C. Boutilier. Sequential optimality and coordination in multiagent systems. In *Proceedings of the Sixteenth International Joint Conferences on Artificial Intelligence (IJCAI-99)*, 1999.
- [4] M. Bowling and M. Veloso. Rational learning of mixed equilibria in stochastic games. In UAI2000 Workshop entitled Beyond MDPs: Representations and Algorithms, June 2000.
- [5] Y-H. Chang and L. P. Kaelbling. Playing is believing: The role of beliefs in multi-agent learning. In *Proceedings of 2001 Conference on Neural Information Processing Systems*, Vancouver, British Columbia, Canada, Dec 3-8 2001.
- [6] C. Claus and C. Boutilier. The dynamics of reinforcement learning in cooperative multiagent systems. In *Proceedings of the Fifteenth National Conference on Artificial Intelligence, AAAI98*, Madison, Wisconsin USA, July 26–30 1997.

- [7] N. Feltovich. Reinforcement-based vs. belief-based learning models in experimental asymmetric-information games. *Econometrica*, 68(3):605–641, 2000.
- [8] N. Frolich and J. Oppenheimer. When is universal contribution best for the group? Characterizing optimality in the prisoner’s dilemma. *Journal of Conflict Resolution*, 40(3):502–516, 1996.
- [9] D. Fudenberg and D. K. Levine. *The Theory of Learning in Games*. MIT Press, 1998.
- [10] D. Fudenberg and D. Levine. *The Theory of Learning in Games*. MIT Press, 1998.
- [11] M. A. Goodrich, J. W. Crandall, and J. L. Stimpson. Neglect tolerant teaming: Issues and dilemmas. In *Proceedings of the 2003 AAAI Spring Symposium on Human Interaction with Autonomous Systems in Complex Environments*, 2003. To appear.
- [12] H. Hamburger. N-person prisoner’s dilemmas. *Journal of Mathematical Sociology*, 3:27–48, 1973.
- [13] J. Hu and M. P. Wellman. Multiagent reinforcement learning: Theoretical framework and an algorithm. In J. Shavlik, editor, *Proceedings of the Fifteenth International Conference on Machine Learning*, pages 242–250, July 1998.
- [14] J. Hu and M. P. Wellman. Experimental results on Q-learning for general-sum stochastic games. In *Proceedings of the Seventeenth International Conference on Machine Learning (ICML-2000)*, San Francisco, July 2000. AAAI Press.
- [15] E. Kalai and E. Lehrer. Rational learning leads to Nash equilibrium. *Econometrica*, 61(5):1019–1045, September 1993.
- [16] R. Karandikar, D. Mookherjee, D. Ray, and F. Vega-Redondo. Evolving aspirations and cooperation. *Journal of Economic Theory*, 80:292–331, 1998.
- [17] H. W. Kuhn. *Classics in Game Theory*. Princeton University Press, 1997. Includes a reprint of Nash’s original paper on the Nash Bargaining solution.
- [18] M. L. Littman and P. Stone. Leading best-response strategies in repeated games. In *Proceedings of 17th International Joint Conference on Artificial Intelligence*, Seattle, WA, USA, Aug 4-10 2001.
- [19] M. L. Littman and P. Stone. A polynomial-time Nash equilibrium algorithm for repeated games. In *Proceedings of 2003 ACM Conference on Electronic Commerce*, San Diego, CA, USA, June 9-12 2003.
- [20] R. D. Luce and H. Raiffa. *Games and Decisions*. John Wiley, New York, 1957.
- [21] Y. Mor and J. S. Rosenschein. Time and the prisoner’s dilemma. In *Proceedings of the First International Conference on Multi-Agent Systems*, pages 276–282, Menlo Park, CA, USA, June 1995. AAAI Press.
- [22] M. Mundhe and S. Sen. Evaluating concurrent reinforcement learners. In *Proceedings of the Fourth International Conference on Multiagent Systems*, pages 421–422, Los Alamitos, CA, 2000. IEEE Press. Poster paper.
- [23] H. A. Simon. *The Sciences of the Artificial*. MIT Press, 3rd edition, 1996.
- [24] J. Stimpson, M. A. Goodrich, and L. Walters. Satisficing learning and cooperation in the prisoner’s dilemma. In *Proceedings of IJCAI 2001*, 2001.
- [25] J. L. Stimpson. Satisficing solutions to a multi-agent social dilemma. Master’s thesis, Brigham Young University, Provo, UT, 84602, USA, 2002.
- [26] W. C. Stirling, M. A. Goodrich, and D. J. Packard. Satisficing equilibria: A non-classical approach to games and decisions. *Autonomous Agents and Multi-Agent Systems Journal*, 2000. To appear.
- [27] Michael Bowling Manuela Veloso. An analysis of stochastic game theory for multiagent reinforcement learning. In *ICML*, 2000. Submitted.