



ELSEVIER

Available online at www.sciencedirect.com

SCIENCE @ DIRECT®

Journal of Economic Behavior & Organization

Vol. xxx (2004) xxx–xxx

 JOURNAL OF
 Economic Behavior
 & Organization

www.elsevier.com/locate/econbase

Learning strategies

Nobuyuki Hanaki^a, Rajiv Sethi^{b,*}, Ido Erev^c, Alexander Peterhansl^b

^a *Earth Institute, Columbia University, USA*

^b *Department of Economics, Barnard College, Columbia University, 3009 Broadway, New York 10027, USA*

^c *Faculty of Industrial Engineering and Management, Technion, Israel*

Received 15 January 2003; received in revised form 12 December 2003; accepted 21 December 2003

Abstract

Adaptive models of learning in experiments typically share two features: initial attractions are given exogenously and players learn about stage-game actions rather than repeated-game strategies. We develop a model of strategy learning with endogenous initial attractions. Learning occurs in two phases. In an initial long-run phase, players explore a complete set of strategies with bounded complexity. Limiting attractions from this phase are the initial attractions in the second, short-run phase, which can be tested against experimental data. Relative to existing models, we can better account for subject behavior in environments where fairness and reciprocity appear to play a significant role.

© 2004 Published by Elsevier B.V.

JEL classification: D83

Keywords: Reinforcement learning; Repeated game strategies

1. Introduction

Within the literature on learning in games, two distinct strands may be identified. One deals with the abstract question of the long-run convergence properties of learning models,

* Corresponding author. Tel.: +1 212 854 5140; fax: +1 212 854 8947.

E-mail addresses: nh85@columbia.edu (N. Hanaki), rs328@columbia.edu (R. Sethi), erev@tx.technion.ac.il (I. Erev), ap11@columbia.edu (A. Peterhansl).

with particular attention paid to the conditions under which learning leads to Nash equilibrium. The second deals with the more empirical task of describing the manner in which human subjects learn in laboratory interactions. The latter class of learning models can be further subdivided into those that are belief-based, such as fictitious play and those based on reinforcement. In belief-based models, subjects use observed histories of opponent actions to predict future play and respond optimally to such beliefs. Reinforcement learning, in contrast, is based on the hypothesis that the propensity to choose an action increases or decreases in response to the payoff experience resulting from the choice of that action. Both belief-based and reinforcement learning models are special cases of experience-weighted attraction (EWA) learning, which allows for the reinforcement not only of actions taken, but also of actions that were *not* taken, based on the imagined payoffs that such actions would have yielded.¹

Adaptive learning models have been reasonably successful in accounting for observed behavior in certain strategic environments, such as games with unique mixed strategy equilibrium and some coordination games, while failing dramatically to replicate human behavior in others. For example, when experimental subjects are paired to play a Prisoner's Dilemma for a finite number of periods under conditions of full information, convergence to mutual cooperation occurs frequently for many payoff configurations. In contrast, fictitious play predicts convergence to mutual defection for all parameter values. Similarly, in the Battle of the Sexes, fixed subject pairs frequently alternate between the two pure-strategy equilibria of the stage-game, thus managing to achieve payoff profiles that are both equitable and efficient. Neither reinforcement learning nor fictitious play can account for this, with both models predicting convergence to the repeated play of one or the other pure-strategy stage-game equilibria.²

One could conceivably account for the disparity between experimental findings and the predictions of learning models by arguing that subjects care not just about their own monetary payoffs, but also about the payoffs obtained by those with whom they interact. Several recent attempts have been made to identify a richer class of preferences that are able to take such interdependencies into account in a manner consistent with experimental behavior.³ From this perspective, payoff functions must be appropriately transformed before learning models can be properly tested or compared. While this is an important and promising direction for research, there is as yet no consensus on the precise manner in which monetary payoffs should be transformed in order to conform to 'social preferences.' Moreover, as we

¹ Fudenberg and Levine (1998) provide a detailed survey of the theoretical literature. The experimental learning literature is vast; see, for instance, Crawford (1995), Cheung and Friedman (1997), Mookherjee and Sopher (1997) and Erev and Roth (1998). EWA learning was developed by Camerer and Ho (1999), who also show that it generalizes both fictitious play and reinforcement learning. Stahl (1999, 2000) has developed a model allowing for the learning of behavioral rules, defined broadly as mappings from games and histories to probability distributions over actions. See also Day (1963) for an early and pioneering analysis of bounded rationality with adaptation, which can be considered a precursor to the modern literature on learning.

² These and other failures of existing adaptive learning models are discussed further in Section 2 below. McKelvey and Palfrey (2001) identify additional weaknesses of standard learning models, such as their insensitivity to variations in information and matching conditions.

³ See, in particular, Fehr and Schmidt (1999), Bolton and Ockenfels (2000), Levine (1998) and Charness and Rabin (2002). Such preference interdependence is predicted by several evolutionary models, including Güith and Yaari (1992), Huck and Oechssler (1999), Gintis (2000) and Sethi and Somanathan (2001).

57 argue below, behavior that appears to be motivated by a concern for fairness and efficiency
58 can in fact be the consequence of an entirely orthodox process of learning in which material
59 payoffs are the driving force.

60 In this paper, we maintain the hypothesis that subjects are motivated primarily by a
61 concern with their own monetary payoffs, but allow for the possibility that subjects learn not
62 just among stage-game actions but also among repeated game strategies. The limitations of
63 action-learning models are well recognized in the literature. Erev and Roth (1998) Erev and
64 Roth, (p. 872), for instance, note that it will “not generally be the case that learning behavior
65 can be analyzed in terms of stage-game actions alone.” Along similar lines, Camerer and
66 Ho (1999) Camerer and Ho, (p. 871) point out that “stage-game strategies are not always
67 the most natural candidates for the strategies that players learn about” and McKelvey and
68 Palfrey (2001) McKelvey and Palfrey, (p. 19) have argued for the development of “strategic
69 learning” models in which players learn not about the performance of actions but rather of
70 strategies (see also, Stahl and Haruvy, 2002). The development of strategic learning models
71 has been inhibited, however, by two potential obstacles. First, the size of the strategy space
72 precludes experimentation with all but a few strategies in any given interaction. In fact, any
73 learning rule can itself be interpreted as a single repeated game strategy. This problem can
74 be overcome, as McKelvey and Palfrey point out, by restricting the complexity of repeated
75 game strategies.⁴ The second difficulty arises from the fact that if players are learning among
76 repeated game strategies, it becomes impossible to compute the hypothetical payoffs that
77 would have been obtained had a *different* strategy been chosen. Hence, even with observable
78 actions and stage-game payoff functions, neither fictitious play nor the general version of
79 experience-weighted attraction learning can be implemented. As McKelvey and Palfrey
80 (2001) McKelvey and Palfrey, (p. 25) observe, “players face an inference problem going
81 from histories to beliefs” about the strategies of their opponents.

82 When learning responds only to payoffs obtained by strategies actually chosen by the
83 subject, *this inference problem does not arise*. A much maligned attribute of reinforcement
84 learning, therefore, turns out to be an advantage in developing models of learning
85 among repeated game strategies. In a straightforward extension of their earlier models of
86 reinforcement learning, Erev and Roth (2001) have studied the Prisoner’s Dilemma, while
87 allowing for players to choose among the two stage-game actions as well as the “tit-for-tat”
88 repeated game strategy. Allowing for the possibility that subjects can learn to reciprocate
89 significantly improves the predictive power of the model, but it does so at a cost: Erev and
90 Roth assume, in effect, that “tit-for-tat” is the only repeated game strategy to have posi-
91 tive probability weight when the process of learning begins. This choice is fundamentally
92 arbitrary and raises the question of why other strategies cannot also have positive initial
93 weight. More generally, one would like a theory of initial attractions that identifies the set
94 of repeated game strategies that experimental subjects explore.

⁴ Instead of drastically limiting the space of strategies *ex ante* and fully exploring this restricted space, one could use genetic algorithms, as in Miller (1996) and Lindgren (1997), to partially explore a much larger space of arbitrarily complex strategies. In such models new strategies enter the population through mutations of existing strategies and learning takes place through changes in the population composition: strategies yielding higher payoffs become more prevalent in the population relative to those yielding lower payoffs. Such social learning allows for the exploration of a large strategy space but is conceptually distinct from the experience-based individual learning with which we are concerned here.

95 Developing such a theory is the principal aim of this paper. Our contribution can be
96 thought of as a re-inforcement based approach to learning over long horizons in a ‘pre-
97 experimental’ phase that determines which repeated game strategies are salient when sub-
98 jects enter the laboratory. It is in this sense a theory of the initial attractions that appear as
99 parameters in standard learning models. In our model learning occurs over a long horizon
100 and begins with positive weight on each repeated game strategy that satisfies a bounded
101 complexity constraint. Specifically, we consider all strategies that can be represented by
102 automata having no more than two states.⁵

103 The model may be described briefly as follows. A large, fixed population is divided
104 into subject pairs. There are two phases of learning. During the first, ‘pre-experimental’
105 phase, subjects engage in a lengthy process of learning among repeated game strategies,
106 while being occasionally re-matched with other members of the population. There is a
107 finite set of simple repeated game strategies from which subjects choose. At the start of
108 the first phase of learning, each of the repeated game strategies has equal attraction and,
109 hence, equal probability of being chosen. Attractions are updated over time as the pay-
110 offs resulting from strategy choices are observed. Subjects maintain their chosen strate-
111 gies for several repetitions of the stage game with the length of this period determined
112 stochastically. Specifically, at each stage, there is some small and constant probability
113 that attractions will be updated and strategy revision will occur. Only strategies that are
114 actually chosen are updated, based on their observed payoff consequences. If strategy re-
115 vision occurs, the (possibly) new strategy is chosen on the basis of updated attractions.
116 There is also a small probability that at any stage, subject pairs are dispersed and indi-
117 viduals are re-matched with other subjects drawn from the population. Over the course
118 of this process some strategies decline in use, while others are observed with greater
119 frequency. The process continues until convergence to a limiting distribution is approx-
120 imated and this ends the first phase of learning. The limiting attractions from the first
121 phase are then used as initial attractions in the second, which consists of a fixed-pair
122 matching for a small number of periods. Learning also occurs in the second phase, but
123 without re-matching. This corresponds to the conditions of an experiment and enables
124 us to compare our results with reported experimental data. This two-phase learning cap-
125 tures the notion that experimental subjects bring whatever they have learned elsewhere into
126 the laboratory. The first “pre-experimental” phase corresponds to the real life experience
127 of players where they learn which strategies work the best in various strategic environ-
128 ments through interacting with many other people and they bring the knowledge from the
129 first phase into the second “experimental” phase, which is comparable to the laboratory
130 setup. We find that several patterns of behavior that are difficult to reconcile with action-
131 learning models, such as cooperation in the Prisoner’s Dilemma and alternation between
132 pure-strategy equilibria in the Battle of the Sexes, emerge as outcomes of our learning
133 procedure.

⁵ This seems unrestrictive in an analysis of 2×2 games and as we show below, promising results can be obtained without considering a larger strategy space. In an evolutionary model of the repeated Prisoners’ Dilemma, Miller considered automata having up to 16 states and found an endogenous decline in complexity with the survival of strategies similar to ‘tit-for-tat.’

134 **2. Two examples**

135 In Arifovic et al. (2002, henchforth AMP), a number of well-known learning models are
 136 presented side-by-side with the experimental results from an earlier paper by McKelvey and
 137 Palfrey. In one of the treatments reported, the setting was a fixed matching of two players
 138 playing a repeated game. The human experiments consisted of 48 subjects paired up for
 139 24 matches. Each match consisted of paired subjects playing a game for 50 rounds. The
 140 subjects in each pairing both saw the complete payoff matrix and observed their opponent’s
 141 choice of action after each round. Each match produced one binary-tuple of data: the average
 142 payoffs of each of two subjects over the course of the match. A variety of learning models,
 143 including fictitious play, reinforcement learning and EWA learning, were then simulated
 144 under literally the same ‘experimental’ conditions, using the initial conditions and parameter
 145 values obtained in prior studies. Data in this form was produced for eight different games,
 146 including the Prisoner’s Dilemma, Chicken, Battle of the Sexes, 2×2 and 3×3 Stag Hunt
 147 games and strategic form versions of Ultimatum bargaining and the Centipede games.
 148 There were systematic deviations between the simulated and the experimental results. The
 149 differences were, in fact, so large that the authors have called for new learning models firmly
 150 rooted in the experimental evidence and for new methodologies for evaluating them.

151 To get a feel for the differences in the outcomes of the learning models versus the out-
 152 comes of the experiments, consider their findings for the Prisoner’s Dilemma and the Battle
 153 of the Sexes (Fig. 1). Although AMP report only the average payoff profile obtained by each
 154 subject pair, it is possible to make some clear inferences about the path of actions chosen. In
 155 the Prisoner’s Dilemma, over half of the human data are tightly grouped around the payoffs
 156 of (8, 8), implying that the subjects have coordinated on the non-equilibrium action profile
 157 (A, A). The rest of the data points imply a mix of actions, such as exploiting a cooperator by
 158 defecting, as well as being exploited when cooperating. Fictitious play immediately con-
 159 verges to the unique Nash equilibrium with payoffs of (2, 2). Most reinforcement learning
 160 data points are scattered in an area relatively close to the Nash equilibrium. Although better
 161 results might be achieved with a recalibration of the models (as indeed we show below),
 162 there are clear qualitative differences between the predictions of the learning models and
 163 the behavior of most human subjects.

164 In the Battle of the Sexes, the majority of the experimental data points are closely scattered
 165 around the payoff profile (12, 12), implying that players coordinated by alternating between
 166 the two pure-strategy stage-game Nash equilibria. Fictitious play converges to one of the

	A	B
A	8 . 8	1 . 9
B	9 . 1	2 . 2

Prisoner’s Dilemma

	A	B
A	18 . 6	3 . 3
B	3 . 3	6 . 18

Battle of the Sexes

Fig. 1. Two games from AMP.

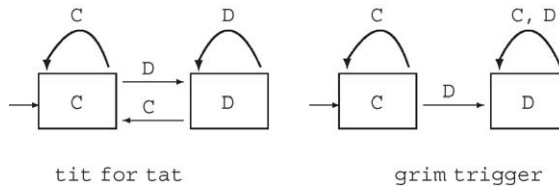


Fig. 2. Two examples of two-state automata.

167 two pure-strategy equilibria, while the data points for reinforcement learning are scattered
 168 between the three equilibria. Similarly, in the 2×2 Stag Hunt and in the game of Chicken,
 169 the majority of experimental subjects coordinate their strategies so as to maximize joint
 170 payoffs, while preserving an equitable distribution, a pattern of behavior that is difficult for
 171 action-learning models to replicate consistently.

172 3. Learning among repeated game strategies

173 Any analysis of learning among repeated game strategies requires some restriction on
 174 the space of available strategies. This is achieved here by restricting the complexity of the
 175 strategies available to players. One way of assessing the complexity of a repeated game
 176 strategy is on the basis of its representation as a finite automaton. The larger the number
 177 of states a strategy requires in automaton representation, the greater its complexity.⁶ This
 178 section starts with a brief description of the manner in which a repeated game strategy can
 179 be represented as a finite automaton. We then proceed to discuss the learning model in some
 180 detail.

181 3.1. Representing repeated game strategies with automata

182 An automaton is described by four components: a *set of states*, an *initial state* that the
 183 automaton occupies at the outset, an *output function* that indicates which action is to be
 184 taken in each particular state and a *transition function* that indicates which state will be
 185 reached in the next period given the current state and the current actions of the opponent.
 186 The current state of an automaton contains all information about the history of play that is
 187 relevant for the execution of the corresponding strategy.

188 The ‘tit-for-tat’ strategy in the Prisoner’s Dilemma can be represented as a two-state
 189 automaton (Fig. 2). The two states in this case are associated with the two available actions,
 190 cooperation and defection. The set of arrows are associated with the opponent’s actions and
 191 represent the transition function. The initial state is cooperation and the automaton stays
 192 in (or returns to) this state each time its opponent cooperates. It enters (or remains in) the
 193 defection state each time its opponent defects.

⁶ This approach to strategic complexity is also utilized, for instance, by Binmore and Samuelson (1992). Chapter 8 of Osborne and Rubinstein (1994) provides a good general introduction to the automaton representation of a repeated game strategy.

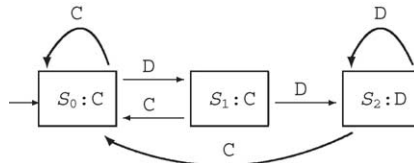


Fig. 3. Tit-for-two-tats.

194 Another strategy that can be represented as a two-state automaton is ‘grim trigger,’
 195 also shown in Fig. 2. Again, the two states are associated with the two available actions,
 196 cooperation and defection. The initial state is again cooperation, but the transition function
 197 is different. A defection by the opponent triggers a move to the defection state, which is
 198 absorbing (the automaton never leaves this state).

199 Each of the above strategies requires a memory of at most one action by the opponent.
 200 A slightly more complex strategy is ‘tit-for-two-tats,’ illustrated in Fig. 3. This strategy
 201 starts with cooperation and defects only if the opponent defects twice in a row. In order to
 202 implement this strategy, a three-state automaton is required. In the figure, the i -th state is
 203 denoted by S_i and the action to be taken in each state follows the colon after the state. The
 204 initial state is S_0 in which the player cooperates. If the opponent defects, state S_1 is reached.
 205 The player still cooperates in this state but remembers that the last action by the opponent
 206 was defection. Cooperation by the opponent when the player is in this state leads to a return
 207 to S_0 . On the other hand, if the opponent defects again, S_2 is reached and the player defects.
 208 Cooperation by the opponent when the player is in S_2 induces a return to S_0 .

209 In this paper, we restrict attention to strategies representable by one- or two-state au-
 210 tomata. When the stage game is 2×2 , this results in a total of 26 possible repeated game
 211 strategies.⁷

212 *3.2. Reinforcement learning among strategies*

213 Consider a population of players $\mathcal{N} = \{1, 2, 3, \dots, N\}$ and a specified symmetric
 214 2×2 stage game. During the first phase of learning, players drawn from the population are
 215 matched pairwise to play the stage game repeatedly. Players use the same strategy across
 216 several periods, but occasionally switch strategies as part of a process of experimentation.
 217 They are also randomly re-matched with other partners from time to time. Let $\rho \in (0, 1)$
 218 represent the probability that a player switches to a (possibly) new strategy at the start of
 219 any given period and $\mu \in (0, \rho)$ the probability that a player is randomly re-matched against
 220 a (possibly) different opponent at the start of any given period.

⁷ There are two possible states for an automaton and two possible actions by the opponent. The transition function thus maps four different possibilities of (own state, the opponent’s action) pairs into the set of two states to be taken in the next period. This generates $2^4 = 16$ cases in total. As each of these can have one of two initial states, we have a total 32 possible automata. Among these, however, four always play the first action and another four always play the second. Elimination of non-unique automata yields a total of 26, of which two are one-state and the rest two-state automata. Appendix B contains a complete listing of these. Appendices are available on the JEBO website.

221 Players have propensities or ‘attractions’ associated with each of their strategies and
 222 these attractions determine the probabilities with which strategies are chosen when players
 223 experiment. At the start of the first phase, all strategies have equal attraction and hence equal
 224 probability of being chosen.⁸ Learning takes place through the evolution of attractions:
 225 prior to updating her strategy, a player evaluates the performance of the strategy she has
 226 been utilizing and updates her attractions accordingly (the precise manner in which this
 227 occurs is described below). Since the learning process is defined over strategies, players
 228 are required to play the stage game a number of times to evaluate their current strategy,
 229 that is, to obtain information on their strategy’s payoff consequences. If ρ is not too large,
 230 meaningful evaluations of repeated game strategies are possible. Notice that unlike the
 231 action learning, players need not update their strategies simultaneously. When players are
 232 re-matched, they also update their strategies.⁹ This process continues until the limiting
 233 distribution of attractions is approximated, at which point the second ‘experimental’ phase
 234 begins. This consists of a fixed number of periods without further re-matching. Learning
 235 occurs also in this phase, building on the attractions generated during the pre-experimental
 236 phase.

237 Let $A_s^i(t)$ denote player i ’s attraction to the strategy $s \in S$ at period t , where $S = \{1, 2,$
 238 $3, \dots, 26\}$ is player i ’s set of 26 strategies. For each player, attractions are updated when
 239 the player updates her strategy. Only the strategy that was chosen at the previous strategy
 240 update is reinforced, as follows. Consider a player who updates her strategy choice at the
 241 start of period t and uses the same strategy $s \in S$ without further updates until the start of
 242 period $t + \tau$. Specifically, suppose that $s^i(t) = s^i(t+1) = \dots = s^i(t + \tau - 1)$, where $s^i(r)$ is the
 243 strategy used by player i in period r . Define the reinforcement value $R^i(t, t + \tau - 1)$ of the
 244 strategy used over the periods $t, \dots, t + \tau - 1$ as the average payoff obtained by player i over
 245 this period:

$$246 \quad R^i(t, t + \tau - 1) = \frac{1}{\tau} \sum_{r=t}^{t+\tau-1} \pi^i(r),$$

247 where $\pi^i(r)$ is the payoff obtained by player i in period r . When strategy revision next occurs
 248 (at the start of period $t + \tau$), player i ’s attraction or propensity for playing strategy s evolves
 249 as a weighted average of its previous value and the reinforcement value:

$$250 \quad A_s^i(t + \tau) = \begin{cases} (1 - \omega)A_s^i(t) + \omega R^i(t, t + \tau - 1) & \text{if } s = s^i(t) = \dots = s^i(t + \tau - 1), \\ A_s(t) & \text{otherwise.} \end{cases} \quad (1)$$

251
 252 Here $\omega \in (0, 1)$ is a weight placed on the reinforcement value, $R(\cdot, \cdot)$, which is the average
 253 payoff the player has obtained from using strategy s since the last strategy update, that is,
 254 between period t and $t + \tau - 1$.¹⁰

⁸ The assumption of equal initial attractions is not critical, since the first phase is long enough to ensure approximate convergence to the limiting distribution of attractions.

⁹ When a player is re-matched, she does not know the previous action played by the new opponent. Rather than assume that the old strategy is retained but enters its initial state; we assume that a strategy revision occurs.

¹⁰ Note that the attraction of each strategy approaches its historical average payoff as the number of updates becomes large.

255 The probability of a player i choosing strategy s , when she updates her strategy in the
 256 beginning of period t , depends on the attractions as follows:

$$257 \quad p_s^i(t) = \frac{e^{\lambda A_s^i(t)}}{\sum_{k \in S} e^{\lambda A_k^i(t)}}. \quad (2)$$

258 The parameter $\lambda \geq 0$ in the logistic transformation represents the extent to which strate-
 259 gies with higher attractions are favored in strategy choice. When $\lambda = 0$, all strategies are
 260 equally likely to be chosen, regardless of their attractions. As λ increases, strategies with
 261 higher attractions become disproportionately more likely to be chosen. In the limiting case
 262 $\lambda \rightarrow \infty$, the strategy with the highest attraction is chosen with probability one.

264 In the long horizon ‘pre-experimental’ phase of learning, the initial attraction, $A_s(0)$, for
 265 all strategies is set equal to the expected payoff given random choice of actions by both
 266 players. As an example, consider the Prisoner’s Dilemma game described in Fig. 1. Here a
 267 player’s initial attraction for each of her 26 strategies is given by $A_1(0) = \dots = A_{26}(0) =$
 268 $\frac{1}{4}(8 + 1 + 9 + 2) = 5$. The pre-experimental phase continues until the limiting distribution
 269 of attractions or probability weights placed on strategies are approximated. This is done
 270 as follows. Let $\bar{p}_s(\cdot)$ be the population average probability weight on strategy s in a given
 271 period and let $\bar{\bar{p}}_s(m)$ be the mean of the population average probability weight on strategy
 272 s over the m -th block of R periods:

$$273 \quad \bar{\bar{p}}_s(m) = \frac{1}{R} \sum_{t=R(m-1)+1}^{Rm} \bar{p}_s(t).$$

274 The convergence criterion employed in the simulation was
 275

$$276 \quad \frac{1}{|S|} \sum_{s \in S} |\bar{\bar{p}}_s(m) - \bar{\bar{p}}_s(m-1)| < \varepsilon$$

277 for several (specifically 20) consecutive m ’s. That is, the pre-experimental phase is termi-
 278 nated if the absolute difference between two consecutive means of the population average
 279 probability weights are, on average, less than ε for a long time.¹¹

280 The key idea behind the ‘pre-experimental’ phase is to eliminate strategies that are
 281 consistently poor performers. As long as ε is chosen to be sufficiently small, most of the
 282 initial strategies are practically eliminated. The set of surviving strategies is found to be
 283 insensitive to changes in ε beyond this point and our results are robust with respect to the
 284 choice of ε in this sense. Since our convergence criterion is based on population average
 285 probability weights, it is possible that such convergence masks offsetting movements by
 286 individuals that are not reflected in the changes in the population mean. For this reason,
 287 we require that individuals inherit their own probability weights (and not the population
 288 average) when they enter the experimental phase.

¹¹ The maximum length of the pre-experimental phase in each of the simulation runs has been set to 500,000 periods. In principle, it is possible to have a simulation run that does not satisfy the convergence criterion before the final period if ε is very small. However, we obtained convergence in all cases, using a value of $\varepsilon = 0.005$ and $R = 1000$.

289 In the second ‘experimental’ phase of learning, we assume that players bring with
 290 them to the laboratory the values of $A_s(\cdot)$ that they have reached at the conclusion of
 291 the first phase. In this latter phase, players are randomly paired to play and learn over the
 292 course of one match with 50 periods (without re-matching). This corresponds to McK-
 293 elvey and Palfrey’s experimental conditions, as reported in AMP. Note that the ‘experi-
 294 mental’ phase can be viewed as a simple continuation of the ‘pre-experimental’ phase,
 295 with the exception that there is no random re-matching of individuals. In this sense, the
 296 experiment itself is part of the string of experiences on the basis of which individuals
 297 learn.

298 There are total of four parameters in this model: the strategy updating rate ρ , the re-
 299 matching probability μ , the weight ω on reinforcement values in attraction updates and the
 300 sensitivity λ of the strategy choice to the attraction level in the logistic transformation. In
 301 addition, the number of players N needs to be large to ensure multiple interactions among
 302 various players in the pre-experimental phase.

303 4. Results

304 We have limited our attention to the four symmetric 2×2 games for which results are
 305 reported in the AMP, namely the 2×2 Stag Hunt, Prisoner’s Dilemma, Chicken and the
 306 Battle of the Sexes.¹² In the first phase of learning, attractions to strategies are ‘initialized’
 307 in anticipation of the second, experimental phase. The final values of $A_s^i(\cdot)$ from the pre-
 308 experimental phase are the initial values $A_s^i(0)$ for the experimental phase. This endogenizes
 309 the initial attractions or initial probability weights that appear as parameters in standard
 310 learning models.

311 We begin by describing results of the pre-experimental phase (the limiting distributions
 312 of probability weights across the 26 strategies) with a focus on strategies that obtain high
 313 limiting weights for a particular set of parameter values. We then proceed to discuss the
 314 results in the experimental phase. The set of parameter values are as follows: the strategy up-
 315 date rate $\rho = 0.05$, the re-matching probability $\mu = 0.02$, weights on the reinforcement values
 316 in attraction updates $\omega = 0.1$ and the sensitivity of the strategy choice to the attraction level
 317 $\lambda = 4$. Among the variety of parameter configurations with which we have experimented,
 318 this set of values provides the highest average performance for the four games considered
 319 here.¹³ Our focus is on qualitative performance, namely the ability of the model to repli-
 320 cate the broad contours of the experimental data with respect to the attainment of fair and
 321 efficient outcomes. Sensitivity of results to changes in parameter values are discussed in
 322 [Appendix C](#).¹⁴

¹² In order to make the game symmetric, the actions for the Column players in the Battle of the Sexes game have been re-labeled as shown in [Figs. 7 and 8](#) below.

¹³ We have experimented with all possible combinations of the following set of parameter values: $\rho = \{0.2, 0.1, 0.05\}$, $\omega = \{0.025, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.5, 0.7\}$, $\mu = \{0.02, 0.01, 0.005\}$ and $\lambda = \{2.5, 3.0, 3.5, 4.0\}$. For all simulations, the population size is kept constant at 1000.

¹⁴ Appendices are available on the JEBO website.

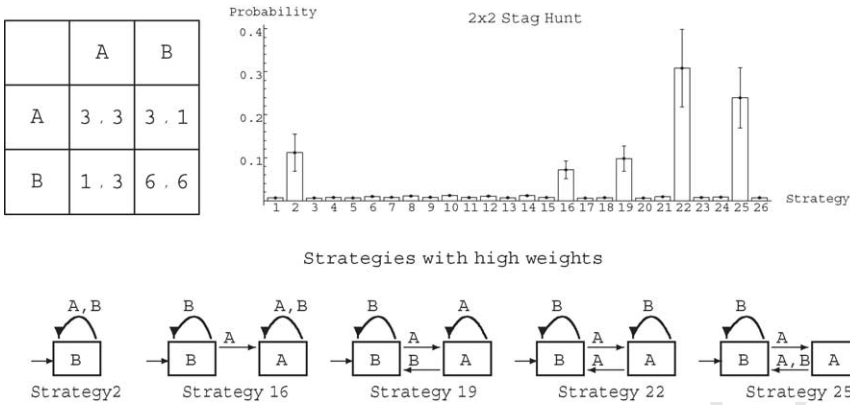


Fig. 4. Approximated limiting probability weights across 26 strategies and strategies with high probability weights for 2 × 2 Stag Hunt game. Strategy indices correspond to those in the Appendix B. Probability weights are averaged over 100 realizations. Error bars in the histogram represent two standard deviations around the mean. Parameter values are $\rho = 0.05$, $\mu = 0.02$, $\omega = 0.1$ and $\lambda = 4.0$.

323 4.1. Pre-experimental phase

324 What are the strategies that simulated players bring with them to the laboratory? In this
 325 section, we discuss the limiting distributions of probability weights in the pre-experimental
 326 phase to answer this question.

327 Figs. 4–7 show the approximate limiting distributions of probability weights for the 26
 328 strategies in each of the four games. (Appendix B contains the complete set of strategies
 329 in automata representation. The strategy indices referred to in the figures as well as in
 330 the text of this section correspond to those in this appendix.) To simplify the discussion,
 331 we focus on strategies with high limiting probability weight. A strategy is said to have a

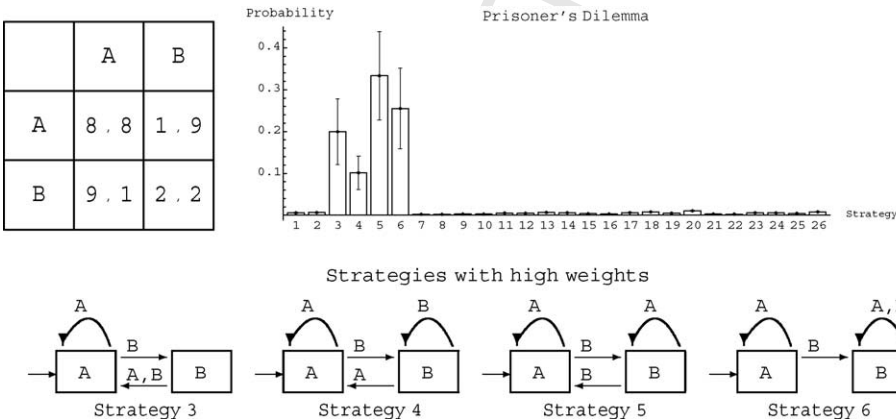


Fig. 5. Approximated limiting probability weights across 26 strategies and strategies with high probability weights for Prisoner's Dilemma.

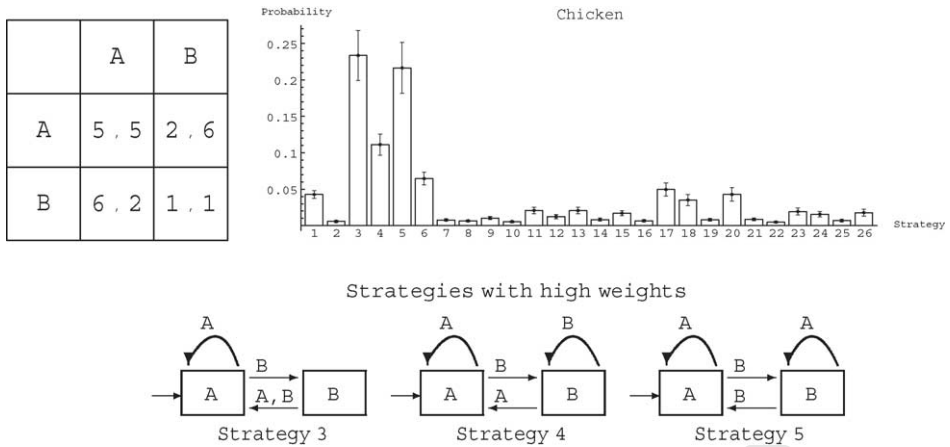


Fig. 6. Approximated limiting probability weights across 26 strategies and strategies with high probability weights for Chicken.

high limiting probability weight if its weight is at least one standard deviation above the average probability weight across all strategies. These strategies are, in a sense, the principal strategies that players “bring to the laboratory” for the experimental phase.

Results for the 2×2 Stag Hunt are shown in Fig. 4. The strategies with high limiting probability weight are ‘always play B,’ ‘grim trigger,’ ‘tit-for-tat,’ ‘punish until the opponent retaliates’ and ‘punish once,’ respectively. The initial state of each of these strategies is B. The first three do not require further explanation, since they have already been discussed in Section 3.1 above. Strategy 22 ‘punishes until the opponent retaliates’: it starts by playing B and stays in this state as long as the opponent also plays B. Once the opponent plays A, however, it switches to playing A. It returns to B only if the opponent

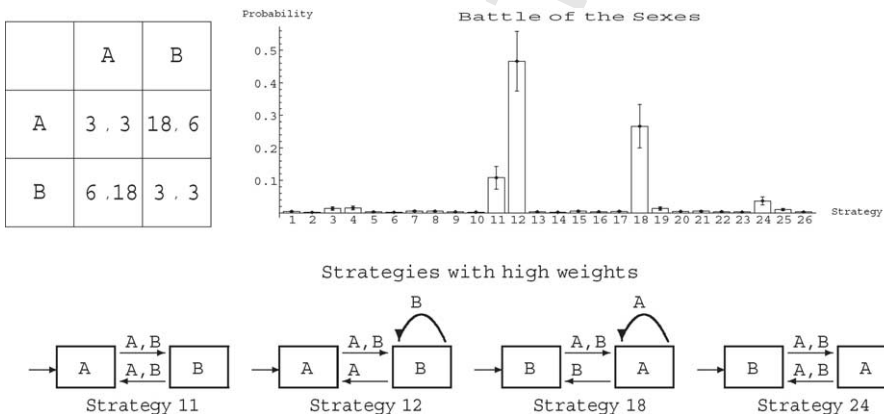


Fig. 7. Approximated limiting probability weights across 26 strategies and strategies with high probability weights for Battle of the Sexes.

342 plays A; otherwise it stays in state A. Strategy 25, which ‘punishes once,’ also starts with
343 action B. It stays in state B unless the opponent plays A. Once the opponent plays A, A,
344 state A is reached for exactly one period, after which the strategy returns to B regard-
345 less of the opponent’s action. Notice that if these five strategies are matched against each
346 other, we will observe all players playing action B forever to achieve the efficient and fair
347 outcome.

348 Fig. 5 shows the outcome of the pre-experimental phase for the Prisoner’s Dilemma.
349 Strategies obtaining high limiting probability weights are ‘punish once,’ ‘tit-for-tat,’ ‘punish
350 until the opponent retaliates’ and ‘grim trigger,’ respectively. The initial state is A for these
351 strategies. As in the case of the 2×2 Stag Hunt game, if these four strategies are played
352 amongst themselves, the observed history of actions will involve mutual cooperation in all
353 periods.

354 Since the Prisoner’s Dilemma has received such widespread attention in economics, the
355 strategies that emerge from the learning process in our model deserve further discussion. The
356 ‘tit-for-tat’ strategy was the winner in two tournaments organized by [Axelrod \(1984\)](#) and
357 has been a subject of extensive study, especially in the context of evolutionary game theory.
358 [Axelrod and Hamilton \(1981\)](#) have shown that ‘tit-for-tat’ is a neutrally stable strategy in the
359 infinitely repeated prisoners’ dilemma with payoffs evaluated according to the limit-of-the-
360 means criterion. This has been interpreted as providing theoretical support for the hypothesis
361 that cooperation sustained through reciprocation is an inevitable outcome of evolutionary
362 process. However, there are a large number of other repeated game strategies that are also
363 neutrally stable and some of them involve mutual defection in most periods. The set of stable
364 strategies can be refined substantially by introducing complexity costs (as in [Binmore and
365 Samuelson](#)) or the possibility of errors in the implementation of strategies (as in [Fudenberg
366 and Maskin, 1990](#)). These refinements result in a prediction of mutual cooperation in the
367 infinitely repeated prisoners’ dilemma, although on the basis of strategies other than ‘tit-
368 for-tat.’ We also find mutual cooperation to be the predicted outcome, although the model
369 considered here is one of finite repetition and bounded complexity. The ‘tit-for-tat’ strategy
370 survives but does not have the highest limiting probability weight: we find ‘punish until the
371 opponent retaliates’ to be the most prolific strategy.

372 Results from the game of Chicken are summarized in [Fig. 6](#). In this game, the strategies
373 ‘punish once,’ ‘tit-for-tat’ and ‘punish until the opponent retaliates’ are the ones with the
374 high probability weights. Again, if these strategies are played only among themselves, we
375 will observe only the efficient and fair outcome as both players continue to play action
376 A. Unlike the two games discussed above, however, there are a few other strategies with
377 non-negligible probability weights. The presence of these strategies will generate outcomes
378 that are not efficient in the experimental phase as we show below.

379 [Fig. 7](#) shows strategies with high weights in the Battle of the Sexes. Strategies 11 and 24
380 mechanically alternate between two states, but have different initial states. Strategy 12 starts
381 by playing action A and regardless of the opponent’s action, switches to state B. It stays in
382 this state if the opponent plays B; otherwise it returns to state A. Strategy 18 is qualitatively
383 identical to strategy 12, with the role of the two actions reversed. Note that strategies 12
384 and 18, when matched with any of the other strategies having significant weights, eventu-
385 ally achieve perfect alternation between the two pure strategy stage-game equilibria. Since
386 these two ‘flexible’ strategies have the highest weights, players in the experimental stage

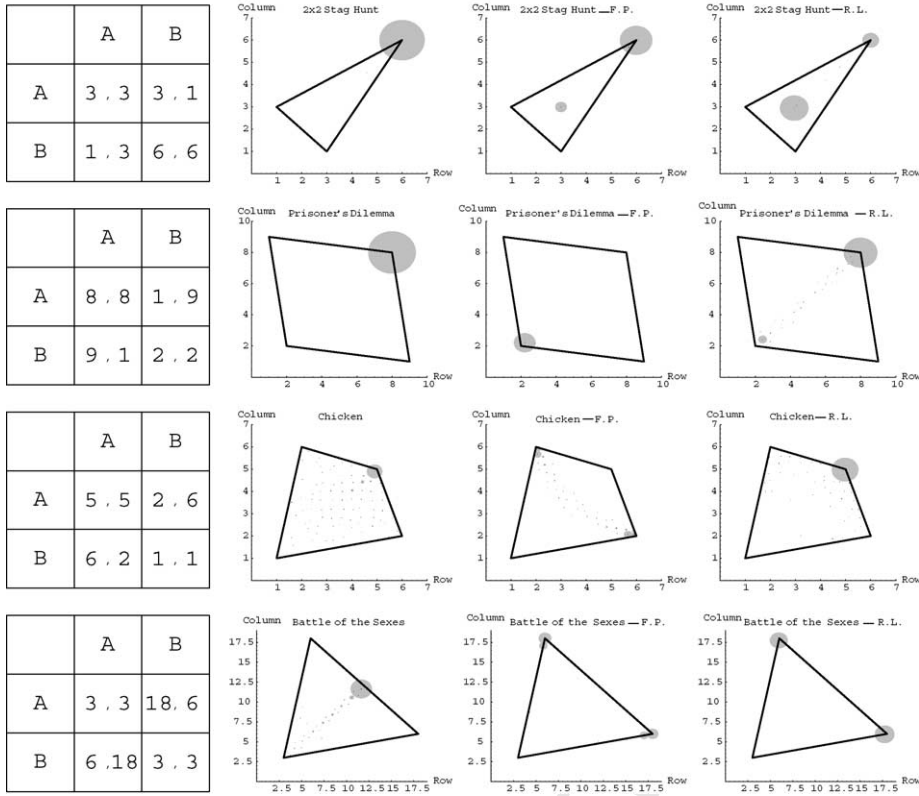


Fig. 8. Comparison of the simulation outcomes among three models of learning: reinforcement learning among repeated game strategies (left), smoothed fictitious play among actions (center) and reinforcement learning among actions (right). $\lambda = 4.0$ for all the models. Other parameters for the first model are set as $\{\rho, \mu, \omega\} = \{0.05, 0.02, 0.1\}$. The polygons in each plots represent the possible per stage average payoff space. Payoff for row (column) players are in the horizontal (vertical) axes. Each point corresponds to payoff profile for a pair and the gray circles around the points represent the relative likelihood of observing the outcome in the center of the circles. Each figure is based on a total of 1000 data points. (The figures for the first model are generated by running two simulations with 1000 players. Each simulation generates 500 data points.) Note that the sizes of circles are only comparable across the three models for the same game and not across the games for the same model.

387 eventually learn to use one of them when initial actions are not coordinated on the equitable
 388 and efficient outcome.

389 Just by examining the strategies with high limiting probability weights from the first
 390 phase of learning, we can expect that efficient and fair outcomes will be observed in the
 391 experimental phase. We now turn to discussion of our results for this phase.

392 *4.2. Experimental phase*

393 Fig. 8 shows the results of simulation runs for the four games we have considered in
 394 the experimental phase. Also presented in the figure are results from two learning mod-

els, smoothed fictitious play and action-reinforcement learning. (See Appendix A for the algorithm used for generating the data for these models.¹⁴)

For three of the four games (2×2 Stag Hunt, Prisoner's Dilemma and Chicken), the efficient outcome is also fair. A majority of the simulated players who learn among repeated game strategies are successful in obtaining such an outcome.¹⁵ They achieve this by repeatedly playing the action profiles $\{B, B\}$ in the 2×2 Stag Hunt and $\{A, A\}$ in the Prisoner's Dilemma, and Chicken. In contrast, fictitious play generates a stage-game Nash equilibrium outcome as the theory predicts. Fictitious play thus generates an outcome that corresponds to the experimental data only in the coordination game (2×2 Stag Hunt), in which the efficient outcome is one of the two pure-strategy stage-game Nash equilibria. Also note that the Pareto superior equilibrium is more likely to be observed than the Pareto inferior one under fictitious play. Reinforcement learning among actions generates an efficient outcome that is not the stage-game Nash equilibrium in Prisoner's Dilemma and Chicken, but it is more likely to result in the Pareto inferior equilibrium in the coordination game. This is an interesting contrast between the two models and deserves further investigation.

In the Battle of the Sexes game, the efficient and fair outcome requires coordinated alternation between two pure-strategy stage-game Nash equilibria. The results for this game are particularly striking. The simulation outcome shows that many of the players successfully learned, out of 26 possible repeated game strategies, to play the strategies that enable them to achieve the efficient and fair outcome. As one can clearly see in the figure, neither of the two action-learning models generates such an outcome. Efficiency and fairness arises in our model in one of two ways. Players may initially adopt strategies that quickly result in convergence to coordinated alternation (if strategy 12 were matched against strategy 11, for instance; see Fig. 8). Alternatively, when initial choices fail to achieve coordination (if both initially adopt strategy 11, for instance) players successfully learn to switch to other strategies. The failure to coordinate eventually induced one of them to experiment with one of the other strategies. After the switch, convergence to efficient alternation occurs rapidly. The probabilistic nature of strategy choice, however, can cause the players to mismatch even after several periods of successful alternation. This results in a fluctuation of payoffs in later periods. The Battle of the Sexes is a striking example of a game in which our approach predicts outcomes that are both consistent with experimental observation and virtually impossible to replicate with action-learning models.

5. Conclusion

We have demonstrated that a simple reinforcement model of learning applied to a restricted set of repeated game strategies can account for the behavior of human subjects in environments, where fairness and reciprocity seem to play a significant role. We have done so without assuming that fairness and reciprocity are primitive concerns. Our results may also be of some interest from the perspective of the problem of equilibrium selection in games. In pure coordination games, where fairness and efficiency are not in conflict, our

¹⁵ The results for Chicken are somewhat weaker than those for the other games.

434 findings predict that learning will converge to the efficient action profile. In the Battle of
 435 the Sexes, where efficient stage-game equilibria are unfair, the model predicts alternation
 436 over time to achieve a profile of average payoffs that is both efficient and fair. In the Pris-
 437 oner's Dilemma, where fairness and efficiency are not in conflict but cannot be attained in
 438 equilibrium, the model predicts convergence to non-equilibrium strategy profiles.

439 One important direction for further research would be to study the feasibility of our
 440 approach in settings of greater complexity, with a larger set of players and stage-game
 441 actions. A potential empirical extension is an analysis of the goodness-of-fit of the model
 442 to the large and varied experimental data that is available. This would require estimation of
 443 the model parameters and out-of-sample comparisons with other learning models.

444 Finally, it would be well worth developing a deeper analytical understanding of the
 445 process by which learning on the basis of material payoffs can result in behavior that
 446 appears to be motivated by fairness and efficiency concerns. A characterization of the class
 447 of games for which the learning dynamics converge to fair and efficient payoff profiles
 448 would be of considerable interest.

449 Acknowledgements

450 We thank Atila Abdulkadiroglu, Alessandra Casella, Peter Dodds, Dale Stahl, Duncan
 451 Watts, seminar participants at Columbia University, Conference participants at USC and an
 452 anonymous referee for their comments and suggestions, and Gueorgi Kossinets and Sibel
 453 Sirakaya for computational advice.

454 Appendix A. Fictitious play and action-reinforcement learning

455 We provide here a brief discussion of two standard learning models. As shown in Camerer
 456 and Ho, both fictitious play and reinforcement learning model can be considered as special
 457 cases of the experience-weighted attraction (EWA) learning model. The following formu-
 458 lation is a simplified version of the EWA model.

459 Let $A_a^i(t)$ be player i 's attraction to the action $a \in S$ at period t , where S is player i 's action
 460 set. For each player, attractions evolve over time as weighted averages of their previous
 461 values and current reinforcement values. Let $a_{-i}(t)$ be the actions chosen by a player's
 462 opponents, denoted by $-i$, at period t . The player's attraction to action a evolves as follows:

$$463 \quad A_a^i(t+1) = (1 - \omega_a^i(t+1))A_a^i(t) + \omega_a^i(t+1)\pi^i(a, a_{-i}(t)). \quad (3)$$

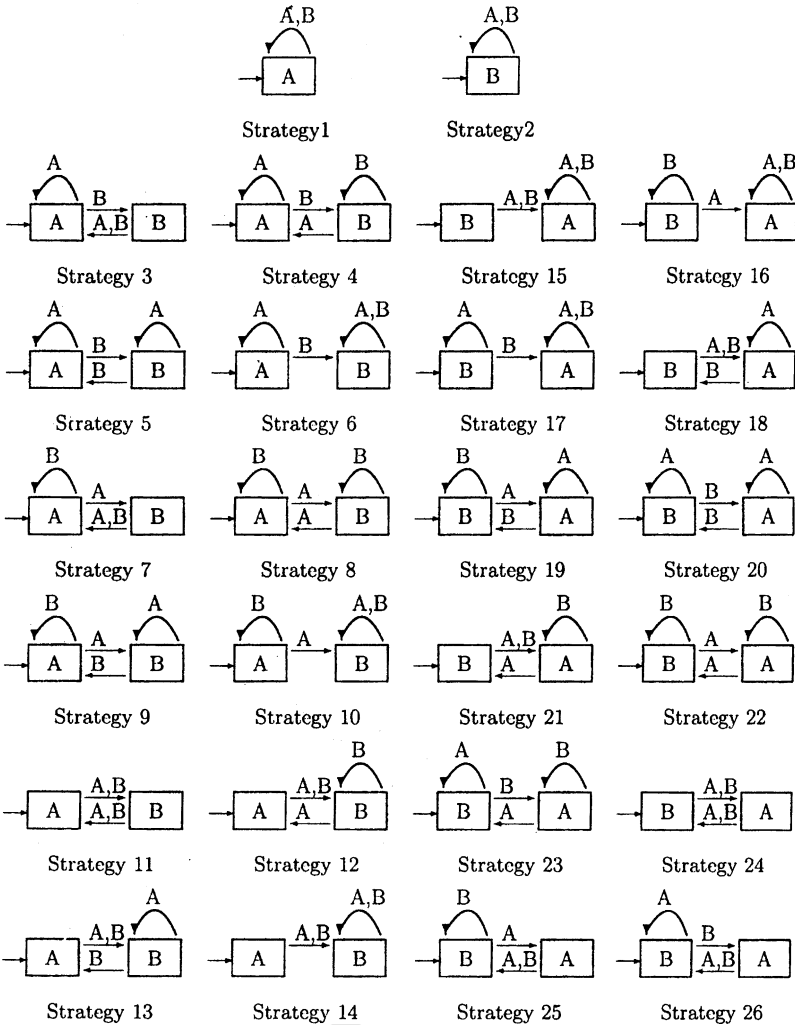
464 It is easy to verify that fictitious play is equivalent to having $\omega_a^i(t+1) = 1/(t+1)$ for all
 465 a . A reinforcement learning model can be obtained by setting

$$466 \quad \omega_a^i(t+1) = \begin{cases} \frac{1}{n_a(t+1)} & \text{if action } a \text{ is chosen in period } t \\ 0 & \text{otherwise} \end{cases}$$

467 where $n_a(\cdot)$ is the total number of times the action a has been chosen since the beginning
 468 of play plus its initial value $n_a(0)$. In the simulation in Section 4.2, we assume no 'pre-

469 experimental' learning for these models, as in the previous literature. We set $n_a(0)$ equal
 470 to 1 and the initial attraction for all actions, $A_a^i(0)$, is set to be the expected payoff given a
 471 random choice by both players. The probability with which each action is chosen is given
 472 by Eq. (2) in Section 3.2.

473 **Appendix B. The complete set of strategies considered**



474 **Appendix C. Sensitivity of results to parameter changes**

475 The model has four parameters: the strategy update rate ρ , the re-matching rate μ , the
 476 weight ω on reinforcement values in the updating of attractions and the sensitivity λ of

Table 1

Population mean per-period average payoff (mean payoff) relative to the payoff in the efficient and fair outcome for each of the game

ρ/μ	0.02	0.01	0.005
2 × 2 Stag Hunt			
0.2	0.999	1.000	0.995
0.1	0.999	1.000	0.999
0.05	0.999	1.000	0.999
Prisoner's Dilemma			
0.2	0.575	1.585	0.602
0.1	1.000	1.000	0.981
0.05	1.000	1.000	0.999
Chicken			
0.2	0.739	0.739	0.742
0.1	0.765	0.769	0.776
0.05	0.868	0.854	0.830
Battle of the Sexes			
0.2	0.671	0.686	0.704
0.1	0.715	0.677	0.694
0.05	0.837	0.836	0.630

λ and ω are set equal to 4.0 and 0.1, respectively.

strategy choice to attractions. (The population size N , as long as it is sufficiently large, has no effect on the results.) We have experimented with all possible combinations of the following set of parameter values: $\rho = \{0.2, 0.1, 0.05\}$, $\omega = \{0.025, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.5, 0.7\}$, $\mu = \{0.02, 0.01, 0.005\}$ and $\lambda = \{2.5, 3.0, 3.5, 4.0\}$. For all simulations, the population size is kept constant at 1000. The two parameters of greatest economic interest are the re-matching rate μ and the frequency of experimentation ρ . We shall discuss here the sensitivity of our results to changes in these two parameters, holding constant $\omega = 0.1$ and $\lambda = 4.0$.

Recall that our measure of performance is qualitative, namely the extent to which the model's predictions fit the broad contours of the experimental data. For the games considered here, experimental subjects appear to coordinate frequently on patterns of play that produce efficient and fair average payoffs over time. One way to track the sensitivity of the model to parameter changes is therefore, to look at deviations of average payoffs from the efficient level. This is done in Table 1, which summarizes the performance of the model for various pairs of ρ and μ . The entries in the table refer to the ratio of the average payoff to the efficient payoff. The parameter pair for which results are reported in the text is in bold.

The results for the 2 × 2 Stag Hunt are relatively insensitive to changes in strategy updating rates and the re-matching probability. We suspect that this is because convergence occurs to a Nash equilibrium of the stage game, as is predicted also by most action-learning models. The other games show different degrees of sensitivity. Results for the Prisoner's Dilemma are quite sensitive to the strategy update rate ρ . If the strategy update rate is high, that is, when $\rho = 0.2$, the mean payoff becomes much lower relative to the efficient outcome. This is because if players update their strategies too frequently, reciprocation strategies such as 'tit-for-tat' lose the opportunity to punish defectors for many periods

in order to discourage others from defecting in the future. Therefore, a high ρ leads to the situation in which players learn to use strategies that involve more defections. The re-matching probability, on the other hand, does not have a strong impact on the result (provided that it remains much lower than the experimentation rate). In the game of Chicken, as in the Prisoner's Dilemma, a higher strategy update rate prevents players from achieving the efficient outcome. The re-matching probability exercises a modest influence in this game, with higher rates of re-matching tending to result in lower payoffs.

The environment in which results are most sensitive to parameter values is the Battle of the Sexes. Results are affected both by the strategy updating rate and the re-matching probability to a greater extent than in the other three cases and these effects are rather complex. When updating is infrequent ($\rho = 0.05$), frequent re-matching is helpful in generating alternation. On the other hand, when updating is itself more frequent, then more frequent re-matching can lead to declines in efficiency. Hence, the learning of successful alternation strategies in this setting depends quite critically on the conditions under which pre-experimental learning occurs. This sensitivity we attribute to the fact that the attainment of fair and efficient outcomes require a pattern of alternation over time that is harder to learn than the repetition of a single action profile.

References

- Arifovic, J., McKelvey, R.D., Pevnitskaya, S., 2002. An initial implementation of the turing tournament to learning in two person games, mimeo. California Institute of Technology.
- Axelrod, R., 1984. *The Evolution of Cooperation*. Basic Books, New York, NY.
- Axelrod, R., Hamilton, W.D., 1981. The evolution of cooperation. *Science* 211, 1390–1396.
- Binmore, K.G., Samuelson, L., 1992. Evolutionary stability in repeated games played by finite automata. *Journal of Economic Theory* 57, 278–305.
- Bolton, G., Ockenfels, A., 2000. *Erc: A theory of equity, reciprocity, and competition*. *American Economic Review* 90, 166–193.
- Camerer, C., Ho, T.-H., 1999. Experience-weighted attraction learning in normal form games. *Econometrica* 7, 827–874.
- Charness, G., Rabin, M., 2002. Understanding social preferences with simple tests. *Quarterly Journal Of Economics* 117, 817–869.
- Cheung, Y.-W., Friedman, D., 1997. Individual learning in normal form games: some laboratory results. *Games and Economic Behavior* 19, 46–76.
- Crawford, V.P., 1995. Adaptive dynamics in coordination games. *Econometrica* 63, 103–143.
- Day, R.H., 1963. *Recursive Programming and Production Response*. North-Holland, Amsterdam.
- Erev, I., Roth, A.E., 1998. Predicting how people play games: reinforcement learning in experimental games with unique, mixed strategy equilibria. *American Economic Review* 88, 848–881.
- Erev, I., Roth, A.E., 2001. Simple reinforcement learning models and reciprocation in the prisoner's dilemma game. In: Gigerenzer, G., Selten, R. (Eds.), *Bounded Rationality: The Adaptive Toolbox*. MIT Press, Cambridge, MA, pp. 215–231.
- Fehr, E., Schmidt, K.M., 1999. A theory of fairness, competition, and cooperation. *The Quarterly Journal of Economics* 114, 817–868.
- Fudenberg, D., Levine, D.K., 1998. *The Theory of Learning in Games*. MIT Press, Cambridge, MA.
- Fudenberg, D., Maskin, E.S., 1990. Evolution and cooperation in noisy repeated games. *American Economic Review Papers and Proceedings* 80, 274–279.
- Gintis, H., 2000. Strong reciprocity and human sociality. *Journal of Theoretical Biology* 206, 169–179.

- 547 Giith, W., Yaari, M., 1992. Explaining reciprocal behavior in simple strategic games: An evolutionary approach.
548 In: Witt, U. (Ed.), *Explaining Process and Change: Approaches to Evolutionary Economics*. University of
549 Michigan Press, Ann Arbor, MI, pp. 23–34.
- 550 Huck, S., Oechssler, J., 1999. The indirect evolutionary approach to explaining fair allocations. *Games and Eco-*
551 *nomic Behavior* 28, 13–24.
- 552 Levine, D.K., 1998. Modeling altruism and spitefulness in experiments. *Review of Economic Dynamics* 1,
553 593–622.
- 554 Lindgren, K., 1997. Evolutionary dynamics in game-theoretic models. In: Arthur, W.B., Durlauf, S.N., Lane, D.A.
555 (Eds.), *The Economy as an Evolving Complex System II*. Perseus Books, Reading, MA, pp. 337–367.
- 556 McKelvey, R.D., Palfrey, T.R., 2001. *Playing in the dark: Information, learning, and coordination in repeated*
557 *games*, mimeo. California Institute of Technology.
- 558 Miller, J.H., 1996. The coevolution of automata in the repeated prisoner's dilemma. *Journal of Economics Behavior*
559 *and Organization* 29, 87–112.
- 560 Mookherjee, D., Sopher, B., 1997. Learning and decision costs in experimental constant sum games. *Games and*
561 *Economic Behavior* 19, 97–132.
- 562 Osborne, M.J., Rubinstein, A., 1994. *A Course in Game Theory*. The MIT Press, Cambridge, MA.
- 563 Sethi, R., Somanathan, E., 2001. Preference evolution and reciprocity. *Journal of Economic Theory* 97, 273–297.
- 564 Stahl, D.O., 1999. Evidence based rules and learning in symmetric normal-form games. *International Journal of*
565 *Game Theory* 28, 111–130.
- 566 Stahl, D.O., 2000. Rule learning in symmetric normal-form games: theory and evidence. *Games and Economics*
567 *Behavior* 32, 105–138.
- 568 Stahl, D.O., Haruvy, E., 2002. Aspiration-based and reciprocity-based rules in learning dynamics for symmetric
569 normal-form games. *Journal of Mathematical Psychology* 46, 531–553.