

Contents lists available at [SciVerse ScienceDirect](http://SciVerse.ScienceDirect.com)

Games and Economic Behavior

www.elsevier.com/locate/geb


Learning in experimental 2×2 games [☆]

 Thorsten Chmura ^a, Sebastian J. Goerg ^{b,c,*}, Reinhard Selten ^d
^a Nottingham University Business School, Nottingham, NG8 1BB, UK

^b Max Planck Institute for Research on Collective Goods, 53113 Bonn, Germany

^c Department of Economics, Florida State University, Tallahassee, FL 32306, USA

^d BonnEconLab, Department of Economics, University of Bonn, 53113 Bonn, Germany

ARTICLE INFO

Article history:

Received 31 March 2011

Available online 30 June 2012

JEL classification:

C72

C91

C92

Keywords:

 2×2 games

Experimental data

Learning

Impulse-matching

Self-tuning EWA

Reinforcement

Action-sampling

ABSTRACT

In this paper, we introduce two new learning models: action-sampling learning and impulse-matching learning. These two models, together with the models of self-tuning EWA and reinforcement learning, are applied to 12 different 2×2 games and their results are compared with the results from experimental data. We test whether the models are capable of replicating the aggregate distribution of behavior, as well as correctly predicting individuals' round-by-round behavior. Our results are two-fold: while the simulations with impulse-matching and action-sampling learning successfully replicate the experimental data on the aggregate level, individual behavior is best described by self-tuning EWA. Nevertheless, impulse-matching learning has the second-highest score for the individual data. In addition, only self-tuning EWA and impulse-matching learning lead to better round-by-round predictions than the aggregate frequencies, which means they adjust their predictions correctly over time.

© 2012 Elsevier Inc. All rights reserved.

1. Introduction

It is well known that rational learning, in the sense of Bayesian updating, leads to the stationary points of the Nash equilibrium (e.g., [Kalai and Lehrer, 1993](#)). But it is also known that actual human behavior does not necessarily converge to Nash equilibrium. In fact, a vast body of literature indicates situations in which standard theory does not perform as a good predictor for subjects' behavior in experiments (e.g., [Brown and Rosenthal, 1990](#); [Erev and Roth, 1998](#)).

A recent publication by [Selten and Chmura \(2008\)](#) documents the predominance of behavioral stationary concepts regarding descriptive power. In the paper, the concepts of impulse-balance equilibrium ([Selten and Chmura, 2008](#)), payoff-sampling equilibrium ([Osborne and Rubinstein, 1998](#)), action-sampling equilibrium ([Selten and Chmura, 2008](#)), and quantal response equilibrium ([McKelvey and Palfrey, 1995](#)) outperform Nash equilibrium in describing the decisions of a population in twelve completely mixed 2×2 games. Moreover, payoff-sampling equilibrium and action-sampling equilibrium perform better than quantal response equilibrium does. In addition, the parameter-free concept of impulse-balance equilibrium performs equally well as the parametric concept of quantal response equilibrium.¹ Furthermore, [Goerg and Selten \(2009\)](#) show that the

[☆] Financial support by the Deutsche Forschungsgemeinschaft and the German–Israeli Foundation for Scientific Research and Development is gratefully acknowledged. We would like to thank Emanuel Castillo and Christian Stracke for excellent research assistance. Christoph Engel, Andreas Nicklisch, and two anonymous referees provided very helpful comments.

* Corresponding author at: Department of Economics, Florida State University, 113 Collegiate Loop, Tallahassee, FL 32306-2180, USA.

E-mail address: goerg@coll.mpg.de (S.J. Goerg).

¹ For further discussions, please refer to [Brunner et al. \(2011\)](#) and [Selten et al. \(2011\)](#).

advantage of impulse-balance equilibrium over Nash equilibrium is not limited to 2×2 games, but also present in cyclic duopoly games.

Presumably stationary behavior is a result of a learning process converging to a stationary distribution of actions for both players, which is, as the above studies demonstrate, not necessarily the Nash equilibrium. Therefore, constructing and testing simple learning models with the predicted stationary states of the better-performing concepts suggests itself. For the two behavioral stationary concepts of action-sampling equilibrium and impulse-balance equilibrium this is quite easy: both yield precise expression for stationary behavior.

The main purpose of this article is to introduce two new learning models, action-sampling learning and impulse-matching learning, which are based on the behavioral reasoning of action-sampling equilibrium and impulse-balance equilibrium, and test them in the environment of twelve repeated 2×2 games with mixed equilibria. Here, the learning rules have to meet two challenges: first, do they reproduce in simulations the aggregate behavior of a human population, and second, can they adequately describe the observed round-by-round behavior of a single individual? The difference between these two approaches is that in the simulations the learning curves of a population of agents is simulated, while for the analyses of the individual round-by-round behavior the actual experimental data from previous rounds is used to predict the next period.

For comparison, we include the models of reinforcement learning (Erev and Roth, 1998) and self-tuning experience-weighted attraction learning (Ho et al., 2007). We decided to compare our results with reinforcement learning, as it was the first model with application of rote learning to economics and it is the most cited learning model in economics. Self-tuning EWA was selected to cover a broader set of different learning variants, as it can describe weighted fictitious play, averaging reinforcement learning, and models in between.²

For the replication of the populations' behavior, we conduct simulations with the four learning models and the twelve 2×2 games experimentally investigated in Selten and Chmura (2008). The simulations replicate the exact situation of the experiments. In each simulation run, eight agents, four deciding as row players and four deciding as column players, are randomly matched in each round over 200 rounds. In each simulation run, one game is played and one learning model is applied. To judge the predictive power of the simulations on the population level, we compare the distribution of choices in the simulation runs with the experimental data from Selten and Chmura (2008).

To investigate how well the learning models predict individuals' round-by-round behavior, we separately evaluate the explanatory power of the learning models for each participant of the 2×2 experiments. For each of the 864 subjects, we compare the actual decision in every round with the decision predicted by the learning model given the subject's history. To judge the power of the learning models, we introduce three benchmarks which all learning models should beat. The first benchmark is the inertia rule, which predicts for each round the same choice as executed in the round before. The second benchmark is a random play with equal probability for each of the two decisions. In addition, if the learning theories describe subjects' behavior correctly over time, their predictions should be more accurate than the observed aggregate frequencies. Thus, as a more demanding benchmark, we include the empirical frequencies as a criterion.

As Erev et al. (2010) state, there are three obstacles for the learning literature: (1) small data sets, (2) problems of overfitting (Salmon, 2001; Hopkins, 2002), and (3) relatively small sets of models. We try to address these issues by (1) using the large data set of Selten and Chmura (2008) with twelve 2×2 games played by 864 subjects; (2) using theories with at most one parameter, adjusting the parameters over all games and applying only non-parametric analysis; and (3) applying and testing four different learning models. In fact, the results presented in this paper are the condensed summary of the analyses with 7 learning models. In addition to the already mentioned models, we introduce and test the concepts of payoff-sampling learning and impulse-balance learning. Because both concepts perform worse than the inertia benchmark and the random play benchmark, we do not cover these two learning models in more detail. More information about these two models can be found in Appendix A. To shed some additional light on the performance of self-tuning EWA, we include a non-parametric version of self-tuning EWA into our analyzes. We will refer to these results in the discussion. Additional information about the omitted concepts as well as the comparison of all 7 learning models can be found in Appendices A, B (Figs. 16 and 17), C (Figs. 18–23) and D (Tables 3–6).

Our results are twofold: our newly introduced models are able to capture the distribution of decisions of the experimental population much better in simulations than self-tuning EWA and reinforcement do, while self-tuning EWA describes the individual data in a much more accurate way. In the simulations, the learning models of impulse-matching learning and action-sampling learning have the smallest distance to the experimental data, while the concepts of self-tuning EWA and reinforcement learning have relatively high distances to the data. On the individual level with the actual individuals' histories, self-tuning EWA and impulse-matching have the highest scores. In addition, these two learning concepts are the only concepts that perform significantly better at describing individual round-by-round behavior than the overall empirical frequencies.

² We decided not to include a pure version of fictitious play in our analyses since a population of fictitious players would converge in the 2×2 games to the Nash equilibrium (Miyasawa, 1961, and Metrick and Polak, 1994), which is clearly outperformed by the stationary concepts of impulse-balance equilibrium and action-sampling equilibrium (Selten and Chmura, 2008). However, with action-sampling learning and self-tuning EWA, variants of fictitious play are included into our analyzes. Furthermore, in this paper we solely focus on learning rules with at most one parameter. Thus, more elaborate versions of fictitious play with additional parameters like the three parameter model by Cheung and Friedman (1997) and the six parameter model of Chen et al. (2011) are ignored.

2. The learning models

In the following, we will introduce impulse-matching learning and action-sampling learning, which are based on the behavioral stationary concepts discussed in Selten and Chmura (2008). In addition to the new learning models, the more established concepts of reinforcement learning (c.p. Erev and Roth, 1998) and self-tuning EWA (Ho et al., 2007) are briefly explained.

Two of the discussed models, namely action-sampling learning and self-tuning EWA, are parametric concepts. In case of action-sampling learning, the parameter is the sample size. Self-tuning EWA is based on the multi-parametric concept of experience-weighted attraction learning (Camerer and Ho, 1999). Self-tuning EWA replaces two of the parameters with numerical values and two with functions. The remaining parameter λ “measures sensitivity of players to attractions” (Camerer and Ho, 1999, p. 835). The version of reinforcement learning examined here does not have a parameter and the initial propensities are not estimated from the data. All investigated learning rules will start with randomization of 0.5 in the first round. Only after all necessary information has been gathered does the corresponding learning rule determine the following decisions.³

The parameters of the parametric concepts are estimated to lead to the best fit over all data and over all games. For more details about the parameter estimation on the aggregate level, refer to the results in Section 4.2; for details on the estimation of the parameters on the individual level, refer to results in Section 5.2.

2.1. Impulse-matching learning

Impulse-matching learning relates to the concepts of impulse-balance equilibrium (Selten et al., 2005 and Selten and Chmura, 2008) and learning direction theory (Selten and Stoecker, 1986 and Selten and Buchta, 1999). After a decision and after the realization of the payoffs, the behavior is adjusted to experience. Selten and Buchta explain the concept by the example of a marksman aiming at a trunk: “If he misses the trunk to the right, he will shift the position of the bow to the left and if he misses the trunk to the left he will shift the position of the bow to the right. The marksman looks at his experience from the last trial and adjusts his behavior [...]” (Selten and Buchta, 1999, p. 86). Impulse-balance equilibrium and impulse-matching learning overcome the limitation of learning direction theory defining a direction only for ordered strategies, e.g., increasing a bid in an auction (cf. Ho et al., 2007) by shifting the probabilities of single actions.

To understand how impulse-matching learning works, suppose that in a period the first of two actions has been chosen and that this action was not the best reply to the action played by the other player. Then the player receives an *impulse* towards the second action. Originally, an impulse was defined as the difference between the payoff the player could have received for his best reply minus the payoff actually received given the decision by the other player in this period. However, the theory of impulse-matching learning is based on another impulse concept. Here, a player always receives an impulse from the action with the lower payoff to the one with the higher payoff. The resulting learning model is similar to the regret-based learning models, which have already been successfully tested by Marchiori and Warglien (2008). The name impulse-matching is due to the fact that this kind of learning leads to probability matching by a player if the probabilities p_1 and $(1 - p_1)$ on the other side are fixed, and the payoffs for the player is one if both players play the strategy with the same number (one or two) and zero otherwise (cf. Estes, 1954).

To incorporate loss aversion, the impulses are not calculated with the original payoffs, but with transformed ones. In games with two pure strategies and a mixed Nash equilibrium, each pure strategy has a minimal payoff and the maximum of the two minimal payoffs is called the pure-strategy maximin. This pure-strategy maximin is the maximal payoff a player can obtain for sure in every round and it forms a natural aspiration level. Amounts below this aspiration level are perceived as losses and amounts above this aspiration level are perceived as gains. In line with prospect theory (Kahneman and Tversky, 1979), losses are counted double in comparison to gains. Thus, gains (the part above the aspiration level) are cut to half for the computation of impulses. Fig. 1 is taken from Selten and Chmura (2008) and illustrates the transformation of the payoffs by the example of game 3.

³ This means, for example, that for impulse-matching learning, impulses into both directions must have been experienced by the subject, and for reinforcement learning, payoff sums for both actions must have been collected.

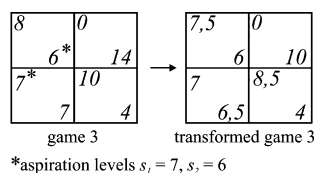


Fig. 1. Example of matrix transformation as given in Selten and Chmura (2008).

Impulse-matching learning can be described as a process in which a subject forms impulse sums. The impulse sum $R_i(t)$ is the sum of all impulses from j towards i experienced up to period $t - 1$. The probabilities for playing action 1 and 2 in period t are proportional to the impulse sums $R_1(t)$ and $R_2(t)$:

$$p_i(t) = \frac{R_i(t)}{R_1(t) + R_2(t)}, \quad \text{for } i = 1, 2. \quad (1)$$

The impulses from action j towards action i in period t are as follows:

$$r_i(t) = \max[0, \pi_i - \pi_j], \quad (2)$$

for $i, j = 1, 2$ and $i \neq j$. Here, π_i is the transformed payoff for action i given the matched agent's decision and π_j the one for action j . Afterwards the impulse sums are updated with the new impulses:

$$R_i(t + 1) = R_i(t) + r_i(t). \quad (3)$$

In the first round, all impulse sums are zero $R_1(1) = R_2(1) = 0$, and until both impulse sums are higher than zero, the probabilities are fixed to $p_1(t) = p_2(t) = 0.5$.

2.2. Action-sampling learning

Action-sampling learning relates to the idea of the action-sampling equilibrium of [Selten and Chmura \(2008\)](#). According to action-sampling equilibrium, a player takes in the stationary state a fixed-size sample of the pure strategies played by the other players in the past and optimizes against this sample. The process of action-sampling learning is a belief-based type of learning, which is very similar to fictitious play. In fact, the model by [Chen et al. \(2011\)](#) is a generalization of action-sampling learning, incorporating inertia, recency, and weighting of the grand mean. Action-sampling learning is much simpler and can be described as a version of fictitious play where only random periods are considered and not the whole data from the history.

In the process of action-sampling learning, the agent randomly takes a sample $A(t)$ with replacement of n earlier actions a_1, \dots, a_n of the other player. Let $\pi_i(a_j)$ be the payoff of action i if the opponent plays action a_j . For $i = 1, 2$ let $P_i(t) = \sum_{j=1}^n \pi_i(a_j)$ be the sum of all payoffs of the player for using her action i against the actions in this sample.

Therefore, in period t , the player chooses her action according to

$$p_i(t) = \begin{cases} 1 & \text{if } P_i(t) > P_j(t), \\ 0.5 & \text{if } P_i(t) = P_j(t), \\ 0 & \text{else,} \end{cases} \quad (4)$$

for $i, j = 1, 2$ and $i \neq j$.

At the beginning, the probabilities are set to $p_1 = p_2 = 0.5$, until both possible actions have been played by the opponent agents.

2.3. Reinforcement learning

The concept of reinforcement learning is one of the oldest and best established learning models in the literature; refer to [Harley \(1981\)](#) for an early application. For experimental economics, it was first formulated and introduced by I. Erev and A.E. Roth ([Erev and Roth, 1998](#), and [Roth and Erev, 1995](#)). In the reinforcement model, a player builds up a payoff sum $B_i(t)$ for each action i according to the following formula:

$$B_i(t + 1) = \begin{cases} B_i(t) + \pi(t) & \text{if action } i \text{ was chosen in } t, \\ B_i(t) & \text{else.} \end{cases} \quad (5)$$

Here $\pi(t)$ is the payoff obtained in period t . After an initial phase, in which both possible actions are used with equal probabilities, the probability of choosing action i in period t is given by:

$$p_i(t) = \frac{B_i(t)}{B_1(t) + B_2(t)}. \quad (6)$$

This model presupposes that all payoffs in a player's payoff matrix are non-negative and at least one payoff in each column and row is positive, a condition fulfilled by all twelve investigated games.⁴ In the first round, the initial payoff sums $B_i(t)$ are zero and the player chooses both possible actions with equal probabilities $p_1 = p_2 = 0.5$. The initial phase ends as soon as both sums are positive, and only from then on, Eq. (6) is applied to determine the probabilities.

⁴ For games with negative payoffs, this approach would not be adequate. To cope with negative payoffs, the model used by [Erev and Roth \(1998\)](#) replaces the payoff $\pi(t)$ in Eq. (1) by $\pi(t) - \pi_{\min}$, where π_{\min} is the smallest possible payoff of the player.

Impulse matching learning and action-sampling learning are both based on behavioral stationary concepts. Reinforcement learning converges in 2×2 constant-sum games to Nash equilibrium, if they are pure or if they are mixed (Beggs, 2005). Therefore, we treat it in the following as the learning concept that corresponds to Nash equilibrium, although this holds only for one half of our games.

2.4. Self-tuning EWA

Self-tuning EWA was introduced by Ho et al. (2007). It is based on the experience-weighted attraction model (Camerer and Ho, 1999), but replaces all but one parameter of this model with functions or fixed values. Of all models discussed in this article, self-tuning EWA is the most sophisticated one because it can capture different types of learning. The decisions are made according to attractions $A_i(t)$ for each strategy i . Attractions are based on the payoffs $\pi(s_i, s^m(t))$ which a subject would have received for playing strategy s_i given the actual decision $s^m(t)$ by the matched player. The attraction updating function depends on an experience weight $N(t)$, a change-detector function $\phi(t)$, and the attention function $\delta(t)$:

$$A_i(t) = \frac{\phi(t)N(t-1)A_i(t-1) + [\delta(t) + (1-\delta(t))I(s_i, s(t))]\pi(s_i, s^m(t))}{N(t)}. \quad (7)$$

Here $I(s_i, s(t))$ is an indicator function equal to 1 for $s(t) = s_i$ and 0 otherwise. An experience weight is applied to each attraction and it is defined as

$$N(t) = N(t-1)\phi(t) + 1, \quad \text{with } N(0) = 1. \quad (8)$$

The change-detector function $\phi(t)$ weights lagged attractions and represents “a player’s perception of how quickly the learning environment is changing” (Ho et al., 2007, p. 182). It is defined as

$$\phi(t) = 1 - \frac{1}{2}S(t) \quad (9)$$

with $S(t)$ being the so called *surprise index*, which measures the deviation of the matched players’ recent decisions from all previous decisions.⁵ $S(t)$ is the quadratic distance between the cumulative history vector $h_k^m(t)$ and the immediate history vector $r_k^m(t)$ for the k strategies of the matched player m . The cumulative history vector gives the relative frequency over all rounds and is defined as

$$h_k^m(t) = \frac{\sum_{\tau=1}^t I(s_k^m, s^m(\tau))}{t}. \quad (10)$$

The immediate history vector gives the relative frequency in the recent rounds. For 2×2 games with mixed equilibria, it is defined as

$$r_k^m(t) = \frac{\sum_{\tau=t-W+1}^t I(s_k^m, s^m(\tau))}{2}. \quad (11)$$

The surprise index is given as:

$$S(t) = \sum_{k=1}^2 (h_k^m(t) - r_k^m(t))^2. \quad (12)$$

The attention function $\delta(t)$ generates a weight for foregone payoffs and turns the attention to strategies which would have yielded higher payoffs. In games with a unique mixed-strategy equilibrium, these payoffs are weighted with $1/W$, with W being the numbers of strategies played in equilibrium. Thus, in our 2×2 games with mixed equilibria, it is set to be $W = 2$.

$$\delta(t) = \begin{cases} \frac{1}{W} & \text{if } \pi(s_j, s_m(t)) \geq \pi(t), \\ 0 & \text{else.} \end{cases} \quad (13)$$

The attention function $\delta(t)$ of self-tuning EWA captures the idea of learning direction theory (Selten and Stoecker, 1986) that subjects have a tendency to move into the direction of the strategy which was ex-post the best response. This is done by shifting the attention and thus the probability towards the strategy with the highest payoff. This is similar to the process of impulse-matching learning. The resulting probability of playing action i in period t , depending on the attractions, is calculated as a logit response function:

$$p_i(t) = \frac{e^{\lambda A_i(t-1)}}{\sum_{j=1}^2 e^{\lambda A_j(t-1)}}. \quad (14)$$

⁵ The experiments were played with random matching and thus no identification of single players is possible. Therefore, we assume that all matched players are perceived as one average player.

Constant sum games			Non-constant sum games						
		L	R			L	R		
Game 1	U	10	8	0	Game 7	10	12	4	22
	D	9	9	10		8	9	9	14
Game 2	U	9	4	0	Game 8	9	7	3	16
	D	6	7	8		5	6	7	11
Game 3	U	8	6	0	Game 9	8	9	3	17
	D	7	7	10		4	7	7	13
Game 4	U	7	4	0	Game 10	7	6	2	13
	D	5	6	9		2	5	6	11
Game 5	U	7	2	0	Game 11	7	4	2	11
	D	4	5	8		1	4	5	10
Game 6	U	7	1	1	Game 12	7	3	3	9
	D	3	5	8		0	3	5	10

The payoffs for the column players are shown in the lower right corner, the payoffs for the row players are shown in the upper left corner. Abbreviations used: L Left, R Right, U Up, D Down

Fig. 2. The twelve 2 × 2 games taken from Selten and Chmura (2008).

Here, λ is the response sensitivity and this parameter must be specified to fit the empirical data. To be consistent with the other models, we have chosen not to estimate any additional values and the simulations start with pure randomization with $p_1 = p_2 = 0.5$.

3. Games and experiments

Our comparison of the investigated learning rules is based on the data of Selten and Chmura (2008). In their study, twelve 2 × 2 games with pure equilibria in mixed strategies were experimentally investigated. To cover a broad set of games, six constant and six non-constant sum games were played. Fig. 2 shows the twelve games used in the experiment. The constant sum games are shown on the left side of the figure, and the non-constant sum games on the right side.

Note that the first six games have the same best-response structure as the second six games and that the concepts of action-sampling equilibrium and Nash equilibrium only depend on this best response structure. Thus, the predictions of Nash equilibrium are the same for the first and the second six games and the same holds true for action-sampling equilibrium. The predictions of Nash equilibrium, action-sampling equilibrium and impulse-balance equilibrium are given in Table 1.

All experiments were run at the BonnEconLab with students mainly majoring in economics or law. The experiment was programmed with RatImage developed by Abbink and Sadrieh (1995). The data was collected in 54 sessions with 16 subjects each. In every session, only one game was played and this game was known by all subjects. The games were played for 200 periods with matching groups consisting out of eight subjects. For each constant sum game, twelve independent matching groups were gathered; for each non-constant sum game, six independent matching groups were gathered. Overall, 864 subjects participated.

The role of the subjects was fixed for the whole experiment, thus four subjects in each matching group decided as column players and the other four as row players throughout the whole experiment. At the beginning of each round, row and column players were randomly matched. After every round, subjects received feedback about the other player's decision, their own payoff, the period number and their own cumulative payoff. Each participant received €5 for showing-up. In addition, the payoffs in the 200 periods were accumulated and transferred into Euro. The exchange rate was €0.016 Cent per

payoff point. An experimental session lasted between 1.5 and 2 hours and the average earning per subject was roughly €24, including show-up fee.

4. Simulating the behavior of populations

In this section, we investigate whether the learning algorithms can replicate in simulations the aggregate distribution of actions generated by human subjects. In the following, we will first introduce our measurement of success for the simulations. Thereafter, we discuss the success of the different learning models in predicting/reproducing the aggregate distribution of behavior.

4.1. Measure of predictive success for the simulations

For this analysis, we conduct simulations keeping everything the same as in the experiment, except that instead of real participants now computer agents interact. Each agent interacts according to her history and to the same learning model over 200 rounds. In each round, eight agents with fixed roles, four deciding as row players and four as column players, are randomly matched and all agents act in accordance to the same learning rule.

After each round, they receive feedback about the matched agent's decision and their payoff. Since none of the learning models makes use of the round number and since the calculation of the cumulated payoff can be done by the agents themselves, this information is not provided to the agents. It is crucial that the agents do not receive more information than the subjects in the experiment did.

All learning models include stochastic elements. To avoid the influence of statistical outliers, 500 simulation runs per game are conducted. In each simulation run, all agents act in accordance with the same learning model. To measure the predictive success on the population level, we will compare the mean frequencies of U and L in the simulations with the mean frequencies obtained in the experiments by means of the quadratic distance. The mean quadratic distance Q is the average quadratic distance over all 12 games and over all 500 simulations. It is defined as

$$Q = \frac{1}{12} \sum_{i=1}^{12} \left(\frac{1}{500} \sum_{n=1}^{500} (s_{in}^L - f_i^L)^2 + (s_{in}^U - f_i^U)^2 \right),$$

with s_{in}^L and s_{in}^U being the frequencies for L and U in game i and simulation run n . Respectively, f_i^L and f_i^U are the mean frequencies for L and U observed in the experiments with game number i . The frequencies of R and D need not be considered in view of $(s_{in}^L - f_i^L)^2 = (s_{in}^R - f_i^R)^2$ and $(s_{in}^U - f_i^U)^2 = (s_{in}^D - f_i^D)^2$. The predictive success of a learning model increases with a decreasing mean quadratic distance, i.e., the smaller the mean quadratic distance is, the better the learning theory fits the experimental data on the aggregate level.

4.2. Parameter estimates

The concepts of action-sampling equilibrium and self-tuning EWA have a parameter which needs to be adjusted to the experimental data. We decided to estimate for each learning model one parameter that minimizes the quadratic distance over all games.⁶ To estimate the optimal parameter, we ran, for each parameter, 500 simulations per game and calculated the mean quadratic distance. Thereafter, the simulations were conducted with the parameter that yielded the smallest quadratic distance.⁷

The parameter of action-sampling learning is the size of the drawn samples. Fig. 3 gives the mean quadratic distances of action-sampling learning for $1 \geq n \geq 15$. A sample size of $n = 12$ leads to the smallest quadratic distance, which is the same sample size that also leads to the smallest distance for the stationary concept (cf. Brunner et al., 2011).

Fig. 4 gives the mean quadratic distances of self-tuning EWA for different lambdas. The left part gives the mean quadratic distances for all tested lambdas between 0 and 10, and the right part gives the quadratic distance for $0.2 < \lambda < 0.3$.

Each point in the graph represents the mean quadratic distance over all twelve games with 500 simulation runs per game with one specific lambda value. The value leading to the smallest quadratic distance is $\lambda = 0.2775$.

⁶ One could fit the parameter of the parametric concepts for each game separately. We believe that this gives an unfair advantage to one-parameter theories over parameter-free ones. This especially holds for the case of 2×2 games, where only two relative frequencies are predicted. Adjusting a parameter separately for each game, so to speak, does half the job. One might use methods to adjust the fit of a theory to the number of parameters used, but this only makes sense if the non-adjusted performance of a model increased in case of parameters being estimated for each game separately. For our simulations, only the quadratic distance of action-sampling learning would benefit from such a procedure. The quadratic distance of self-tuning EWA (0.0805 vs. 0.0786) would change only slightly and this adjustment would not influence the relative ranking of quadratic distances. Therefore, we decided to estimate only one parameter for each theory.

⁷ To speed up this procedure, round-by-round data were only saved for the simulations with the final parameter.

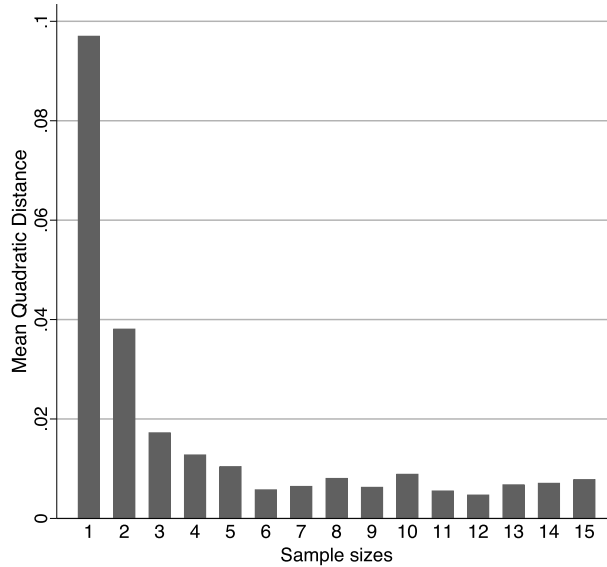


Fig. 3. Quadratic distances of action-sampling learning.

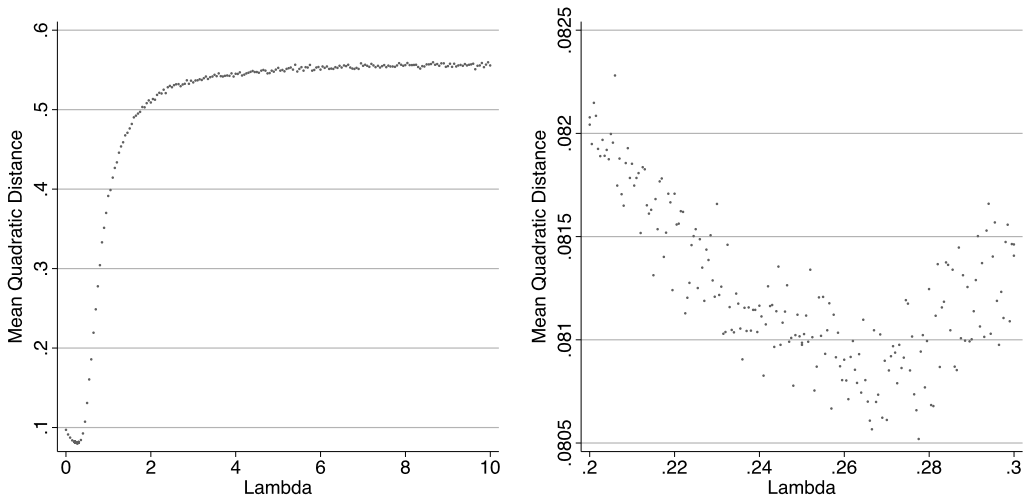


Fig. 4. Quadratic distances of self-tuning EWA for different lambdas. Each point represents the mean quadratic distance over 500 simulations per game. Left figure for $0 \leq \lambda \leq 10$ and right figure for $0.2 \leq \lambda \leq 0.3$.

4.3. Relative frequencies

Table 1 gives the observed mean frequencies for each learning type, mean frequencies predicted by the stationary concepts and the observed frequencies in the experiments. For the experimental games 1 to 6, the mean frequencies observed in a game are based on the observed frequencies in twelve independent matching groups; for games 7 to 12, they are based on the observed frequencies in six independent matching groups. Each matching group consists of eight subjects. For each learning type and game, the mean is based on 500 simulation runs, which produced 500 independent matching groups per game. Each matching group consists of eight agents. Fig. 5 gives the typical development of probabilities over time in game 7 for each of the learning types. This figure in combination with Table 1 already reveals some differences between the learning rules.⁸

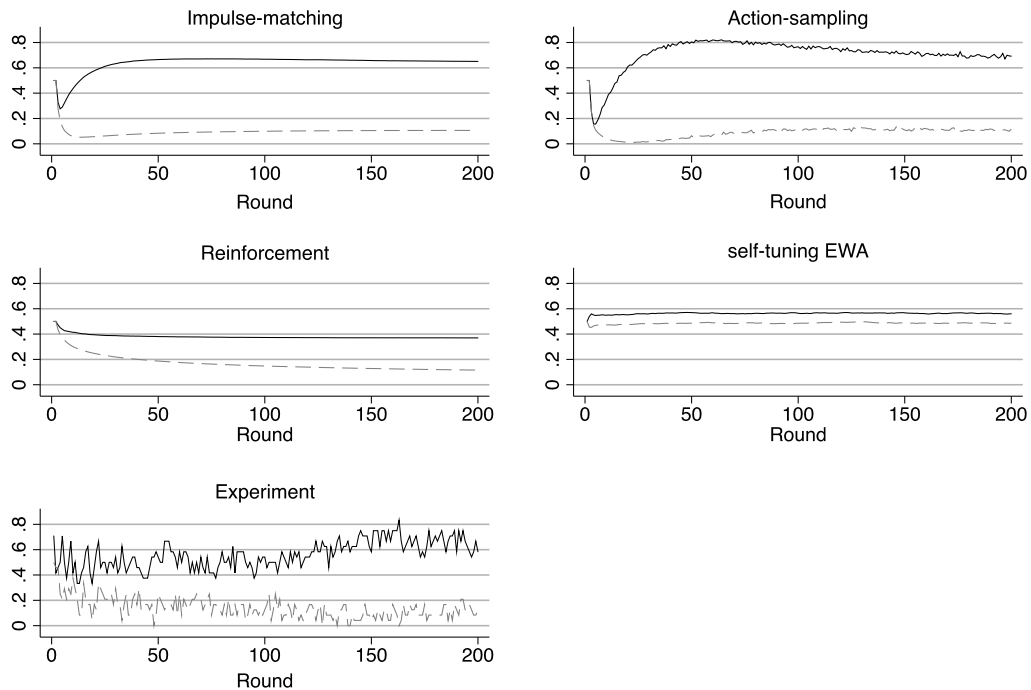
It is surprising that self-tuning EWA yields relative frequencies very near to 0.5 for each of the twelve games. This is probably connected to the fact that, in our simulations, the whole population is of the same type and agents try to adjust to an inaccurate history, and by this process generate a new inaccurate history for themselves and the matched agents. Estimating the free parameter of this model jointly for all games is not a reason for this behavior. If we estimate the

⁸ The course of probabilities for all games is given in Appendix C.1.

Table 1

Relative frequencies for playing left and up in the simulations, predicted by the stationary concepts and observed in the experiments for up and left.

		Impulse- matching learning	Action- sampling learning	Reinforce- ment learning	Self-tuning EWA learning	Impulse- balance equilibrium	Action- sampling equilibrium	Nash equilibrium	Selten and Chmura data
Game 1	L	0.574	0.710	0.345	0.499	0.580	0.705	0.909	0.690
	U	0.063	0.095	0.121	0.499	0.068	0.090	0.091	0.079
Game 2	L	0.495	0.571	0.333	0.477	0.491	0.584	0.727	0.527
	U	0.169	0.193	0.161	0.502	0.172	0.193	0.182	0.217
Game 3	L	0.770	0.763	0.503	0.541	0.765	0.774	0.909	0.793
	U	0.157	0.211	0.128	0.492	0.161	0.208	0.273	0.163
Game 4	L	0.714	0.711	0.587	0.548	0.710	0.719	0.818	0.736
	U	0.259	0.295	0.190	0.494	0.259	0.302	0.364	0.286
Game 5	L	0.632	0.639	0.566	0.524	0.628	0.643	0.727	0.664
	U	0.296	0.323	0.241	0.495	0.297	0.329	0.364	0.327
Game 6	L	0.602	0.596	0.666	0.527	0.600	0.596	0.636	0.596
	U	0.400	0.422	0.265	0.497	0.400	0.426	0.455	0.445
Game 7	L	0.637	0.709	0.380	0.564	0.634	0.705	0.909	0.564
	U	0.098	0.094	0.170	0.485	0.104	0.090	0.091	0.141
Game 8	L	0.563	0.572	0.396	0.540	0.561	0.584	0.727	0.586
	U	0.258	0.193	0.217	0.494	0.258	0.193	0.182	0.250
Game 9	L	0.767	0.762	0.525	0.600	0.764	0.774	0.909	0.827
	U	0.185	0.212	0.165	0.489	0.188	0.208	0.273	0.254
Game 10	L	0.726	0.711	0.640	0.587	0.724	0.719	0.818	0.699
	U	0.303	0.295	0.219	0.487	0.304	0.302	0.364	0.366
Game 11	L	0.648	0.640	0.609	0.572	0.646	0.643	0.727	0.652
	U	0.354	0.324	0.289	0.492	0.354	0.329	0.364	0.331
Game 12	L	0.605	0.596	0.560	0.578	0.604	0.596	0.636	0.604
	U	0.466	0.422	0.342	0.494	0.466	0.426	0.455	0.439

**Fig. 5.** Mean probabilities for left and right in the simulations runs and the experiment for game 7. The mean probability for left is given in black and the mean probability for up is given in gray.

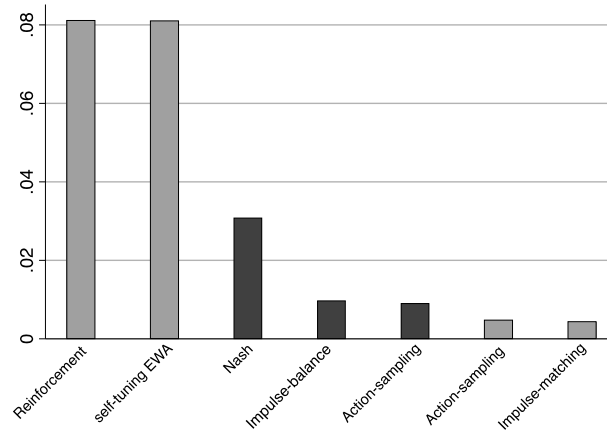


Fig. 6. Mean quadratic distances of the stationary concepts and the learning models to the observed behavior (stationary concepts in dark bars and learning models in light bars).

optimal λ for each game separately, the resulting mean quadratic distance to the data improves only marginally (0.081 vs. 0.079) and observed relative frequencies do not change much.

Impulse-matching learning and action-sampling learning are quite close to their stationary counterparts after 200 periods. The quadratic distances between impulse-matching learning and impulse-balance equilibrium, as well as the one between action-sampling learning and action-sampling equilibrium, are smaller than 0.001. If we treat reinforcement learning as the learning counterpart to Nash equilibrium, the difference in quadratic distances is 0.158. Self-tuning EWA has the highest distances towards all stationary concepts. This closeness results in high correlations between the frequencies of the simulations and corresponding stationary concepts. For impulse-matching learning and action-sampling learning, this is true for both players (pairwise correlation with $r > 0.9$ and $p < 0.01$ for row and column players), and for reinforcement only for row players (pairwise correlation with $r > 0.9$ and $p < 0.01$). Correlations between observed frequencies from the experiments and the simulations with action-sampling learning and impulse-matching learning are high (pairwise correlation with $r > 0.8$ and $p < 0.01$ for both players), but lower than the ones with the stationary concepts. The frequencies of reinforcement learning are correlated with the ones of the row player in the experiment (pairwise correlation with $r > 0.8$ and $p < 0.01$), but not for the ones of the column players. The frequencies of self-tuning EWA are not significantly correlated with the empirical data for either players.

4.4. Overall performance

Fig. 6 gives the mean of the quadratic distance between the experiment and simulations over all games and rounds for self-tuning EWA learning, reinforcement learning, action-sample learning, and impulse-matching learning. In addition, the figure gives the mean quadratic distances between of the stationary counterparts (if existing) and the data.⁹

We first turn our attention to the comparison of the simulations. The figure reveals a clear order of explanatory power. The order from worst to best (highest quadratic distance to lowest quadratic distance) is as follows: reinforcement learning, self-tuning EWA learning, action-sampling learning, and impulse-matching learning. Because of the high number of observations (6000 per learning type), the order given by Fig. 6 is statistically robust (for all $p < 0.01$ Fisher–Pitman permutation test for paired replicates). The difference between self-tuning EWA and reinforcement is very small and irrelevant. However, the similarity between the two quadratic distances does not mean that both theories make similar predictions. This can be seen, for example, in Table 1 and in Fig. 5. The figure demonstrates that the concepts of self-tuning EWA and reinforcement fail to describe the aggregate behavior in the 2×2 experiments, in contrast to the other concepts. The quadratic distance of self-tuning EWA is 18 times higher than the one of impulse-matching learning.

The quadratic distances of reinforcement learning and self-tuning EWA learning are significantly bigger not only over all games, but also for the subsets of constant sum games and non-constant sum games. However, reinforcement performs better in constant sum games than self-tuning EWA does, while self-tuning EWA performs better in non-constant sum games (both $p < 0.01$ Fisher–Pitman permutation test for paired replicates). While the quadratic distance of impulse-matching is stable in constant and non-constant sum games, the one of action-sampling learning is smaller in constant sum games. Thus, action-sampling learning performs significantly better in constant sum games, and impulse-matching learning performs significantly better in non-constant sum games (both $p < 0.01$ Fisher–Pitman permutation test for paired replicates).

Comparing the stationary concepts with the learning models reveals that self-tuning EWA and reinforcement learning are not only outperformed by impulse-matching learning and action-sampling learning, but by all stationary concepts.

⁹ The mean quadratic distances of the stationary concepts are either taken from Selten and Chmura (2008) or from Brunner et al. (2011). There were some flaws in the paper by Selten and Chmura (2008). For a detailed discussion, refer to Brunner et al. (2011) and Selten et al. (2011).

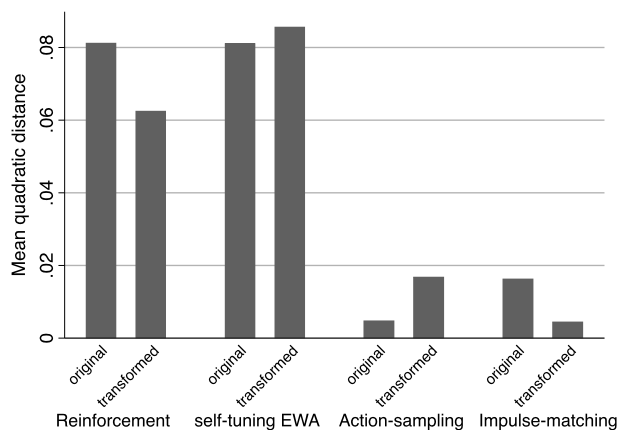


Fig. 7. Mean quadratic distance in original and transformed games.

In contrast, the learning models of action-sampling and impulse-matching perform very well. Both learning models have higher predictive success than the other learning models and additionally a higher predictive success than all stationary concepts.

4.5. Original vs. transformed games

The concept of impulse-matching learning is applied to the transformed game rather than the original one. This transformation is an essential part of impulse-matching learning and impulse-balance equilibrium (Selten and Chmura, 2008 and Goerg and Selten, 2009), because both concepts involve a fixed loss-aversion. Losses with respect to the pure-strategy maximin are counted double. While double counting of losses with respect to the pure-strategy maximin is an essential part of impulse-matching learning, it is ignored by the other concepts (reinforcement learning, action-sampling learning, self-tuning EWA). This raises the question whether the good performance of impulse-matching learning is an artifact of the incorporation of loss aversion. To investigate this point, we apply all learning models to the transformed and to the original matrices.

Fig. 7 shows the overall mean quadratic distances for self-tuning EWA learning, reinforcement learning, payoff-sampling learning, impulse-balance learning, action-sampling learning, and impulse-matching learning applied to the original games and to the transformed games, which are again based on 500 simulation runs per game and learning model.

It can be seen that impulse-matching learning and reinforcement learning perform better when applied to the transformed games, whereas self-tuning EWA learning and action-sampling learning do less well. While the improvement of impulse-matching learning in transformed games is expected, the benefit of applying reinforcement learning to transformed games is unexpected. This improvement is substantial, in the original game the quadratic distance is nearly 1.3 times higher than in the transformed ones.

The theory of Roth and Erev (1995) applies a transformation of the original game by replacing the payoff of a player by its difference to the minimal value in her matrix. The transformation used here is different since it involves double weights for losses with respect to the pure-strategy maximin. However, in Selten and Chmura (2008), no improvement of the predictive power of the Nash equilibrium was observed when applied to the transformed game rather than to the original one. It is interesting that the picture looks different for the simulations over 200 rounds with reinforcement learning, although it corresponds very much with Nash equilibrium (Beggs, 2005).

Although reinforcement learning improves when applied to the transformed matrix, it still performs significantly worse than impulse-matching learning. Therefore, and because of self-tuning EWA and action-sampling learning performing worse in the transformed matrices, we can conclude that the good performance of impulse-matching learning is not driven by the transition to the transformed matrices alone.

4.6. Changes over time

Learning processes are always dependent on time and history, and therefore it is of interest to check whether our above results remain stable over time. To check stability of the order of explanatory power over time, we compare the first hundred periods with the second hundred periods. Fig. 8 gives the mean quadratic distances for periods 1–100 (left) and 101–200 (right) for the six learning models. The basis of the comparison is always the observed mean frequencies for the corresponding rounds (either round 1–100 or 101–200) in the experiments.

It is easy to recognize that in the second half of the simulation runs the explanatory power of self-tuning EWA and reinforcement learning decreases significantly, while the one of impulse-matching learning improves significantly (all Fisher–Pitman permutation test for paired replicates $p < 0.01$). The concept of action-sampling learning is rather stable over time

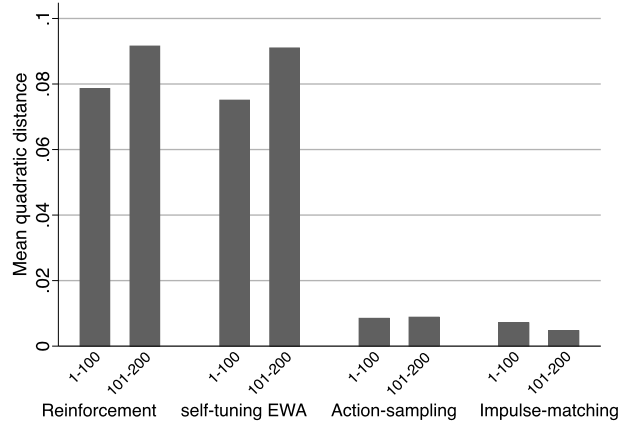


Fig. 8. Mean quadratic distance over time.

and the statistically significant disimprovement of action-sampling learning ($p < 0.01$) is economically negligible, with an increase of the quadratic distances of only 0.0004.

The ranking of concepts by mean quadratic distances is stable over time, the overall ranking for round 1–200 is the same in rounds 1–100 and 101–200.

5. Individual round-by-round predictions

In this part, we investigate how well the learning rules describe the individual behavior of the subjects in the 2×2 experiments. To judge the performance on the individual level, we compare the individual decisions in every round of the experiment with the predicted decisions or predicted probability by the learning rule, given the history of the subject.

5.1. Measure of predictive success for the individual round-by-round behavior

To measure the predictive success of the learning theories describing the behavior of a single individual, we apply the quadratic scoring rule on each of the 864 subjects for each learning rule. The quadratic scoring rule was first introduced by Brier (1950) in the context of weather forecasting. The rationale behind the quadratic scoring rule is that for each round a score is determined, which evaluates the nearness of the predicted probability distribution to the observed outcome.

In Selten (1998), the quadratic scoring rule is axiomatically characterized. The characterizing properties of the quadratic scoring rule, as described in Selten (1998), are: symmetry, elongational invariance, incentive compatibility, and neutrality. Symmetry means that the score of a theory must not depend on the numbering of the decision alternatives. Elongational invariance assures that the score of a theory is not influenced by adding or leaving an alternative which is predicted with a probability of zero. Incentive compatibility requires that predicting the actual probabilities yields the highest score. Finally, neutrality means that in the comparison of two theories, of which one is right in the sense that it predicts the actual probabilities, and the other is wrong, the score for the right theory does not depend on which of the two theories is the right one. This means that the score does not prejudice one of the theories depending on the location of the theory in the space of probability distributions.

We apply the quadratic scoring rule to measure the predictive success of a theory for every period and subject and then calculate the mean over subjects, rounds, and games. Accordingly, a score depending on the predicted probabilities and the actually observed action is computed. In order to compute the score, the observation is interpreted as a frequency distribution, where for the chosen action the relative frequency is one, and for the action not chosen it is zero.

The quadratic score $q(t)$ of a learning theory for subject choosing action i in period t is given as¹⁰:

$$q(t) = 2p_i(t) - p_i^2(t) - (1 - p_i(t))^2.$$

Here $p_i(t)$ is the predicted probability of the learning theory. The predicted probability of the learning theory is calculated by applying the theory's learning algorithm to the whole playing history of this player. If no history yielding a positive number smaller than 1 for $p_i(t)$ is available, the player randomizes with $p_i(t) = 0.5$. This rule provides an initial phase. As soon as both probabilities are positive, they will remain positive forever.

If a player decides completely in line with the prediction of the theory, he receives a score of 1; if he decides in complete contrast to the prediction of the theory, he receives a score of -1 . The mean score \bar{q} is given as the mean of $q(t)$ over all

¹⁰ If the decision maker has n choices it is defined as: $q(t) = 2p_i(t) - \sum_{j=1}^n p_j^2(t)$. The formula in the text holds for the special case of $n = 2$.

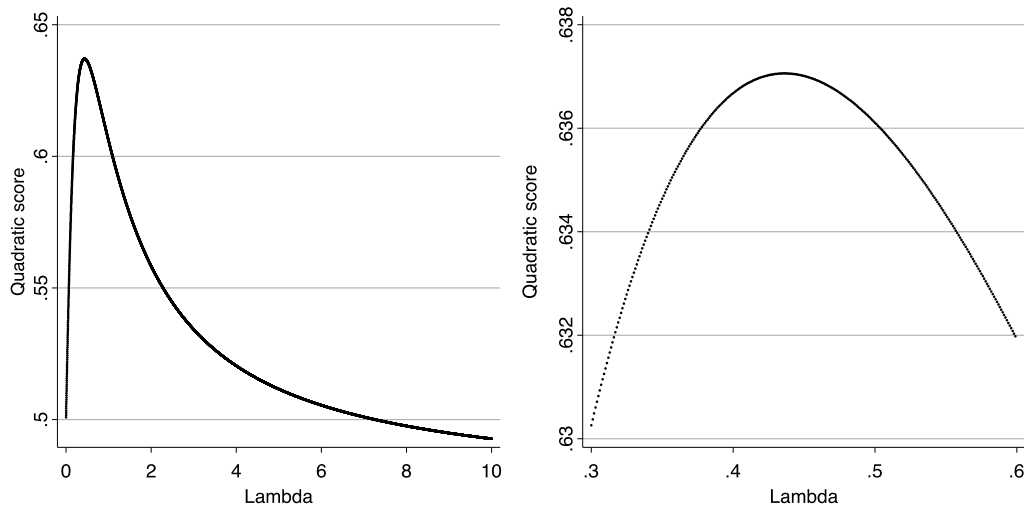


Fig. 9. Quadratic scores of self-tuning EWA for different lambdas, each point represents the mean quadratic score for all 864 subjects in left figure for $0 \leq \lambda \leq 10$ and right figure for $0.3 \leq \lambda \leq 0.6$.

200 rounds and all 864 subjects. Of course, \bar{q} must be in the closed interval between -1 and $+1$. Thus, in contrast to our measurement for the success on the aggregate level, the success of a theory on the individual level increases with the score.

The concept of action-sampling learning always yields a probability of 1, 0 or 0.5 for one of the possible actions. Which action is chosen depends on the randomly drawn sample. Therefore we calculate the probability of drawing a sample that commands playing action 1 or action 2 as the predictions of this concept.

In addition to the investigated learning rules we introduce three benchmarks. The first one is a heuristic which we call the inertia rule. This rule commands to “do exactly the same as in the preceding round”. This does not apply to the first period in which both possible actions are chosen with equal probabilities. The player is required to repeat the decision of the preceding period even if he deviated from this rule in the past. Obviously, the inertia rule is not a serious decision rule, but it serves as a benchmark that every learning rule should beat. The second benchmark is the score an agent would receive if he decided randomly between the two actions with $p = 0.5$. In this case, the score would be 0.5 and again every learning rule should beat this benchmark. The third benchmark is the aggregated observed frequencies taken from the experiments. If a learning theory adequately describes the adjustments over time, it should yield higher scores than these stationary probabilities, which do not depend on this information.

5.2. Parameter estimates

On the individual level, we calculated for each parametric learning theory one parameter, which leads to the highest mean quadratic score over all 864 subjects. Fig. 9 gives the mean quadratic scores for the parameter lambda of self-tuning EWA between 0 and 10 (left side) and for lambda between 0.3 and 0.6 (right side). The highest mean quadratic score is reached for $\lambda = 0.436$.

Fig. 10 gives the mean quadratic scores of action-sampling for sample sizes between 1 and 15. The optimal sample size for action-sampling is $n = 6$. Note that the optimal sample size of action-sampling learning for the performance on the aggregate level ($n = 12$) also performs very well on the individual level; it leads to the second-highest mean quadratic score.

5.2.1. Overall mean quadratic scores

Fig. 11 gives the mean quadratic scores in the 108 independent observations for each learning model. The randomization benchmark is included as a horizontal line. The figure reveals a clear order of predictive success, from best to worst: self-tuning EWA, impulse-matching learning, empirical frequencies, reinforcement learning, action-sampling learning, randomization benchmark, and inertia benchmark.

Applying a two-sided permutation test for the pairwise comparison of the mean scores over all independent observations reveals that the order given by the graph is statistically robust. All pairwise comparisons between two learning models over all games are at least significant on the 1% level. Table 2 gives all test results over all games (top), over the constant-sum games (middle), and over the non-constant sum games (bottom). This ranking is robust over all games, as well as for the subsets of constant sum games and non-constant sum games.

All reported learning models perform significantly better than the inertia and randomization benchmarks. However, only impulse-matching learning and self-tuning EWA perform significantly better than the aggregated empirical frequencies in describing the individual round-by-round behavior. No significant difference between reinforcement learning and the empirical frequencies are observed, and action-sampling learning performs even worse than the empirical frequencies.

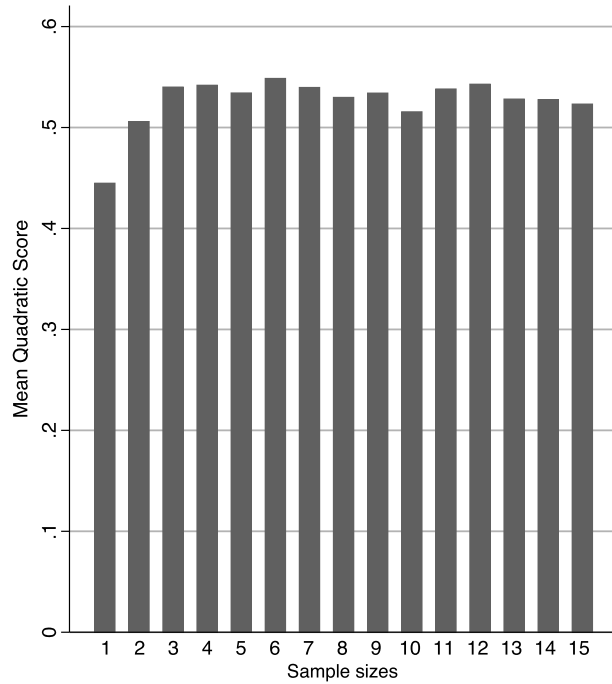


Fig. 10. Quadratic scores of action-sampling learning for different sample sizes.

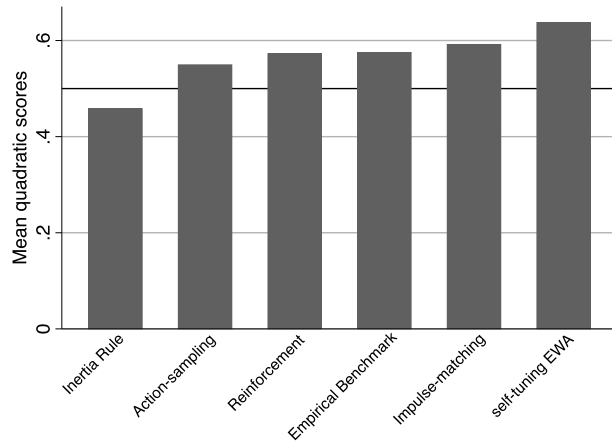


Fig. 11. Mean quadratic scores, over all 108 independent observations. The solid line gives the random rule benchmark.

Although the inertia benchmark performs rather badly in describing subjects' behavior, it does not imply that low inertia rates are observed. On the contrary, in line with Erev and Haruvy (2005) and Erev et al. (2010), very high inertia rates are observed. In 74% of all cases, subjects stick to their previous decision. On the other side, this means that, in 26% of all cases, inertia predicts exactly the opposite of observed behavior and receives the lowest score of all models (−1). Nevertheless, 13% of the subjects are best described by the inertia rule. Overall, roughly 8% of the subjects are best described by reinforcement and action-sampling learning, 23% by impulse-matching learning, and the majority of 47% is best described by self-tuning EWA.¹¹

5.3. Mean quadratic scores over time

To conclude our analysis, we now take a look at the quadratic scores over time. Fig. 12 gives the mean quadratic scores for rounds 1–100 and rounds 101–200. In the first and the second half of the experiments, the same order of

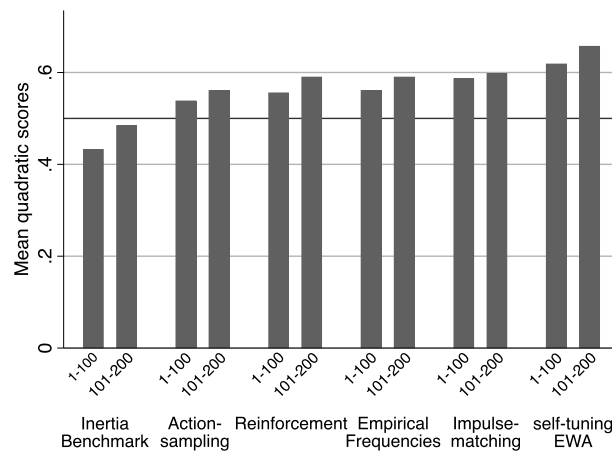
¹¹ Please note that this calculation of proportion is problematic: it depends on the number and the performance of included learning concepts. We prefer the mean quadratic score as it does not depend on the competing learning concepts, but report these fractions for completeness.

Table 2

Two-sided Significances in Favor of Row Concepts, Monte-Carlo approximation of the two-sided Fisher–Pitman permutation test for paired replicates.

	Impulse matching	Empirical frequencies	Reinforcement	Action sampling	Random benchmark	Inertia benchmark
Self-tuning EWA learning	1%	1%	1%	1%	1%	1%
Impulse matching learning		1% 1% 5%	1% 1% 10%	1% 1% 1%	1% 1% 1%	1% 1% 1%
Empirical frequencies			n.s. n.s. n.s.	1% 1% 1%	1% 1% 1%	1% 1% 1%
Reinforcement learning				1% 1% 10%	1% 1% 5%	1% 1% 1%
Action sampling learning					1% 1% 5%	1% 1% 1%
Randomization benchmark						1% 1% 1%

Notes: Above: all 108 experiments; Middle: 72 constant-sum game experiments; Below: 36 non-constant sum game experiments.

**Fig. 12.** Mean quadratic scores in the first and second half of the experiments.

success is present as over all rounds. For rounds 1–100, all differences between the scores, except for the one of reinforcement learning and the empirical benchmark, are highly significant (all Fisher–Pitman permutation test for paired replicates $p < 0.01$). As for the overall comparison, no significant difference between the mean quadratic scores of reinforcement learning, and the empirical benchmark is observed. Although the order remains the same over time, in the second half of the experiments the differences between impulse-matching learning, reinforcement learning and the empirical benchmark decrease. The difference between reinforcement learning and impulse-matching learning is no longer significant, while impulse-matching learning still performs significantly better than the empirical benchmark.

Fig. 13 gives the development of the mean quadratic scores per round over time. Impulse-matching learning has the fastest increase of scores in the very early rounds. In the first 10 rounds, impulse-matching learning performs better than self-tuning EWA. The performance of self-tuning EWA increases continuously over time, leading to higher scores per round after round 10, and after 25 rounds, the overall quadratic score of self-tuning EWA is above the one of impulse-matching. The score of reinforcement learning also increases continuously over time, but at a slower rate, approaching the performance of impulse-matching learning only in the last 50 rounds. In addition, the graph reveals the increase of inertia over the rounds. Over time, the score of the inertia benchmark approaches the score of the pure randomization benchmark (0.5). In the last periods, inertia has a higher predictive power than randomization, but still performs worse than the learning concepts.

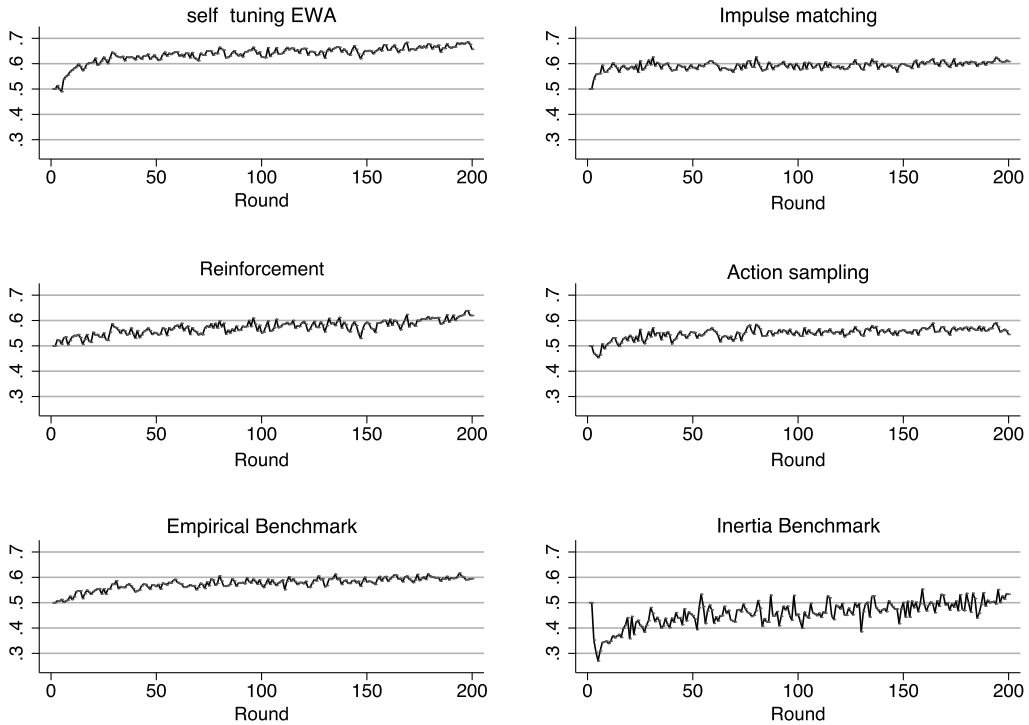


Fig. 13. Mean quadratic scores per round over time for the learning models and benchmarks.

6. Summary and discussion

In this article, the models of impulse-matching learning and action-sampling learning have been introduced. Together with reinforcement learning and self-tuning EWA, they were applied and tested in the environment of 12 repeated 2×2 games.

The newly-introduced learning models are based on the behavioral reasoning of action-sampling equilibrium and impulse-balance equilibrium, which had been successfully tested in experimental 2×2 games by Selten and Chmura (2008). Therefore the experimental dataset obtained by Selten and Chmura (2008) was used as a testbed for the learning models. The experimental data comprises aggregate and individual behavior in 12 completely mixed 2×2 games, 6 constant sum games with 12 independent subject groups each, and 6 non-constant sum games with 6 independent subject groups each. Each subject group consists of eight participants being randomly matched over 200 periods.

The learning models had to prove whether they could replicate the aggregate behavior of the experimental population and whether they could explain the individual round-by-round behavior of single subjects. For the comparison with the aggregate behavior of the population, 500 simulation runs per game and learning model were conducted. As in the experiment, 200 rounds with random matching and four agents deciding as row players and four agents as column players were simulated. Our measure of predictive power for the aggregate population is the quadratic distance between observed relative frequencies in simulation runs and the mean frequencies observed in the experiments. For the comparison with the individuals' behavior, the models were applied to the history of each participant. Then the actual decisions of every round were compared with the predictions of the learning models given the actual subject's history. For each subject and round, a quadratic score, a measurement for the accuracy of a prediction, was calculated and averaged over rounds, subjects, and games.

For our comparisons with the aggregate distribution of the experimental population and the individual round-by-round behavior, we can conclude two main results:

Main Result 1. *The models of impulse-matching learning and action-sampling learning are able to replicate with simulations the aggregate behavior of the experimental population.*

The comparison of the four models yields the following order of predictive success from best to worst: impulse-matching learning, action-sampling learning, reinforcement learning, self-tuning EWA learning. Due to the high number of simulation runs, this order is statistically robust, all pairwise comparisons are at least significant on the 1% level.

The predominance of the new models, impulse-matching learning and action-sampling learning, over the established models of reinforcement learning and self-tuning EWA is stable over time and across the different game types (constant

sum and non-constant sum games). A further interesting result is that for reinforcement learning the quadratic distance to the data is about 22% lower if applied to the transformed matrixes instead to the original ones.

Main Result 2. *On the individual level, self-tuning EWA outperforms all other learning concepts in predicting individuals' round-by-round behavior.*

Overall, the models of action-sampling learning, reinforcement learning, impulse-matching learning, and self-tuning EWA perform better than simple randomization with 0.5 and the inertia benchmark does. But only the models of self-tuning EWA and impulse-matching learning perform significantly better in describing round-by-round behavior than the aggregated frequencies from the experiment. The mean quadratic score of self-tuning EWA is significantly above the ones of all other investigated learning concepts, and impulse-matching learning has the second-highest score.

The good performance of self-tuning EWA on the individual level is remarkable, and not the results of the free parameter. A non-parametric version of self-tuning EWA¹² yields an only marginal lower quadratic score (0.624 vs. 0.637) and still performs significantly better than impulse-matching learning. However, on the aggregate level, the distance of the simulated populations to the experimental data increases even further (0.079 vs. 0.111), and therefore impulse-matching learning performs significantly better in this domain (both Fisher–Pitman permutation test for paired replicates $p < 0.01$). In addition, our data shows that over time inertia increases, and therefore adding an inertia component to the learning models might improve their predictive power. For example, the model by Chen et al. (2011) provides a multi-parameter generalization of action-sampling learning that takes inertia and recency into account.

We conclude that impulse-matching learning produces good results across fields of applications, while self-tuning EWA organizes existing individual data exceptionally well. Our results suggest that if one is interested in the aggregate behavior of a population in a certain 2×2 game without any prior information upfront (i.e., no information for parameter estimates), simulations with impulse-matching learning result in an extremely good approximation. In addition, impulse-matching learning provides good predictions for round-by-round behavior, but in this domain it is clearly outperformed by self-tuning EWA. Obviously, self-tuning EWA interprets actual history very well, while it fails to generate accurate behavior in the simulations. Interestingly, only the two concepts of impulse-matching learning and self-tuning EWA have components of regret-based learning and only these two concepts outperform the empirical frequency benchmark.

Appendix A. Additional learning models

A.1. Impulse-balance learning

The algorithm for impulse-balance learning is very similar to the one of impulse-matching learning. Only the calculation of impulses differs, therefore Eq. (A.1) replaces Eq. (2) of impulse-matching learning. All other equations ((1) and (3)) remain the same. In contrast to impulse-matching, a player receives only actual impulses. This means the player does not receive an impulse if his action was a best reply against the other player's decision. Thus, for impulse-balance learning, the impulses from action j towards action i in period t are as follows:

$$r_i(t) = \begin{cases} \max[0, \pi_i - \pi_j] & \text{if the chosen action is } j, \\ 0 & \text{else,} \end{cases} \quad (\text{A.1})$$

for $i, j = 1, 2$ and $i \neq j$. Again, π_i is the transformed payoff for action i given the matched agent's decision and π_j the one for action j . Afterwards, the impulse sums are updated with the new impulses. In the first round, all impulse sums are zero $R_1(1) = R_2(1) = 0$ and until both impulse sums are higher than zero the probabilities are fixed to $p_1(t) = p_2(t) = 0.5$.

In fact, both learning rules are so similar that they lead to the same stationary points in 2×2 games. To illustrate this, we take a look at the structure of the investigated experimental 2×2 games, as introduced by Selten and Chmura (2008).

	L	R
U	$a_L + c_L$ b_U	a_R $b_U + d_U$
D	a_L $d_D + d_D$	$a_R + c_R$ b_D

Fig. 14. The structure of the experimental 2×2 games.

Fig. 14 shows the transformed payoffs, the payoffs for the column players are shown in the lower right corner, and the payoff for the row players are shown in the upper left corner. The following equations must be fulfilled: $a_L, a_R, b_U, b_D \geq 0$ and $c_L, c_R, d_U, d_D > 0$. In the following, p_U and p_D are the probabilities of the row player for U and D and q_L and q_R are

¹² A non-parametric self-tuning EWA can be obtained by replacing the equation for calculating the probabilities with $p_i(t) = \frac{A_i(t-1)}{\sum_{j=1}^2 A_j(t-1)}$. We thank an anonymous referee for this suggestion.

the probabilities for L and R by the column player. In the following, we will only look at the row player; the behavior in equilibrium of the column player is calculated analogously.

In case of impulse-balance equilibrium, the expected impulses for each of both strategies must be the same. Here, the row player receives only an impulse towards U for the proportion of plays in which he would choose down (given by p_D) and the other player at the same time would have chosen L (given by q_L). Therefore the expected impulse for U is given by $p_D q_L c_L$. Applying the same reasoning leads to $p_U q_R c_R$ as the expected impulse for D of the row player. Thus the impulse-balance equation, which must be fulfilled in equilibrium, is given as:

$$p_D q_L c_L = p_U q_R c_R.$$

In case of impulse-matching equilibrium, the row player always receives an impulse of c_L towards U if the column player plays L. The column player does so with a probability of q_L . In addition the row player always receives an impulse of c_R towards D if the column player chooses R. The column player plays R with a probability of q_R . Impulse-matching equilibrium is reached if the ratio of the two probabilities of U and D is the same as the ratio of expected impulses for U and D.

$$\frac{p_U}{p_D} = \frac{q_L c_L}{q_R c_R}.$$

By transforming we obtain the impulse-balance equation of impulse-balance equilibrium:

$$p_D q_L c_L = p_U q_R c_R.$$

Therefore, impulse-matching equilibrium and impulse-balance equilibrium have the same mixed stationary points in the case of the described 2×2 games. However, for other types of games, both concepts do not necessarily lead to the same stationary points.

A.2. Payoff-sampling learning

Payoff-sampling learning relates to the stationary concept of Osborne and Rubinstein (1998), which was first applied to experimental data in Selten and Chmura (2008). The behavioral explanation of the stationary concept is that a player chooses her action after sampling each alternative an equal number of times, picking the action that yields the highest payoff.

To implement this behavior, payoff-sampling learning is based on samples from earlier periods. The samples are randomly drawn with replacement and fixed sample sizes. The agent draws two samples ($s_1(t), s_2(t)$) of earlier payoffs, one sample with payoffs from rounds in which she chose action 1, and one with payoffs from rounds in which she chose action 2. In the following, $S_1(t)$ and $S_2(t)$ denote the payoff sums in $s_1(t)$ and $s_2(t)$, respectively.

After the drawing of the samples, the cumulated payoffs $S_1(t)$ and $S_2(t)$ are calculated and the action with the higher cumulated payoff is played, if there is one. If the samples of both possible actions have the same cumulated payoff, the agent randomizes with $p_1 = p_2 = 0.5$.

$$p_i(t) = \begin{cases} 1 & \text{if } S_i(t) > S_j(t), \\ 0.5 & \text{if } S_i(t) = S_j(t), \\ 0 & \text{else,} \end{cases} \tag{A.2}$$

for $i, j = 1, 2$ and $i \neq j$.

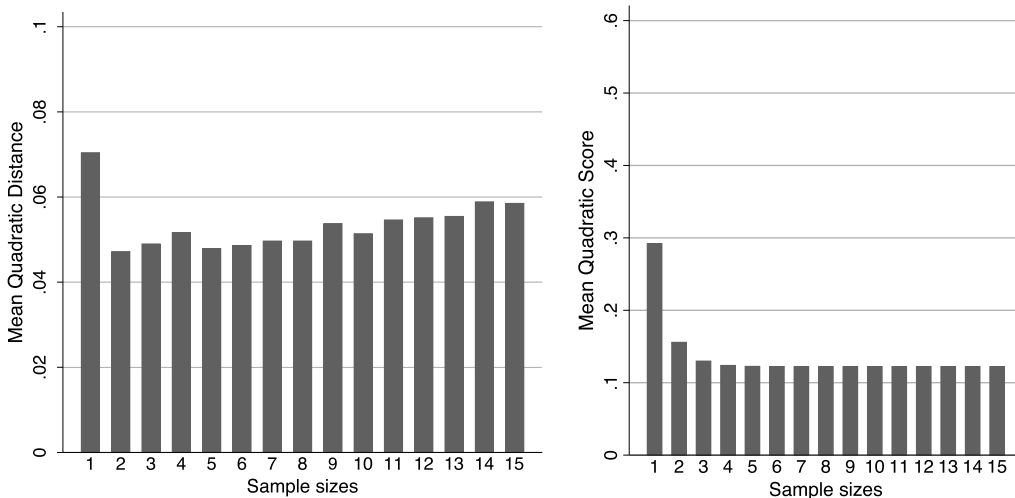


Fig. 15. Figure gives the mean quadratic distance (left) and the mean quadratic score (right) of payoff-sampling for different sample sizes.

As before, $p_i(t)$ is the probability of playing action i in period t . At the beginning and until positive payoffs for each action have been obtained at least once, the agent chooses both actions with equal probabilities, i.e., $p_1 = p_2 = 0.5$.

The optimal sample sizes are calculated analogously to the ones of action-sampling learning. Fig. 15 gives the mean quadratic distances and the mean quadratic scores for different sample sizes. For the quadratic distance, the optimal sample size is $n = 2$, and for the quadratic scores, it is $n = 1$.

A.3. Parameter-free self-tuning EWA

One of the referees suggested to test additionally a parameter-free version of self-tuning EWA and we are very grateful for this idea. This allows a better comparison between the similar action choice mechanisms of impulse-matching learning (based on impulses) and self-tuning EWA (based on attractions) without the bias of an estimated parameter. Therefore, the probability of playing action i in period t depending on the attractions is not calculated as a logit response function with the parameter λ . Instead, the probabilities are calculated with the relative attractions. Thus, Eq. (14) is replaced by Eq. (A.3):

$$p_i(t) = \frac{A_i(t-1)}{\sum_{j=1}^2 A_j(t-1)}. \quad (\text{A.3})$$

Everything else stays unchanged.

Appendix B. Comparison of all learning models

B.1. Simulations

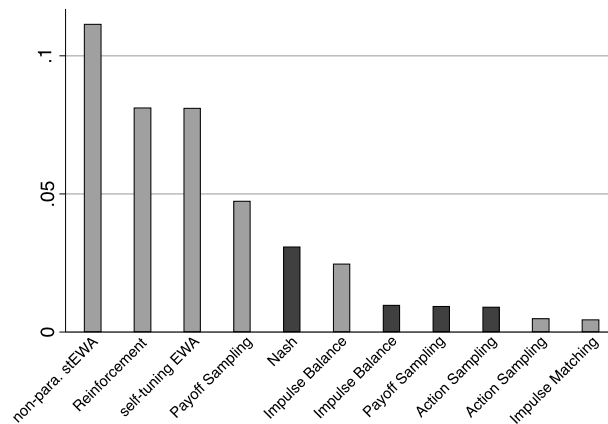


Fig. 16. Mean quadratic distance for all learning models and stationary concepts (dark grey).

B.2. Individual round-by-round behavior

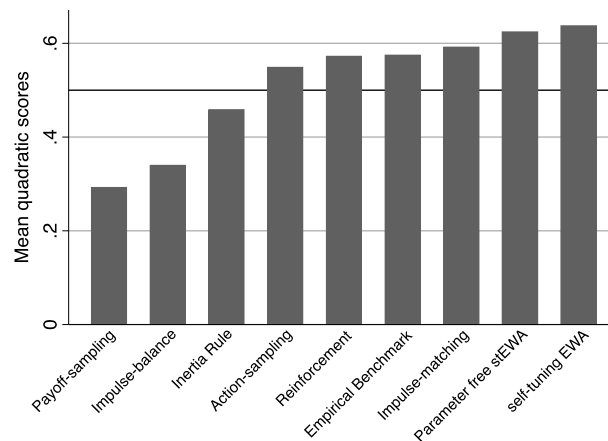


Fig. 17. Mean quadratic scores for all learning types.

Appendix C. Mean probabilities of simulations over time

C.1. Main learning models

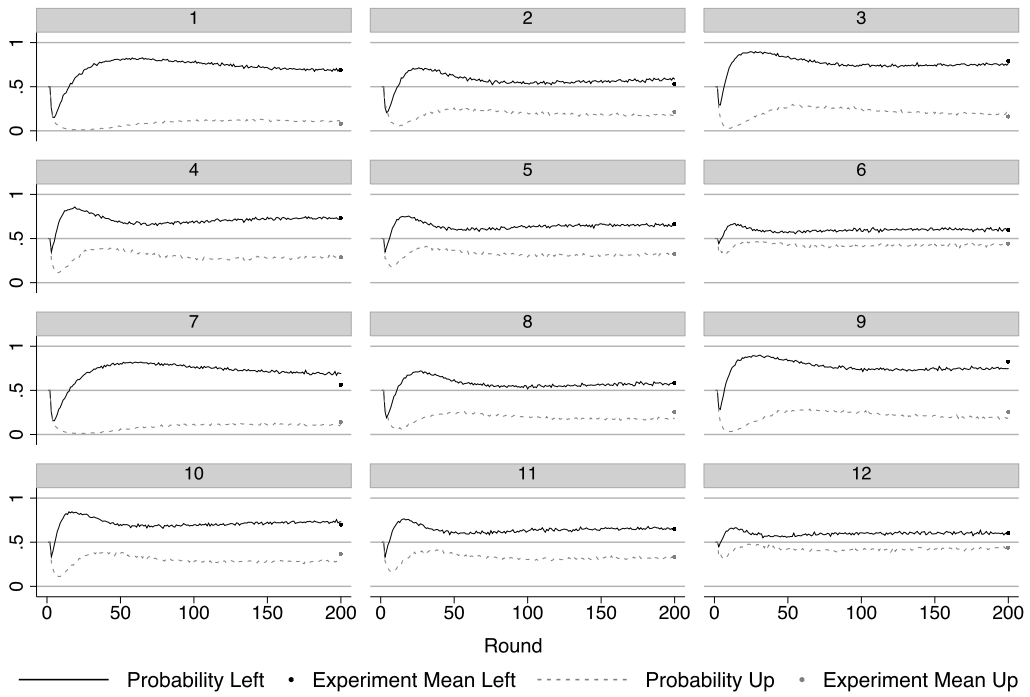


Fig. 18. Mean probabilities of *action-sampling learning* for left and up over 500 simulations runs per game, and frequencies in the experiments over 200 rounds.

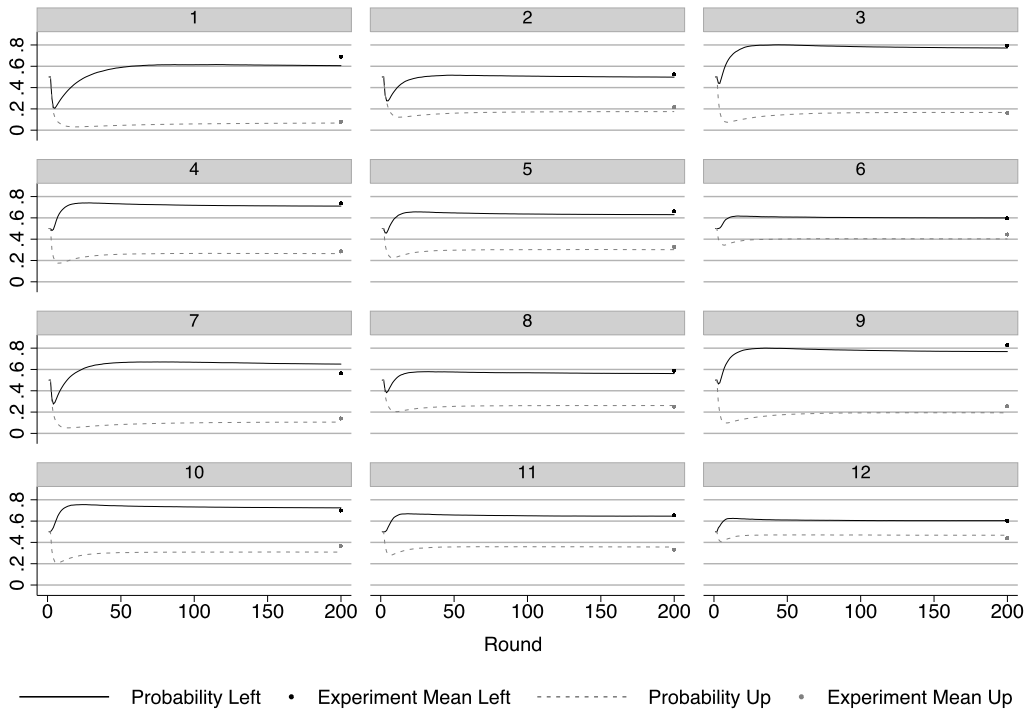


Fig. 19. Mean probabilities of *impulse-matching learning* for left and up over 500 simulations runs per game, as well as mean frequencies in the experiments over 200 rounds.

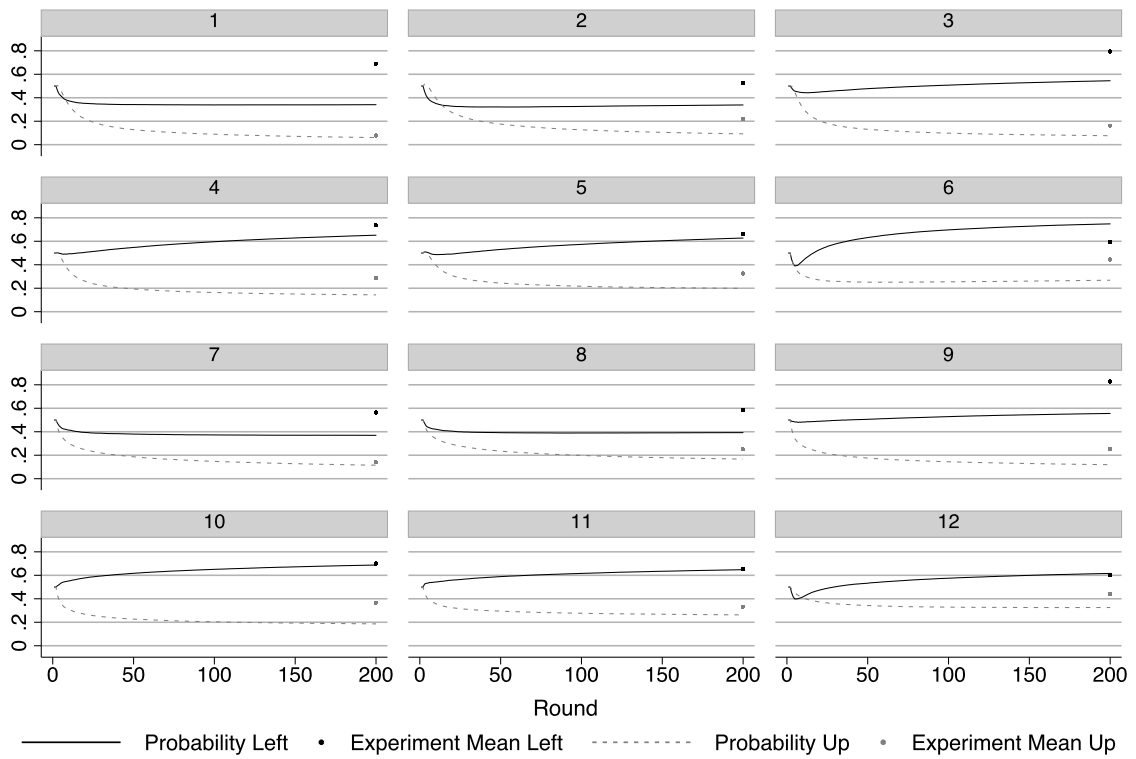


Fig. 20. Mean probabilities of reinforcement learning for left and up over 500 simulations runs per game, as well as mean frequencies in the experiments over 200 rounds.

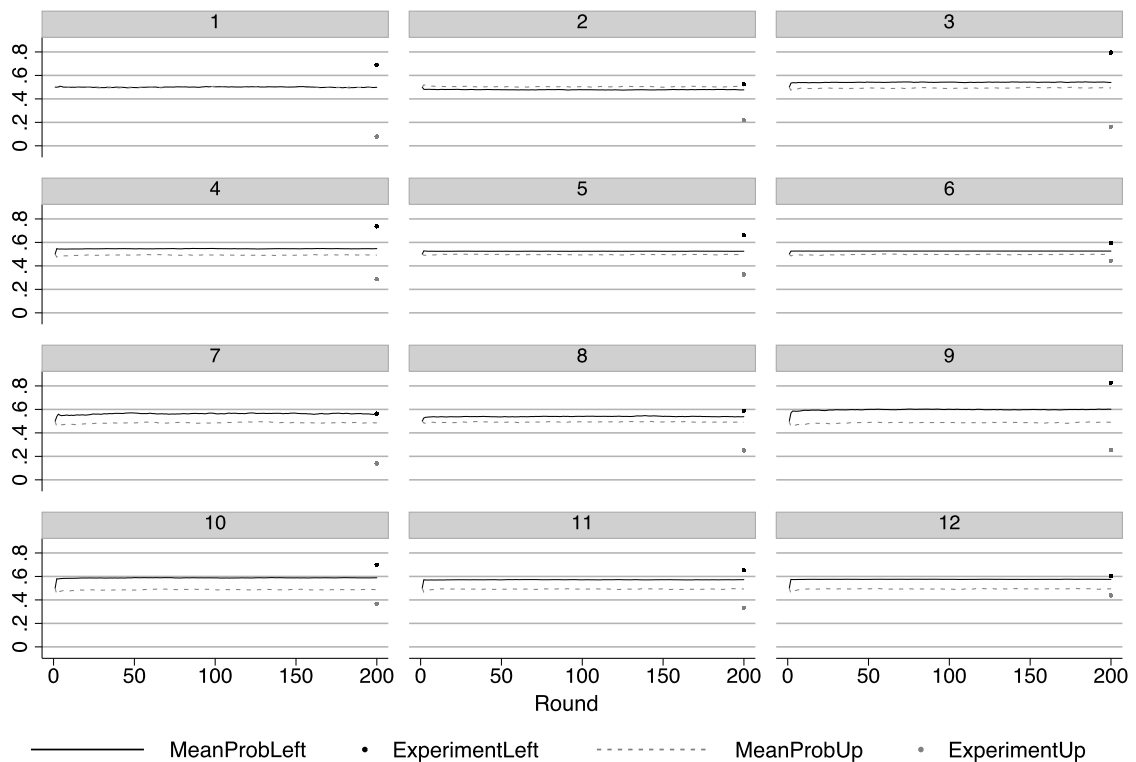


Fig. 21. Mean probabilities of self-tuning EWA learning for left and up over 500 simulations runs per game, as well as mean frequencies in the experiments over 200 rounds.

C.2. Excluded learning models

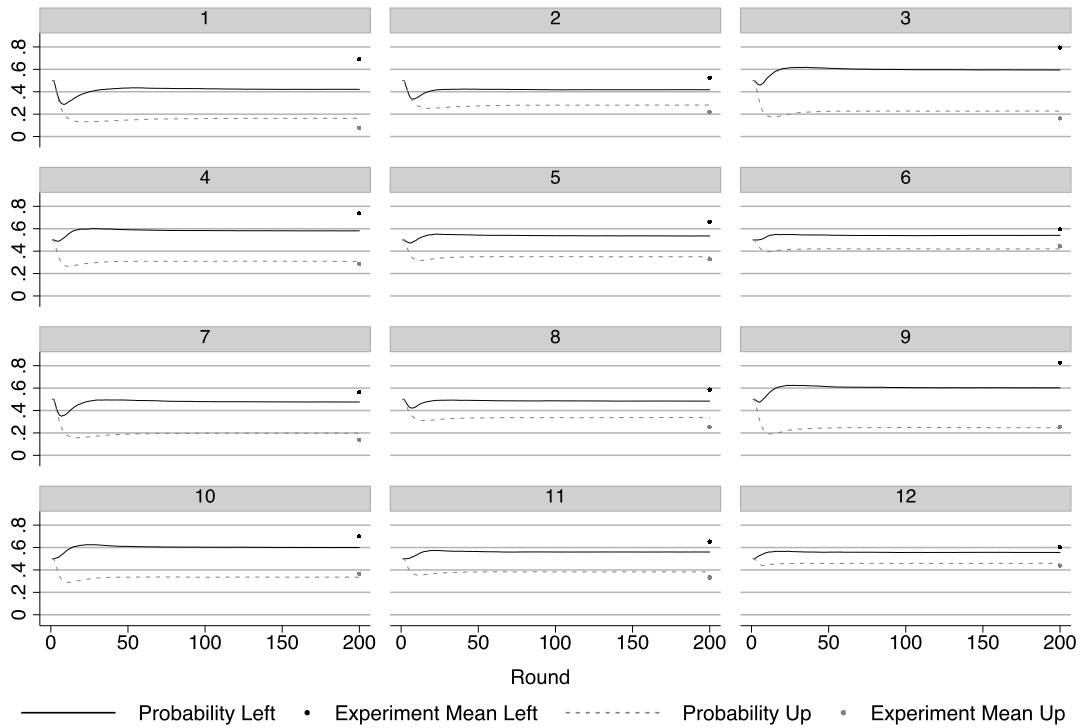


Fig. 22. Mean probabilities of *impulse-balance learning* for left and up over 500 simulations runs per game, as well as mean frequencies in the experiments over 200 rounds.

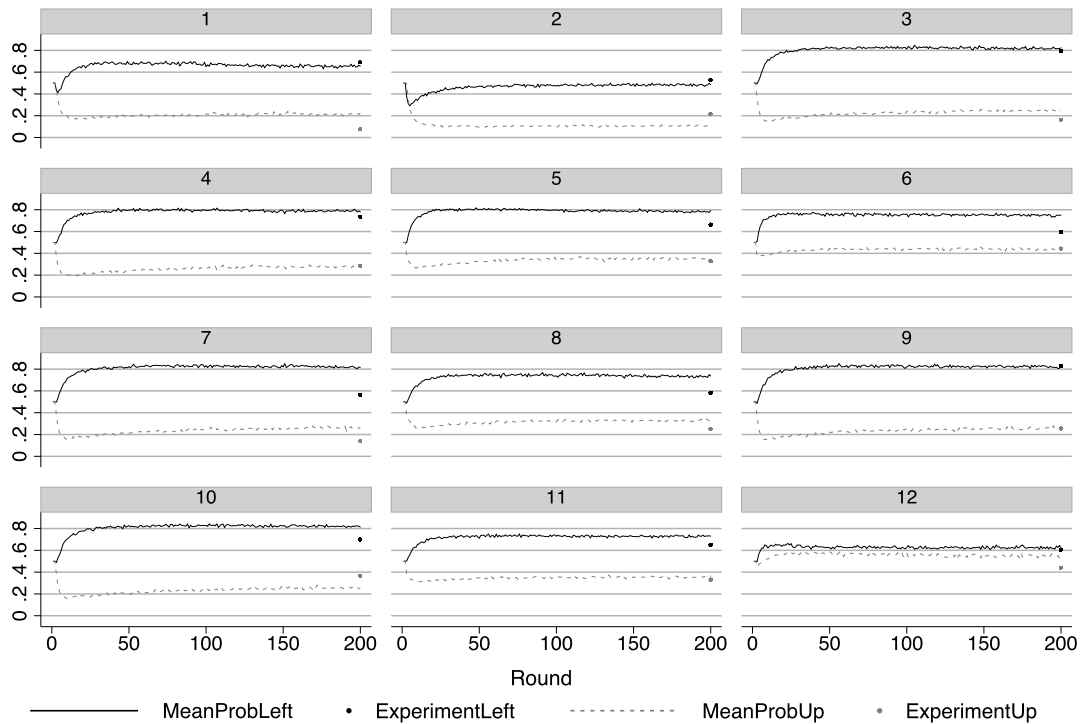


Fig. 23. Mean probabilities of *payoff-sampling learning* for left and up over 500 simulations runs per game, as well as mean frequencies in the experiments over 200 rounds.

- Harley, Calvin B., 1981. Learning the evolutionarily stable strategy. *J. Theor. Biol.* 89 (4), 611–633.
- Ho, Teck H., Camerer, Colin, Chong, Juin-Kuan, 2007. Self-tuning experience-weighted attraction learning in games. *J. Econ. Theory* 133, 177–198.
- Hopkins, E., 2002. Two competing models of how people learn in games. *Econometrica* 2002 (70), 2141–2166.
- Kahneman, Daniel, Tversky, Amos, 1979. Prospect theory: An analysis of decision under risk. *Econometrica* 47 (2), 263–291.
- Kalai, Ehud, Lehrer, Ehud, 1993. Rational learning leads to Nash equilibrium. *Econometrica* 61 (5), 1019–1045.
- Marchiori, Davide, Warglien, Massimo, 2008. Predicting human interactive learning by regret-driven neural networks. *Science* 319, 1111–1113.
- McKelvey, Richard D., Palfrey, Thomas R., 1995. Quantal response equilibria for normal form games. *Games Econ. Behav.* 10 (1), 6–38.
- Metrick, Andrew I., Polak, Ben, 1994. Fictitious play in 2×2 games: A geometric proof of convergence. *Econ. Theory* 4, 923–933.
- Miyasawa, K., 1961. On the Convergence of the Learning Process in a 2×2 Non-Zero-Sum Game. Economic Research Program. Princeton University, Research Memorandum, 33.
- Osborne, Martin J., Rubinstein, Ariel, 1998. Games with procedurally rational players. *Amer. Econ. Rev.* 88 (4), 834–847.
- Roth, Alvin E., Erev, Ido, 1995. Learning in extensive-form games: Experimental data and simple dynamic models in the intermediate term. *Games Econ. Behav.* 8, 164–212.
- Salmon, Timothy C., 2001. An evaluation of econometric models of adaptive learning. *Econometrica* 69 (6), 1597–1628.
- Selten, Reinhard, 1998. Axiomatic characterization of the quadratic scoring rule. *Exper. Econ.* 1, 43–62.
- Selten, Reinhard, Abbink, Klaus, Cox, Ricarda, 2005. Learning direction theory and the winner's curse. *Exper. Econ.* 8, 5–20.
- Selten, Reinhard, Buchta, Joachim, 1999. Experimental sealed bid first price auctions with directly observed bid functions. In: Budescu, David, Erev, Ido, Zwick, Rami (Eds.), *Games and Human Behavior: Essays in the Honor of Amnon Rapoport*. Lawrence Associates, Mahwah, NJ, pp. 79–104.
- Selten, Reinhard, Chmura, Thorsten, 2008. Stationary concepts for experimental 2×2 -games. *Amer. Econ. Rev.* 98 (3), 938–966.
- Selten, Reinhard, Chmura, Thorsten, Goerg, Sebastian J., 2011. Stationary concepts for experimental 2×2 -games: Reply. *Amer. Econ. Rev.* 101 (2), 1041–1044.
- Selten, Reinhard, Stoecker, Rolf, 1986. End behavior in sequences of finite prisoner's dilemma supergames. *J. Econ. Behav. Organ.* 7, 47–70.