



Should I remember more than you? Best responses to factored strategies

René Levínský¹ · Abraham Neyman² · Miroslav Zelený³

Accepted: 26 August 2020
© Springer-Verlag GmbH Germany, part of Springer Nature 2020

Abstract

In this paper we offer a new, unifying approach to modeling strategies of bounded complexity. In our model, the strategy of a player in a game does not directly map the set H of histories to the set of her actions. Instead, the player's perception of H is represented by a map $\varphi: H \rightarrow X$, where X reflects the "cognitive complexity" of the player, and the strategy chooses its mixed action at history h as a function of $\varphi(h)$. In this case we say that φ is a factor of a strategy and that the strategy is φ -factored. Stationary strategies, strategies played by finite automata, and strategies with bounded recall are the most prominent examples of factored strategies in multistage games. A factor φ is recursive if its value at history h' that follows history h is a function of $\varphi(h)$ and the incremental information $h' \setminus h$. For example, in a repeated game with perfect monitoring, a factor φ is recursive if its value $\varphi(a_1, \dots, a_t)$ on a finite string of action profiles (a_1, \dots, a_t) is a function of $\varphi(a_1, \dots, a_{t-1})$ and a_t . We prove that in a discounted infinitely repeated game and (more generally) in a stochastic game with finitely many actions and perfect monitoring, if the factor φ is recursive, then for every profile of φ -factored strategies there is a pure φ -factored strategy that is a best reply, and if the stochastic game has finitely many states and actions and the factor φ has a finite range then there is a pure φ -factored strategy that is a best reply in all the discounted games with a sufficiently large discount factor.

Keywords Bounded rationality · Factored strategies · Bounded recall strategies · Finite automata

René Levínský: Research was supported by GAČR Grant 17-19672S and by the Ministry of Education, Youth and Sports of the Czech Republic through the project *SHARE – CZ* + (CZ.02.1.01/0.0/0.0/16_013/0001740), Abraham Neyman: Research was supported in part by Israel Science Foundation grant 1123/06. Miroslav Zelený: Research was supported by MSM research project 0021620839 financed by MSMT and GAČR Grant 17-19672S.

Extended author information available on the last page of the article

22 1 Introduction

23 There are two widely studied approaches to modeling strategies of bounded complex-
 24 ity in repeated games. One approach is to assume that players have stationary bounded
 25 recall Aumann (1981), Lehrer (1988), Aumann and Sorin (1989). Such players have
 26 imperfect consciousness of the actual stage of the game, and their actions in the cur-
 27 rent stage game rely only on the signals, e.g., players' actions, of the last k rounds.
 28 The strategies of such players are called stationary bounded recall strategies, hence-
 29 forth, SBR strategies. The other approach is to assume that players have stationary
 30 finite memory Neyman (1985), Rubinstein (1986), Abreu and Rubinstein (1988), Ben-
 31 Porath (1993). Such players track a finite summary of the history of play and update
 32 their summaries as a function of the current memory and the new input. The strategies
 33 of such players are implementable by finite automata (Moore machines).

34 Both approaches provide a measure of the complexity of the strategy. Under the
 35 bounded recall approach, the complexity of a strategy is described by the "depth of
 36 recall" and under the finite memory approach the complexity of a strategy is measured
 37 by the minimal number of states the automaton must have to play the given strategy.

38 In this paper, the question raised by Kalai (1990), "What information system (size
 39 and structure) should a player maintain when playing a strategic game?," is examined
 40 in the context of strategies of bounded complexity. Explicitly, we study the minimal
 41 complexity of a strategy that is the best response to a strategy of a given complexity.
 42 Abreu and Rubinstein (1988) shows that for every finite automaton A_1 of Player 1,
 43 there exists a finite automaton A_2 of Player 2 whose number of states is less than or
 44 equal to the number of states of A_1 , such that A_2 maximizes its own payoff in the
 45 discounted repeated game against A_1 . Here, we address this question in the broader
 46 context of the newly defined concept of factored strategies.

47 In addition, we study this question not only in the classic repeated game context,
 48 namely, in supergames with finitely many stage actions and perfect monitoring, but
 49 also in the stochastic game context.

50 Let us also remark that relationships between the minimal size of an automaton that
 51 implements a strategy and the complexity of the strategy were investigated by Kalai
 52 and Stanford (1988) who studied the complexity of strategies forming equilibria.

53 A strategy is any map from the player's information sets to the set of feasible
 54 actions. In the repeated game model with perfect monitoring the information sets are
 55 the possible histories, i.e., the finite strings of action profiles.

56 Under our bounded rationality approach, the player may not be cognitively capable
 57 of distinguishing two distinct histories. Hence, the player can base her actions only
 58 on some equivalence classes of indistinguishable histories.¹

59 Let φ be the map that assigns to any possible history h its equivalence class of
 60 indistinguishable histories x . Hence, φ is a function from the set H of all possible
 61 histories onto the set X of equivalence classes of indistinguishable histories. Naturally,
 62 we are interested in cases where the function φ is not injective. A strategy whose action
 63 at any history h is a function of $\varphi(h)$ is called a φ -factored strategy.

¹ Partitions of the space of histories have already been used for other purposes in the theory of Markovian strategies Maskin and Tirole (2001).

Examples of factored strategies in the repeated game are SBR strategies and strategies that are implementable by an automaton. For example, let us illustrate that a stationary bounded k -recall strategy is a factored strategy: it does not distinguish between two different histories h and h' that are identical in the last k coordinates; hence, its choice of an action depends on the image of the factor φ that maps a history to its last k action profiles.

The question that arises is, what are the game models and the factors φ such that for any φ -factored strategy of Player 1, Player 2 has a best reply that is a φ -factored strategy?

We show an example of a discounted repeated game, a factor φ , and a profile of φ -factored strategies for which no φ -factored strategy is a best reply. Hence, even in the discounted repeated game model an additional assumption on the factor φ is needed in order to derive a positive result.

The additional assumption is recursiveness. A factor in the repeated game is *recursive* if its value on any finite string of actions (a_1, \dots, a_t) is a function of a_t and its value on (a_1, \dots, a_{t-1}) . Note that the factor of a SBR strategy and of a strategy that is implementable by an automaton (with deterministic transitions) is recursive.

Our main result asserts that if the factor φ is recursive, then, for any profile of φ -factored strategies in the discounted repeated game with finitely many stage actions, there is a best reply that is a pure φ -factored strategy. This result holds not only in the classic repeated game context but also in the stochastic game context.

Our main result generalizes analogous results of Blackwell (1962) and Derman (1965) for a Markov decision process (MDP), and the proof relies on the same ideas as those used in the proofs of Blackwell (1962) and Derman (1965).

All relevant notions will be defined and discussed in the next section. Section 3 introduces the concept of factored strategies and presents examples. Section 4 contains the main results and their proofs. In Sect. 5 we connect our results with well-known theorems of Derman and Blackwell, Sect. 6 concludes.

2 The game models

If X is a finite or countable set (or a measurable space), then $\Delta(X)$ denotes the set of all probabilities on X . For ease of exposition we focus on multistage two-person games with perfect monitoring. The results of the present paper on the best reply can be extended (from two-person games) to n -person games ($n > 2$) by considering Players 1, 2, \dots , $n - 1$ as a single player.

2.1 Supergames

We start by recalling the model of the two-person supergame with finite action sets. Let $G = \langle A_1, A_2, u_1, u_2 \rangle$ be a strategic game, where A_i is the nonempty finite set of actions of player $i \in \{1, 2\}$ and $u_i: A_1 \times A_2 \rightarrow \mathbb{R}$ is the payoff function of player i . The corresponding supergame G^∞ is played as follows. In each period $t \in \mathbb{N} = \{1, 2, 3, \dots\}$, Players 1 and 2 make simultaneous and independent moves

104 $a_t^i \in A_i, i = 1, 2$. A *play* of the supergame is a sequence of action profiles $(a_t)_{t=1}^\infty$
 105 with $a_t = (a_t^1, a_t^2) \in A := A_1 \times A_2$. A play $(a_t)_{t=1}^\infty$ defines a stream $(u_i(a_t))_{t=1}^\infty$ of
 106 payoffs to player i .

107 We assume perfect monitoring. Hence, a *pure strategy for player i in the supergame*
 108 G^∞ is a mapping $\sigma: A^* \rightarrow A_i$, where A^* denotes the set of all finite sequences of
 109 elements from A (including the empty one). A pure strategy σ plays in the t -th round
 110 the action $\sigma(a_1, \dots, a_{t-1})$, where $(a_1, \dots, a_{t-1}) \in A^{t-1}$ is the sequence of actions
 111 that have already been played.

112 A *behavioral strategy for player i in the supergame G^∞* is a mapping $\sigma: A^* \rightarrow$
 113 $\Delta(A_i)$. A behavioral strategy σ plays in the t -th round an action $a_t^i \in A_i$ with the
 114 probability $\sigma(a_1, \dots, a_{t-1})(a_t^i)$, where $(a_1, \dots, a_{t-1}) \in A^{t-1}$ is the sequence of
 115 actions that have already been played. Pure strategies can be viewed as a special case
 116 of behavioral strategies by identifying A_i with the Dirac measures on A_i . This point
 117 of view will be used throughout the paper.

118 2.2 Supergames with time-dependent stage games

119 The previous concept can be generalized as follows. Let $\Gamma = \{(A_1(t), A_2(t),$
 120 $u_1(t), u_2(t))\}$ be a sequence of stage games. The game Γ is played as follows. In
 121 each period $t \in \mathbb{N}$, Players 1 and 2 make simultaneous and independent moves
 122 $a_t^i \in A_i(t), i = 1, 2$. A *play* of Γ is a sequence of action profiles $(a_t)_{t=1}^\infty$ with
 123 $a_t = (a_t^1, a_t^2) \in A(t) := A_1(t) \times A_2(t)$. A play $(a_t)_{t=1}^\infty$ defines a stream $(u_i(t)(a_t))_{t=1}^\infty$
 124 of payoffs to player i . The pure and behavioral strategies of player i in Γ are defined
 125 as a family of mappings $\sigma_t, t \geq 1$, from $\times_{s < t} A(s)$ to $\Delta(A_i(t))$.

126 2.3 Stochastic games

127 A two-person *stochastic game* with finite action sets is a 5-tuple $\Gamma = \langle S, A, u, p, \mu \rangle$,
 128 where

- 129 • S is a nonempty set, called the *state space*;
- 130 • $A(z) = A_1(z) \times A_2(z)$ is a cartesian product of two nonempty finite sets. It is the
 131 set of *action pairs*: for every state $z \in S$, $A_i(z)$ is the finite set of feasible actions
 132 of player $i \in \{1, 2\}$ in state z ;
- 133 • $u = (u_1, u_2)$ is a function from the set of pairs (z, a) with $z \in S$ and $a \in A(z)$ to
 134 \mathbb{R}^2 . It is the *payoff function*: $u_i(z, a)$ is the payoff to player i when the action pair
 135 $a \in A(z)$ is played in state $z \in S$;
- 136 • p is a *transition function*: for each state $z \in S$ and each action pair $a \in A(z)$,
 137 $p(\cdot | z, a) \in \Delta(S)$ is the distribution of the next state; i.e., $p(z' | z, a)$ is the
 138 probability of moving to state z' if the players played the action pair a in state z ;
- 139 • $\mu \in \Delta(S)$ is the distribution of the initial state.

140 A *play* of the stochastic game Γ is a sequence of states and action pairs
 141 $(z_1, a_1, \dots, z_t, a_t, \dots)$ with $a_t \in A(z_t)$. The set of all plays in Γ is denoted by H^∞ ,
 142 i.e.,

$$143 \quad H^\infty = \{(z_1, a_1, \dots, z_t, a_t, \dots) : z_t \in S \text{ and } a_t \in A(z_t) \text{ for every } t \in \mathbb{N}\}.$$

144 Additional measurability conditions, which are not stated here, are assumed in the
145 case where the state space is a measurable space.

146 Consider a two-person stochastic game with perfect monitoring. In such a game,
147 each player observes the sequence $(z_1, a_1, \dots, a_{t-1}, z_t)$ of past states and action pairs
148 before choosing her action in state t . Hence, a *pure strategy of player i* specifies her
149 action $a_t^i \in A_i(z_t)$ as a function of the sequence $(z_1, a_1, \dots, a_{t-1}, z_t)$ of past states
150 and action pairs. Similarly, a *behavioral strategy of player i* specifies the probability
151 of her action $a_t^i \in A_i(z_t)$ as a function of the sequence $(z_1, a_1, \dots, a_{t-1}, z_t)$ of past
152 states and action pairs.

153 A pair σ^1 and σ^2 of behavioral strategies of Players 1 and 2, respectively, defines a
154 probability distribution P_{σ^1, σ^2} on the space of plays H^∞ as follows. The distribution
155 of z_1 is μ . The conditional probability that $a_t = (a^1, a^2) \in A(z_t)$, given z_1, a_1, \dots, z_t ,
156 is

$$157 \quad \sigma^1(z_1, a_1, \dots, z_t)(a^1) \cdot \sigma^2(z_1, a_1, \dots, z_t)(a^2).$$

158 The conditional distribution of z_{t+1} , given $z_1, a_1, \dots, z_t, a_t$, is $p(\cdot | z_t, a_t)$.

159 Similarly, a finite initial segment h of a play of the stochastic game along a pair σ^1
160 and σ^2 of behavioral strategies (of Players 1 and 2, respectively) defines a probability
161 distribution $P_{(\sigma^1, \sigma^2|h)}$ on the set $H(h)$ of all plays that extend the finite sequence h .
162 The expectation with respect to the probability distribution $P_{(\sigma^1, \sigma^2|h)}$ is denoted by
163 $E_{(\sigma^1, \sigma^2|h)}$.

164 2.3.1 Payoffs

165 Let σ^1 and σ^2 be strategies of Players 1 and 2, respectively, in the stochastic game
166 Γ and let $h = (z_1, a_1, \dots, z_\ell)$, or $h = (z_1, a_1, \dots, a_{\ell-1})$, be a finite history. The
167 subgame that follows the initial history h is denoted by $(\Gamma | h)$.

168 The β -discounted payoff to player i , $0 \leq \beta < 1$, in $(\Gamma | h)$ is defined by

$$169 \quad V_\beta^i(\sigma^1, \sigma^2 | h) = E_{(\sigma^1, \sigma^2|h)} \left(\sum_{t=\ell}^{\infty} \beta^{t-\ell} u_i(z_t, a_t) \right).$$

170 2.3.2 Subclasses of stochastic games

171 A *Markov decision process* (henceforth *MDP*) is a one-person stochastic game; hence,
172 it is a special case of a stochastic game.

173 A *supergame* is a stochastic game with a single state; hence, it is a special case of
174 a stochastic game.

175 Similarly, a *supergame with time-dependent stage games* can be viewed as a stochas-
176 tic game whose state space is the set of positive integers \mathbb{N} and the transitions are
177 deterministic: $\mathbb{N} \ni t \mapsto t + 1$.

178 Therefore, the β -discounted payoff is well defined for a supergame, as well as for
179 a supergame with time-dependent stage games as long as the stage payoffs are either
180 bounded or grow at a subexponential rate in t .

181 Therefore, results for stochastic games with finitely many states have direct conse-
 182 quences for supergames and results for stochastic games with countably many states
 183 have direct consequences for supergames with time-dependent stage games.

184 3 Factored strategies

185 In this section we define the concept of a factored strategy. This concept plays a central
 186 role in the results of the paper. For simplicity, we start with the introduction of the
 187 concept in the supergame model, and thereafter continue with the definition in the
 188 more general model of a stochastic game.

189 3.1 Factored strategies in supergames

190 Let F denote the set of all information sets. In a supergame with perfect monitoring
 191 F is the set of all finite histories A^* . Let X be a set and let φ be a mapping from F
 192 to X . We say that a behavioral strategy σ is a φ -factored strategy for player i in the
 193 supergame G^∞ if there is an action function $\omega: X \rightarrow \Delta(A_i)$ such that $\sigma = \omega \circ \varphi$.

194 The factor φ is called *recursive* if there is a function $g: X \times A \rightarrow X$ such that
 195 $\varphi(a_1, \dots, a_l) = g(\varphi(a_1, \dots, a_{l-1}), a_l)$. Roughly speaking, if a player's cognition is
 196 captured by a recursive factor then she never learns something she previously ignored.

197 We continue by recalling the concepts of time-independent bounded recall strate-
 198 gies and of (their superset of) strategies that are implementable by time-independent
 199 automata with finitely many states. Such strategies are examples of factored strategies
 200 with a recursive factor.

201 Let $k \in \mathbb{N}$. A behavioral k -SBR strategy for player i in the supergame G^∞ is defined
 202 by a pair (e, ω) , where $e = (e_1, e_2, \dots, e_k) \in A^k$ and a mapping $\omega: A^k \rightarrow \Delta(A_i)$.
 203 The strategy that is defined by (e, ω) is played as follows. If moves $a_1, \dots, a_l \in A$ have
 204 been played, then player i takes the sequence b of the last k elements of the sequence
 205 $(e_1, \dots, e_k, a_1, \dots, a_l)$, and his $(l + 1)$ -th move is $a \in A_i$ with the (conditional)
 206 probability $\omega(b)(a)$. A pure k -SBR strategy for player i in the supergame G^∞ is
 207 defined in a similar way.

208 Defining $\varphi(a_1, \dots, a_l)$ by the last k elements of the finite sequence of action profiles

$$209 \quad (e_1, \dots, e_k, a_1, \dots, a_l),$$

210 we see that the k -SBR strategy σ defined above obeys $\sigma = \omega \circ \varphi$, and the factor φ
 211 is recursive; thus, σ is a factored strategy with a recursive factor φ that has a finite
 212 range.

213 We say that a (pure) behavioral strategy σ is a (pure) behavioral SBR strategy if σ
 214 is a (pure) behavioral k -SBR strategy for some $k \in \mathbb{N}$.

215 A behavioral automaton (for Player 1 in the supergame G^∞) is a quadruple
 216 $\langle M, m^*, \alpha, \tau \rangle$, where M is a nonempty set (the state space), $m^* \in M$ is the ini-
 217 tial state, $\alpha: M \rightarrow \Delta(A_1)$ is a probabilistic action function, and $\tau: M \times A \rightarrow M$
 218 is a transition function. A k -state behavioral automaton is a behavioral automa-

219 ton where the set M has k elements. A behavioral automaton $\langle M, m^*, \alpha, \tau \rangle$ defines
 220 a behavioral strategy σ^1 (for Player 1) inductively: $m_1 = m^*$, $\sigma^1(\emptyset) = \alpha(m_1)$,
 221 $\sigma^1(a_1, \dots, a_{t-1}) = \alpha(m_t)$, where $m_t = \tau(m_{t-1}, a_{t-1})$. A behavioral automaton
 222 $\langle M, m^*, \alpha, \tau \rangle$ defines a factored strategy with a recursive factor φ , where $X = M$,
 223 $\varphi(\emptyset) = m^*$, $\varphi(a_1, \dots, a_t) = \tau(\varphi(a_1, \dots, a_{t-1}), a_t)$, and $\omega = \alpha$. A k -state (deter-
 224 ministic) automaton is defined by replacing $\Delta(A_1)$ (in the definition of α) by A_1 .

225 3.2 Factored strategies in stochastic games

226 A factor φ in a stochastic game keeps track of the current state as well as of some
 227 statistics of the past history. That is, φ maps the set $F := S \times (A \times S)^*$ of all information
 228 sets in the stochastic game to a set X and the map φ obeys $\varphi(z_1, a_1, \dots, z_t) \neq$
 229 $\varphi(z'_1, a'_1, \dots, z'_t)$ whenever $z_t \neq z'_t$. Equivalently, X is the disjoint union of sets X_z
 230 and φ is a mapping from F to X such that $\varphi(z_1, a_1, \dots, z_t) \in X_{z_t}$.

231 A behavioral strategy σ is a φ -factored strategy of player i in the stochastic game
 232 if there is an action function $\omega: X \rightarrow \Delta(A_i)$ such that

$$233 \quad \sigma(z_1, a_1, \dots, a_{t-1}, z_t) = \omega(\varphi(z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t)).$$

234 In a stochastic game, a factor φ is called *recursive* if there is a function $g: X \times A \times$
 235 $S \rightarrow X$ such that

$$236 \quad \varphi(z_1, a_1, \dots, z_t, a_t, z_{t+1}) = g(\varphi(z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t), a_t, z_{t+1}).$$

237 3.3 Time-dependent factored strategies

238 A time-dependent factor φ in a multistage game keeps track of the stage number as
 239 well as of some statistic of the past history. That is, if h and h' are two information sets
 240 at two distinct stages t and t' , then $\varphi(h) \neq \varphi(h')$. Hence, a time-dependent factor φ
 241 is a factor φ that maps the set of histories of a multistage game into the disjoint union
 242 $X = \bigcup_{t \in \mathbb{N}} X_t$, where a history up to the play of the t -th stage is mapped into X_t .

243 A *time-dependent factored strategy* for player i is a factored strategy with respect
 244 to a time-dependent factor φ .

245 We continue by recalling the concepts of a time-dependent action automaton, a
 246 time-dependent transition automaton, and a time-dependent (action and transition)
 247 automaton. A strategy that is defined by a time-dependent automaton is a time-
 248 dependent factored strategy.

249 A *time-dependent action automaton* is defined by replacing the action function α
 250 by a sequence of action functions α_t , $t \geq 1$, where α_t defines the action at stage t [see,
 251 e.g., (Neyman 1997)]. Similarly, a *time-dependent transition automaton* is obtained by
 252 replacing the (stationary) transition function τ with a sequence of time-dependent
 253 transitions τ_t , $t \geq 1$, where τ_t defines the transition at stage t . Finally, a *time-*
 254 *dependent (action and transition) automaton* in the supergame G^∞ is a quadruple
 255 $\langle M, m^*, (\alpha_t)_{t=1}^\infty, (\tau_t)_{t=1}^\infty \rangle$, where M is a nonempty set (the state space), $m^* \in M$ is the

256 initial state, $\alpha_t: M \rightarrow \Delta(A_t)$ is a probabilistic action function, and $\tau_t: M \times A \rightarrow M$
 257 is a deterministic transition function.

258 A time-dependent automaton $\langle M, m^*, (\alpha_t)_{t=1}^\infty, (\tau_t)_{t=1}^\infty \rangle$ defines a behavioral strategy
 259 σ^1 (for Player 1) inductively: $m_1 = m^*, \sigma^1(\emptyset) = \alpha_1(m_1)$, and $\sigma^1(a_1, \dots, a_{t-1}) =$
 260 $\alpha_t(m_t)$, where $m_t = \tau_t(m_{t-1}, a_{t-1})$.

261 Note that a time-dependent automaton $\langle M, m^*, (\alpha_t)_{t=1}^\infty, (\tau_t)_{t=1}^\infty \rangle$ defines the same
 262 strategy as the automaton $\langle M \times \mathbb{N}, m^{**}, \alpha, \tau \rangle$ with $m^{**} = (m^*, 1)$, $\alpha(m, t) = \alpha_t(m)$,
 263 and $\tau((m, t), a) = (\tau_t(m, a), t + 1)$. Therefore, the corresponding strategy is a
 264 recursive φ -factored strategy, where $\varphi: A^* \rightarrow M \times \mathbb{N}$ is given by $\varphi(a_1, \dots, a_t) =$
 265 $\tau(\varphi(a_1, \dots, a_{t-1}), a_t)$ and $\omega = \alpha$.

266 We recall the concept of a strategy with time-dependent recall. Let $k: \mathbb{N} \rightarrow \mathbb{N}$
 267 be a function with $k(t) < t$ for every $t \in \mathbb{N}$. A *behavioral* (respectively, *pure*) k -
 268 *SBR strategy* in the supergame G^∞ is defined analogously to the above case but the
 269 action at stage t depends only on the last $k(t)$ stage actions [see, e.g., (Neyman and
 270 Okada 2009)]. Let σ be such a strategy. Setting $\varphi(a_1, \dots, a_t) = (a_{t-k(t)+1}, \dots, a_t)$,
 271 we easily see that σ is φ -factored.

272 We now define the concept of a time-dependent strategy with time-dependent recall.
 273 Let $k: \mathbb{N} \rightarrow \mathbb{N}$ be a function with $k(t) < t$ for every $t \in \mathbb{N}$. A *behavioral* (respectively,
 274 *pure*) *time-dependent k-BR strategy* in the supergame G^∞ is defined analogously to
 275 the above case but the action at stage t depends only on t and the last $k(t)$ stage actions.
 276 Let σ be such a strategy. Setting $\varphi(a_1, \dots, a_t) = (t, (a_{t-k(t)+1}, \dots, a_t))$, we easily
 277 see that σ is φ -factored. Moreover, if $k(t + 1) \leq k(t) + 1$ for every $t \in \mathbb{N}$, then φ is
 278 recursive.

279 4 The Main results

280 In this section we state and prove the main results in the stochastic game context. In a
 281 later section we will spell out the implications of the main results in the special cases
 282 of a supergame, a supergame with time-dependent stage games, and an MDP.

283 It should be noted that the main result is analogous to, and implies, known results
 284 in the theory of Markov decision processes, and the proof follows the same ideas as
 285 the proofs of the analogous result for an MDP. However, there is a slight difference,
 286 on which we will comment in Remark 5.4.

287 The first result, Theorem 4.1, gives conditions on a stochastic game such that for
 288 every recursive factor φ , a φ -factored strategy of Player 1, and a discount factor β ,
 289 Player 2 has a φ -factored strategy that is a best reply in the β -discounted game.

290 **Theorem 4.1** *Let $\Gamma = \langle S, A, u, p, \mu \rangle$ be a two-person stochastic game with countably*
 291 *many states,² finitely many actions in each state, and a bounded payoff function u_2 .*
 292 *Let σ^1 be a φ -factored behavioral strategy of Player 1 in Γ .*

293 *If φ is recursive, then, for every $\beta \in (0, 1)$, there exists a φ -factored pure strategy*
 294 *σ^2 of Player 2, such that for every behavioral strategy ρ of Player 2 in Γ and every*

² The assumption of countably many states is not essential here. The result holds also for stochastic games with a continuum of states and some measurability assumptions. For simplicity, and in particular to avoid the need to spell out the needed measurability assumptions, we assume countably many states.

295 *initial history* $h = (z_1, a_1, \dots, z_l)$ we have

$$296 \quad V_\beta^2(\sigma^1, \sigma^2 | h) \geq V_\beta^2(\sigma^1, \rho | h).$$

297 The second result, Theorem 4.2, gives conditions on a stochastic game such that for
298 every recursive factor φ and a φ -factored strategy of Player 1, Player 2 has a φ -factored
299 strategy that is a best reply in each one of the β -discounted games with a sufficiently
300 large discount factor $\beta < 1$.

301 **Theorem 4.2** *Let $\Gamma = \langle S, A, u, p, \mu \rangle$ be a two-person stochastic game with finitely
302 many actions in each state. Let φ be a recursive factor with a finite range³ and let σ^1
303 be a φ -factored behavioral strategy of Player 1 in Γ .*

304 *Then, there is a φ -factored pure strategy σ^2 of Player 2 and a discount factor
305 $\beta_0 \in (0, 1)$, such that for every behavioral strategy ρ (of Player 2 in Γ) and every
306 initial history $h = (z_1, a_1, \dots, z_l)$ we have, for every $\beta \in [\beta_0, 1)$,*

$$307 \quad V_\beta^2(\sigma^1, \sigma^2 | h) \geq V_\beta^2(\sigma^1, \rho | h).$$

308 The next result, Theorem 4.3, gives conditions on a stochastic game such that for
309 every recursive factor φ and a φ -factored strategy of Player 1, Player 2 has a φ -factored
310 strategy that is (a) a best reply in each one of the β -discounted games with a sufficiently
311 large discount factor $\beta < 1$ as well as in the limiting-average payoff game, and (b) an
312 approximate (within a possible error $\varepsilon > 0$) best reply in each of the sufficiently long
313 finite-horizon stochastic games.

314 **Theorem 4.3** *Let $\Gamma = \langle S, A, u, p, \mu \rangle$ be a two-person stochastic game with finitely
315 many actions in each state. Let φ be a recursive factor with a finite range and let σ^1
316 be a φ -factored behavioral strategy of Player 1 in Γ .*

317 *Let σ^2 be a φ -factored strategy and $\beta_0 \in (0, 1)$, such that for every behavioral strat-
318 egy ρ (of Player 2 in Γ) and every finite history h , $V_\beta^2(\sigma^1, \sigma^2 | h) \geq V_\beta^2(\sigma^1, \rho | h)$
319 for every $\beta \in [\beta_0, 1)$.⁴ Then, there is a function $N: (0, \infty) \rightarrow \mathbb{N}$, such that for every
320 strategy ρ of Player 2 in Γ and every initial history $h = (z_1, a_1, \dots, z_l)$ we have*

$$321 \quad E_{(\sigma^1, \sigma^2 | h)} \left(\liminf_{n \rightarrow \infty} \frac{1}{n} \sum_{t=l}^{l+n-1} u_2(z_t, a_t) \right) \geq E_{(\sigma^1, \rho | h)} \left(\limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{t=l}^{l+n-1} u_2(z_t, a_t) \right), \quad (1)$$

322 *and*

$$323 \quad E_{(\sigma^1, \sigma^2 | h)} \left(\frac{1}{n} \sum_{t=l}^{l+n-1} u_2(z_t, a_t) \right) \geq E_{(\sigma^1, \rho | h)} \left(\frac{1}{n} \sum_{t=l}^{l+n-1} u_2(z_t, a_t) \right) - \varepsilon \quad \forall n \geq N(\varepsilon). \quad (2)$$

³ The existence of a factor with a finite range implies that the stochastic game has finitely many states.

⁴ The existence of such a strategy σ^2 is guaranteed by Theorem 4.2.

324 **Proof of Theorem 4.1**

325 Fix $\beta \in (0, 1)$. Let σ^1 be a strategy of Player 1 in Γ and let $h = (z_1, a_1, \dots, z_t)$ be a
 326 finite history. Define $V_\beta^2(\sigma^1 | h)$ by

$$327 \quad V_\beta^2(\sigma^1 | h) = \sup_{\sigma^2} V_\beta^2(\sigma^1, \sigma^2 | h),$$

328 where the supremum is taken over all strategies of Player 2. Note that a strategy
 329 τ of Player 2 is a best reply to σ^1 in the β -discounted game $(\Gamma | h')$ whenever
 330 $V_\beta^2(\sigma^1, \tau | h') = V_\beta^2(\sigma^1 | h')$.

331 Assume that τ is a strategy such that for every finite history $h = (z_1, a_1, \dots, z_t)$,

$$332 \quad E_{(\sigma^1, \tau | h)} \left(u_2(z_t, a_t) + \beta V_\beta^2(\sigma^1 | h, a_t, z_{t+1}) \right) = V_\beta^2(\sigma^1 | h). \quad (3)$$

333 Fix a finite history $h' = (z_1, \dots, a_{t'-1}, z_{t'})$ of the stochastic game. For $t \geq t'$ we
 334 denote by U_t the expectation w.r.t. $P_{(\sigma^1, \tau | h')}$ of $V_\beta^2(\sigma^1 | h)$ where h is the finite history
 335 $h = (h', \dots, z_t)$ up to the play at stage t . By taking the expectation w.r.t. $P_{(\sigma^1, \tau | h')}$ in
 336 equality (3), we deduce that

$$337 \quad E_{(\sigma^1, \tau | h')} u_2(z_t, a_t) = U_t - \beta U_{t+1}.$$

338 Multiplying this inequality by $\beta^{t-t'}$ and summing over all $t \geq t'$, we conclude that

$$339 \quad E_{(\sigma^1, \tau | h')} \left(\sum_{t=t'}^{\infty} \beta^{t-t'} u_2(z_t, a_t) \right) = \sum_{t=t'}^{\infty} \beta^{t-t'} U_t - \sum_{t=t'+1}^{\infty} \beta^{t-t'} U_t \\ = U_{t'} = V_\beta^2(\sigma^1 | h').$$

340 Hence, τ is a best reply to σ^1 in the β -discounted game $(\Gamma | h)$.

341 Let φ be a recursive factor and σ^1 be a φ -factored strategy for Player 1 in the
 342 stochastic game Γ . Thus, φ maps the histories (z_1, \dots, z_t) to a set X , and there exist
 343 functions $\omega: X \rightarrow \Delta(A_1)$ and $g: X \times A \times S \rightarrow X$ such that

$$344 \quad \sigma^1(z_1, a_1, \dots, z_t) = \omega(\varphi(z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t)) \quad \text{and} \\ 345 \quad \varphi(z_1, a_1, \dots, z_t, a_t, z_{t+1}) = g(\varphi(z_1, a_1, \dots, z_{t-1}, a_{t-1}, z_t), a_t, z_{t+1}).$$

346 As σ^1 is φ -factored and φ is recursive, $V_\beta^2(\sigma^1 | h) = V_\beta^2(\sigma^1 | h')$ whenever
 347 $\varphi(h) = \varphi(h')$. Let τ be a pure strategy of Player 2 such that $\tau(h) = \tau(h')$ whenever
 348 $\varphi(h) = \varphi(h')$ and $\tau(h)$ maximizes the left-hand side of (3). The existence of such
 349 a strategy τ follows from the facts that φ is recursive and that there are only finitely

350 many possible actions for Player 2 in each state. Hence, τ is a φ -factored strategy. For
 351 every $h = (z_1, a_1, \dots, z_t)$ and every strategy σ^2 we have

$$\begin{aligned}
 E_{(\sigma^1, \tau|h)}(u_2(z_t, a_t) + \beta V_\beta^2(\sigma^1 | h, a_t, z_{t+1})) &\geq E_{(\sigma^1, \sigma^2|h)}(u_2(z_t, a_t) + \beta V_\beta^2(\sigma^1 | h, a_t, z_{t+1})) \\
 &\geq E_{(\sigma^1, \sigma^2|h)}(u_2(z_t, a_t) + \beta V_\beta^2(\sigma^1, \sigma^2 | h, a_t, z_{t+1})) \\
 &= V_\beta^2(\sigma^1, \sigma^2 | h).
 \end{aligned}$$

353 This implies that

$$E_{(\sigma^1, \tau|h)}(u_2(z_t, a_t) + \beta V_\beta^2(\sigma^1 | h, a_t, z_{t+1})) \geq V_\beta^2(\sigma^1 | h). \quad (4)$$

355 Since the opposite inequality in (4) clearly holds we conclude that τ satisfies (3).
 356 Therefore τ is a best reply to σ^1 in the β -discounted games $(\Gamma | h)$. This completes
 357 the proof of Theorem 4.1.

358 Proof of Theorem 4.2

359 In the forthcoming lemmas we fix a stochastic game Γ with finitely many states and
 360 actions and assume that φ is a recursive factor with a finite range X . Let H denote the
 361 set of all finite histories $h = (z_1, a_1, \dots, z_t)$ of the stochastic game.

362 Let \mathbb{R}^X be the finite-dimensional space of all real-valued functions defined on X
 363 with the maximum norm, and let $\ell_\infty(H)$ be the Banach space of all bounded real-
 364 valued functions defined on H with the supremum norm. Let σ^1 and σ^2 be strategies
 365 of Player 1 and Player 2, respectively. Define the map $T_\beta: \ell_\infty(H) \rightarrow \ell_\infty(H)$ by

$$(T_\beta V)(h) = E_{(\sigma^1, \sigma^2|h)}(u_2(z(h), a_t) + \beta V(h, a_t, z_{t+1})),$$

367 where $h = (z_1, a_1, \dots, z_t)$ is a finite history. Then, T_β is a strict contraction ($\|T_\beta V -$
 368 $T_\beta V'\| \leq \beta \|V - V'\|$) and hence has a unique fixed point $V_\beta \in \ell_\infty(H)$. As

$$V_\beta^2(\sigma^1, \sigma^2 | h) = E_{(\sigma^1, \sigma^2|h)}(u_2(z(h), a_t) + \beta V_\beta^2(\sigma^1, \sigma^2 | h, a_t, z_{t+1})),$$

370 we get $V_\beta(h) = V_\beta^2(\sigma^1, \sigma^2 | h)$.

371 Assume that σ^1 and σ^2 are φ -factored strategies of Player 1 and Player 2, respec-
 372 tively. As φ is recursive, the map $T: \mathbb{R}^X \rightarrow \mathbb{R}^X$ that is defined by

$$(TU)(x) = E_{(\sigma^1, \sigma^2|h)}(u_2(z(h), a_t) + \beta U(\varphi(h, a_t, z_{t+1}))),$$

374 where $h = (z_1, a_1, \dots, z_t)$ is any finite history with $\varphi(h) = x$, is well defined. In
 375 addition, T is a strict contraction ($\|TU - TU'\| \leq \beta \|U - U'\|$), and hence has a
 376 unique fixed point $U_\beta \in \mathbb{R}^X$.

377 Note that $V \in \ell_\infty(H)$ that is defined by $V(h) = U_\beta(\varphi(h))$ is a fixed point of T_β .
 378 Hence, $V_\beta^2(\sigma^1, \sigma^2 | h)$ equals $U_\beta(\varphi(h))$ and is a function of $\beta, \sigma^1, \sigma^2$, and $\varphi(h)$.

379 **Lemma 4.4** For every finite history h and φ -factored strategies σ^1 of Player 1 and σ^2
 380 of Player 2, $V_\beta^2(\sigma^1, \sigma^2 | h)$ is a rational function of $\beta \in (0, 1)$.

381 **Proof** Recall that a factor φ in a stochastic game determines the current state. The
 382 state determined by $x \in X$ is denoted by $z(x)$. Therefore, if $h = (z_1, a_1, \dots, z_t)$ and
 383 $\varphi(h) = x$ we have $z(x) = z_t$. Fix a pair of φ -factored strategies $\sigma^1 = \omega^1 \circ \varphi$ and
 384 $\sigma^2 = \omega^2 \circ \varphi$. Define the vector $u \in \mathbb{R}^X$ and the $X \times X$ transition matrix Q by

$$385 \quad u(x) = E_{(\sigma^1, \sigma^2 | h)} u_2(z_t, a_t)$$

386 and

$$387 \quad Q(x, x') = P_{(\sigma^1, \sigma^2 | h)}(\varphi(h, a_t, z_{t+1}) = x'),$$

388 where $h = (z_1, \dots, z_t)$ with $\varphi(h) = x$. As the factor φ is recursive and σ^1 and σ^2 are
 389 φ -factored, u and Q are well defined. Note that the vector $U \in \mathbb{R}^X$ that is given by

$$390 \quad U(x) = \sum_{x' \in X} \left(u(x') \cdot \sum_{j=0}^{\infty} \beta^j Q^j(x, x') \right), \quad x \in X,$$

391 is a fixed point of T . As Q is a transition matrix, $\sum_{j=0}^{\infty} \beta^j Q^j = (I - \beta Q)^{-1}$ and

$$392 \quad U(x) = \sum_{x' \in X} u(x') \cdot (I - \beta Q)^{-1}(x, x').$$

393 Hence,

$$394 \quad V_\beta \left(\sigma^1, \sigma^2 | h \right) = U_\beta(\varphi(h)) = U(\varphi(h)) = \sum_{x' \in X} u(x') (I - \beta Q)^{-1}(\varphi(h), x').$$

395 Each one of the entries of $(I - \beta Q)^{-1}$ is a rational function of $\beta \in (0, 1)$. Hence,
 396 $V_\beta^2(\sigma^1, \sigma^2 | h)$ is a rational function of $\beta \in (0, 1)$.

397 **Lemma 4.5** For every $x \in X$ and a φ -factored strategy σ^1 of Player 1 there exists
 398 $\beta_1(x) \in (0, 1)$ such that if $\varphi(h) = x$ then $\beta \mapsto V_\beta^2(\sigma^1 | h)$ is a rational function of
 399 $\beta \in [\beta_1(x), 1)$.

400 **Proof** Fix a finite history h and a φ -factored strategy σ^1 of Player 1. By Theorem 4.1,
 401 for every $\beta \in (0, 1)$, there is a pure φ -factored strategy τ of Player 2 such that
 402 $V_\beta^2(\sigma^1 | h) = V_\beta^2(\sigma^1, \tau | h)$. There are finitely many pure φ -factored strategies
 403 τ of Player 2 and for each such strategy τ , $V_\beta^2(\sigma^1, \tau | h)$ is a rational function of
 404 $\beta \in (0, 1)$ by Lemma 4.4. Hence, $V_\beta^2(\sigma^1 | h)$ is the maximal value of finitely many

405 rational functions of β . Two distinct rational functions coincide in at most finitely
 406 many points. Therefore, there is $\beta_*(h) \in (0, 1)$ such that $V_\beta^2(\sigma^1 | h)$ is a rational
 407 function of $\beta \in (\beta_*(h), 1]$.

408 The range of φ is finite; therefore, there is $\beta_1(x) \in (0, 1)$ such that, for every $h \in H$,
 409 $\beta \mapsto V_\beta^2(\sigma^1 | h)$ is a rational function of $\beta \in [\beta_1, 1)$.

410 Let $\sigma^1 = \omega^1 \circ \varphi$ be a φ -factored strategy, where $\varphi: H \rightarrow X$ is a recursive factor
 411 with a finite range. For every $\beta \in (0, 1)$, $x \in X$, and $a^2 \in A_2(z(x))$, we define

$$412 \quad M(\beta, x, a^2) = \sum_{a^1 \in A_1(z(x))} \left(\omega^1(x)(a^1) \cdot \left[u_2(z(x), (a^1, a^2)) + \beta \cdot \sum_{z \in S} (p(z(x), (a^1, a^2))(z) \cdot V_\beta^2(\sigma^1 | h, (a^1, a^2), z)) \right] \right)$$

413 where $h \in H$ satisfies $\varphi(h) = x$. Since σ^1 is φ -factored and φ is recursive, the function
 414 M is well defined. Using Lemma 4.5 we see that the function $\beta \mapsto M(\beta, x, a^2)$ is
 415 a rational function of $\beta \in [\beta_1, 1)$, where β_1 is taken from Lemma 4.5 and x, a^2
 416 are fixed. Hence, there is $\beta_0(x) \in (0, 1)$ and $\omega^2(x) \in A_2(z(x))$ such that, for every
 417 $\beta \in [\beta_0(x), 1)$, we have

$$418 \quad M(\beta, x, \omega^2(x)) = \max_{a^2 \in A_2(z(x))} M(\beta, x, a^2).$$

419 By definition we also get $M(\beta, x, \omega^2(x)) = V_\beta^2(\sigma^1 | h)$ for every $x \in X, h \in H$ with
 420 $\varphi(h) = x$, and $\beta \in [\beta_0(x), 1)$. By setting $\beta_0 = \max_{x \in X} \beta_0(x)$, for every $x \in X$ and
 421 $\beta \in [\beta_0, 1)$, we have

$$422 \quad M(\beta, x, \omega^2(x)) = \max_{a^2 \in A_2(z(x))} M(\beta, x, a^2).$$

423 Therefore, for every discount factor $\beta \in [\beta_0, 1)$ and $h \in H$, we have $\omega^2 \circ \varphi$ is a best
 424 reply to σ^1 in the β -discounted games $(\Gamma | h)$.

425 Proof of Theorem 4.3

426 Let $\varphi: H \rightarrow X$ be a recursive factor with a finite range and σ^1 be a φ -factored strategy
 427 of Player 1. Let σ^2 be a φ -factored strategy of Player 2 and $\beta_0 \in (0, 1)$ be a discount
 428 factor such that for every discount factor $\beta \in [\beta_0, 1)$ and history h , σ^2 is a best reply to
 429 σ^1 in the β -discounted game $(\Gamma | h)$; that is, $V_\beta^2(\sigma^1, \sigma^2 | h) \geq V_\beta^2(\sigma^1, \rho | h)$ for every
 430 strategy ρ of Player 2. The existence of such σ^2 and β_0 follows from Theorem 4.2.

431 Set $x_t = \varphi(z_1, a_1, \dots, z_t)$ and $v_\lambda(x) := \lambda V_{1-\lambda}^2(\sigma^1, \sigma^2 | h)$ whenever $\varphi(h) = x$. As
 432 σ^1 and σ^2 are φ -factored and φ is a recursive factor, $v_\lambda(x)$ is well defined. The function
 433 $(0, 1] \ni \lambda \mapsto v_\lambda(x)$ is a bounded rational function. Hence, it is a Lipschitz function of
 434 λ . Let K be a sufficiently large positive constant such that $|v_\lambda(x) - v_{\lambda'}(x)| \leq K|\lambda - \lambda'|$
 435 for every $x \in X$. Note that $v_\infty(x) := \lim_{\lambda \rightarrow 0^+} v_\lambda(x)$ exists.

436 For $1 > \beta \geq \beta_0$, $V_\beta^2(\sigma^1, \sigma^2 | h) = V_\beta^2(\sigma^1 | h)$, and hence, for $h = (z)$ and for
 437 every $0 < \lambda \leq 1 - \beta_0$, we have

$$438 \quad E_{(\sigma^1, \sigma^2 | z)}(\lambda u_2(z_t, a_t) + (1 - \lambda)v_\lambda(x_{t+1}) | \mathcal{H}_t) \geq v_\lambda(x_t), \quad (5)$$

439 where \mathcal{H}_t is the algebra of events defined by histories up to stage t . For any strategy
 440 ρ of Player 2 and every $0 < \lambda \leq 1 - \beta_0$, we have

$$441 \quad E_{(\sigma^1, \rho | z)}(\lambda u_2(z_t, a_t) + (1 - \lambda)v_\lambda(x_{t+1}) | \mathcal{H}_t) \leq v_\lambda(x_t). \quad (6)$$

442 The continuation of the proof uses the proof of Mertens and Neyman (1981). Let
 443 A be the largest absolute value of a payoff appearing in the game matrices. Note
 444 that inequality (5) is identical to inequality (2.1) of Mertens and Neyman (1981), and
 445 inequality (6) is a dual version of inequality (2.1) of Mertens and Neyman (1981).
 446 The integrable function $\psi : (0, 1] \rightarrow \mathbb{R}$ that satisfies inequality (2.2) of Mertens and
 447 Neyman (1981) is simply the constant function K .

448 As in the proof of Mertens and Neyman (1981), for every $\varepsilon > 0$ and $1 > \beta^* \geq \beta_0$,
 449 there are

- 450 • a sequence $(\beta_t)_{t=1}^\infty$, where β_t is a function of the history (z_1, a_1, \dots, z_t) with
- 451 $\beta_t \geq \beta^*$ and
- 452 • a constant s_0 ,

453 such that the following conditions are satisfied.

- 454 (a) We denote $\lambda_t := 1 - \beta_t$. Then we have $\lim_{t \rightarrow \infty} \lambda_t = 0$ $P_{(\sigma^1, \sigma^2 | z)}$ -a.e.
- 455 (b) Let I be the indicator function. Then $E_{(\sigma^1, \sigma^2 | z)}(\sum_{t=1}^\infty I(\beta_t = \beta^*)) \leq \frac{2A}{\varepsilon(1-\beta^*)}$.
- 456 (c) We have

$$457 \quad E_{(\sigma^1, \sigma^2 | z)}(v_{\lambda_t}(x_{t+1})) \geq v_\infty(\varphi(z)) - 3\varepsilon.$$

- 458 (d) There exists a function Y_∞ such that $v_{\lambda_t}(x_{t+1}) \rightarrow Y_\infty$ $P_{(\sigma^1, \sigma^2 | z)}$ -a.e., and
- 459 $E_{(\sigma^1, \sigma^2 | z)}(Y_\infty | z) \geq v_\infty(\varphi(z)) - 2\varepsilon$.
- 460 (e) For every n we have

$$461 \quad \sum_{t < n} u_2(z_t, a_t) \geq \sum_{t < n} v_{\lambda_t}(x_{t+1}) - s_0 - 2A \sum_{t < n} I(\beta_t = \beta^*) - 4n\varepsilon.$$

462 Therefore, we have

$$463 \quad E_{(\sigma^1, \sigma^2 | z)}\left(\liminf_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n u_2(z_t, a_t)\right) \geq E_{(\sigma^1, \sigma^2 | z)}\left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n v_{\lambda_t}(x_{t+1})\right) - 4\varepsilon$$

$$= E_{(\sigma^1, \sigma^2 | z)}(Y_\infty | z) - 4\varepsilon \geq v_\infty(\varphi(z)) - 6\varepsilon.$$

464 Further, for n sufficiently large such that $\frac{s_0}{n} + \frac{2A}{n} \frac{2A}{\varepsilon(1-\beta^*)} < \varepsilon$, we have

$$465 \quad E_{(\sigma^1, \sigma^2|z)} \left(\frac{1}{n} \sum_{t=1}^n u_2(z_t, a_t) \right) \geq E_{(\sigma^1, \sigma^2|z)} \left(\frac{1}{n} \sum_{t=1}^n v_{\lambda_t}(z_{t+1}) \right) - 4\varepsilon \\ \geq v_\infty(\varphi(z)) - 5\varepsilon.$$

466 Similarly, if ρ is a strategy of player 2, then

$$467 \quad E_{(\sigma^1, \rho|z)} \left(\limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n u_2(z_t, a_t) \right) \leq E_{(\sigma^1, \sigma^2|z)} \left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n v_{\lambda_t}(x_{t+1}) \right) + 4\varepsilon \\ = E_{(\sigma^1, \sigma^2|z)}(Y_\infty | z) + 4\varepsilon \leq v_\infty(\varphi(z)) + 6\varepsilon,$$

468 and for n sufficiently large such that $\frac{s_0}{n} + \frac{2A}{n} \frac{2A}{\varepsilon(1-\beta^*)} < \varepsilon$,

$$469 \quad E_{(\sigma^1, \rho|z)} \left(\frac{1}{n} \sum_{t=1}^n u_2(z_t, a_t) \right) \leq E_{(\sigma^1, \sigma^2|z)} \left(\frac{1}{n} \sum_{t=1}^n v_{\lambda_t}(x_{t+1}) \right) + 4\varepsilon \\ \leq v_\infty(\varphi(z)) + 5\varepsilon.$$

470 The last four displayed inequalities prove (1) and (2) for $h = z$. The proof of inequalities (1) and (2) for a general finite history $h = (z_1, a_1, \dots, z_{t'})$ is the same: just replace, in the above inequalities, $P_{(\sigma^1, \sigma^2|z)}$ by $P_{(\sigma^1, \sigma^2|h)}$, $P_{(\sigma^1, \rho|z)}$ by $P_{(\sigma^1, \rho|h)}$, n by $t' + n$, and $t \in \mathbb{N}$ by $t \geq t'$.

474 **Remark 4.6** An alternative proof of Theorem 4.3 that does not rely on the Mertens–Neyman proof (Mertens and Neyman (1981)) follows the proof of parts 3) and 4) of [Neyman (2003) Proposition 3].

477 5 Corollaries of the main result

478 The next corollary follows from the facts that a supergame and a supergame with time-dependent stage games are special cases of a stochastic game and that a behavioral k -SBR, a behavioral time-dependent $k(t)$ -BR strategy, a behavioral k -state automaton strategy, and a behavioral k -state time-dependent automaton strategy are factored strategies with the corresponding natural recursive factors. Note that the results for supergames and strategies with a finite history can be found in Press and Dyson (2012).

484 **Corollary 5.1** Fix a discount factor $\beta \in (0, 1)$. Let an infinite-stage game be either a supergame, or a time-dependent supergame, or a stochastic game with countably many states. If the stage games have finitely many actions, then the β -discounted infinite-stage game obeys the following best-reply properties.

488 (a) Let k be a nonnegative integer. Then, for every behavioral k -SBR strategy σ^1 , there is a pure k -SBR strategy σ^2 that is a best reply of Player 2.

- 490 (b) Let $k : \mathbb{N} \rightarrow \mathbb{N}$ be a function with $k(t + 1) \leq k(t) + 1$ for every $t \in \mathbb{N}$. Then, for
 491 every behavioral time-dependent k -BR strategy σ^1 , there is a pure k -BR strategy
 492 σ^2 that is a best reply.
- 493 (c) Let k be a positive integer. Then, for every strategy σ^1 that is defined by a k -state
 494 behavioral automaton, there is a best reply σ^2 of Player 2 that is defined by a
 495 deterministic k -state automaton.
- 496 (d) Let k be a positive integer. Then, for every strategy σ^1 that is defined by a k -state
 497 time-dependent automaton, there is a best reply σ^2 that is defined by a deterministic
 498 k -state time-dependent automaton.

499 The next two theorems are well-known results in the theory of Markov decision
 500 processes. The first is analogous to Theorem 4.1 and the second is analogous to The-
 501 orem 4.2. As an MDP is a one-person stochastic game, these results follow from our
 502 main result. An MDP is defined by a tuple $\langle S, A, u, p, \mu \rangle$, where S is the state space,
 503 A is the set of actions of the decision maker (DM), u is the payoff to the DM, p is the
 504 transition probability, and μ is the probability of the initial state. If σ is a strategy of
 505 DM, then $v_\beta(\sigma)$ denotes the β -discounted payoff of DM as a function of the strategy
 506 σ .

507 **Theorem 5.2** (Derman 1965) Let $\langle S, A, u, p, \mu \rangle$ be an MDP with countably many
 508 states, finitely many actions in each state, and a bounded reward function. Then for
 509 each $\beta \in (0, 1)$ there is a stationary pure strategy σ such that, for every strategy ρ ,
 510 we have $v_\beta(\sigma) \geq v_\beta(\rho)$.

511 **Theorem 5.3** (Blackwell 1962) Let $\langle S, A, u, p, \mu \rangle$ be an MDP with finitely many
 512 states and actions. Then there is a stationary pure strategy σ and $\beta_0 \in (0, 1)$ such
 513 that, for every strategy ρ and for every $\beta \in [\beta_0, 1)$, we have $v_\beta(\sigma) \geq v_\beta(\rho)$.

514 **Remark 5.4** Theorem 4.1 and Theorem 4.2 are analogous to Theorem 5.2 and Theorem
 515 5.3 respectively. A natural approach is to try to deduce Theorems 4.1 and 4.2 from
 516 these theorems by assigning to a φ -factored strategy σ^1 in a two-person stochastic
 517 game with countably many states and finitely many actions the following auxiliary
 518 MDP. The state space of the auxiliary MDP is X , where X is the image of the recursive
 519 factor φ . The action sets of DM are those of Player 2 in the stochastic game. The stage
 520 payoffs are defined by

$$521 \quad u(x, a^2) = \sum_{a^1 \in A_1} \omega^1(x)(a^1) \cdot u_2(a^1, a^2), \quad x \in X \text{ and } a^2 \in A_2,$$

522 where $\omega^1 \circ \varphi$ is the φ -factored strategy of Player 1, and the transitions are defined by

$$523 \quad p(x, a^2)(x') = \sum_{a^1 \in A_1} \omega^1(x)(a^1) \cdot p(x, a^1, a^2)(x'), \quad x, x' \in X \text{ and } a^2 \in A_2.$$

524 By Theorem 5.2 the decision maker has a pure stationary optimal strategy σ in the
 525 auxiliary β -discounted MDP, and a stationary strategy in the auxiliary MDP maps
 526 naturally to a φ -factored strategy in the stochastic game.

Hence, it is tempting to argue that this proves that this φ -factored strategy in the stochastic game is a best reply of Player 2 to the φ -factored strategy σ^1 of Player 1.

However, in the stochastic game, a strategy of Player 2 may depend also on the past actions of Player 1. In the mapping of the stochastic game to the auxiliary MDP the history of past actions of Player 1 disappears. Hence, the φ -factored strategy σ is a best reply among all strategies that ignore past actions of Player 1, but, at least formally, one cannot argue that it is a best reply among all strategies of Player 2.

6 Remarks

6.1 The role of recursiveness

Recursiveness of φ is not a necessary condition in the main results, even in the model of the supergame. For example, consider a supergame, where each player has at least two actions, and let φ be a nonrecursive factor that is independent of the actions of Player 2. Then, for every action function $\omega^1: X \rightarrow \Delta(A_1)$ of Player 1, if $\omega^2: X \rightarrow \Delta(A_2)$ is the best reply to ω^1 , i.e., $u_2(\omega^1(x), \omega^2(x)) \geq u_2(\omega^1(x), a^2)$ for every $x \in X, a^2 \in A_2$, then the φ -factored strategy $\omega^2 \circ \varphi$ of Player 2 is a best reply to the φ -factored strategy $\omega^1 \circ \varphi$ of Player 1 in any β -discounted game $(G | h)$.

However, if the factor φ is nonrecursive and obeys an additional condition that is spelled out below, then the conclusion of Theorem 4.1 does not hold even in the model of the discounted supergame. We illustrate this remark in the model of the supergame where each player has at least two actions. Without loss of generality we can assume that $A_i = \{0, 1\}$ and that the factor $\varphi: H \rightarrow X$ obeys $\varphi(h) = \varphi(h')$, $\varphi(h, (1, 1)) \neq \varphi(h', (1, 1))$ for some histories h, h' . The added condition is

$$\{\varphi(h), \varphi(h, (1, 1)), \varphi(h', (1, 0))\} \cap \{\varphi(h', (1, 1)), \varphi(h, (1, 0))\} = \emptyset. \quad (7)$$

Then we can find

- a payoff function $u_2: A_1 \times A_2 \rightarrow [0, 1]$ and
- a pure action function $\omega^1: X \rightarrow A_1$,

such that any strategy τ^2 of Player 2 that is a best reply to $\omega \circ \varphi$ in both of the β -discounted games, $(G | h)$ and $(G | h')$, obeys $\tau^2(h) \neq \tau^2(h')$; thus τ^2 is not φ -factored.

Using (7) we can define a φ -factored strategy $\omega^1 \circ \varphi$ of Player 1 by setting

$$\begin{aligned} \omega^1(\varphi(h)) = \omega^1(\varphi(h')) = \omega^1(\varphi(h, (1, 1))) = \omega^1(\varphi(h', (1, 0))) = 1 \quad \text{and} \\ \omega^1(\varphi(h', (1, 1))) = \omega^1(\varphi(h, (1, 0))) = 0. \end{aligned}$$

Let the payoff function u_2 of Player 2 be given by

$$u_2(1, \cdot) = 1 \text{ and } u_2(0, \cdot) = 0.$$

561 We will deal with the supgame whose stage game G is $\langle A_1, A_2, u_1, u_2 \rangle$. Fix a
 562 sufficiently small discount factor $\beta \in (0, 1)$ such that $\beta > \beta^2/(1 - \beta)$. Consider the
 563 supgame $(G | h)$, i.e., the supgame that follows history h . If Player 2 plays action 1
 564 in the first stage of $(G | h)$, then she guarantees that her (unnormalized) β -discounted
 565 payoff is least $1 + \beta$. If Player 2 plays action 0 in the first stage of the supgame
 566 $(G | h)$, then her (unnormalized) β -discounted payoff is at most $1 + \beta^2/(1 - \beta)$.
 567 Therefore, any strategy of Player 2 in which action 1 is played in the first stage of
 568 $(G | h)$ dominates any strategy of Player 2 in which action 0 is played in the first stage
 569 of $(G | h)$. Therefore, if σ^2 is a best reply of Player 2 to $\omega^1 \circ \varphi$ in the β -discounted
 570 game $(G | h)$, then $\sigma^2(h) = 1$.

571 Similarly, consider the supgame $(G | h')$. If Player 2 plays action 0 in the first
 572 stage of $(G | h')$, then she guarantees that her (unnormalized) β -discounted payoff
 573 is at least $1 + \beta$. If Player 2 plays action 1 in the first stage of $(G | h')$, then her
 574 (unnormalized) β -discounted payoff is at most $1 + \beta^2/(1 - \beta)$. Therefore any strategy
 575 of Player 2 in which action 0 is played in the first stage of $(G | h')$ dominates any
 576 strategy of Player 2 in which action 1 is played in the first stage of $(G | h')$. Therefore,
 577 if σ^2 is a best reply of Player 2 to $\omega^1 \circ \varphi$ in the β -discounted game $(G | h')$, then
 578 $\sigma^2(h') = 0$.

579 Hence, if σ^2 is a best reply of Player 2 to $\omega^1 \circ \varphi$ in each one of the β -discounted
 580 games $(G | h)$ and $(G | h')$, then $\sigma^2(h) \neq \sigma^2(h')$ and, therefore, as $\varphi(h) = \varphi(h')$,
 581 σ^2 is not φ -factored.

582 6.2 Compact action spaces

583 A natural extension of our model is to consider players with compact action sets A_i . In
 584 this extension, there arises a new problem not found in games with finite action profiles,
 585 namely, the existence of a best reply to a given strategy σ . Consider, for example, the
 586 following two-player supgame, where the set of actions of each player is the interval
 587 $[0, 1]$ and the stage payoff of Player 2 is (at any time) given by $u_2(a^1, a^2) = a^1 + a^2$.
 588 Now, suppose that Player 1 plays the 1-SBR strategy given by

$$589 \quad \sigma^1(a_1, \dots, a_{t-1}) = \begin{cases} 1, & \text{if } a_{t-1}^2 < 1 \text{ and } t > 1, \\ 0, & \text{otherwise,} \end{cases} \quad e_1 = (0, 0).$$

590 The strategy σ^1 is a factor-based strategy with a recursive factor. Indeed, we set
 591 $X = \{B, C\}$, $\omega(B) = 1$ and $\omega(C) = 0$, and $\varphi(a_1, \dots, a_{t-1}) = B$ if $a_{t-1}^2 < 1$ and
 592 $t > 1$, $\varphi(a_1, \dots, a_{t-1}) = C$. Then we have $\sigma^1 = \omega \circ \varphi$. However, in the β -discounted
 593 game there is no φ -factored best reply, and any φ -factored reply is dominated by
 594 (another) φ -factored reply.

595 Of course, there does not exist any general best reply to σ^1 . The difficulty stems
 596 from the fact that the factor φ is not continuous. However, in analogy to, and using
 597 the tools of, say, Maitra (1968), one can generalize Theorem 4.1 to the setup with a
 598 continuous factor.

599 **6.3 Public vs. private strategies**

600 A natural recursive factor φ arises in the theory of repeated games with imperfect
 601 monitoring. This natural factor φ maps any finite history of the game to the finite
 602 history x of the public signals, and we can identify the φ -factored strategies with
 603 the so-called *public strategies* [see, e.g., (Radner et al. 1986)]. By contrast, a *private*
 604 *strategy* [see, e.g., (Kandori and Obara 2006)] is a strategy where the current action
 605 depends on the history of public signals (i.e., on elements of X), as well as on the
 606 history of private signals (e.g., past own actions).


607 Our question at the outset of this paper can then be reformulated to “If my opponent
 608 is using a public strategy, when can I benefit by using my (additional) private signals?”;
 609 in other words, “Can a private strategy fare better than any public strategy against a
 610 public strategy?” The answer is that in the discounted repeated game with imperfect
 611 monitoring (finitely many stage actions and countably many public signals) one does
 612 not profit from the additional private signals since the factor φ is obviously recursive.

613 **References**

- 614 Abreu D, Rubinstein A (1988) The structure of Nash equilibrium in repeated games with finite automata.
 615 *Econometrica* 56(6):1259–1281
- 616 Aumann RJ (1981) Survey of repeated games. In: Essays in game theory and mathematical economics in
 617 Honour of Oskar Morgenstern, pp 11–42. Wissenschaftsverlag, Bibliographisches Institut
- 618 Aumann RJ, Sorin S (1989) Cooperation and bounded recall. *Games Econ Behavior* 1(1):5–39
- 619 Ben-Porath E (1993) Repeated games with finite automata. *J Econ Theory* 59:17–32
- 620 Blackwell D (1962) Discrete dynamic programming. *Ann Math Stat* 33:719–726
- 621 Derman C (1965) Markovian sequential control processes: denumerable state spaces. *J Math Anal Appl*
 622 10:295–302
- 623 Kalai E (1990) Bounded rationality and strategic complexity in repeated games. In: Ichiishi T, Neyman A,
 624 Tauman Y (eds) *Game theory and applications*. Academic Press, San Diego, pp 131–157
- 625 Kalai E, Stanford W (1988) Finite rationality and interpersonal complexity in repeated games. *Econometrica*
 626 56(2):397–410
- 627 Kandori M, Obara I (2006) Efficiency in repeated games revisited: the role of private strategies. *Econometrica*
 628 74(2):499–519
- 629 Lehrer E (1988) Repeated games with stationary bounded recall strategies. *J Econ Theory* 46(1):130–144
- 630 Maitra A (1968) Discounted dynamic programming on compact metric spaces. *Sankhyā Indian J Stat Ser*
 631 *A* 30(2):211–216
- 632 Maskin E, Tirole J (2001) Markov perfect equilibrium. I. Observable actions. *J Econ Theory* 100(2):191–219
- 633 Mertens J-F, Neyman A (1981) Stochastic games. *Internat J Game Theory* 10(2):53–66
- 634 Neyman A (1997) Cooperation, repetition, and automata. In: Hart S, Mas-Colell A (eds) *Cooperation: game*
 635 *theoretic approaches*, volume 155 of NATO ASI Series F, pp 233–255
- 636 Neyman A (1985) Bounded complexity justifies cooperation in the finitely repeated prisoners’ dilemma.
 637 *Econ Lett* 19:227–229
- 638 Neyman A (2003) From Markov chains to stochastic games. In: Neyman A, Sorin S (eds) *Stochastic games*
 639 *and applications*. Springer, Dordrecht, pp 9–25
- 640 Neyman A, Okada D (2009) Growth of strategy sets, entropy, and nonstationary bounded recall. *Games*
 641 *Econ Behavior* 66(1):404–425
- 642 Press WH, Dyson FJ (2012) Iterated prisoner’s dilemma contains strategies that dominate any evolutionary
 643 opponent. *PNAS* 109(26):10409–10413
- 644 Radner R, Myerson R, Maskin E (1986) An example of a repeated partnership game with discounting and
 645 with uniformly inefficient equilibria. *Rev Econ Stud* 53(172):59–69
- 646 Rubinstein A (1986) Finite automata play the repeated prisoner’s dilemma. *J Econ Theory* 39(1):83–96

647 **Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps
648 and institutional affiliations.

649 Affiliations

650 **René Levínský¹  · Abraham Neyman² · Miroslav Zelený³**

651 ✉ René Levínský
652 rene.levinsky@cerge-ei.cz

653 Abraham Neyman
654 aneyman@math.huji.ac.il

655 Miroslav Zelený
656 zeleny@karlin.mff.cuni.cz

657 ¹ Economics Institute of the Czech Academy of Sciences, Politických vězňů 7, 111 21 Praha 1,
658 Czech Republic

659 ² Institute of Mathematics and Center for the Study of Rationality, The Hebrew University of
660 Jerusalem, Givat Ram, 91904 Jerusalem, Israel

661 ³ Charles University, Faculty of Mathematics and Physics, Sokolovská 83, 186 75 Praha 8, Czech
662 Republic