



# Learning and self-confirming long-run biases

P. Battigalli <sup>a,\*</sup>, A. Francetich <sup>b</sup>, G. Lanzani <sup>c</sup>, M. Marinacci <sup>a</sup>

<sup>a</sup> Department of Decision Sciences and IGIER, Università Bocconi, Via Röntgen, 1, 20136 Milano, Italy

<sup>b</sup> School of Business, University of Washington Bothell, 18115 Campus Way NE, Bothell, WA 98011, USA

<sup>c</sup> Department of Economics, Massachusetts Institute of Technology, 50 Memorial Drive, Cambridge, MA 02139, USA

Received 29 March 2017; final version received 13 July 2019; accepted 23 July 2019

Available online 30 July 2019

---

## Abstract

We consider an ambiguity averse, sophisticated decision maker facing a recurrent decision problem where information is generated endogenously. In this context, we study self-confirming actions as the outcome of a process of active experimentation. We provide inter alia a learning foundation for self-confirming equilibrium with model uncertainty (Battigalli et al., 2015), and we analyze the impact of changes in ambiguity attitudes on convergence to self-confirming equilibria. We identify conditions under which the set of self-confirming equilibrium actions is invariant to changes in ambiguity attitudes, and yet ambiguity aversion may affect the dynamics. Indeed, we argue that ambiguity aversion tends to stifle experimentation, increasing the likelihood that the decision maker gets stuck into suboptimal “certainty traps.”

© 2019 Elsevier Inc. All rights reserved.

*JEL classification:* C72; D81; D83

*Keywords:* Learning; Stochastic control; Ambiguity aversion; Self-confirming equilibrium

---

## 1. Introduction

We study the dynamic behavior of a decision maker (DM) who faces a recurrent decision problem in which the actions he selects depend on the information endogenously gathered through

---

\* Corresponding author.

E-mail addresses: [pierpaolo.battigalli@unibocconi.it](mailto:pierpaolo.battigalli@unibocconi.it) (P. Battigalli), [aletich@uw.edu](mailto:aletich@uw.edu) (A. Francetich), [lanzani@mit.edu](mailto:lanzani@mit.edu) (G. Lanzani), [massimo.marinacci@unibocconi.it](mailto:massimo.marinacci@unibocconi.it) (M. Marinacci).

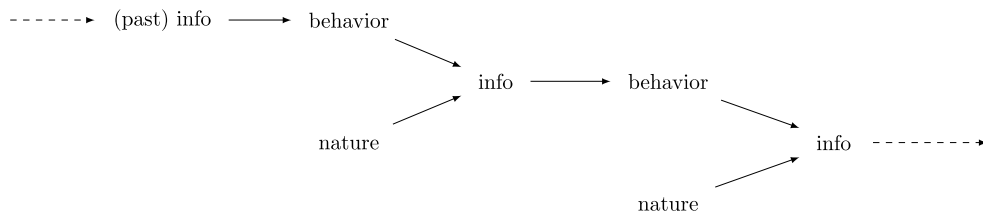


Fig. 1. Timeline.

his past behavior as, for example, in multiarmed bandit problems (cf. Gittins, 1989). We diagram the flow of actions and information in Fig. 1.

Our DM is ambiguity averse, finitely patient, and uncertain about the stochastic process of states of nature. In this setting, there are three crucial elements of our analysis. First, the process of states is governed by an unknown objective probability model (e.g., the composition of an urn). Second, the uncertainty of the DM about the objective model is represented through a subjective probability measure, a belief, which is updated according to information feedback. Each period, the DM evaluates the possible actions (given his updated belief) according to a dynamic version of the smooth ambiguity criterion of Klibanoff et al. (2005), which separates ambiguity attitudes (a personal trait) from the evolving perception of ambiguity, and allows for a Bayesian analysis of learning. Third, the DM uses a rational strategy given his prior belief.

It is essential to understand the meaning of the term “rational” in our setting. An uncertainty averse DM may have *dynamically-inconsistent preferences* (cf. Example 5). While we allow for such reversal of preferences, we assume that the DM is sophisticated in the sense that he formulates a *dynamically-consistent strategy*, that is, a strategy that satisfies the one-deviation (or intrapersonal equilibrium) property: There is no instance where the DM has an incentive to choose an action different from the one prescribed by the given strategy. In a finite-horizon model, this is equivalent to folding-back planning. But we cannot rely on folding back, because we focus on infinite-horizon models to study the limit properties of behavior and beliefs, and to exploit the ensuing stationarity of the dynamic decision problem.

We study how steady-state actions arise from an active experimentation process, providing a novel convergence result. Specifically, we show that the stochastic process of beliefs and actions converges with probability 1 to a random limit action-belief pair. This random limit pair satisfies almost surely the following *self-confirming equilibrium* conditions: The limit action maximizes the one-period value given the limit belief, and the limit belief assigns probability 1 to the set of probability models that are observationally equivalent to the true one given the limit action.

Therefore, even if the DM cares about the future, the limit action-belief pair must be a self-confirming equilibrium of the one-period game repeatedly played against nature. Since the belief may only partially identify the true model (nature’s “behavior strategy”), such limit behavior may be very different from the “Nash” (or “rational expectations”) equilibrium, in which the DM plays the objective best reply.<sup>1</sup>

Since we assume that the state process is exogenous, that is, the DM’s actions cannot influence the probabilities of states in future periods, our framework cannot model long-run interactions with a fixed set of co-players. However, our exogeneity assumption is justified within the scenario

<sup>1</sup> Our definition of self-confirming equilibrium (also called “conjectural equilibrium”) is broader than the one of Fudenberg and Levine (1993, 1998). See the discussion in Battigalli et al. (2015).

of large population games. Indeed, our setup can represent the point of view of a DM who plays a game recurrently with other agents independently drawn from large, statistically stable, populations. Hence, the DM recognizes to be unable to influence the evolution of the environment with his actions. The probability models describe the distribution of behaviors in the co-players' populations. Experimentation is valuable to the DM since a better understanding of the correct distribution may allow him to select a better strategy in the following periods. In particular, our setting is consistent with a steady-state learning environment *à la* Fudenberg and He (2018), where individual agents learn through their life, but the population's statistics are constant.

Under this interpretation, we provide a learning foundation for self-confirming equilibrium with model uncertainty (Battigalli et al., 2015, henceforth BCMM). Specifically, the random limit pair corresponds to the “smooth” self-confirming equilibrium concept of BCMM since the limit action is a myopic best response, and the evidence generated by the limit action and the steady-state distribution of opponents' actions confirms the limit belief. BCMM prove that higher ambiguity aversion yields a larger set of self-confirming equilibrium actions. Intuitively, the reason is that a self-confirming equilibrium action is “tested,” hence it yields known risks (objective probabilities of consequences), whereas deviations yield unknown risks that are the less attractive the higher the aversion to ambiguity. Since we show that self-confirming equilibrium emerges as the long-run outcome of an active experimentation and learning process, the comparative statics result of BCMM implies that higher ambiguity aversion reduces the predictability of long-run behavior.

We provide special conditions under which the BCMM theorem holds as an invariance result: The set of self-confirming equilibrium actions does not depend on ambiguity attitudes. Nonetheless, ambiguity aversion may still affect the dynamics. Specifically, we argue that ambiguity aversion tends to stifle experimentation, increasing the likelihood that the DM gets stuck into suboptimal “certainty traps.” The intuition is as follows. Suppose that the DM can only learn from observing his realized payoffs. The actions perceived as ambiguous, that is, those with uncertain distributions of payoffs, are those from which the DM expects to learn. If instead an action is perceived as unambiguous, the DM expects to have the same belief before and after choosing it, i.e., he does not expect to learn from it. Hence, ambiguity aversion biases the DM toward “exploitation” and against “exploration.”

**Related literature** We point out that there is a formal connection between the concept of SCE and the literature on active learning (or “stochastic control”), and in particular the seminal work by Easley and Kiefer (1988, henceforth EK). The working paper version<sup>2</sup> provides a translation between our setup and the active learning setup. Our paper departs from EK in two fundamental aspects. First, we allow for non-neutral ambiguity attitudes and dynamically inconsistent preferences. Second, EK requires the DM to assign positive subjective probability to (every neighborhood of) the correct model, whereas our sufficient condition for convergence to an SCE allows for misspecified beliefs.<sup>3</sup>

Our definition of self-confirming equilibrium is related to the notion of subjective equilibrium of Kalai and Lehrer (1993 and 1995, henceforth KL). Relatively minor details aside, there are two key differences between the two concepts. First, KL define and analyze subjective equilibrium as the rest point of a process of updated beliefs about the path of play in a supergame. Such

<sup>2</sup> Available as IGIER w.p. 588.

<sup>3</sup> See Section 3. Moreover, the working paper version considers the more general case of non i.i.d. state generating process.

beliefs can be interpreted either as “subjective averages” of probability models, or as subjective Dirac beliefs over probability models. Focusing on such beliefs is without loss of generality under subjective expected utility maximization, but not under non-neutral ambiguity attitudes (see Sections 2.1 and 4). Second, since in KL the set of interacting players is fixed once and for all, their analysis concerns the convergence to a steady state of beliefs about supergame behavior. Our analysis, instead, is consistent with steady-state learning in a population game scenario; thus, we obtain convergence to an equilibrium of the one-period game.

Our results on ambiguity aversion and experimentation are consistent with the findings in Li (2019) and Anderson (2012). Li (2019) characterizes the optimal experimentation strategy under ambiguity aversion in an independent K-armed bandit problem. Aside from focusing on this specific case, the key difference with our paper is that Li (2019) models ambiguity aversion following the two-stage multiple-prior model of Marinacci (2002), while we employ the smooth ambiguity criterion of Klibanoff et al. (2005). As a result, the comparative-statics analysis in Li (2019) considers the impact of changes in the perception of ambiguity, while ours studies the effect of changes in ambiguity attitudes. Moreover, Li (2019) uses a recursive version of the maxmin expected-utility criterion and is thus able to employ standard dynamic programming techniques. Such a recursive representation is precluded in our setting. Unlike Li (2019) and our paper, Anderson (2012) derives the predictions of his model under the implicit assumption that the decision maker can commit ex-ante to any strategy. However, the empirical evidence he presents is consistent with the theoretical predictions of our model.

**Outline** The paper is structured as follows. Section 2 presents the static and dynamic decision framework, as well as preliminary concepts. Section 3 describes the endogenous information process. Section 4 describes the DM’s intertemporal preferences. Section 5 analyzes self-confirming equilibrium, rational strategies, and presents our results on convergence to SCE. Section 6 presents our comparative dynamics results with respect to changes in ambiguity attitudes. Finally, Section 7 briefly relates our analysis to the literature on learning in games and concludes. Proofs are relegated to the appendix. We refer to the working paper version for the complete derivation of the results presented in the examples.

## 2. Framework

### 2.1. Static environment

Let  $S$  be a finite space of states of nature and let  $M$  be a finite outcome space. We consider a control setup where a finite set  $A$  of actions (or controls) is available to the DM, and actions and states translate into outcomes through a feedback function  $f : A \times S \rightarrow M$ .<sup>4</sup> We assume that outcomes are observable, while states are not directly observable. Thus, unless the feedback function is injective given the chosen action, inference on the states is partial. The quadruple  $(A, S, M, f)$  is the basic structure of the decision problem.

Given a probability measure  $\theta$  on  $S$ , an action  $a$  induces a pushforward measure over outcomes via the function  $F : A \times \Delta(S) \rightarrow \Delta(M)$  defined by:

$$\forall m \in M, F(a, \theta)(m) = \sum_{s \in f_a^{-1}(m)} \theta(s),$$

<sup>4</sup> We endow all finite sets with the discrete topology.

where  $f_a := f(a, \cdot)$  is the section of  $f$  at  $a$ . Our first maintained assumption is that the DM is an Expected Utility Maximizer with respect to these (objective) lotteries.

**Assumption 1** (*Expected utility on lotteries*). There exists a utility function  $u : A \times M \rightarrow \mathbb{R}$  such that, for every objective probability measure  $\theta$  on  $S$  and pair of actions  $a', a'' \in A$ , the DM prefers the (objective) distribution over outcomes induced by  $a'$  to the one induced by  $a''$  if and only if:

$$\sum_{m \in M} u(a', m) F(a', \theta)(m) \geq \sum_{m \in M} u(a'', m) F(a'', \theta)(m).$$

Under Assumption 1, we define the expected payoff over outcomes as:

$$R(a, \theta) := \mathbb{E}_{F(a, \theta)}[u_a],$$

where  $u_a := u(a, \cdot)$  is the section of  $u$  at  $a$ .

Let  $\Theta \subseteq \Delta(S)$  be a compact set of probability measures on  $S$ . These measures, which we call *models*, represent the structural (often physical) information available to the DM, with  $\bar{\theta} \in \Theta$  denoting the objectively true model.<sup>5</sup> We identify  $\Theta$  with a subset of the simplex of dimension  $|S| - 1$  and endow it with the Borel  $\sigma$ -algebra  $\mathcal{B}(\Theta)$ . Incompleteness of information is captured by the non-singleton nature of  $\Theta$ . Under model uncertainty (cf. Marinacci, 2015), the DM ranks actions according to the *smooth ambiguity* criterion of Klibanoff et al. (2005):

$$\bar{V}(a, \mu) := \phi^{-1} \left( \int_{\Theta} \phi(R(a, \theta)) \mu(d\theta) \right), \tag{1}$$

where  $\mu$  is a prior probability measure on  $(\Theta, \mathcal{B}(\Theta))$ , and  $\phi : [\min_{a, \theta} R(a, \theta), \max_{a, \theta} R(a, \theta)] \rightarrow \mathbb{R}$  is a strictly increasing and continuous function that describes attitudes towards ambiguity.<sup>6</sup> In particular, a concave  $\phi$  captures ambiguity aversion, while a linear  $\phi$  (e.g., the identity function) corresponds to the classical subjective expected utility criterion (Cerrei-a-Vioglio et al., 2013b)<sup>7</sup>:

$$\bar{V}(a, \mu) = \int_{\Theta} R(a, \theta) \mu(d\theta) = R(a, \theta_{\mu}),$$

where  $\theta_{\mu} \in \Theta$  is the predictive probability given by  $\theta_{\mu}(E) := \int_{\Theta} \theta(E) \mu(d\theta)$  for all  $E \subseteq S$ . Finally, note that:

- (i) When the support of  $\mu$ ,  $\text{supp } \mu$ , is a singleton  $\{\theta\}$ , criterion (1) reduces to the expected payoff criterion  $R(a, \theta)$ ;

<sup>5</sup> For example, if  $S = \{b, g\}$  is the set of possible colors of the ball drawn from an urn of 90 balls that are either black or green, then  $\Theta = \{\theta \in \Delta(S) : \theta(b) = i/90 = 1 - \theta(g), i \in \{0, \dots, 90\}\}$ .

<sup>6</sup> See Theorem 6 in Cerrei-a-Vioglio et al. (2013a) for an axiomatization. In our setting, the domain of  $\phi$  is well defined by finiteness of  $A$  and continuity of  $R(a, \theta)$  with respect to  $\theta$ .

<sup>7</sup> To map our decision criterion into theirs, let their space of consequences be  $A \times M$  and identify each action  $a$  with the act  $g(s) = (a, f(a, s))$ .

- (ii) The limit case of criterion (1) as ambiguity aversion increases is a version of the maximin criterion  $\min_{\theta \in \text{supp } \mu} R(a, \theta)$  of Gilboa and Schmeidler (1989); see Proposition 3 in Klibanoff et al. (2005).

The static decision problem can be summarized by:

$$\Gamma = (A, S, M, \Theta, f, u, \phi, \mu). \tag{2}$$

### 2.2. Dynamic environment

**Notation** For every finite set  $Z$ , we let  $Z^t = \prod_{\tau=1}^t Z$  and  $Z^\infty = \prod_{\tau=1}^\infty Z$ .<sup>8</sup> We endow the space  $Z^\infty$  with the Borel  $\sigma$ -algebra,  $\mathcal{B}(Z^\infty)$ , corresponding to the product topology on  $Z^\infty$ ; this is the same as the  $\sigma$ -algebra generated by the elementary cylinders  $\{z_1\} \times \dots \times \{z_t\} \times Z^\infty$  (see, e.g., Proposition 1.3 in Folland, 2013). We denote by  $z^t = (z_1, \dots, z_t) \in Z^t$  both the histories and the elementary cylinders that they identify through the following map:

$$(z_1, \dots, z_t) \mapsto \{z_1\} \times \dots \times \{z_t\} \times Z^\infty.$$

We denote by  $z^\infty = (z_1, \dots, z_t, \dots)$  a generic element of  $Z^\infty$ .

**Environment** Given  $S$ , let  $(S^\infty, \mathcal{B}(S^\infty))$  be the measurable space on which a coordinate state process  $(s_1, s_2, \dots)$  is defined, with  $s_t : S^\infty \rightarrow S$  for each  $t$ .<sup>9</sup> We will use the less demanding notation  $\mathbf{s}^\infty$  for the state process describing the exogenous uncertainty in the decision problem. Its realizations are denoted by  $s^\infty \in S^\infty$ . Similarly, we write  $\mathbf{s}^t = (s_1, \dots, s_t)$  with realization  $s^t$  for finite state processes.

For a generic stochastic process  $(z_1, z_2, \dots)$  defined on  $(S^\infty, \mathcal{B}(S^\infty))$ , we denote by  $\sigma(\mathbf{z}^t)$  the  $\sigma$ -algebra generated by the random variables  $\mathbf{z}_1, \dots, \mathbf{z}_t$ , namely, by the process up to time  $t$ ;  $\sigma(\mathbf{z}^0)$  denotes the trivial  $\sigma$ -algebra.

Finally, for every  $\theta \in \Theta$  we define  $p_\theta \in \Delta(S^\infty)$  as the unique i.i.d. extension on  $\mathcal{B}(S^\infty)$  of the measure given on all elementary cylinders by:

$$p_\theta(s^t \times S^\infty) = \prod_{\tau=1}^t \theta(s_\tau)$$

for every  $t$  and every  $s^t$  in  $S^t$ .<sup>10</sup>

**Actions and outcomes** We describe the DM’s choices as a sequence  $(a_t) \in A^\infty$  that consists of an action  $a_t$  for each period  $t$ . At each such  $t$ , there is a time-independent feedback function  $f : A \times S \rightarrow M$ , where  $f(a_t, s_t)$  is the outcome that the DM observes ex-post (i.e., after the decision) at the end of period  $t$  if he chooses action  $a_t$  and state  $s_t$  obtains.

<sup>8</sup> Unless otherwise stated, it is understood that  $t$  is an element of  $\mathbb{N}$ , the set of natural numbers. We use the terms “time” and “period” interchangeably to refer to  $t$ .

<sup>9</sup> We use boldface letters for random variables and normal letters for realizations.

<sup>10</sup> Existence and uniqueness are guaranteed by the Kolmogorov Extension Theorem.

**Information feedback** In a dynamic setting, the outcome that the DM observes provides feedback about past states, which is a source of “endogenous” (choice dependent) information.<sup>11</sup> Its relevance is peculiar to the dynamic setting and will play a key role in the paper. By selecting action  $a_t \in A$  at time  $t$ , the DM observes ex-post the outcome  $m_t = f(a_t, s_t)$  if state  $s_t$  realizes. Thus, a DM who selects action  $a_t$  and observes outcome  $m_t$  ex-post knows that the realized state  $s_t$  belongs to the set  $\{s \in S : f(a_t, s) = m_t\} = f_{a_t}^{-1}(m_t)$ .

In general, ex-post information about the state is typically endogenous; that is, the partition of the state space  $S$  induced by outcomes,

$$\left\{ f_a^{-1}(m) : m \in M \right\},$$

may depend on the choice of action  $a$ . If the DM receives the same information about states regardless of his action, namely, if:

$$\forall a, a' \in A, \left\{ f_a^{-1}(m) : m \in M \right\} = \left\{ f_{a'}^{-1}(m) : m \in M \right\},$$

we say that feedback satisfies *own-action independence*. In particular, there is *perfect feedback* when the DM observes the realized state  $s_t$  ex-post; that is, if  $f_a$  is injective for each  $a \in A$ .

Actions and outcomes are remembered: At each period  $t > 1$ , the ex-ante endogenous information—that is, the endogenous information gathered prior to the period- $t$  decision—is given by the history of outcomes  $m^{t-1} = (m_1, \dots, m_{t-1})$  that obtained in the previous periods as a result of the history of actions  $a^{t-1} = (a_1, \dots, a_{t-1})$  and states  $s^{t-1} = (s_1, \dots, s_{t-1})$ .

**Example 1 (Prelude).** Consider an urn that contains black ( $B$ ), green ( $G$ ), and yellow ( $Y$ ) balls. At each time  $t$ , the DM is asked to bet 1 euro on the color of the ball that will be drawn from the urn; therefore, the possible bets are black ( $b$ ), green ( $g$ ), and yellow ( $y$ ). Suppose that the DM is told ex-ante that one-third of the balls are black (and that the only possible colors are  $B$ ,  $G$ , and  $Y$ ), as in the classical Ellsberg’s paradox. That is, the set of posited models is  $\Theta = \{\theta \in \Delta(\{B, G, Y\}) : \theta(B) = 1/3\}$ . Ex post, after the draw, he only learns the result of his bet, namely, whether or not he wins 1 euro. Here,  $S = \{B, G, Y\}$ ,  $A = \{b, g, y\}$ , and  $M = \{0, 1\}$ . The feedback function is described in the following table:

$f$	$B$	$Y$	$G$
$b$	1	0	0
$y$	0	1	0
$g$	0	0	1

Therefore, we have:

$$\begin{aligned} f_b^{-1}(1) &= \{B\}, & f_b^{-1}(0) &= \{Y, G\}, \\ f_y^{-1}(1) &= \{Y\}, & f_y^{-1}(0) &= \{B, G\}, \\ f_g^{-1}(1) &= \{G\}, & f_g^{-1}(0) &= \{B, Y\}. \end{aligned}$$

Note that own-action independence is violated: Ex post, betting on  $b$  yields the partition  $\{\{B\}, \{Y, G\}\}$  of  $S$ , while the bets on  $y$  and  $g$  respectively yield the partitions  $\{\{Y\}, \{B, G\}\}$  and  $\{\{G\}, \{B, Y\}\}$ . ▲

<sup>11</sup> We refer to the working paper version of this paper for a more general setting that separates outcomes and feedback.

**Example 2 (Two-arm bandit).** There are two urns,  $I$  and  $II$ , with black and green balls. The DM chooses an urn, say  $k$ , and wins 1 euro if the ball drawn from urn  $k$  is green ( $G_k$ , good outcome from urn  $k$ ) and zero if it is black ( $B_k$ , bad outcome from urn  $k$ ). The outcome for the chosen urn is observed ex-post. Here,  $S = \{B_I B_{II}, B_I G_{II}, G_I B_{II}, G_I G_{II}\}$ ,  $A = \{I, II\}$ , and  $M = \{0, 1\}$ . The following table describes the feedback function:

$f$	$B_I B_{II}$	$B_I G_{II}$	$G_I B_{II}$	$G_I G_{II}$
$I$	0	0	1	1
$II$	0	1	0	1

Therefore:

$$f_I^{-1}(1) = \{G_I B_{II}, G_I G_{II}\}, \quad f_I^{-1}(0) = \{B_I B_{II}, B_I G_{II}\},$$

$$f_{II}^{-1}(1) = \{B_I G_{II}, G_I G_{II}\}, \quad f_{II}^{-1}(0) = \{B_I B_{II}, G_I B_{II}\}.$$

Own-action independence of feedback is once again violated.  $\blacktriangle$

### 2.3. Strategies and information

**Strategies** At each period  $t$ , the overall ex-ante information available to the DM is given by the histories of actions and outcomes,  $a^{t-1}$  and  $m^{t-1}$ . The ex-ante information history  $h_t$  at time  $t$  is given by:

$$h_1 = (a^0, m^0); \quad \forall t > 1, h_t = (a^{t-1}, m^{t-1}) = (h_{t-1}, a_{t-1}, m_{t-1}),$$

where  $(a^0, m^0)$  represents null data. Hence, the ex-ante information history space  $H_{t+1}$  at the beginning of period  $t + 1$ , determined by information about previous periods, is:

$$H_{t+1} = \{(a^t, m^t) \in A^t \times M^t : \exists s^t \in S^t, \forall k \in \{1, \dots, t\}, m_k = f(a_k, s_k)\}.$$

By definition,  $H_1 = \{(a^0, m^0)\}$ .

Strategies specify an action for each possible information history. Thus, they are modeled as sequences  $\alpha = (\alpha_t)$  of time- $t$  strategies, with  $\alpha_t : H_t \rightarrow A$  for each  $t$ . Since  $H_1 = \{(a^0, m^0)\}$  is a singleton, the first element in the sequence,  $\alpha_1$ , prescribes a non-contingent action. Sometimes it is useful to refer to the strategy  $\alpha|h_t$  that behaves as specified by  $h_t$  at the information sets preceding  $h_t$ , and coincides with  $\alpha$  elsewhere.

**Information and strategies** A state process  $\mathbf{s}^\infty$  and a strategy  $\alpha = (\alpha_t)$  recursively induce an action process  $(\mathbf{a}_t^\alpha)$ , an outcome process  $(\mathbf{m}_t^\alpha)$ , and an information process  $\mathbf{h}^\alpha = (\mathbf{h}_t^\alpha)$ , as follows:

- (i)  $\mathbf{h}_1^\alpha = (a^0, m^0)$  and  $\mathbf{a}_1^\alpha = \alpha_1(a^0, m^0)$ ;
- (ii)  $\mathbf{m}_1^\alpha = f(\mathbf{a}_1^\alpha, \mathbf{s}_1)$ ;
- (iii)  $\mathbf{h}_{t+1}^\alpha = (\mathbf{h}_t^\alpha, \mathbf{a}_t^\alpha, \mathbf{m}_t^\alpha)$ ,  $\mathbf{a}_{t+1}^\alpha = \alpha_{t+1}(\mathbf{h}_{t+1}^\alpha)$ , and  $\mathbf{m}_{t+1}^\alpha = f(\mathbf{a}_{t+1}^\alpha, \mathbf{s}_{t+1})$  for each  $t$ .

In words, at each period  $t$ , an action  $a_t$  is selected according to  $\alpha_t$  based on the information history  $h_t = (h_{t-1}, a_{t-1}, m_{t-1})$ . In turn, its execution generates an outcome  $m_t$  that the DM may consider in subsequent periods. Note that  $\alpha_1$  prescribes only one action,  $\alpha_1(a^0, m^0)$ , which, together with realization  $s_1$  of  $\mathbf{s}_1$ , initializes the recursion by sending outcome  $m_1$ .



The sequence of  $\sigma$ -algebras  $(\sigma(\mathbf{h}_t^\alpha))$  on  $S^\infty$  generated by the information process  $(\mathbf{h}_t^\alpha)$  is a filtration that describes the information structure generated and used by strategy  $\alpha$ . Since feedback will typically not be perfect, this filtration is coarser than the one generated by the state process  $s^\infty$ ; that is,  $\sigma(\mathbf{h}_t^\alpha) \subseteq \sigma(s^{t-1})$  for each  $t > 1$ . For this reason, without loss of generality, we can regard  $\mathbf{h}_t^\alpha$  as well as  $\mathbf{a}_t^\alpha$  and  $\mathbf{m}_{t-1}^\alpha$  as functions defined on  $S^{t-1}$ .<sup>12</sup>

Each finite history  $h_t = (a^{t-1}, m^{t-1})$  corresponds to the cylinder:

$$I(h_t) = \{s^\infty \in S^\infty : \forall \tau \in \{1, \dots, t-1\}, f(a_\tau, s_\tau) = m_\tau\} \in \sigma(s^{t-1}).$$

This is the information about the realized sequence of states revealed by  $h_t$ .

Since states are not directly observed, we can focus on processes  $(\mathbf{a}_t^\alpha)$ ,  $(\mathbf{m}_t^\alpha)$ , and  $(\mathbf{h}_t^\alpha)$ , keeping the underlying parameterized probability space  $(S^\infty, \mathcal{B}(S^\infty), p_\theta)$  in the background. We write events in terms of the processes observable by the DM. In particular,

$$[\mathbf{h}_{t+1}^\alpha = (a^t, m^t)] = \begin{cases} I(a^t, m^t), & \text{if } \forall \tau \in \{1, \dots, t\}, \alpha_\tau(a^{\tau-1}, m^{\tau-1}) = a_\tau, \\ \emptyset, & \text{otherwise.} \end{cases}$$

**Example 3 (Act I).** Assume that only bets on either black or yellow are possible, not on green. As a result, we now have  $A = \{b, y\}$  and the table in the Prelude becomes:

$f$	$B$	$Y$	$G$
$b$	1	0	0
$y$	0	1	0

Throughout we will consider two strategies, denoted by  $\alpha^{NE}$  (No Experimentation) and  $\alpha^E$  (Experimentation). Strategy  $\alpha^{NE}$  dictates betting on black forever. Strategy  $\alpha^E$  dictates experimenting with yellow in period 1, and, from period 2 onwards, the action prescribed is constant but depends on the result of this experimentation: If a success is observed in period 1,  $y$  is chosen; otherwise  $b$  is chosen every period thereafter. Formally:

**Strategy  $\alpha^{NE}$ :** For each  $h_t = (a^{t-1}, m^{t-1})$ ,

$$\alpha_t^{NE}(h_t) = \begin{cases} b & \text{if } t = 1, \\ y & \text{if } t > 1, \text{ and } (y, 1) \in \{(a_1, m_1), \dots, (a_{t-1}, m_{t-1})\}, \\ b & \text{if } t > 1, \text{ and } (y, 1) \notin \{(a_1, m_1), \dots, (a_{t-1}, m_{t-1})\}. \end{cases}$$

(Of course, to assess deviations, the strategy must specify actions to be taken at histories that the strategy itself excludes, such as what to do after having bet on yellow.)

By always betting on black, the DM cannot observe the relative frequencies of  $Y$  and  $G$ . In particular, for each period  $t$  and state history  $s^t$ ,

$$\begin{aligned} \mathbf{a}_t^{\alpha^{NE}}(s^{t-1}) &= b, \\ \mathbf{m}_t^{\alpha^{NE}}(s^t) &= \begin{cases} 1 & \text{if } s_t = B, \\ 0 & \text{if } s_t \in \{Y, G\}, \end{cases} \\ \mathbf{h}_{t+1}^{\alpha^{NE}}(s^t) &= \begin{cases} (\mathbf{h}_t^{\alpha^{NE}}(s^{t-1}), b, 1) & \text{if } s_t = B, \\ (\mathbf{h}_t^{\alpha^{NE}}(s^{t-1}), b, 0) & \text{if } s_t \in \{Y, G\}. \end{cases} \end{aligned}$$

<sup>12</sup> Recall that  $\sigma(s^0)$  is the trivial  $\sigma$ -algebra.

**Strategy  $\alpha^E$ :** For each  $h_t = (a^{t-1}, m^{t-1})$ ,

$$\alpha_t^E(h_t) = \begin{cases} y & \text{if } t = 1, \\ y & \text{if } t > 1, \text{ and } (y, 1) \in \{(a_1, m_1), \dots, (a_{t-1}, m_{t-1})\}, \\ b & \text{if } t > 1, \text{ and } (y, 1) \notin \{(a_1, m_1), \dots, (a_{t-1}, m_{t-1})\}. \end{cases}$$

The only difference between this strategy and  $\alpha^{NE}$  is the action chosen in the first period. Next we describe the induced processes of actions and outcomes:

$$\begin{aligned} \mathbf{a}_1^{\alpha^E} &= y, \\ \mathbf{m}_1^{\alpha^E}(s_1) &= \begin{cases} 1 & \text{if } s_1 = Y, \\ 0 & \text{if } s_1 \in \{B, G\}, \end{cases} \\ \mathbf{h}_2^{\alpha^E}(s_1) &= \begin{cases} (y, 1) & \text{if } s_1 = Y, \\ (y, 0) & \text{if } s_1 \in \{B, G\}, \end{cases} \end{aligned}$$

and, for each  $t > 1$  and  $s^t$ ,

$$\begin{aligned} \mathbf{a}_t^{\alpha^E}(s^{t-1}) &= \begin{cases} y & \text{if } s_1 = Y, \\ b & \text{else,} \end{cases} \\ \mathbf{m}_t^{\alpha^E}(s^t) &= \begin{cases} 1 & \text{if } s_1 = Y \text{ and } s_t = Y, \\ 1 & \text{if } s_1 \in \{B, G\} \text{ and } s_t = B, \\ 0 & \text{else,} \end{cases} \\ \mathbf{h}_{t+1}^{\alpha^E}(s^t) &= \begin{cases} (\mathbf{h}_t^{\alpha^E}(s^{t-1}), y, 1) & \text{if } s_1 = Y \text{ and } s_t = Y, \\ (\mathbf{h}_t^{\alpha^E}(s^{t-1}), b, 1) & \text{if } s_1 \in \{B, G\} \text{ and } s_t = B, \\ (\mathbf{h}_t^{\alpha^E}(s^{t-1}), y, 0) & \text{if } s_1 = Y \text{ and } s_t \in \{B, G\}, \\ (\mathbf{h}_t^{\alpha^E}(s^{t-1}), b, 0) & \text{if } s_1 \in \{B, G\} \text{ and } s_t \in \{Y, G\}. \quad \blacktriangle \end{cases} \end{aligned}$$

### 3. Models and learning

#### 3.1. Distributions and updating

**Prior, predictive, and posterior probabilities** A probability measure  $\mu : \mathcal{B}(\Theta) \rightarrow [0, 1]$  is called a prior. A prior induces a predictive distribution  $p_\mu \in \Delta(S^\infty)$  defined by  $p_\mu(B) = \int_\Theta p_\theta(B) \mu(d\theta)$  for all  $B \in \mathcal{B}(S^\infty)$ . We make the following maintained assumption about the prior of the DM.

**Assumption 2.** The prior  $\mu$  has finite support and satisfies the following “no-surprise” property:

$$\forall t \in \mathbb{N}, \forall h_t \in H_t \quad p_{\bar{\theta}}(I(h_t)) > 0 \Rightarrow p_\mu(I(h_t)) > 0. \tag{3}$$

In words,  $p_{\bar{\theta}}$ -almost surely, the DM is not surprised by what he observes.<sup>13</sup> A sufficient condition for (3) is to have  $\mu(\bar{\theta}) > 0$ . A weaker sufficient condition for (3), one which allows for

<sup>13</sup> Our results would still hold without the “no-surprise” assumption if we endow the DM with a conditional probability system describing how he updates, or revises, his beliefs over models (see the working paper version of Battigalli et al., 2019).

misspecification, is the existence of a model  $\hat{\theta} \in \text{supp } \mu$  such that  $\bar{\theta}$  is absolutely continuous with respect to  $\hat{\theta}$ .<sup>14</sup> For example, it is sufficient that the prior assigns positive probability to a full support model. With this, the posterior of the DM after every information history  $h_t$  with positive objective probability (i.e., with  $p_{\bar{\theta}}(I(h_t)) > 0$ ) is given by Bayes rule:

$$\forall \theta \in \Theta \quad \mu(\theta|h_t) = \frac{p_{\theta}(I(h_t)) \mu(\theta)}{p_{\mu}(I(h_t))}.$$

Our use of Bayes rule for an ambiguity sensitive DM may seem surprising since Bayesian updating has a revealed-preference foundation in the Subjective Expected Utility axioms for preferences over strategies. However, if we do not allow for commitment, we cannot rely on a revealed-preference approach to justify Bayesian updating in this setting: Strategies cannot be chosen, only actions can be chosen. This impossibility of commitment to a particular strategy is critical in a context where preferences are allowed to be dynamically inconsistent.

By sticking to Bayesian updating, we can preserve the separation between ambiguity attitudes and the perception of ambiguity of the static KMM decision criterion in this dynamic setting. This separation is lost if we consider dynamically consistent rules for updating beliefs.<sup>15</sup> Our approach complements the analysis of the dynamic choices of a KMM decision maker by Hanany et al. (2019). In their work, they maintain dynamic consistency of the preferences of the DM but consider a different updating rule for beliefs.

**Observationally equivalent models** Given a strategy  $\alpha$  and a probability distribution  $p_{\theta}$ , let  $p_{\theta}^{\alpha} : \sigma(\mathbf{h}^{\alpha}) \rightarrow [0, 1]$  denote the restriction of  $p_{\theta}$  to the  $\sigma$ -algebra generated by the  $\alpha$ -observable events:  $\sigma(\mathbf{h}^{\alpha}) = \sigma(\cup_t \sigma(\mathbf{h}_t^{\alpha}))$ . With this, we define the  $\sigma(\mathbf{h}_t^{\alpha})$ -measurable correspondence (random set) representing the collection of models that are deemed possible and that, conditional on  $\mathbf{h}_t^{\alpha}(s^{t-1})$ , are *observationally equivalent* to the true model  $\bar{\theta}$  under  $\alpha$  and prior  $\mu$ . Formally,<sup>16</sup>

$$\Theta_t^{\alpha, \mu}(s^{\infty}) = \left\{ \theta \in \text{supp } \mu \left( \cdot | \mathbf{h}_t^{\alpha}(s^{t-1}) \right) : p_{\theta}^{\alpha} \left( \cdot | \mathbf{h}_t^{\alpha}(s^{t-1}) \right) = p_{\bar{\theta}}^{\alpha} \left( \cdot | \mathbf{h}_t^{\alpha}(s^{t-1}) \right) \right\}.$$

Note that, for some  $s^{\infty}$ , the set  $\Theta_t^{\alpha, \mu}(s^{\infty})$  may be empty if  $\bar{\theta} \notin \text{supp } \mu$ .

The next lemma establishes a monotonicity property of this correspondence. We introduce the following abuse of notation/terminology: When a property holds  $p_{\theta}$ -almost surely, we will simply say that it holds  $\theta$ -almost surely ( $\theta$ -a.s. for short).

**Lemma 1.** *For every true model  $\bar{\theta} \in \Theta$  and every period  $t$ ,  $\Theta_t^{\alpha, \mu} \subseteq \Theta_{t+1}^{\alpha, \mu}$   $\bar{\theta}$ -a.s.*

The intuition behind the lemma is as follows. The set  $\Theta_t^{\alpha, \mu}$  may contain models that disagree with  $\bar{\theta}$  on the relative probabilities of past events (up to  $t - 1$ ), but that agree with  $\bar{\theta}$  on the relative probabilities of future events (from  $t$ ). Almost surely, every model that agrees with  $\bar{\theta}$  on future events conditional on information up to  $t - 1$  also agrees on future events conditional on information up to  $t$ . Act II of our running Example shows that the inclusion can be strict.

It follows from the lemma that,  $\bar{\theta}$ -a.s.,

$$\Theta_1^{\alpha, \mu} := \left\{ \theta \in \text{supp } \mu : p_{\theta}^{\alpha} = p_{\bar{\theta}}^{\alpha} \right\} \subseteq \Theta_t^{\alpha, \mu}$$

<sup>14</sup> That is,  $\text{supp } \bar{\theta} \subseteq \text{supp } \hat{\theta}$ . Under perfect feedback, this condition is equivalent to (3).

<sup>15</sup> See Hanany and Klibanoff (2009), further discussed in Section 7, and Maccheroni et al. (2006).

<sup>16</sup> It is actually enough to require  $p_{\theta}^{\alpha}(E | \mathbf{h}_t^{\alpha}(s^{t-1})) = p_{\bar{\theta}}^{\alpha}(E | \mathbf{h}_t^{\alpha}(s^{t-1}))$  for all  $E \in \cup_{t \geq 1} \sigma(\mathbf{h}_t^{\alpha})$ . That is, observational equivalence is determined by the  $\alpha$ -observable events.

for every  $t$ . Set  $\Theta_1^{\alpha,\mu}$  represents the irreducible model uncertainty that, when  $\bar{\theta}$  is the true model, the DM faces if he plays  $\alpha$  and holds belief  $\mu$ .<sup>17</sup> When  $\Theta_1^{\alpha,\mu} = \text{supp } \mu$ , such uncertainty and strategy do not allow any learning, as all the models that the DM initially deems possible are  $\alpha$ -observationally equivalent to the true model. The opposite is true when  $\Theta_1^{\alpha,\mu} = \{\bar{\theta}\}$ , since in this case the DM will assign probability arbitrarily close to 1 to the true model as he accumulates observations.

In what follows, we will often study properties of a triple  $(\alpha, \mu, \bar{\theta})$  where  $\alpha$  is the strategy carried out by the DM,  $\mu$  is his prior belief over models at period 0, and  $\bar{\theta}$  is the correct model. We want to understand the behavior and learning of a DM who follows strategy  $\alpha$  when his prior is  $\mu$  and the true model is  $\bar{\theta}$ . Therefore, we have a particular interest in statements that hold  $\bar{\theta}$ -a.s., that is, that are almost surely true for the correct model. The notion of observationally equivalent models motivates the following definition.

**Definition 1.** A triple  $(\alpha, \mu, \bar{\theta})$  is consistent at time  $t$  if  $\Theta_t^{\alpha,\mu} \neq \emptyset \bar{\theta}$ -a.s.

In words, a triple  $(\alpha, \mu, \bar{\theta})$  is consistent at time  $t$  if, conditional on the available information  $\mathbf{h}_t^\alpha$ , at least one model deemed possible is  $\alpha$ -observationally equivalent to the true model.<sup>18</sup>

Let:

$$\sigma_{\geq t}(\mathbf{h}^\alpha(s^{t-1})) = \sigma \left( \cup_{\tau \geq 1} \left\{ B \in \sigma(\mathbf{h}_{t+\tau}^\alpha) : B \subseteq I(\mathbf{h}_t^\alpha(s^{t-1})) \right\} \right)$$

denote the sigma-algebra of  $\alpha$ -observable events from date  $t$  onwards given  $s^{t-1}$ . Then:

$$\Theta_t^{\alpha,\mu}(s^\infty) = \left\{ \theta \in \text{supp } \mu \left( \cdot | \mathbf{h}_t^\alpha(s^{t-1}) \right) : \forall E \in \sigma_{\geq t}(\mathbf{h}^\alpha(s^{t-1})), p_\theta^\alpha(E) = p_{\bar{\theta}}^\alpha(E) \right\}.$$

Hence,  $(\alpha, \mu, \bar{\theta})$  is consistent at  $t$  if, for  $\bar{\theta}$ -almost every  $s^{t-1}$ , there exists some  $\theta \in \text{supp } \mu(\cdot | \mathbf{h}_t^\alpha(s^{t-1}))$  such that  $p_\theta^\alpha(E) = p_{\bar{\theta}}^\alpha(E)$  for all  $E \in \sigma_{\geq t}(\mathbf{h}^\alpha(s^{t-1}))$ . Of course, an obvious sufficient condition for consistency is that  $\mu(\bar{\theta}) > 0$ .

In view of Lemma 1, for a triple  $(\alpha, \mu, \bar{\theta})$  it is easier to meet the condition for consistency as  $t$  gets larger. Let  $T(\alpha, \mu, \bar{\theta})$  denote the infimum of the set of  $t$  at which  $(\alpha, \mu, \bar{\theta})$  is consistent.<sup>19</sup> If  $T(\alpha, \mu, \bar{\theta}) = T < \infty$ , we say that the triple is consistent from period  $T$ ; if the triple is consistent for some period  $T' < \infty$ , we say that it is eventually consistent. We begin by showing that, under our consistency assumption, beliefs converge almost surely.

**Lemma 2.** If  $(\alpha, \mu, \bar{\theta})$  is eventually consistent, then the process  $(\mu(\cdot | \mathbf{h}_t^\alpha))$  converges  $\bar{\theta}$ -a.s.

For the next result, note that convergence of beliefs along a path does not imply convergence of actions. That is, more than one action can be played infinitely often along a path.

<sup>17</sup> In this work, we use the term “belief” to denote the probability assessment over (stochastic) models. Using the terminology of Marinacci (2015), this belief represents how the DM addresses epistemic uncertainty, whereas models capture the (perceived) physical uncertainty.

<sup>18</sup> The word “consistent” may remind the reader of the consistency criterion imposed in Arrow and Green (1973). However, theirs is an “existence of equilibrium condition” requiring that, given any DM’s action and true model, there exists a subjective model conceivable by the DM that is observationally equivalent to the actual one.

<sup>19</sup> We set  $T(\alpha, \mu, \bar{\theta}) = \infty$  if no such  $t$  exists.

**Proposition 1.** *If  $(\alpha, \mu, \bar{\theta})$  is eventually consistent, then*

$$\lim_{t \rightarrow \infty} \mu(\{\theta \in \Theta : F(\bar{a}, \theta) = F(\bar{a}, \bar{\theta})\} | \mathbf{h}_t^\alpha(s^\infty)) = 1 \tag{4}$$

for  $\bar{\theta}$ -almost every  $s^\infty$  and every action  $\bar{a}$  played infinitely often by  $\alpha$  on path  $s^\infty$ .

In words, a triple  $(\alpha, \mu, \bar{\theta})$  that is eventually consistent allows the DM to learn the  $\alpha$ -observable implications of the true model  $\bar{\theta}$  in the long run.<sup>20</sup>

Under perfect feedback, the true model is asymptotically identified. This is the classical result of Doob (1949).

**Corollary 1.** *If  $(\alpha, \mu, \bar{\theta})$  is eventually consistent, then, under perfect feedback,*

$$\mu(\bar{\theta} | \mathbf{h}_t^\alpha) \rightarrow 1 \quad \bar{\theta}\text{-a.s.}$$

**Example 4 (Act II).** Suppose that the DM:

1. knows that 1/3 of the balls are black (and so all his models  $\theta$  are such that  $\theta(B) = 1/3$ );
2. has a 3-point prior  $\mu$  with  $\text{supp } \mu = \{\theta^Y, \theta^{uni}, \theta^G\}$  and believes it is equally likely that the true model is either  $\theta^Y$  (with  $\theta^Y(Y) = 2/3$ ), the uniform model  $\theta^{uni}$ , or  $\theta^G$  (with  $\theta^G(G) = 2/3$ ):

Marginals	<i>B</i>	<i>Y</i>	<i>G</i>
$\theta^Y$	$\frac{1}{3}$	$\frac{2}{3}$	0
$\theta^{uni}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
$\theta^G$	$\frac{1}{3}$	0	$\frac{2}{3}$

Prior	$\theta^Y$	$\theta^{uni}$	$\theta^G$
$\mu$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

By requiring to always bet on the color with known proportion, strategy  $\alpha^{NE}$  does not allow the DM to learn anything. Formally,

$$\forall s^t \in S^t, \quad \mu(\cdot | \mathbf{h}_{t+1}^{\alpha^{NE}}(s^t)) = \mu.$$

Here,  $T(\alpha^{NE}, \mu, \bar{\theta}) = 1$  and  $\Theta_1^{\alpha^{NE}, \mu} = \text{supp } \mu$  regardless of what  $\bar{\theta}$  is; strategy  $\alpha^{NE}$  only allows partial identification. We also have that  $T(\alpha^E, \mu, \bar{\theta}) = 1$  for  $\alpha^E$  if  $\bar{\theta} \in \text{supp } \mu$ , although  $\alpha^E$  may allow identification of the true model. After a success in period 1, we have  $\text{supp } \mu(\cdot | (y, 1)) = \{\theta^Y, \theta^{uni}\}$  and:

$$\Theta_2^{\alpha^E, \mu}(s^\infty) = \begin{cases} \{\bar{\theta}\} & \text{if } \mathbf{h}_2^{\alpha^E}(s_1) = (y, 1), \\ \{\theta^Y, \theta^{uni}, \theta^G\} & \text{if } \mathbf{h}_2^{\alpha^E}(s_1) = (y, 0), \end{cases}$$

<sup>20</sup> In this respect we differ from the literature on learning with a misspecified prior, and in particular from Esponda and Pouzo (2016). The standard result of this literature is that, starting from a misspecified prior, beliefs about the model and its observable implications may converge to a distribution of observables that is different from the one implied by the model. Instead, although our notion of consistency allows for misspecification, by Proposition 1 beliefs about observables converge to the distribution implied by the model.

Notice that the inclusion in Lemma 1 can be strict. Strategy  $\alpha^E$  experiments with  $y$  at  $t = 1$ , but reverts to  $b$  if a failure is observed, that is, if  $s_1 \in \{B, G\}$ . Given  $\alpha^E$ , model  $\theta^Y$  is the only one that predicts success at  $t = 1$  with probability  $2/3$ , whereas all models predict success with probability  $1/3$  (betting on  $B$ ) from  $t = 2$  if a failure is observed at  $t = 1$ . Therefore, if  $\bar{\theta} = \theta^Y$ ,

$$\Theta_1^{\alpha^E, \mu}(s^\infty) = \{\theta^Y\} \subset \{\theta^Y, \theta^{uni}, \theta^G\} = \Theta_2^{\alpha^E, \mu}(s^\infty)$$

for  $s_1 = B, G$ .

By Proposition 1,

$$\mu(\cdot | \mathbf{h}_t^{\alpha^E}) \rightarrow \begin{cases} \delta_{\bar{\theta}} & \text{if } \mathbf{h}_2^{\alpha^E} = (y, 1), \\ \mu(\cdot | (y, 0)) & \text{if } \mathbf{h}_2^{\alpha^E} = (y, 0), \end{cases}$$

where  $\delta_{\bar{\theta}}$  denotes the Dirac measure on  $\bar{\theta}$ . If experimentation succeeds, the true model is asymptotically learned. Otherwise, if  $h_2 = (y, 0)$ , posterior beliefs attain their limit value as early as the second period, and the DM remains in the dark.  $\blacktriangle$

#### 4. Value

We posit that, in the absence of model uncertainty, the DM ranks alternative strategies according to the standard Discounted Expected Utility criterion. Let  $R : A \times \Theta \rightarrow \mathbb{R}$  denote the (objective) expected reward function defined in Section 2.1.

**Assumption 1'** (*Discounted expected utility on lotteries*). There exists a constant  $\delta \in [0, 1)$  such that, for every objective probability measure  $\theta$ , history  $h_t$  with  $p_\theta(I(h_t)) > 0$ , and strategies  $\alpha$  and  $\beta$ , the DM prefers  $\alpha$  to  $\beta$  if and only if:

$$\sum_{\tau=t}^{\infty} \delta^{\tau-t} \sum_{h_\tau \in H_\tau} R(\alpha(h_\tau), \theta) p_\theta([\mathbf{h}_\tau^{\alpha|h_t} = h_\tau] | h_t) \\ \geq \sum_{\tau=t}^{\infty} \delta^{\tau-t} \sum_{h_\tau \in H_\tau} R(\beta(h_\tau), \theta) p_\theta([\mathbf{h}_\tau^{\beta|h_t} = h_\tau] | h_t).$$

Under model uncertainty, we postulate a dynamic version of smooth ambiguity preferences. Let  $\phi$  be the function capturing ambiguity attitudes introduced in Section 2.1. Given prior  $\mu$ , if history  $h_t$  with  $p_\mu(I(h_t)) > 0$  is observed the DM ranks strategy  $\alpha$  according to the present value of the continuation stream of utility certainty equivalents:<sup>21</sup>

$$V(\alpha, \mu | h_t) := \sum_{\tau=t}^{\infty} \delta^{\tau-t} \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{h_\tau \in H_\tau} R(\alpha(h_\tau), \theta) p_\theta(h_\tau | h_t) \right) \mu(d\theta | h_t) \right). \quad (5)$$

This criterion evaluates uncertain one-period outcomes according to the smooth ambiguity model and then aggregates their (utility) certainty equivalents over time through discounting. Therefore, (utility) smoothing over time is irrelevant. Indeed, when the DM evaluates two continuation streams of utility certainty equivalents, he is interested only in their discounted sum, not on their

<sup>21</sup> We abbreviate  $([\mathbf{h}_\tau^{\alpha|h_t} = h_\tau] | h_t)$  as  $(h_\tau | h_t)$ .

variability over time. Note that the value only depends on the continuation strategy induced by  $\alpha$  starting from  $h_t$ .

In particular, we obtain:

- (i)  $V(\alpha, \mu | h_t) = \sum_{\tau=t}^{\infty} \delta^{\tau-t} (\sum_{h_{\tau} \in H_{\tau}} R(\alpha(h_{\tau}), \sum \mu(\theta | h_t) \theta) \int_{\Theta} p_{\theta}(h_{\tau} | h_t) \mu(d\theta | h_t))$   
(up to an affine transformation) when  $\phi$  is linear;
- (ii)  $V(\alpha, \mu | h_t) = \sum_{\tau=t}^{\infty} \delta^{\tau-t} \sum_{h_{\tau} \in H_{\tau}} R(\alpha(h_{\tau}), \theta) p_{\theta}(h_{\tau} | h_t)$  when  $\text{supp } \mu = \{\theta\}$ .

Our analysis relies on the continuity of  $V$  with respect to the prior  $\mu$ .

**Lemma 3.** *For every strategy  $\alpha$  and history  $h_t$ , the functional  $V(\alpha, \cdot | h_t)$  is continuous on the set of priors  $\mu$  such that  $p_{\mu}(I(h_t)) > 0$ .*

With this lemma, we obtain the following additional corollary to Proposition 1.

**Corollary 2.** *If  $(\alpha, \mu, \bar{\theta})$  is eventually consistent, then,  $\bar{\theta}$ -a.s.,*

$$|V(\alpha, \mu | \mathbf{h}_t^{\alpha}) - V(\alpha, \delta_{\bar{\theta}} | \mathbf{h}_t^{\alpha})| \rightarrow 0.$$

This corollary tells us that, in a consistent triple, the strategy becomes unambiguous on path in terms of value. Of course, this result does not imply that  $\mu(\cdot | \mathbf{h}_t^{\alpha}) \rightarrow \delta_{\bar{\theta}}$ , only that the present value of the strategy that is used converges to the true value. In particular, even in the limit, alternative strategies may entail unknown outcome distributions. Therefore, this corollary is the dynamic version of the observation made in the static analysis of BCMM that only equilibrium strategies have to be unambiguous.

Note that, except for the benchmark case of ambiguity neutrality, this time-additive value *does not admit a recursive formulation*. This is related to the well-known dynamic inconsistency of decision makers with non-neutral attitudes toward ambiguity. For this reason, we are precluded from employing many of the standard dynamic programming results. We provide an example of these inconsistencies in our setting.

**Example 5 (Dynamic inconsistency).** Consider a modified version of our running example. There are only two periods. Only bets on either black or yellow are possible, not on green. However, by paying a small cost  $\varepsilon$ , it is also possible to bet on black and to observe the color of the selected ball (action  $bo$ ). Thus, the outcome has two components:  $m = (m_1, m_2)$ , where  $m_1$  is the monetary payoff and  $m_2$  is the color of the drawn ball under action  $bo$ , and null information (denoted by  $*$ ) under actions  $b$  and  $y$ . Finally, we normalize payoffs as  $u(a, (m_1, m_2)) = m_1$ . The feedback function  $f$  is thus described by the following table:

$f$	$B$	$Y$	$G$
$b$	$1, *$	$0, *$	$0, *$
$y$	$0, *$	$1, *$	$0, *$
$bo$	$1 - \varepsilon, B$	$0 - \varepsilon, Y$	$0 - \varepsilon, G$

Suppose that the DM:

1. knows that  $1/3$  of the balls are black (and so all her models  $\theta$  are such that  $\theta(B) = 1/3$ );
2. believes it is equally likely that the true model is either  $\hat{\theta}^Y$  or  $\hat{\theta}^G$ ;

Marginals	B	Y	G
$\hat{\theta}^Y$	$\frac{1}{3}$	$\frac{5}{12}$	$\frac{1}{4}$
$\hat{\theta}^G$	$\frac{1}{3}$	$\frac{1}{4}$	$\frac{5}{12}$

Prior	$\hat{\theta}^Y$	$\hat{\theta}^G$
$\mu$	$\frac{1}{2}$	$\frac{1}{2}$

Let  $\phi(u) = -e^{-10u}$ . Then, the ex-ante optimal strategy if  $\varepsilon$  is sufficiently small is:

Strategy  $\beta$ : “Bet on black observing the color at  $t = 1$ . For  $t = 2$ , given yellow in the first period, bet on yellow, otherwise bet on black.”<sup>22</sup> The ex-ante value of strategy  $\beta$  is<sup>23</sup>:

$$V(\beta, \mu) = 0.\bar{3} - \varepsilon + \delta 0.3364.$$

However,  $\beta$  does not satisfy the one-deviation property, i.e., it is not incentive compatible. Indeed, after having observed yellow, the DM prefers to bet on black. The posterior belief after having chosen *bo* and having observed yellow is:

Posterior	$\hat{\theta}^Y$	$\hat{\theta}^G$
$\mu(\cdot   (bo, Y))$	$\frac{5}{8}$	$\frac{3}{8}$

Hence,

$$V(\beta, \mu | (bo, Y)) \cong 0.3207 < \frac{1}{3},$$

where  $1/3$  is the value of betting on black.

This is a typical example of dynamically-inconsistent preferences. At period 0, for sufficiently small  $\varepsilon$ , the DM would want to commit to conditioning his behavior on the observed draw. In particular, he would like to choose *y* if the draw in the first period is *Y*, that is, after history  $(bo, Y)$ . Indeed, even if betting on yellow leads to ambiguous consequences, the DM is confident that with high probability, if  $\hat{\theta}^G$  is the true model, *Y* will not be the first-period draw. Therefore, even under model  $\hat{\theta}^G$ , this strategy presents a moderately high expected value. However, after having observed  $(bo, Y)$ , even if the posterior probability of  $\hat{\theta}^G$  is lower, the DM considers the consequences of choosing action *y* to be too ambiguous. Indeed, the expected value under model  $\hat{\theta}^G$ ,  $1/4$ , is quite small. Therefore, since the DM is highly ambiguity averse, he will select *b*.

Moreover, it can be shown that the strategy “always bet on black” has a lower ex-ante value,  $(1 + \delta)/3$ , but satisfies the one-deviation property. A sophisticated DM will not pay the cost  $\varepsilon$  anticipating that he will not condition his behavior on the observed outcome, even if this conditioning is ex-ante optimal. ▲

#### 4.1. Stationary strategies

A strategy  $\alpha$  is *stationary* if, given the prior  $\mu$ , it depends on history only through the induced posterior belief; that is, for all  $t, t' \in \mathbb{N}$  and for every two histories  $h_t \in H_t, h'_{t'} \in H_{t'}$  such that  $p_\mu(I(h_t))$  and  $p_\mu(I(h'_{t'}))$  are strictly positive:

<sup>22</sup> Note that this is not a proper strategy since it does not assign an action to every information history. In particular, it does not assign an action to personal histories ruled out by the strategy itself. However, the specification of the actions selected at those information histories is irrelevant in determining ex-ante optimality.

<sup>23</sup> We refer to the working paper version for the computations used to obtain formulas in the examples.



$$\mu(\cdot|h_t) = \mu(\cdot|h'_t) \Rightarrow \alpha(h_t) = \alpha(h'_t).$$

Note that stationarity is a property of the pair  $(\alpha, \mu)$  of strategy and prior.

The following result, obtained using a coupling argument, shows that when the DM uses a stationary strategy, the function  $V$  depends on the history only through belief updating.

**Lemma 4.** *For all  $\alpha$  and  $\mu$  satisfying stationarity, if  $h_t$  and  $h'_t$  are two histories such that  $p_\mu(I(h_t))$  and  $p_\mu(I(h'_t))$  are strictly positive, then  $\mu(\cdot|h_t) = \mu(\cdot|h'_t)$  implies  $V(\alpha, \mu|h_t) = V(\alpha, \mu|h'_t)$ .*

Under ambiguity neutrality, Hinderer (1970)<sup>24</sup> proves the existence of an optimal stationary strategy for arbitrary beliefs. Therefore, it is without loss of generality to focus on stationary strategies. We cannot adopt this approach “as is” because it relies on a notion of global optimality that may violate incentive compatibility under dynamically inconsistent preferences, as shown in Example 5. However, Proposition 2 in the next section provides a partially analogous result for the notion of rationality (intra-personal equilibrium) used in this paper. Given this result, we focus on stationary strategies in the sequel. Therefore, with an abuse of notation, we will often regard  $\alpha$  as a function of beliefs over probability models.<sup>25</sup> More precisely, given prior  $\mu$ , we call *belief-range of  $\mu$*  the set of beliefs that the DM may hold with positive probability under  $p_\mu$ :

$$\{\mu(\cdot|h_t) : t \in \mathbb{N}, h_t \in H_t, p_\mu(I(h_t)) > 0\}.$$

Then for all  $\nu = \mu(\cdot|h_t)$  in the belief range of  $\mu$ ,  $\alpha(\nu)$  is defined as:

$$\alpha(\nu) = \alpha(h_t).$$

## 5. Self-confirming equilibrium and learning

### 5.1. Self-confirming equilibrium

A self-confirming equilibrium (SCE) is a state in which the DM chooses a subjectively optimal action according to the static smooth criterion and, given this choice, his posterior belief coincides with his prior belief. Formally:

**Definition 2.** A triple  $(a^*, \mu^*, \bar{\theta}) \in A \times \Delta(\Theta) \times \Theta$  is a *self-confirming equilibrium* (SCE) if:

- (i)  $\mu^*(\{\theta \in \Theta : F(a^*, \theta) = F(a^*, \bar{\theta})\}) = 1$ ;
- (ii)  $a^* \in \arg \max_{a \in A} \bar{V}(a, \mu^*)$ .

We say that  $a^*$  is an SCE action if it is part of an SCE  $(a^*, \mu^*, \bar{\theta})$ , and that  $a^*$  is a Nash equilibrium action if  $a^* \in \arg \max_{a \in A} R(a, \bar{\theta})$ .

The second condition says that  $a^*$  is a (myopic, or one-period) best response to  $\mu^*$  given the ambiguity attitude determined by  $\phi$ . The first condition is a self-confirming property adapted to

<sup>24</sup> Building on the results of Dynkin (1965).

<sup>25</sup> We refer the interested reader to the working paper version for more general results for non-stationary strategies in a non i.i.d. setting.

the static framework. It is equivalent to requiring that,  $\bar{\theta}$ -a.s.,  $\mu^*(\cdot|a^*, f(a^*, \mathbf{s})) = \mu^*(\cdot)$ , because all models in the support of  $\mu^*$  yield the same distribution of outcomes as the true model  $\bar{\theta}$  given action  $a^*$ . We can interpret this condition as follows. The distribution of outcomes that the DM “observes” in the long run if he always plays  $a^*$  is precisely what he expects it to be. In this sense,  $a^*$  is *unambiguous* for  $\mu^*$ ; since payoffs are observable, the self-confirming property implies that the expected distribution of payoffs coincides with the one implied by the true model  $\bar{\theta}$ .

**Remark 1.** If  $(a^*, \mu^*, \bar{\theta})$  is an SCE,  $R(a^*, \theta) = \mathbb{E}_{F(a^*, \theta)}[u_{a^*}]$  is constant over  $\text{supp} \mu^*$ . ▲

### 5.2. Rational learning dynamics

When the DM faces a recurrent choice problem, the notion of SCE characterizes behavior and beliefs after the latter have “converged.” In other words, the data provided by the equilibrium strategy do not lead to any further updating because the models that the DM deems possible in an SCE cannot be distinguished from each other or from the true model.

In dynamic settings, we are interested not only in behavior after beliefs have reached a steady state, but also in rational behavior as the DM is learning from the data. To define our notion of rationality, we introduce the following notation. For any information history  $h_t$  and action  $a$ , let  $\alpha/(h_t, a)$  denote the continuation strategy that selects  $a$  at  $h_t$  and behaves as  $\alpha$  thereafter.<sup>26</sup>

**Definition 3.** A pair  $(\alpha, \mu)$  is *rational* if it satisfies stationarity and if

$$\forall t \in \mathbb{N}, \forall h_t \in H_t, \forall a \in A \quad p_\mu(I(h_t)) > 0 \Rightarrow V(\alpha, \mu | h_t) \geq V(\alpha/(h_t, a), \mu | h_t).$$

Besides stationarity, this condition is the one-deviation property, which says that—for every information history  $h_t$  that the DM deems reachable with positive probability—action  $\alpha(\mu(\cdot | h_t))$  maximizes the continuation value conditional on  $h_t$  given that  $\alpha$  is expected to apply in the future. The motivation is the following: Strategy  $\alpha$  is a plan formulated by a sophisticated DM who understands his sequential incentives. In each period  $t$ , the DM only controls the action in that period, and therefore we require that he maximizes his value with respect to what he can control, given the predicted behavior of his “future selves,” i.e., his continuation strategy. If the time horizon is finite, then this condition is equivalent to folding-back planning: In all periods, the DM predicts to behave in the last period according to his last-period contingent incentives; given such last-period prediction, the DM predicts to behave in the second-to-last period according to his second-to-last period contingent incentives, and so on. When the DM is ambiguity neutral, the one-deviation principle implies that if  $(\alpha, \mu)$  is rational, then strategy  $\alpha$  is subjectively optimal given  $\mu$ .

A strategy  $\alpha$  is said to be *rational given  $\mu$*  if the pair  $(\alpha, \mu)$  is rational. The proposition below establishes the existence of rational strategies given any prior.

**Proposition 2.** For every prior  $\mu$  there exists a strategy  $\alpha$  that is rational given  $\mu$ .

We illustrate the concept of rationality in our running example.

<sup>26</sup> Formally, the value function maps strategies into real numbers. However, we have already pointed out that the value at a history  $h_t$  depends only on the continuation of the posited strategy. Therefore, the expression  $V(\alpha/(h_t, a), \mu | h_t)$  is well defined.

**Example 6 (Act III).** We normalize payoffs as  $u(a, 0) = 0$  and  $u(a, 1) = 1$ . Outcomes are thus the bets' payoffs. Moreover, we assume that  $\phi(u) = -e^{-\lambda u}$ , so that higher (absolute) ambiguity aversion corresponds to higher  $\lambda$  (see Klibanoff et al., 2005).

Suppose that the DM features the prior  $\mu$  presented in Act II. We consider the same strategies  $\alpha^{NE}$  and  $\alpha^E$  analyzed there. The former strategy involves no experimentation as it recommends always betting on black, the color with the known proportion. Thus, the value of this strategy is independent of histories and beliefs, and it is given by<sup>27</sup>:

$$V(\alpha^{NE}, \mu|h_t) = \frac{\frac{1}{3}}{1 - \delta}.$$

The latter strategy recommends betting on  $y$  at  $t = 1$  and then switching to  $b$  permanently if and only if this first bet is unsuccessful. If the DM chooses  $y$ , the outcomes are informative about the distribution, and he updates his belief. Recall by Act II that the DM has a uniform 3-point prior  $\mu$  with  $\text{supp } \mu = \{\theta^Y, \theta^{uni}, \theta^G\}$ . If we denote  $\mu(\cdot|h_t) := (\mu(\theta^{uni}|h_t), \mu(\theta^Y|h_t), \mu(\theta^G|h_t))$ , the posterior is:

$$\mu(\cdot|(y, 1)) = \left(\frac{1}{3}, \frac{2}{3}, 0\right)$$

if the outcome is 1 (success), and:

$$\mu(\cdot|(y, 0)) = \left(\frac{1}{3}, \frac{1}{6}, \frac{1}{2}\right)$$

otherwise.

After the first period, strategy  $\alpha^E$  recommends a fixed action. Thus, the continuation value is equal to  $\frac{1}{1-\delta}$  times the single-period expected payoff. Specifically, for any history  $h_t$  ( $t > 1$ ) that induces belief  $\mu(\cdot|h_t)$ , the continuation-value after a success in period 1 is

$$V(\alpha^E, \mu|h_t) = \frac{\phi^{-1}(\mu(\theta^{uni}|h_t)\phi(\frac{1}{3}) + \mu(\theta^Y|h_t)\phi(\frac{2}{3}))}{1 - \delta}.$$

For the initial (empty) history, we have:

$$V(\alpha^E, \mu|h_1) = \phi^{-1}\left(\frac{1}{3}\phi\left(\frac{1}{3}\right) + \frac{1}{3}\phi\left(\frac{2}{3}\right) + \frac{1}{3}\phi(0)\right) + \frac{\delta}{1 - \delta}\phi^{-1}\left(\frac{1}{3}\phi\left(\frac{1}{3}\right) + \frac{1}{3}\phi\left(\frac{5}{9}\right) + \frac{1}{3}\phi\left(\frac{1}{3}\right)\right).$$

Two forces affect the option value of experimentation: ambiguity aversion (the higher the value of  $\lambda$ , the lower the value of experimentation) and patience (the higher the value of  $\delta$ , the higher the value of experimentation). Given this, strategy  $\alpha^{NE}$  is preferred if either  $\delta = 0$  or  $\lambda$  is high enough given  $\delta > 0$ ; if so, the pair  $(\alpha^{NE}, \mu)$  is rational. As for strategy  $\alpha^E$ , if  $\delta$  is sufficiently high and  $\lambda$  is low enough, e.g.,  $\lambda = 1$  and  $\delta = 0.39$ , strategy  $\alpha^E$  satisfies the one-deviation property at  $(a^0, m^0)$ . However, because of experimentation, we need to consider two different contingencies.

<sup>27</sup> This holds only for histories allowed by the strategy, namely on path.

1. If experimentation is successful (i.e.,  $s_1 = Y$ ), the DM learns that model  $\theta^G$  is false and updates his belief from  $(1/3, 1/3, 1/3)$  to  $(1/3, 2/3, 0)$ . Moreover, at every information history, Bayesian updating implies that the posterior will be of the form  $(1 - k, k, 0)$ , with  $k \in (0, 1)$ . At this point, the strategy recommends sticking to  $y$ . It can be checked that this recommendation is better than trying out  $b$  once before switching to  $y$  thereupon. In other words, it satisfies the one-deviation property: For all  $\delta \in (0, 1)$  and all  $\lambda > 0$ , the value at an information history  $h_t = ((y, 1), \dots)$  with  $\mu(\cdot|h_t) = (1 - k, k, 0)$  satisfies:

$$\begin{aligned} V(\alpha^E, \mu|h_t) &= \frac{\phi^{-1}((1 - k)\phi(\theta^{uni}(Y)) + k\phi(\theta^Y(Y)))}{1 - \delta} \\ &= \frac{\phi^{-1}((1 - k)\phi(\frac{1}{3}) + k\phi(\frac{2}{3}))}{1 - \delta} \\ &> \frac{1}{3} + \delta \frac{\phi^{-1}((1 - k)\phi(\frac{1}{3}) + k\phi(\frac{2}{3}))}{1 - \delta} \\ &= V(\alpha^E/(h_t, b), \mu|h_t). \end{aligned}$$

2. If experimentation is unsuccessful (i.e.,  $s_1 \in \{B, G\}$ ), the posterior lowers the weight of model  $\theta^Y$  relative to models  $\theta^{uni}$  and  $\theta^G$ , so that  $p_\mu(Y | (y, 0)) < p_\mu(B | (y, 0)) = 1/3$ . Thereupon, strategy  $\alpha^E$  recommends switching (and sticking) to black, so that the continuation value is the same as that under  $\alpha^{NE}$ . Moreover, since betting on black does not lead to any further updating, it is enough to check the inequality with second-period beliefs. For sufficiently small  $\delta$ , or for sufficiently high  $\lambda$ ,

$$V(\alpha^E, \mu|h_2) = \frac{1}{3} \frac{1}{1 - \delta} > V(\alpha^E/(h_1, y), \mu|h_1).$$

In particular, this inequality holds with  $\lambda = 1$  and  $\delta = 0.39$ , and we have already argued that  $(\alpha^E, \mu)$  satisfies the one-deviation property at the initial history; therefore,  $(\alpha^E, \mu)$  is rational.  $\blacktriangle$

### 5.3. Convergence to SCE

We are interested in studying the limit behavior of a rational DM. In particular, we investigate the conditions that imply convergence to SCE. Building on Lemma 2, we provide a learning foundation to the SCE concept of BCMM.

The stochastic process of actions and beliefs  $(\mathbf{a}_t^\alpha, \mu(\cdot|h_t^\alpha))$  converges to an SCE if, for  $\bar{\theta}$ -almost every  $s^\infty$ , beliefs converge to a limit  $\mu_{s^\infty}^\alpha$  and there exists a finite time  $t$  such that  $(\mathbf{a}_\tau^\alpha(s^{\tau-1}), \mu_{s^\infty}^\alpha, \bar{\theta})$  forms an SCE for all  $\tau \geq t$ . Note that the tail sequence of actions  $(\mathbf{a}_\tau^\alpha(s^{\tau-1}))_{\tau \geq t}$  is not required to be constant, but each action in the tail is a one-period best reply to the limit belief, which is confirmed given said action.

**Proposition 3.** Assume that  $(\alpha, \mu)$  is rational. If  $(\alpha, \mu, \bar{\theta})$  is eventually consistent, then the stochastic process of actions and beliefs  $(\mathbf{a}_t^\alpha, \mu(\cdot|h_t^\alpha))$  converges to an SCE.

The intuition is as follows. Since the action set  $A$  is finite, after a certain amount of time, every action chosen by  $\alpha$  is played infinitely often. Under the stated assumptions, beliefs converge almost surely to a random limit  $\mu_{s^\infty}^\alpha$ . Thus, each action chosen by  $\alpha$  in the long run must be a myopic best reply to the limit belief. This holds because the updated beliefs converge and the

value of experimentation vanishes for actions played infinitely often. Proposition 1 implies that, for  $\bar{\theta}$ -almost every  $s^\infty$ , the limit belief  $\mu_{s^\infty}^\alpha$  assigns probability 1 to the models that induce the same probabilities over consequences as  $\bar{\theta}$  given the actions played in the long run. Therefore, for every action  $a^*$  chosen by  $\alpha$  in the long run,  $(a^*, \mu_{s^\infty}^\alpha, \bar{\theta})$  must be an SCE.

As noted above, the realized sequence of actions  $(\mathbf{a}_t^\alpha(s^{t-1}))$  need not converge unless there is a unique myopic best reply to the limit belief  $\mu_{s^\infty}^\alpha$ . If the myopic best reply is a unique action  $a^*$ , the action sequence  $(\mathbf{a}_t^\alpha(s^{t-1}))$  is eventually constant at  $a^*$  and  $(a^*, \mu_{s^\infty}^\alpha, \bar{\theta})$  is an SCE. Moreover, after a finite time, the agent chooses an action that maximizes the one-period value given the current belief (and not only limit one); that is, exploration (experimentation) becomes irrelevant, all that matters is exploitation. Formally:

**Proposition 4.** *Assume that  $(\alpha, \mu)$  is rational and that  $(\mathbf{a}_t^\alpha, \mu(\cdot|\mathbf{h}_t^\alpha))$  converges to an SCE on path  $s^\infty$ . If:*

$$\arg \max_{a \in A} \phi^{-1} \left( \int_{\Theta} \phi(R(a, \theta)) \mu_{s^\infty}^\alpha(d\theta) \right) = \{a^*\}$$

for some  $a^* \in A$ , there exists some  $\tau$  such that, for all  $t \geq \tau$ ,

$$\mathbf{a}_t^\alpha(s^{t-1}) = a^* \in \arg \max_{a \in A} \phi^{-1} \left( \int_{\Theta} \phi(R(a, \theta)) \mu(d\theta|\mathbf{h}_t^\alpha(s^{t-1})) \right).$$

It is important to stress that this convergence implies neither that the limit belief is the Dirac measure supported by the correct model, nor that the limit action is the objective myopic best reply. However, the limit pairs of beliefs and actions almost surely satisfy the standard properties of stochastic limits in the (expected utility) stochastic control literature. Indeed, the realization  $(a^*, \mu_{s^\infty}^\alpha)$  features the following:

- (Confirmed Beliefs):  $\mu_{s^\infty}^\alpha$  assigns probability 1 to the models that are observationally equivalent to the true  $\bar{\theta}$  given  $a^*$  (see, Proposition 1);
- (Subjective Myopic Best Reply): Even if the discount factor is strictly positive, the agent maximizes his one-period value. That is, exploitation prevails over exploration.

In contrast, in a Nash equilibrium, beliefs are correct (i.e.,  $\mu = \delta_{\bar{\theta}}$ ) and the action played is an objective myopic best reply. A sufficient condition for convergence to an SCE where the action played is objectively optimal is to have own-action independence.

**Corollary 3.** *Under own-action independence, if  $(\alpha, \mu, \bar{\theta})$  is eventually consistent and  $(\alpha, \mu)$  is rational, then the stochastic action process  $(\mathbf{a}_t^\alpha)$  converges to a Nash equilibrium action.*

Note that own-action independence guarantees convergence to a Nash equilibrium under observable payoffs, a maintained assumption in this work. If we relax this hypothesis, the stronger condition of perfect feedback is needed.

Our running example illustrates how the true model may remain unidentified in the limit.

**Example 7 (Act IV).** Consider the strategy  $\alpha^E$  of the previous acts. Again, recall that the DM has a uniform 3-point prior  $\mu$  with  $\text{supp } \mu = \{\theta^Y, \theta^{uni}, \theta^G\}$ . In Act II, we show that  $(\alpha^E, \mu, \bar{\theta})$

is consistent from period 1, whereas, in Act III, we have proved that with parameters  $\lambda = 1$  and  $\delta = 0.39$ ,  $(\alpha^E, \mu)$  is rational. We can show how our convergence result obtains in this specific case. Suppose that  $\bar{\theta} = \theta^Y$ . By Lemma 2, beliefs converge. In particular:

$$\mu_{s^\infty}^{\alpha^E} = \left( \mu_{s^\infty}^{\alpha^E}(\theta^{uni}), \mu_{s^\infty}^{\alpha^E}(\theta^Y), \mu_{s^\infty}^{\alpha^E}(\theta^G) \right) = \begin{cases} (\frac{1}{3}, \frac{1}{6}, \frac{1}{2}) & \text{if } s_1 \in \{B, G\}, \\ (0, 1, 0) & \text{if } s_1 = Y. \end{cases}$$

If experimentation is unsuccessful, the posterior of  $\mu$  lowers the weight of model  $\theta^Y$  relative to models  $\theta^{uni}$  and  $\theta^G$ ; thereupon, strategy  $\alpha^E$  recommends switching (and sticking) to black, so there is no additional updating. On the other hand, if the experimentation is successful, strategy  $\alpha^E$  prescribes sticking to yellow thereupon, and then the correct model  $\theta^Y$  is asymptotically identified.

If  $s_1 \in \{B, G\}$ , for every  $t > 1$ ,

$$a_t^{\alpha^E}(s^{t-1}) = b,$$

and  $(b, (1/3, 1/6, 1/2), \theta^Y)$  is the SCE that obtains in the limit. Note that in this case the DM will end up choosing an objectively sub-optimal action.

If  $s_1 = Y$ ,

$$a_t^{\alpha^E}(s^{t-1}) = y,$$

for every  $t > 1$ , and  $(y, (0, 1, 0), \theta^Y)$  is the SCE that obtains in the limit. It is immediate to see that these actions maximize one-period value for limit beliefs and that the distribution of probabilities over outcomes confirms them.

Finally, consider strategy  $\alpha^{NE}$ . In Act III, we argue that  $(\alpha^{NE}, \mu)$  is rational if the DM is sufficiently ambiguity averse. In this case, regardless of the correct model  $\bar{\theta} \in \Theta = \{\theta^Y, \theta^{uni}, \theta^G\}$ , we have almost sure convergence to an SCE from period 1. Indeed, the DM sticks to black from the first period onwards, and black is the myopic best reply to the confirmed prior  $\mu = (1/3, 1/3, 1/3)$ . However, note that if the correct model is  $\theta^Y$ , betting on black is objectively sub-optimal. ▲

## 6. Comparative dynamics for changes in ambiguity aversion

### 6.1. Certainty traps: the general case

Act III of our running example suggests that as ambiguity aversion increases, experimentation becomes less attractive. In this section we formalize this intuition. Actions that induce the same probabilities of payoffs under all the models that the DM deems possible are appealing under ambiguity aversion, but generate no new information about the underlying probability model. To obtain evidence on the correct model, the DM has to choose an action that will potentially induce a different probability measure over payoffs under the different models he deems possible—that is, he has to experiment. An ambiguity averse DM is inclined to avoid such ambiguous actions.

Fix an arbitrary belief  $\nu$ ; we say that action  $a$  is  $\nu$ -unambiguous if  $F(a, \theta) = F(a, \theta')$  for every  $\theta, \theta' \in \text{supp } \nu$ . In words, an action  $a$  is unambiguous given the DM's beliefs if all models entertained by the DM assign the same probabilities to outcomes given  $a$ . Otherwise, we say that  $a$  is  $\nu$ -ambiguous. The next proposition shows that, if a strategy  $\alpha$  is rational given  $\mu$  for an ambiguity-neutral DM and prescribes unambiguous actions, then every strategy  $\beta$  that is rational given  $\mu$  for a strictly ambiguity averse DM must also prescribe unambiguous actions. In other

words, if experimentation is not rational for an ambiguity neutral DM, then it cannot be rational for an ambiguity averse DM with the same beliefs. For notational simplicity, in this section we will assume that the utility function  $u$  is injective.

**Proposition 5.** *Assume that  $(\alpha, \mu)$  is rational under ambiguity neutrality and  $(\beta, \mu)$  is rational under strict ambiguity aversion. For every belief  $v$  in the belief-range of  $\mu$ , if  $\alpha(v)$  is  $v$ -unambiguous, then the same holds for  $\beta(v)$ .*

In what follows, we restrict our attention to cases where there is a *unique  $\mu$ -ambiguous action*.<sup>28</sup> Although this assumption is restrictive, it encompasses interesting stochastic control problems such as two-armed bandits with a safe arm. Moreover, this restriction parallels the one needed for important comparative statics results in the case of choice under risk.<sup>29</sup>

The next proposition establishes that, given a prior  $\mu$ , the sequence of actions and posteriors induced by a strategy  $\beta$  that is rational under strict ambiguity aversion will almost surely converge faster to an SCE (given the true model) than the sequence of actions and posterior beliefs corresponding to a strategy  $\alpha$  that is rational under ambiguity neutrality. The intuition is as follows. Under the stated assumptions, the decision problem amounts to deciding how long to experiment, choosing the unique ambiguous action, say  $a^*$ . Assuming that convergence occurs, there are two possibilities: Either the DM never stops experimenting, or he stops at some finite time  $t$ , choosing the best unambiguous action thereafter. In the first case, convergence to an SCE occurs (typically) at infinity. In the second case, it occurs at the stopping time. By Proposition 5, if an ambiguity neutral DM prefers to stop at time  $t$ , then a strictly ambiguity averse DM prefers to stop as well. Alternatively, if an ambiguity neutral DM chooses  $a^*$  forever, an ambiguity averse DM either does the same, or starts choosing the best unambiguous action from some period  $t$  onwards and is henceforth “trapped” in a self-confirming equilibrium.

**Proposition 6.** *Assume that (i) there is a unique  $\mu$ -ambiguous action  $a^*$ , (ii)  $(\alpha, \mu)$  is rational under ambiguity neutrality and  $(\beta, \mu)$  is rational under strict ambiguity aversion, and (iii)  $(\alpha, \mu, \bar{\theta})$  is eventually consistent. Then,  $\bar{\theta}$ -a.s., the action-belief process  $(\beta(\mu(\cdot|\mathbf{h}_t^\beta)), \mu(\cdot|\mathbf{h}_t^\beta))$  converges to an SCE at least as fast as the action-belief process  $(\alpha(\mu(\cdot|\mathbf{h}_t^\alpha)), \mu(\cdot|\mathbf{h}_t^\alpha))$ .*

The next proposition shows that an ambiguity averse DM is less likely than an ambiguity neutral DM to eventually play the Nash equilibrium (i.e., objectively optimal) action.

**Proposition 7.** *Assume that (i) there is a unique  $\mu$ -ambiguous action  $a^*$ , (ii)  $(\alpha, \mu)$  and  $(\beta, \mu)$  are, respectively, rational under ambiguity neutrality and strict ambiguity aversion, and (iii)  $(\alpha, \mu, \bar{\theta})$  and  $(\beta, \mu, \bar{\theta})$  are eventually consistent. Then,  $\bar{\theta}$ -a.s., if  $(\mathbf{a}_t^\beta)$  converges to a Nash Equilibrium action, so does  $(\mathbf{a}_t^\alpha)$ .*

We illustrate the previous results in the classical setup introduced by Rothschild (1974) where a monopolist trades-off exploration of the demand curve for his good against exploitation using the price that (subjectively)-maximizes one period profit.

<sup>28</sup> We refer to Battigalli et al. (2019) for an in-depth analysis of the role of this assumption.

<sup>29</sup> Specifically, the seminal contribution by Arrow (1971) shows that, when there is a unique risky asset, the amount invested in such an asset is decreasing in the risk aversion of the DM. Additional reasons to focus on this case are provided by Proposition 8.

**Example 8.** A monopolist who is uncertain about the demand for its product faces a new customer each period. The cost of producing a unit is  $c > 0$ . The monopolist can charge either a low price,  $p_L$ , or a high price,  $p_H$ , where  $p_H > p_L > c \geq 0$ . Each new customer has a reservation price in  $\{p_H, p_L, 0\}$ , and he buys the product if and only if his reservation price is weakly larger than the ask price. If the price is set to  $p_i$  and a sale is made, the monopolist’s profit is  $p_i - c$ ; otherwise, the profit is 0. Thus, we are considering build-to-order production. Here,  $A = \{p_L, p_H\}$ ,  $M = \{0, 1\}$ ,  $S = \{p_H, p_L, 0\}$ ,

$$f(a, s) = \begin{cases} 1 & a \leq s, \\ 0 & a > s, \end{cases}$$

and:

Payoff $u(p, m)$	$m = 1$	$m = 0$
$p = p_L$	$p_L - c$	0
$p = p_H$	$p_H - c$	0

Suppose that according to his prior  $\mu$ , the monopolist believes that  $\theta = (\theta(p_H), \theta(p_L), \theta(0))$  is either  $\theta_1 := (0.8, 0.1, 0.1)$ , or  $\theta_2 := (0, 0.9, 0.1)$ , and that these two models are equally likely. Here, the monopolist is certain that if he posts the low price he sells with probability 0.9, but he is uncertain about the selling probability at a high price, which is—therefore—a  $\mu$ -ambiguous action.

Moreover, suppose that the correct model is  $\theta_1$  and that  $0.9(p_L - c) < 0.8(p_H - c)$ . Our previous results imply that an ambiguity averse monopolist will stop experimentation with the high price earlier and that he will be more likely to be trapped in the (objectively) suboptimal SCE where he posts the low price (the  $\mu$ -unambiguous action). ▲

At first sight, the previous results may seem surprising. Indeed, if a DM is ambiguity averse, why does he not experiment more, so as to eliminate (or reduce) the uncertainty about the true model? This reasoning tacitly relies on a different notion of experimentation. It is true that, typically, an ambiguity averse DM is willing to pay more to eliminate model uncertainty (see e.g., Theorem 2 in Anderson, 2012). However, in our active-learning setting, the DM cannot simply buy information about the true model; learning happens only when actions with ambiguous probabilities of consequences are chosen. Since an ambiguity averse DM dislikes those actions, he will end up resolving less ambiguity than his ambiguity neutral counterpart.

The results above relate to the findings in Anderson (2012). On the theoretical side, his Theorem 1 for two-armed bandits with a safe arm strictly relates to our Proposition 5. The main difference is that Anderson (2012) implicitly assumes the possibility to commit to a strategy. Indeed, the Gittins indices used in that paper characterize the *ex-ante* optimal strategy for a decision maker, that is, the strategy that maximizes the value at the initial history. However, like us, he assumes the DM performs Bayesian updating, a feature that paired with ambiguity aversion induces dynamically inconsistent preferences.

On the experimental side, the theoretical predictions of our model are consistent with the findings presented in Anderson (2012): The behavior of a Subjective Expected Utility maximizer cannot explain joint data about willingness to pay for information about the stochastic process characterizing the ambiguous arm and the amount of experimentation that is performed. In particular, the resulting experimentation is too low, which is the prediction of our model under ambiguity aversion.



Another reason to focus on the case of a unique ambiguous action is to illustrate how our analysis adds to BCMM. Indeed, the following proposition shows that, when there is a unique ambiguous action, the set of SCE actions is *invariant* with respect to the (positive) degree of ambiguity aversion captured by the (concave) function  $\phi$ . Yet, as Propositions 6 and 7 show, ambiguity aversion has “dynamic” effects on the persistence of experimentation and the distribution of long-run outcomes.

**Proposition 8.** *Assume that there is a unique  $\mu$ -ambiguous action  $a^*$ . For every concave and strictly increasing  $\phi, \phi'$  and every action  $\bar{a}$ , if  $(\bar{a}, \mu, \theta)$  is an SCE under ambiguity attitudes  $\phi$ , then, for some belief  $\mu'$ ,  $(\bar{a}, \mu', \theta)$  is an SCE under ambiguity attitudes  $\phi'$ .*

Our running example illustrates. As the unique ambiguous action is to bet on yellow, Proposition 8 implies that the monotonicity result of BCMM (that the SCE set is weakly increasing in the degree of ambiguity aversion) holds vacuously: The set of equilibrium actions is not affected by ambiguity attitudes. Nonetheless, if the true model is  $\theta_Y$ , the example shows that beliefs converge to the true model with positive probability under ambiguity neutrality, while the process of actions and beliefs is trapped in a non-Nash SCE under ambiguity aversion.

The next example illustrates the relevance of the assumption of a unique ambiguous action.

**Example 9 (Multiple ambiguous actions).** There are four possible states,  $S = \{g, \bar{g}, b, \bar{b}\}$ , and two possible models in  $\Theta$ , the good model  $\theta_g$  and the bad model  $\theta_b$  defined as follows:

$$\theta_g(g) = 0.9 = \theta_b(b) \text{ and } \theta_g(\bar{g}) = 0.1 = \theta_b(\bar{b}).$$

The DM has three actions: He can bet aggressively (action  $a$ ), bet conservatively (action  $c$ ), or not bet at all (action  $n$ ). The feedback received by the DM is his monetary payoff. We also assume risk-neutrality: For all  $\bar{a} \in A, m \in M, u(a, m) = m$ . Feedback and payoffs are summarized in the following table:

$f$	$g$	$\bar{g}$	$b$	$\bar{b}$
$a$	10	0	0	10
$c$	5	5	4	4
$n$	4.2	4.2	4.2	4.2

Therefore,

$$R(a, \theta_g) = 9, R(a, \theta_b) = 1, R(c, \theta_g) = 5 \text{ and } R(c, \theta_b) = 4.$$

For simplicity, suppose  $\delta = 0$  and  $\mu(\theta_g) = \mu(\theta_b) = 1/2$ . In this case, in the first period, an ambiguity neutral DM bets aggressively ( $a$ ):

$$\frac{R(a, \theta_g) + R(a, \theta_b)}{2} = 5 > 4.5 = \frac{R(c, \theta_g) + R(c, \theta_b)}{2}.$$

Similarly, one can check that a DM with intermediate ambiguity attitudes (i.e.,  $\lambda = 1$ ) and with the same belief  $\mu$  bets conservatively. Now, suppose that the true model is  $\theta_g$ . First note that, since  $c$  perfectly reveals the model, the ambiguity averse DM discovers the true model at the end of the first period and starts to bet aggressively from the second period. In other words, convergence to Nash equilibrium (since  $a$  is the objectively optimal action under  $\theta_g$ ) happens in one period and with probability 1. Now, consider the ambiguity neutral DM. He bets on  $a$  in the

first period, and with probability  $0.1 = \theta_g(\bar{g})$  he receives message  $m = 0$  and updates his beliefs to:

$$\mu(\theta_g | (a, 0)) = 0.1.$$

With this, from the second period onwards, he stops betting (i.e., he chooses  $n$ ) and remains in the dark. Thus, there is at least probability 0.1 that the ambiguity neutral DM will be trapped in a non-Nash SCE. ▲

In this example there is a misalignment between the most informative action ( $c$ ) and the most ambiguous one ( $a$ ). Therefore, by avoiding the most ambiguous action, an ambiguity averse DM quickly gathers information about the true model. We believe that this kind of misalignment is unlikely to arise in applications. Still, we conjecture that the results of this section can be extended to the case of multiple ambiguous actions by first giving an adequate definition of an ambiguity order (see, e.g., Jewitt and Mukerji, 2017), and then imposing a condition of comonotonicity between informativeness and ambiguity of actions.<sup>30</sup>

### 6.2. Certainty traps: myopic decision makers

Sharper versions of our comparative dynamics results can be provided when the DM is myopic, i.e., when  $\delta = 0$ . This follows from the fact that, in the present framework where each period consists of a one-stage decision problem, a myopic DM is not vulnerable to dynamic inconsistencies. Despite this simplification, the behavior over time of a myopic DM evolves in interesting ways as he gathers information about the true stochastic process. Indeed, several models of learning in games use the assumption of myopic players (see, e.g., Fudenberg and Kreps, 1995, and Fudenberg and Levine, 1998). We show that when the DM is myopic our comparative statics results hold for the entire spectrum of ambiguity attitudes.

For the rest of this section, we assume that the value is given by:

$$V(\alpha, \mu | h_t) := \phi^{-1} \left( \int_{\Theta} \phi(R(\alpha(h_t), \theta)) \mu(d\theta | h_t) \right).$$

We let  $\phi'$  be a strictly increasing and concave transformation of  $\phi$ , i.e., we assume that the DM with ambiguity attitudes  $\phi'$  is *strictly more ambiguity averse* than the one with ambiguity attitudes  $\phi$ .

**Proposition 9.** *Let  $(\alpha, \mu)$  be rational under ambiguity attitudes  $\phi$  and let  $(\beta, \mu)$  be rational under ambiguity attitudes  $\phi'$ . Then, for every belief  $\nu$  in the belief-range of  $\mu$ , if  $\alpha(\nu)$  is  $\nu$ -unambiguous, the same holds for  $\beta(\nu)$ .*

Similarly, the speed of convergence to the SCE is monotone for the entire spectrum of ambiguity attitudes.

**Proposition 10.** *Assume that there is a unique  $\mu$ -ambiguous action  $a^*$ ; let  $(\alpha, \mu)$  be rational under ambiguity attitudes  $\phi$  and let  $(\beta, \mu)$  be rational under ambiguity attitudes  $\phi'$ ; furthermore, assume that  $(\alpha, \mu, \bar{\theta})$  is eventually consistent. Then,  $\bar{\theta}$ -a.s., the action-belief process*

<sup>30</sup> We can prove this for the case of a myopic DM ( $\delta = 0$ ).

$(\beta(\mu(\cdot|\mathbf{h}_t^\beta)), \mu(\cdot|\mathbf{h}_t^\beta))$  converges to an SCE at least as fast as the action-belief process  $(\alpha(\mu(\cdot|\mathbf{h}_t^\alpha)), \mu(\cdot|\mathbf{h}_t^\alpha))$ .

As a consequence, the probability of converging to a Nash equilibrium is decreasing in ambiguity aversion as well.

**Proposition 11.** *Assume that there is a unique  $\mu$ -ambiguous action  $a^*$ . Let  $(\alpha, \mu)$  and  $(\beta, \mu)$  be rational under ambiguity attitudes  $\phi$  and  $\phi'$ , respectively, and let  $(\alpha, \mu, \bar{\theta})$  and  $(\beta, \mu, \bar{\theta})$  be eventually consistent. Then,  $\bar{\theta}$ -a.s., if  $(\mathbf{a}_t^\beta)$  converges to a Nash Equilibrium action, so does  $(\mathbf{a}_t^\alpha)$ .*

## 7. Concluding remarks

The concept of self-confirming equilibrium with standard expected utility maximizing agents has been given a rigorous learning foundation. We note that the literature on stochastic control problems implicitly addresses this issue, showing that the behavior and beliefs of an ambiguity neutral agent, who faces an unknown i.i.d. process of states affecting the outcome of his actions, almost surely converges to what we call an SCE.<sup>31</sup> As for games against other agents, convergence cannot be taken for granted; but if it occurs, the limit point must be an SCE (e.g., Fudenberg and Levine, 1993, and Fudenberg and Kreps, 1995).

This learning foundation cannot be mechanically applied to the case of non-neutral ambiguity attitudes. Ambiguity averse agents typically have dynamically inconsistent preferences over strategies, and dynamic inconsistency prevents us from applying standard dynamic programming techniques. Given such difficulties, to derive results and insights about convergence to SCE under ambiguity aversion, we focus on the case of repeated play against nature, assuming that the decision maker is sophisticated and thus takes future incentives into account as he chooses actions in earlier periods. With this, we obtain a result of convergence to SCE under ambiguity aversion (Proposition 3).

We point out that, in several interesting problems, the set of SCE actions is independent of ambiguity attitudes (Proposition 8). Yet, ambiguity aversion affects the dynamics: Higher ambiguity aversion tends to decrease experimentation and therefore makes convergence to Nash equilibrium (best reply to the correct model) less likely. In particular, we show that ambiguity aversion may make it more likely that the agent falls into a suboptimal “certainty trap” whereby he keeps choosing an unambiguous action from which he cannot learn, which leads him to settle faster and prevents him from identifying the objectively optimal action (Propositions 5 through 11).

The adoption of the smooth ambiguity model allows us to separate ambiguity attitudes, a personal feature of the DM, from the perception of ambiguity, which depends on his beliefs about statistical models and therefore changes as new evidence accumulates. We model the process of updating such beliefs in a standard Bayesian fashion, as we regard the chain rule, hence Bayesian updating, as part of rational cognition (cf. Section 3 and Battigalli et al., 2019). Yet, the decision-theoretic literature does not take a clear stand on whether ambiguity averse players should update beliefs according to the standard rules of conditional probabilities (see, for example, Epstein and Schneider, 2007; Hanany and Klibanoff, 2009). We remark that one may conduct an analysis

<sup>31</sup> See Easley and Kiefer (1988) and Section 5 of the working paper version of this article for a detailed analysis of the connection with their work.

similar to ours by considering a DM who uses the Hanany and Klibanoff (2009) updating rule for beliefs. In particular, we can prove that every SCE is stable according to their updating rule. Therefore, the main message of Section 6 that “ambiguity aversion makes convergence to a non-Nash equilibrium more likely” still holds.<sup>32</sup>

We can give a game-theoretic interpretation of our analysis within a population-game scenario. In this setting, the DM recognizes to be unable to influence the actions of future co-players. Nevertheless, experimentation is valuable for him, since a better understanding of the correct distribution of behaviors in co-players’ populations may allow him to select a better strategy in the following periods. Fudenberg and Levine (1993) put forward a model of this kind. The main difference with their work is that they consider an overlapping generations model with finitely lived agents. Since we assume an infinite horizon, we have to slightly modify our model by introducing a constant probability of death to embed our analysis in an overlapping generations model, as in Fudenberg and He (2018).<sup>33</sup>

**Acknowledgments**

We thank Federico Bobbio, Roberto Corrao, Carlo Cusumano, Nicodemo De Vito, Ignacio Esponda, Francesco Fabbri, Drew Fudenberg, Filippo Massari, David Ruiz Gomez, Ran Spiegler, Muhamet Yildiz, and anonymous referees for useful comments and suggestions. Special thanks go to Simone Cerreia-Vioglio and Fabio Maccheroni with whom this project started. Financial support from the European Research Council (advanced grant 324219), Guido Cazzavillan Scholarship, and the AXA Research Fund is gratefully acknowledged.

**Appendix A. Proofs and related material**

For our proofs, it is often convenient to use the notation  $R(a, \theta) = \sum_{s \in S} r(a, s) \theta(s)$ , where  $r : A \times S \rightarrow \mathbb{R}$  is the payoff (or reward) function  $r(a, s) := (u_a \circ f_a)(s)$ .

**Lemma 5.** *If we endow  $\Delta(\Theta)$  with the topology of weak convergence of measures then, for every  $a \in A$ , the functional  $\bar{V}(a, \cdot)$  is continuous.*

**Proof.** Note that  $R(a, \cdot)$  is a bounded function, since  $|R(a, \cdot)| \leq \max_{s \in S} |r(a, s)|$ . Moreover, it is an affine function on a finite dimensional space, so it is continuous. Thus,  $\int_{\Theta} \phi(R(a, \theta))(\cdot) (d\theta)$  is continuous. Since  $\phi$  is strictly increasing and continuous on the interval

$$\left[ \min_{s \in S} r(a, s), \max_{s \in S} r(a, s) \right],$$

$\phi^{-1}$  is continuous as well. Since  $\bar{V}(a, \cdot) = \phi^{-1} \circ \int_{\Theta} \phi(R(a, \theta))(\cdot) (d\theta)$ , the result follows.  $\square$

*A.1. Models and learning*

Instrumental for the following proofs is the correspondence  $\iota_t^\alpha : S^{t-1} \rightarrow \sigma(\mathbf{h}_t^\alpha)$  defined by:

$$\iota_t^\alpha \left( s^{t-1} \right) := I \left( \mathbf{h}_t^\alpha \left( s^{t-1} \right) \right) \times S^\infty = \left\{ \bar{s}^\infty \in S^\infty : \mathbf{h}_t^\alpha \left( \bar{s}^{t-1} \right) = \mathbf{h}_t^\alpha \left( s^{t-1} \right) \right\}. \tag{6}$$

<sup>32</sup> A formal statement and proof are available by request.

<sup>33</sup> See also Blanchard (1985).

We can regard  $\iota_t^\alpha$  as the identification correspondence determined by  $\alpha$  at time  $t$ . This correspondence models the information about state histories which is available ex-ante at time  $t$  to a DM who is acting according to strategy  $\alpha$ . Clearly,  $s^{t-1} \times S^\infty \in \iota_t^\alpha (s^{t-1})$ , and so the correspondence induces a partition of  $S^{t-1}$ . We have *perfect (state) identification* under  $\alpha$  when  $\iota_t^\alpha (s^{t-1}) = \{s^{t-1}\} \times S^\infty$  for each  $s^{t-1}$  and each  $t > 1$ ; in this case, the DM knows the actual past history  $s^{t-1}$ . Otherwise, we have partial identification. This dependence on  $\alpha$  of the identification correspondence plays a key role in our results. Of course, there is no such dependence under own-action independence of feedback, in which case we can write  $\iota_t (s^{t-1})$ ; in particular, under perfect feedback,  $\iota_t (s^{t-1}) = \{s^{t-1}\} \times S^\infty$ .

**Proof of Lemma 1.** Fix  $s^t$  with  $p_{\bar{\theta}} (s^t) > 0$ . Note that  $s^t \in \iota_{t+1}^\alpha (s^t) \subseteq \iota_t^\alpha (s^{t-1})$ ; thus,  $p_{\bar{\theta}} (\iota_t^\alpha (s^{t-1})) \geq p_{\bar{\theta}} (\iota_{t+1}^\alpha (s^t)) \geq p_{\bar{\theta}} (s^t) > 0$ . Let  $\theta \in \Theta_t^{\alpha, \mu} (S^\infty)$ ; By definition  $\theta \in \text{supp } \mu (\cdot | \mathbf{h}_t^\alpha (s^{t-1}))$ , and so  $p_\theta (\iota_t^\alpha (s^{t-1})) > 0$ . We want to show that  $\theta \in \Theta_{t+1}^{\alpha, \mu} (S^\infty)$ . To this end, notice that if  $E \cap \iota_{t+1}^\alpha (s^t) = \emptyset$ ,  $p_\theta^\alpha (E | \mathbf{h}_{t+1}^\alpha (s^t)) = p_{\bar{\theta}}^\alpha (E | \mathbf{h}_{t+1}^\alpha (s^t)) = 0$ . Therefore, it is enough to show that  $p_\theta^\alpha (\cdot | \mathbf{h}_{t+1}^\alpha (s^t)) = p_{\bar{\theta}}^\alpha (\cdot | \mathbf{h}_{t+1}^\alpha (s^t))$  for  $E \in \sigma(\mathbf{h}^\alpha)$  with  $E \subseteq \iota_{t+1}^\alpha (s^t)$ . Fix such an  $E$ . Since  $\theta \in \Theta_t^{\alpha, \mu} (S^\infty)$ , then:

$$\frac{p_\theta (E)}{p_\theta (\iota_t^\alpha (s^{t-1}))} = \frac{p_{\bar{\theta}} (E)}{p_{\bar{\theta}} (\iota_t^\alpha (s^{t-1}))}; \quad \frac{p_\theta (\iota_{t+1}^\alpha (s^t))}{p_\theta (\iota_t^\alpha (s^{t-1}))} = \frac{p_{\bar{\theta}} (\iota_{t+1}^\alpha (s^t))}{p_{\bar{\theta}} (\iota_t^\alpha (s^{t-1}))}.$$

The second equality implies  $p_\theta (\iota_{t+1}^\alpha (s^t)) > 0$ . Since:

$$\frac{p_\theta (E)}{p_\theta (\iota_{t+1}^\alpha (s^t))} \frac{p_\theta (\iota_{t+1}^\alpha (s^t))}{p_\theta (\iota_t^\alpha (s^{t-1}))} = \frac{p_{\bar{\theta}} (E)}{p_{\bar{\theta}} (\iota_{t+1}^\alpha (s^t))} \frac{p_{\bar{\theta}} (\iota_{t+1}^\alpha (s^t))}{p_{\bar{\theta}} (\iota_t^\alpha (s^{t-1}))},$$

it follows that:

$$p_\theta^\alpha (E | \mathbf{h}_{t+1}^\alpha (s^t)) = \frac{p_\theta (E)}{p_\theta (\iota_{t+1}^\alpha (s^t))} = \frac{p_{\bar{\theta}} (E)}{p_{\bar{\theta}} (\iota_{t+1}^\alpha (s^t))} = p_{\bar{\theta}}^\alpha (E | \mathbf{h}_{t+1}^\alpha (s^t)).$$

Hence,  $p_\theta^\alpha (\cdot | \mathbf{h}_{t+1}^\alpha (s^t)) = p_{\bar{\theta}}^\alpha (\cdot | \mathbf{h}_{t+1}^\alpha (s^t))$ . Since  $\theta \in \text{supp } \mu (\cdot | \mathbf{h}_t^\alpha (s^{t-1}))$  and  $p_\theta (\iota_{t+1}^\alpha (s^t)) > 0$ , it follows that  $\theta \in \text{supp } \mu (\cdot | \mathbf{h}_{t+1}^\alpha (s^t))$ , hence  $\theta \in \Theta_{t+1}^{\alpha, \mu} (S^\infty)$ .  $\square$

**Lemma 6.** For every  $\hat{\theta} \in \Theta$ , the process  $(\mu (\hat{\theta} | \mathbf{h}_t^\alpha))$  is a uniformly bounded martingale in  $(S^\infty, \mathcal{B}(S^\infty), p_\mu)$ .

**Proof.** Uniform boundedness is immediate from the fact that the process is  $[0, 1]$ -valued. For every  $t$ , we want to show that, for all  $k \in \mathbb{R}$  such that  $p_\mu ([\mu (\hat{\theta} | \mathbf{h}_{t-1}^\alpha) = k]) > 0$  we have:

$$\mathbb{E}_{p_\mu} \left[ \mu (\hat{\theta} | \mathbf{h}_t^\alpha) \mid \mu (\hat{\theta} | \mathbf{h}_{t-1}^\alpha) = k \right] = k.$$

Let  $h_{t-1}$  be an arbitrary element of  $H_{t-1}$  such that  $\mu (\hat{\theta} | h_{t-1}) = k$  and  $p_\mu (h_{t-1}) > 0$ . By the Law of Iterated Expectations, it is enough to prove that:

$$\mathbb{E}_{p_\mu} \left[ \mu (\hat{\theta} | \mathbf{h}_t^\alpha) \mid \mathbf{h}_{t-1}^\alpha = h_{t-1} \right] = k.$$

Recall that the Bayes map yields:

$$\Delta(\Theta) \times A \times M \rightarrow \Delta(\Theta)$$

$$(\mu, a, m) \mapsto B(\mu, a, m)(\theta) = \frac{F(a, \theta)(m) \mu(\theta)}{\sum_{\theta' \in \text{supp } \mu} F(a, \theta')(m) \mu(\theta')}$$

for each  $m$  deemed possible according to  $\mu$  given action  $a$ , that is, each  $m$  such that the denominator is positive. Define:

$$M(h_{t-1}) = \left\{ m \in M : \sum_{\theta' \in \text{supp } \mu(\cdot|h_{t-1})} F(a, \theta')(m) \mu(\theta'|h_{t-1}) > 0 \right\}.$$

With this,

$$\begin{aligned} & \mathbb{E}_{p_\mu} \left[ \mu(\hat{\theta} | \mathbf{h}_t^\alpha) | \mathbf{h}_{t-1}^\alpha = h_{t-1} \right] \\ &= \sum_{m \in M(h_{t-1})} p_\mu \left[ \mathbf{m}_{t-1}^\alpha = m | \mathbf{h}_{t-1}^\alpha = h_{t-1} \right] B(\mu(\cdot|h_{t-1}), \alpha(h_{t-1}), m)(\hat{\theta}) \\ &= \sum_{m \in M(h_{t-1})} \sum_{\theta \in \text{supp } \mu} \frac{\mu(\theta) p_\theta \left( \left\{ s^\infty \in I(h_{t-1}) : s_t \in f_{\alpha(h_{t-1})}^{-1}(m) \right\} \right)}{\sum_{\theta' \in \text{supp } \mu} \mu(\theta') p_{\theta'}(I(h_{t-1}))} \\ & \quad \times B(\mu(\cdot|h_{t-1}), \alpha(h_{t-1}), m)(\hat{\theta}) \\ &= \sum_{m \in M(h_{t-1})} \sum_{\theta \in \text{supp } \mu} \frac{\mu(\theta) p_\theta(I(h_{t-1})) F(\alpha(h_{t-1}), \theta)(m)}{\sum_{\theta' \in \text{supp } \mu} \mu(\theta') p_{\theta'}(I(h_{t-1}))} \\ & \quad \times B(\mu(\cdot|h_{t-1}), \alpha(h_{t-1}), m)(\hat{\theta}) \\ &= \sum_{m \in M(h_{t-1})} \sum_{\theta \in \text{supp } \mu} F(\alpha(h_{t-1}), \theta)(m) \mu(\theta|h_{t-1}) \\ & \quad \times \frac{k F(\alpha(h_{t-1}), \hat{\theta})(m)}{\sum_{\theta' \in \text{supp } \mu(\cdot|h_{t-1})} \mu(\theta'|h_{t-1}) F(\alpha(h_{t-1}), \theta')(m)} \\ &= \sum_{m \in M(h_{t-1})} k F(\alpha(h_{t-1}), \hat{\theta})(m) \frac{\sum_{\theta \in \text{supp } \mu} \mu(\theta|h_{t-1}) F(\alpha(h_{t-1}), \theta)(m)}{\sum_{\theta' \in \text{supp } \mu(\cdot|h_{t-1})} \mu(\theta'|h_{t-1}) F(\alpha(h_{t-1}), \theta')(m)} \\ &= \sum_{m \in M(h_{t-1})} k F(\alpha(h_{t-1}), \hat{\theta})(m) = k, \end{aligned}$$

where the first equality comes from the definition of expected value, the second by the definition of  $p_\mu$ , the third from the fact that the environment is i.i.d., the fourth from the definition of Bayes map, and the fifth from rearranging the terms.

For the last equality notice that either  $\mu(\hat{\theta}|h_{t-1}) = k = 0$ , and the equality is trivial, or  $\mu(\hat{\theta}|h_{t-1}) > 0$ , and therefore  $m \notin M(h_{t-1})$  implies that  $F(\alpha(h_{t-1}), \hat{\theta})(m) = 0$ .  $\square$

**Lemma 7.** For every  $\hat{\theta} \in \text{supp } \mu$ , the process  $(\mu(\cdot|h_t^\alpha))_{t \in \mathbb{N}_0}$  converges  $\hat{\theta}$ -a.s. to a random limit  $\mu_{s^\infty}^\alpha$ .

**Proof.** By Lemma 6, the stochastic process  $(\mu(\cdot|\mathbf{h}_t^\alpha))_{t \in \mathbb{N}_0}$  is a uniformly-bounded martingale in  $(S^\infty, \mathcal{B}(S^\infty), p_\mu)$ . By the Martingale Convergence Theorem (Billingsley, 2012, Theorem 35.5), the limit random variable  $\mu_{s^\infty}^\alpha$  exists  $p_\mu$ -almost surely. This means that there exists a set  $E \in \mathcal{B}(S^\infty)$  such that  $p_\mu(E) = 1$ , so  $p_\mu(S^\infty \setminus E) = 0$ , and:

$$\lim_{t \rightarrow +\infty} \mu(\cdot|\mathbf{h}_t^\alpha(s^\infty)) = \mu_{s^\infty}^\alpha$$

for every  $s^\infty \in E$ . Note that, since  $\mu(\hat{\theta}) > 0$ ,

$$\sum_{\theta \in \text{supp } \mu} p_\theta(S^\infty \setminus E) \mu(\theta) = p_\mu(S^\infty \setminus E) = 0$$

implies that  $p_{\hat{\theta}}(S^\infty \setminus E) = 0$  and so  $p_{\hat{\theta}}(E) = 1$ .  $\square$

**Proof of Lemma 2.** Let  $T(\alpha, \mu, \bar{\theta}) = T$ . If  $\bar{\theta} \in \text{supp } \mu$  the result follows immediately from Lemma 7. If  $\bar{\theta} \notin \text{supp } \mu$ , by consistency the set  $O := \{s^\infty : \Theta_T^{\alpha, \mu}(s^\infty) \neq \emptyset\}$  has  $p_{\bar{\theta}}$ -probability 1. Define the set  $E^* := \{s^\infty : \exists v \in \Delta(\Theta) : \lim_{t \rightarrow +\infty} \mu(\cdot|\mathbf{h}_t^\alpha(s^\infty)) = v\}$ . For every  $s^\infty \in O$ , we can find some  $\theta(s^\infty) \in \text{supp } \mu(\cdot|\mathbf{h}_T^\alpha(s^{T-1}))$  such that  $p_{\bar{\theta}}(A) = p_{\theta(s^\infty)}(A)$  for every  $A \in \sigma_{\geq T}(\mathbf{h}_T^\alpha(s^{T-1}))$ . Repeating the arguments of Lemmata 6 and 7 with the probability space  $(I(\mathbf{h}_T^\alpha(s^{T-1})), \sigma_{\geq T}(\mathbf{h}_T^\alpha(s^{T-1})), \mu(\cdot|\mathbf{h}_T^\alpha(s^{T-1})))$  in place of  $(S^\infty, \mathcal{B}(S^\infty), \mu)$ , and  $p_{\theta(s^\infty)}(\cdot|\mathbf{h}_T^\alpha(s^{T-1}))$  in place of  $p_{\hat{\theta}}$ , we obtain that there exists a set  $E \in \sigma_{\geq T}(\mathbf{h}_T^\alpha(s^{T-1}))$  such that  $p_{\theta(s^\infty)}(E|\mathbf{h}_T^\alpha(s^{T-1})) = 1$ , and:

$$\forall s^\infty \in E, \exists v \in \Delta(\Theta) : \lim_{t \rightarrow +\infty} \mu(\cdot|\mathbf{h}_t^\alpha(s^\infty)) = v.$$

Moreover,

$$p_{\bar{\theta}}(E^*|\mathbf{h}_T^\alpha(s^{T-1})) \geq p_{\bar{\theta}}(E|\mathbf{h}_T^\alpha(s^{T-1})) = p_{\theta(s^\infty)}(E|\mathbf{h}_T^\alpha(s^{T-1})) = 1,$$

and therefore  $p_{\bar{\theta}}(E^*|\mathbf{h}_T^\alpha(s^{T-1})) = 1$ . Since  $s^\infty$  was an arbitrary element of the set  $O$  with  $p_{\bar{\theta}}(O) = 1$ , by the Law of Iterated Expectations (see, e.g., 9.7i in Williams, 1991)  $p_{\bar{\theta}}(E^*) = 1$ .  $\square$

For every path  $s^\infty$ , denote by  $\mathbf{a}_\infty^\alpha(s^\infty)$  the set of actions played infinitely often (i.o.) under strategy  $\alpha$  along this path. Also for every  $\theta \in \Theta$  let  $E_\theta^\alpha$  denote the event that an action and a state that occur i.o. also occur jointly i.o. More formally,

$$E_\theta^\alpha = \left\{ s^\infty : \forall (a, \bar{s}) \in \mathbf{a}_\infty^\alpha(s^\infty) \times \text{supp } \hat{\theta}, (\alpha(\mathbf{h}_{t-1}^\alpha(s^\infty)), s_t) = (a, \bar{s}) \text{ i.o.} \right\}.$$

**Lemma 8.** For every  $\hat{\theta} \in \text{supp } \mu$ , we have  $p_{\hat{\theta}}(E_\theta^\alpha) = 1$ .

**Proof.** For every  $(a, \bar{s}) \in \mathbf{a}_\infty^\alpha(s^\infty) \times \text{supp } \hat{\theta}$ , denote by  $E(a, \bar{s}, n) \subseteq S^\infty$  the set of sequences  $s^\infty$  such that  $a \in \mathbf{a}_\infty^\alpha(s^\infty)$  but  $(\alpha(\mathbf{h}_{t-1}^\alpha(s^\infty)), s_t) \neq (a, \bar{s})$  for every  $t \geq n$ . We have:

$$S^\infty \setminus E_\theta^\alpha \subseteq \bigcup_{a \in \mathbf{a}_\infty^\alpha(s^\infty)} \bigcup_{\bar{s} \in \text{supp } \hat{\theta}} \bigcup_{n \in \mathbb{N}} E(a, \bar{s}, n).$$

In turn, for every  $k \in \mathbb{N}$ ,

$$E(a, \bar{s}, n) \subseteq \bigcap_{j=1}^k \left\{ s^\infty : \exists s_j \in S \setminus \{\bar{s}\}, t_j \geq n, \left( \alpha \left( \mathbf{h}_{t_j-1}^\alpha (s^\infty) \right), s_{t_j} \right) = (a, s_j) \right\},$$

and so:

$$\begin{aligned} p_{\hat{\theta}}(E(a, \bar{s}, n)) &\leq \times_{j=1}^k p_{\hat{\theta}} \left( \left\{ s^\infty : \exists s_j \in S \setminus \{\bar{s}\}, t_j \geq n, \left( \alpha \left( \mathbf{h}_{t_j-1}^\alpha (s^\infty) \right), s_{t_j} \right) = (a, s_j) \right\} \right) \\ &\leq \times_{j=1}^k \hat{\theta}(S \setminus \{\bar{s}\}) \\ &= (\hat{\theta}(S \setminus \{\bar{s}\}))^k. \end{aligned}$$

Since  $\bar{s} \in \text{supp } \hat{\theta}$ , this inequality implies that  $p_{\hat{\theta}}(E(a, \bar{s}, n)) = 0$ . It follows that  $p_{\hat{\theta}}(S^\infty \setminus E_{\hat{\theta}}^\alpha) = 0$ , or  $p_{\hat{\theta}}(E_{\hat{\theta}}^\alpha) = 1$ .  $\square$

**Lemma 9.** For every  $\hat{\theta} \in \text{supp } \mu$ , we have,  $\hat{\theta}$ -a.s.,

$$\forall a \in \mathbf{a}_\infty^\alpha (s^\infty), \mu_{s^\infty}^\alpha \left( \left\{ \theta \in \Theta : F(a, \theta) = F(a, \hat{\theta}) \right\} \right) = 1.$$

**Proof.** Fix any  $\hat{\theta} \in \text{supp } \mu$ . Let  $E$  be as in the proof of Lemma 7. Define the set  $\overline{E}_{\hat{\theta}}^\alpha := E \cap E_{\hat{\theta}}^\alpha$ , and fix a sample path  $s^\infty \in \overline{E}_{\hat{\theta}}^\alpha$ . Suppose by way of contradiction that there is some  $a \in \mathbf{a}_\infty^\alpha (s^\infty)$  and some  $m \in M$  such that, for some  $\theta \in \text{supp } \mu_{s^\infty}^\alpha$ , we have  $F(a, \theta)(m) \neq F(a, \hat{\theta})(m)$ . This implies that, for any  $\bar{s} \in f_a^{-1}(m)$ , we have  $B(\mu_{s^\infty}^\alpha, a, f(a, \bar{s})) \neq \mu_{s^\infty}^\alpha$ ; thus, we can find some  $\varepsilon > 0$  such that  $\|B(\mu_{s^\infty}^\alpha, a, f(a, \bar{s})) - \mu_{s^\infty}^\alpha\| = 2\varepsilon$ . By continuity of the Bayes map at  $(\mu_{s^\infty}^\alpha, a, f(a, \bar{s}))$ , there exists some  $\delta > 0$  such that  $\|\mu(\cdot | \mathbf{h}_{t-1}^\alpha (s^\infty)) - \mu_{s^\infty}^\alpha\| < \delta$  implies<sup>34</sup>:

$$\|B(\mu(\cdot | \mathbf{h}_{t-1}^\alpha (s^\infty)), a, f(a, \bar{s})) - B(\mu_{s^\infty}^\alpha, a, f(a, \bar{s}))\| < \varepsilon.$$

Therefore,

$$\begin{aligned} 2\varepsilon &= \|B(\mu_{s^\infty}^\alpha, a, f(a, \bar{s})) - \mu_{s^\infty}^\alpha\| \\ &\leq \|B(\mu(\cdot | \mathbf{h}_{t-1}^\alpha (s^\infty)), a, f(a, \bar{s})) - \mu_{s^\infty}^\alpha\| \\ &\quad + \|B(\mu(\cdot | \mathbf{h}_{t-1}^\alpha (s^\infty)), a, f(a, \bar{s})) - B(\mu_{s^\infty}^\alpha, a, f(a, \bar{s}))\| \\ &< \|B(\mu(\cdot | \mathbf{h}_{t-1}^\alpha (s^\infty)), a, f(a, \bar{s})) - \mu_{s^\infty}^\alpha\| + \varepsilon, \end{aligned}$$

so  $\|B(\mu(\cdot | \mathbf{h}_{t-1}^\alpha (s^\infty)), a, f(a, \bar{s})) - \mu_{s^\infty}^\alpha\| > \varepsilon$ . Invoking Lemma 7 and since  $s^\infty \in \overline{E}_{\hat{\theta}}^\alpha$ , there exists a sequence of dates  $(t_n)_{n \in \mathbb{N}}$  such that, for every  $n$ ,

$$\left( \alpha \left( \mathbf{h}_{t_n-1}^\alpha (s^\infty) \right), s_{t_n} \right) = (a, \bar{s}) \text{ and } \|\mu(\cdot | \mathbf{h}_{t_n-1}^\alpha (s^\infty)) - \mu_{s^\infty}^\alpha\| < \delta,$$

and so:

$$\begin{aligned} &\|\mu(\cdot | \mathbf{h}_{t_n}^\alpha (s^\infty)) - \mu_{s^\infty}^\alpha\| \\ &= \|B(\mu(\cdot | \mathbf{h}_{t_n-1}^\alpha (s^\infty)), \alpha(\mathbf{h}_{t_n-1}^\alpha (s^\infty)), f(\alpha(\mathbf{h}_{t_n-1}^\alpha (s^\infty)), s_{t_n})) - \mu_{s^\infty}^\alpha\| \\ &= \|B(\mu(\cdot | \mathbf{h}_{t_n-1}^\alpha (s^\infty)), a, f(a, \bar{s})) - \mu_{s^\infty}^\alpha\| > \varepsilon. \end{aligned}$$

<sup>34</sup> Notice that the Bayes map is continuous at all  $(v, a, m)$  such that  $\sum_\theta F(a, \theta)(m) v(\theta) > 0$ . Since  $s^\infty \in \overline{E}_{\hat{\theta}}^\alpha$ , this is the case for  $(\mu_{s^\infty}^\alpha, a, f(a, \bar{s}))$ .



This contradicts Lemma 7. Since by Lemmata 7 and 8  $p_{\hat{\theta}}(\overline{E_{\theta}^{\alpha}}) = 1$ , the result follows.  $\square$

**Proof of Proposition 1.** The result is obtained from Lemma 9 using the same argument used to derive Lemma 2 from Lemma 7, replacing the set  $E^*$  with:

$$E^* := \{s^{\infty} : \forall \bar{a} \in \mathbf{a}_{\infty}^{\alpha}(s^{\infty}), \mu_{s^{\infty}}^{\alpha}(\{\theta \in \Theta : F(\bar{a}, \theta) = F(\bar{a}, \bar{\theta})\}) = 1\}.$$

Further details are omitted.  $\square$

**Proof of Corollary 1.** Since under perfect feedback

$$\forall \bar{a} \in A \quad \{\theta \in \Theta : F(\bar{a}, \theta) = F(\bar{a}, \bar{\theta})\} = \bar{\theta}$$

and  $A$  is finite, the statement follows from Proposition 1.  $\square$

A.2. Value

In the rest of the Appendix, we will make use of the fact that the dependence on the state in the value can be made explicit. Indeed, by definition we have:

$$\sum_{h_{\tau} \in H_{\tau}} R(\alpha(h_{\tau}), \theta) p_{\theta}([\mathbf{h}_{\tau}^{\alpha} = h_{\tau}] | h_t) = \sum_{s^{\tau} \in S^{\tau}} r(\mathbf{a}_{\tau}^{\alpha}(s^{\tau-1}), s_{\tau}) p_{\theta}(s^{\tau} | h_t),$$

so that the value function can be expressed as:

$$V(\alpha, \mu | h_t) = \sum_{\tau=t}^{\infty} \delta^{\tau-t} \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s^{\tau} \in S^{\tau}} r(\mathbf{a}_{\tau}^{\alpha}(s^{\tau-1}), s_{\tau}) p_{\theta}(s^{\tau} | h_t) \right) \mu(d\theta | h_t) \right).$$

**Proof of Lemma 3.** It is immediate to see that the map:

$$W_{\bar{t}} : \theta \mapsto \phi \left( \sum_{s^{\bar{t}} \in S^{\bar{t}}} r(\mathbf{a}_{\bar{t}}^{\alpha}(s^{\bar{t}-1}), s_{\bar{t}}) p_{\theta}(s^{\bar{t}} | h_t) \right)$$

is continuous and bounded by  $\max_{(a,s) \in A \times S} |\phi(r(a, s))|$ . Moreover, as argued in the proof of Lemma 5,  $\phi^{-1}$  is continuous. Since the space of measures endowed with the topology of weak convergence is metrizable, it is enough to establish sequential continuity. By continuity of Bayesian updating with respect to positive probability events and  $p_{\mu}(I(h_t)) > 0$ ,

$$\mu_n \rightarrow \mu \Rightarrow \mu_n(\cdot | h_t) \rightarrow \mu(\cdot | h_t).$$

Therefore, by definition of weak convergence of measures and continuity of  $W_{\bar{t}}$  and  $\phi^{-1}$ :

$$\mu_n(\cdot | h_t) \rightarrow \mu(\cdot | h_t) \Rightarrow \phi^{-1} \left( \int_{\Theta} W_{\bar{t}}(\theta) \mu_n(d\theta | h_t) \right) \rightarrow \phi^{-1} \left( \int_{\Theta} W_{\bar{t}}(\theta) \mu(d\theta | h_t) \right).$$

Let  $\varepsilon > 0$ . Since  $\delta < 1$  and  $W_{\tau}$  is bounded, there exists a  $T$  such that for every  $\mu(\cdot | h_t)$ ,

$$\left| \sum_{\tau=T}^{\infty} \delta^{\tau-t} \phi^{-1} \left( \int_{\Theta} W_{\tau}(\theta) \mu(d\theta | h_t) \right) \right| < \varepsilon.$$

But then, let  $n$  be such that for every  $\tau \leq T$ ,

$$|\phi^{-1}(W_{\tau}(\theta) \mu_n(d\theta | h_{\tau})) - \phi^{-1}(W_{\tau}(\theta) \mu(d\theta | h_{\tau}))| < \varepsilon.$$

It follows that:

$$\begin{aligned} & |V(\alpha, \mu | h_t) - V(\alpha, \mu_n | h_t)| \\ &= \left| \sum_{\tau=t}^{\infty} \delta^{\tau-t} \phi^{-1} \left( \int_{\Theta} W_{\tau}(\theta) \mu(d\theta | h_{\tau}) \right) - \sum_{\tau=t}^{\infty} \delta^{\tau-t} \phi^{-1} \left( \int_{\Theta} W_{\tau}(\theta) \mu_n(d\theta | h_{\tau}) \right) \right| \\ &\leq \left| \sum_{\tau=t}^T \delta^{\tau-t} \phi^{-1} \left( \int_{\Theta} W_{\tau}(\theta) \mu(d\theta | h_{\tau}) \right) - \sum_{\tau=t}^T \delta^{\tau-t} \phi^{-1} \left( \int_{\Theta} W_{\tau}(\theta) \mu_n(d\theta | h_{\tau}) \right) \right| \\ &\quad + \left| \sum_{\tau=T+1}^{\infty} \delta^{\tau-t} \phi^{-1} \left( \int_{\Theta} W_{\tau}(\theta) \mu(d\theta | h_{\tau}) \right) \right. \\ &\quad \left. - \sum_{\tau=T+1}^{\infty} \delta^{\tau-t} \phi^{-1} \left( \int_{\Theta} W_{\tau}(\theta) \mu_n(d\theta | h_{\tau}) \right) \right| \\ &< (T + 2)\varepsilon. \end{aligned}$$

Since  $\varepsilon$  has been chosen arbitrarily, the result is proved.  $\square$

**Proof of Corollary 2.** By Proposition 1 and Lemma 2, the set:

$$E^* = \left\{ s^{\infty} : \exists v \in \Delta(\Theta) : \lim_{t \rightarrow +\infty} \mu(\cdot | \mathbf{h}_t^{\alpha}(s^{\infty})) = v, v(\{\theta \in \Theta : F(a, \theta) = F(a, \bar{\theta})\}) = 1 \right\}$$

has  $\bar{\theta}$ -probability 1. For each  $s^{\infty} \in E^*$ , consider:

$$\begin{aligned} \lim_{t \rightarrow \infty} V(\alpha, \mu | \mathbf{h}_t^{\alpha}(s^{\infty})) &= \lim_{t \rightarrow \infty} \sum_{\tau=t}^{\infty} \delta^{\tau-t} \phi^{-1} \left( \int_{\Theta} W_{\tau}(\theta) \mu_{s^{\infty}}^{\alpha}(d\theta) \right) \\ &= \lim_{t \rightarrow \infty} V(\alpha, \delta_{\bar{\theta}} | \mathbf{h}_t^{\alpha}(s^{\infty})), \end{aligned}$$

where the first equality follows from Lemma 3, and the second equality follows from the definition of  $E^*$ .  $\square$

**Proof of Lemma 4.** First note that since  $p_{\mu}(I(h_t))$  and  $p_{\mu}(I(h'_t))$  are strictly positive, then:

$$\frac{\mu(\theta) p_{\theta}(I(h'_t))}{p_{\mu}(I(h'_t))} = \mu(\theta | h'_t) = \mu(\theta | h_t) = \frac{\mu(\theta) p_{\theta}(I(h_t))}{p_{\mu}(I(h_t))}.$$

In particular,

$$\mu(\theta | h'_t) = \mu(\theta | h_t) > 0 \Rightarrow p_{\theta}(I(h_t)) > 0, \text{ and } p_{\theta}(I(h'_t)) > 0.$$

That is, the models in the support of the  $\mu(\cdot | h'_t) = \mu(\cdot | h_t)$  assign positive probability to the two conditioning events. In turn, this implies that  $p_{\theta}(\cdot | h_t)$  and  $p_{\theta}(\cdot | h'_t)$  are obtained by Bayes rule. Hence we have:

$$V(\alpha, \mu | h_t) = \sum_{\tau=t}^{\infty} \delta^{\tau-t} \phi^{-1} \left( \int_{\text{supp } \mu(\cdot | h_t)} \phi \left( \sum_{s^\tau \in S^\tau} r \left( \mathbf{a}_\tau^\alpha \left( s^{\tau-1} \right), s_\tau \right) p_\theta \left( s^\tau | h_t \right) \right) \mu \left( d\theta | h_t \right) \right).$$

To show our result, we will prove that for every  $n$  in  $\mathbb{N}_0$ ,

$$\begin{aligned} & \phi^{-1} \left( \int_{\text{supp } \mu(\cdot | h_t)} \phi \left( \sum_{s^{t+n} \in S^{t+n}} r \left( \mathbf{a}_{t+n}^\alpha \left( s^{t+n-1} \right), s_{t+n} \right) p_\theta \left( s^{t+n} | h_t \right) \right) \mu \left( d\theta | h_t \right) \right) \\ &= \phi^{-1} \left( \int_{\text{supp } \mu(\cdot | h'_t)} \phi \left( \sum_{s^{t'+n} \in S^{t'+n}} r \left( \mathbf{a}_{t'+n}^\alpha \left( s^{t'+n-1} \right), s_{t'+n} \right) p_\theta \left( s^{t'+n} | h'_t \right) \right) \mu \left( d\theta | h'_t \right) \right). \end{aligned}$$

Since  $V(\alpha, \mu | h_t)$  and  $V(\alpha, \mu | h'_t)$  are defined as the discounted sum from  $n = 0$  to infinity of, respectively, the first and second line above, the statement will follow.

Let  $n \in \mathbb{N}_0$ ,  $\theta \in \text{supp } \mu(\cdot | h_t) = \text{supp } \mu(\cdot | h'_t)$ , and  $(k_0, \dots, k_n) \in S^{n+1}$  such that  $\theta(k_i) \neq 0$  for every  $i \in \{0, \dots, n\}$ . Define:

$$K(k_0, \dots, k_n) := \{s^{t+n} \in I(h_t) | s_t = k_0, \dots, s_{t+n} = k_n\}$$

and

$$K'(k_0, \dots, k_n) := \{s^{t'+n} \in I(h'_t) | s_{t'} = k_0, \dots, s_{t'+n} = k_n\}.$$

By definition of  $p_\theta$ , for every  $s^{t+n} \in K(k_0, \dots, k_n)$  and every  $s^{t'+n} \in K'(k_0, \dots, k_n)$ ,

$$p_\theta(s^{t+n} | h_t) = \prod_{i=0}^n \theta(k_i) = p_\theta(s^{t'+n} | h'_t).$$

To ease notation, fix  $(k_0, \dots, k_n)$  momentarily and let  $K = K(k_0, \dots, k_n)$  and  $K' = K'(k_0, \dots, k_n)$ . We show that

$$r \left( \mathbf{a}_{t+n}^\alpha \left( s^{t+n-1} \right), s_{t+n} \right)$$

is constant on  $K$ . Indeed, we prove by way of induction that for every  $j \in \{0, \dots, n\}$ ,  $\mathbf{a}_{t+j}^\alpha(s^{t+j-1})$  is constant on  $K$ . Since for every  $s^{t+n} \in K$  we have  $\mathbf{h}_t^\alpha(s^{t-1}) = h_t$ ,

$$\mathbf{a}_t^\alpha(s^{t-1}) = \alpha(\mu(\cdot | h_t)).$$

Suppose by way of induction that the statement holds for  $j' \leq j$ . Thus, for every  $s^{t+n} \in K$  we have:

$$\begin{aligned} \mathbf{h}_{t+j}^\alpha(s^{t+j-1}) &= \left( h_t, \mathbf{a}_t^\alpha(s^{t-1}), f \left( \mathbf{a}_t^\alpha(s^{t-1}), s_t \right), \dots, \mathbf{a}_{t+j-1}^\alpha(s^{t+j-2}), \right. \\ & \quad \left. f \left( \mathbf{a}_{t+j-1}^\alpha(s^{t+j-2}), s_{t+j-1} \right) \right), \end{aligned}$$

which, by definition of  $K$  and by the inductive hypothesis, is constant on  $K$ . It follows that:

$$\mathbf{a}_{t+j}^\alpha(s^{t+j-1}) = \alpha \left( \mu \left( \cdot | \mathbf{h}_{t+j}^\alpha(s^{t+j-1}) \right) \right)$$

is constant on  $K$ . Since  $s_{t+n} = k_n$  for every  $s^{t+n} \in K$ , we have shown that  $r(\mathbf{a}_{t+n}^\alpha(s^{t+n-1}), s_{t+n})$  is also constant on  $K$ . A similar argument shows that  $r(\mathbf{a}'_{t+n}^\alpha(s^{t'+n-1}), s'_{t+n})$  is constant on  $K'$ . Moreover, for every  $s^{t+n}$  in  $K$  and  $s^{t'+n}$  in  $K'$ , and for every  $j$  in  $\{0, \dots, n\}$ ,

$$\mu(\cdot | \mathbf{h}_{t+j}^\alpha(s^{t+j-1})) = \mu(\cdot | \mathbf{h}'_{t+j}^\alpha(s^{t'+j-1})).$$

We prove this equality by induction on  $j$ . By hypothesis, it is true for  $j = 0$ . Let  $j \in \{1, \dots, n\}$  and suppose that it is true for  $j - 1$ . This implies that:

$$\begin{aligned} \mathbf{a}_{t+j-1}^\alpha(s^{t+j-2}) &= \alpha\left(\mu\left(\cdot | \mathbf{h}_{t+j-1}^\alpha(s^{t+j-2})\right)\right) \\ &= \alpha\left(\mu\left(\cdot | \mathbf{h}'_{t+j-1}^\alpha(s^{t'+j-2})\right)\right) \\ &= \mathbf{a}'_{t'+j-1}^\alpha(s^{t'+j-2}). \end{aligned}$$

Therefore:

$$\begin{aligned} &\mu\left(\theta | \mathbf{h}_{t+j}^\alpha(s^{t+j-1})\right) \\ &= \frac{\mu\left(\theta | \mathbf{h}_{t+j-1}^\alpha(s^{t+j-2})\right) F\left(\mathbf{a}_{t+j-1}^\alpha(s^{t+j-2}), \theta\right) \left(f\left(\mathbf{a}_{t+j-1}^\alpha(s^{t+j-2}), k_{j-1}\right)\right)}{F\left(\mathbf{a}_{t+j-1}^\alpha(s^{t+j-2}), \theta_{\mu\left(\cdot | \mathbf{h}_{t+j-1}^\alpha(s^{t+j-2})\right)}\right) \left(f\left(\mathbf{a}_{t+j-1}^\alpha(s^{t+j-2}), k_{j-1}\right)\right)} \\ &= \frac{\mu\left(\theta | \mathbf{h}'_{t'+j-1}^\alpha(s^{t'+j-2})\right) F\left(\mathbf{a}'_{t'+j-1}^\alpha(s^{t'+j-2}), \theta\right) \left(f\left(\mathbf{a}'_{t'+j-1}^\alpha(s^{t'+j-2}), k_{j-1}\right)\right)}{F\left(\mathbf{a}'_{t'+j-1}^\alpha(s^{t'+j-2}), \theta_{\mu\left(\cdot | \mathbf{h}'_{t'+j-1}^\alpha(s^{t'+j-2})\right)}\right) \left(f\left(\mathbf{a}'_{t'+j-1}^\alpha(s^{t'+j-2}), k_{j-1}\right)\right)} \\ &= \mu\left(\theta | \mathbf{h}'_{t'+j-1}^\alpha(s^{t'+j-1})\right). \end{aligned}$$

This in turn implies that, for every  $s^{t+n}$  in  $K$  and  $s^{t'+n}$  in  $K'$ ,

$$\begin{aligned} r\left(\mathbf{a}_{t+n}^\alpha(s^{t+n-1}), s_{t+n}\right) &= r\left(\alpha\left(\mu\left(\cdot | \mathbf{h}_{t+n}^\alpha(s^{t+n-1})\right)\right), k_n\right) \\ &= r\left(\alpha\left(\mu\left(\cdot | \mathbf{h}'_{t+n}^\alpha(s^{t'+n-1})\right)\right), k_n\right) \\ &= r\left(\mathbf{a}'_{t'+n}^\alpha(s^{t'+n-1}), s'_{t+n}\right). \end{aligned} \tag{7}$$

Now, we restart to explicitly highlight the dependence on  $(k_0, \dots, k_n)$  of  $K$ . Moreover, for every  $n \in \mathbb{N}_0$  and for every  $(k_1, \dots, k_n) \in S^n$ , let  $r(k_0, \dots, k_n) = r(\mathbf{a}_{t+n}^\alpha(s^{t+n-1}), s_{t+n}) = r(\mathbf{a}'_{t'+n}^\alpha(s^{t'+n-1}), s'_{t+n})$ , where  $s^{t+n} \in K(k_0, \dots, k_n)$  and  $s^{t'+n} \in K'(k_0, \dots, k_n)$ . By (7), this quantity is well defined. We have:

$$\begin{aligned} &\sum_{s^{t+n} \in S^{t+n}} r\left(\mathbf{a}_{t+n}^\alpha(s^{t+n-1}), s_{t+n}\right) p_\theta(s^{t+n} | h_t) \\ &= \sum_{s^{t+n} \in S^{t+n}} r\left(\mathbf{a}_{t+n}^\alpha(s^{t+n-1}), s_{t+n}\right) \prod_{i=0}^n \theta(k_i) \\ &= \sum_{s^{t'+n} \in S^{t'+n}} r\left(\mathbf{a}'_{t'+n}^\alpha(s^{t'+n-1}), s'_{t+n}\right) p_\theta(s^{t'+n} | h'_{t'}). \end{aligned}$$

Finally, since we have  $\mu(\cdot | h_t) = \mu(\cdot | h'_t)$ , this implies that:

$$\begin{aligned} & \phi^{-1} \left( \int_{\text{supp } \mu(\cdot | h_t)} \phi \left( \sum_{s^{t+n} \in S^{t+n}} r \left( \mathbf{a}_{t+n}^\alpha (s^{t+n-1}), s_{t+n} \right) p_\theta(s^{t+n} | h_t) \right) \mu(d\theta | h_t) \right) \\ &= \phi^{-1} \left( \int_{\text{supp } \mu(\cdot | h'_t)} \phi \left( \sum_{s^{t'+n} \in S^{t'+n}} r \left( \mathbf{a}_{t'+n}^\alpha (s^{t'+n-1}), s_{t'+n} \right) p_\theta(s^{t'+n} | h'_t) \right) \mu(d\theta | h'_t) \right) \end{aligned}$$

and the thesis follows.  $\square$

**Proof of Proposition 2.** The result follows by considering the subgame perfect equilibrium of an ancillary game. We describe the structure of this game:

- The set of players is  $P = \mathbb{N} \cup \{0\}$ . That is, players are the different periods and Nature;
- The set of actions available to each player  $t \in \mathbb{N}$  is  $A$ , the set of actions available to Nature is  $S$ ;
- The timing is as follows: at period  $t \in \mathbb{N}$ , player  $t$  chooses an action. After player  $t$  has chosen his action, Nature chooses  $s_t \in S$  accordingly to a (stationary) uniform distribution over the states;
- Every player observes the sequence of actions played and states realized in the previous periods.

Notice that the set of subgames where a player  $t \in \mathbb{N}$  is going to move is one to one with the set  $H_t$ . Therefore, a strategy profile  $\sigma$  of the ancillary game is mapped in an obvious way into a strategy of our single-agent decision problem.

- Finally, we specify the payoff of player  $t$  when the strategy profile of the single-agent decision problem induced by  $\sigma$  is  $\alpha$  and the sequence of state realized is  $s^\infty$ :

$$U_t(\alpha, s^\infty) = V(\alpha, \mu(\cdot | \mathbf{h}_t^\alpha(s^t)) | \mathbf{h}_t^\alpha(s^t))$$

where  $V$  is the value in the single-agent decision problem.

The ancillary game thus obtained is a perfect information game, and by Lemma 4 it satisfies all the assumptions of Theorem 4 in Hellwig and Leininger (1987). Therefore it admits a subgame-perfect equilibrium in pure strategies  $\alpha$ . Moreover, by Theorem 4 in Hellwig and Leininger (1987) and the coupling argument of Lemma 4 every player with the same belief plays the same action, and by definition of the payoff function of each player,  $(\alpha, \mu)$  is rational.  $\square$

### A.3. Convergence to SCE

By Lemma 4, we know that if  $h_t$  and  $h'_t$  are two histories such that  $p_\mu(I(h_t))$  and  $p_\mu(I(h'_t))$  are strictly positive, then  $\mu(\cdot | h_t) = \mu(\cdot | h'_t)$  implies  $V(\alpha, \mu | h_t) = V(\alpha, \mu | h'_t)$ . Therefore, the value function at a particular history depends on the history only through the updated beliefs. In particular, whenever  $p_\mu(I(h_t)) > 0$ ,

$$V(\alpha, \mu|h_t) = V(\alpha, \mu(\cdot|h_t)|h_1).$$

Therefore, to ease notation, in the following proofs we will use the ancillary function  $\hat{V}$  mapping stationary strategies and beliefs into real numbers:

$$\hat{V}(\alpha, \nu) := V(\alpha, \nu|h_1).$$

By the previous argument,  $\hat{V}(\alpha, \nu) = V(\alpha, \mu|h_t)$  whenever  $\nu = \mu(\cdot|h_t)$ .

**Lemma 10.** *Let  $(\alpha, \mu, \bar{\theta})$  be such that:*

1.  $\mu(\{\theta \in \Theta : p_\theta^\alpha = p_{\bar{\theta}}^\alpha\}) = 1;$
2. *For every action  $a$ , period  $t$ , and information history  $h_t$ ,*

$$p_\mu(I(h_t)) > 0 \Rightarrow V(\alpha, \mu|h_t) \geq V(\alpha/(h_t, a), \mu|h_t).$$

*Then,  $(\alpha(\mu), \mu, \bar{\theta})$  is an SCE.*

**Proof.** It is immediate to see that condition (1) implies condition (i) of SCE. Now, we show that  $(\alpha, \mu, \bar{\theta})$  features myopic best reply on path, that is,  $(\alpha(\mu), \mu, \bar{\theta})$  satisfies property (ii) of SCE. By way of contradiction, suppose there is  $a \in A$  such that:

$$\phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a, s)\theta(s) \right) \mu(d\theta) \right) > \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(\alpha(\mu), s)\theta(s) \right) \mu(d\theta) \right).$$

By property 2, it must be the case that  $\hat{V}(\alpha/a, \mu) \leq \hat{V}(\alpha, \mu)$ , where  $\alpha/a$  is the strategy that prescribes  $a$  in the first period to come and coincide with  $\alpha$  otherwise. However, we have:

$$\begin{aligned} & \hat{V}(\alpha, \mu) \\ &= \sum_{s \in S} r(\alpha(\mu), s)\bar{\theta}(s) + \delta \hat{V}(\alpha, \mu) \\ &\leq \sum_{s \in S} r(\alpha(\mu), s)\bar{\theta}(s) + \delta \min_{m: F(a, \theta_\mu)(m) > 0} \hat{V}(\alpha, \mu(\cdot|(a, m))) \\ &< \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a, s)\theta(s) \right) \mu(d\theta) \right) + \\ &\delta \min_{m: F(a, \theta_\mu)(m) > 0} \left( \sum_{\tau=2}^{\infty} \delta^{\tau-2} \phi^{-1} \right. \\ &\times \left. \left( \int_{\Theta} \phi \left( \sum_{s^\tau \in S^\tau} r(\mathbf{a}_\tau^\alpha(s^{\tau-1}), s_\tau) p_\theta(s^\tau|(a, m)) \right) \mu(d\theta|(a, m)) \right) \right) \\ &= \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a, s)\theta(s) \right) \mu(d\theta) \right) + \\ &\delta \left( \sum_{\tau=2}^{\infty} \delta^{\tau-2} \phi^{-1} \left( \min_{m: F(a, \theta_\mu)(m) > 0} \right. \right. \end{aligned}$$

$$\begin{aligned}
 & \times \left( \int_{\Theta} \phi \left( \sum_{s^\tau \in S^\tau} r \left( \mathbf{a}_\tau^\alpha \left( s^{\tau-1} \right), s_\tau \right) p_\theta \left( s^\tau | (a, m) \right) \right) \mu \left( d\theta | (a, m) \right) \right) \Big) \\
 & \leq \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a, s) \theta(s) \right) \mu(d\theta) \right) + \\
 & \delta \left( \sum_{\tau=2}^{\infty} \delta^{\tau-2} \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s^\tau \in S^\tau} r \left( \mathbf{a}_\tau^\alpha \left( s^{\tau-1} \right), s_\tau \right) p_\theta \left( s^\tau | (a, m) \right) \right) \mathbb{E}_{p_\mu} [\mu(d\theta | (a, m))] \right) \right) \\
 & = \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a, s) \theta(s) \right) \mu(d\theta) \right) + \\
 & \delta \left( \sum_{\tau=2}^{\infty} \delta^{\tau-2} \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s^\tau \in S^\tau} r \left( \mathbf{a}_\tau^\alpha \left( s^{\tau-1} \right), s_\tau \right) p_\theta \left( s^\tau | (a, m) \right) \right) \mu(d\theta) \right) \right) \\
 & = \sum_{\tau=1}^{\infty} \delta^{\tau-2} \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s^\tau \in S^\tau} r \left( \mathbf{a}_\tau^{\alpha/a} \left( s^{\tau-1} \right), s_\tau \right) p_\theta \left( s^\tau | (a, m) \right) \right) \mu(d\theta) \right) \\
 & = V(\alpha/a, \mu),
 \end{aligned}$$

where the first equality comes from property 1, the strict inequality comes from hypothesis, the second equality comes from the fact that  $\phi$  is strictly increasing, the third equality by Lemma 6, and the fourth and fifth equalities by the definition of  $\alpha/a$ . Note that we will be done as soon as we prove the first weak inequality, that is:

$$\hat{V}(\alpha, \mu) \leq \min_{m: F(a, \theta_\mu)(m) > 0} \hat{V}(\alpha, \mu(\cdot | (a, m))).$$

Indeed, it would follow that  $\hat{V}(\alpha, \mu) < \hat{V}(\alpha/a, \mu)$ , a contradiction with the fact that  $(\alpha, \mu, \bar{\theta})$  satisfies 2.

Suppose that there exists  $m$  such that  $F(a, \theta_\mu)(m) > 0$  with  $\hat{V}(\alpha, \mu(\cdot | (a, m))) < \hat{V}(\alpha, \mu)$ . The fact that  $F(a, \theta_\mu)(m) > 0$  implies that  $\theta_\mu(I(a, m)) > 0$ . On the other hand, by property 1,  $\mu(\{\theta \in \Theta : p_\theta^\alpha = p_{\bar{\theta}}^\alpha\}) = 1$ , and in particular:

$$\mu(\{\theta \in \Theta : F(\alpha(\mu), \theta)(m) = F(\alpha(\mu), \bar{\theta})(m)\}) = 1.$$

Then, let  $B = \{\theta \in \Theta : F(\alpha(\mu), \theta) = F(\alpha(\mu), \bar{\theta})\}$ . By Bayes rule, we have that:

$$\mu(B| (a, m)) = \mu(B) = 1.$$

But then, it follows that:

$$\begin{aligned}
 & \hat{V}(\alpha, \mu(\cdot | (a, m))) \\
 & < (1 - \delta) \hat{V}(\alpha, \mu) + \delta \hat{V}(\alpha, \mu(\cdot | (a, m))) \\
 & = \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(\alpha(\mu), s) \theta(s) \right) \mu(d\theta) \right) + \delta \hat{V}(\alpha, \mu(\cdot | (a, m)))
 \end{aligned}$$

$$\begin{aligned}
 &= \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(\alpha(\mu), s) \theta(s) \right) \mu(d\theta|(a, m)) \right) + \delta \hat{V}(\alpha, \mu(\cdot|(a, m))) \\
 &= \sum_{\tau=1}^{\infty} \delta^{\tau-1} \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s^\tau \in S^\tau} r(\mathbf{a}_\tau^{\alpha/\alpha(\mu)}(s^{\tau-1}), s_\tau) p_\theta(s^\tau) \right) \mu(d\theta|(a, m)) \right) \\
 &= \hat{V}(\alpha/\alpha(\mu), \mu(\cdot|(a, m))).
 \end{aligned}$$

This contradicts the fact that  $(\alpha, \mu, \bar{\theta})$  satisfies property 2.  $\square$

**Proof of Proposition 3.** Let  $(\alpha, \mu, \bar{\theta})$  be consistent from period  $T$ . First, we have that the hypotheses of Lemma 7 are satisfied, so let  $E$  be as in the corresponding proof. Define

$$\hat{E} = E \bigcap_{t \in \mathbb{N}} \left\{ s^\infty : \bar{\theta} \left( l_t^\alpha(s^{t-1}) \right) > 0 \right\}.$$

By Lemma 3, the value (4) is continuous in beliefs  $\mu$ . Fix  $s^\infty \in \hat{E}$ ; for every  $a$  in  $A$ ,

$$\lim_{t \rightarrow \infty} \hat{V}(\alpha/a, \mu(\cdot|\mathbf{h}_t^\alpha(s^{t-1}))) = \hat{V}(\alpha/a, \mu_{s^\infty}^\alpha).$$

Let  $A_\infty := \arg \max_{a \in A} \hat{V}(\alpha/a, \mu_{s^\infty}^\alpha)$ . Note that, in general, our definition of rationality does not require that  $\alpha(\mu_{s^\infty}^\alpha) \in A_\infty$ . Indeed, if there is no  $h_t$  such that  $p_\mu(I(h_t)) > 0$  and  $\mu(\cdot|h_t) = \mu_{s^\infty}^\alpha$ , then  $\alpha(\mu_{s^\infty}^\alpha)$  does not need to satisfy the one-deviation property. Since  $s^\infty \in \hat{E}$ , it follows that  $p_{\bar{\theta}}(I(\mathbf{h}_t^\alpha(s^{t-1}))) > 0$  for every finite  $t$ . By Assumption 2,  $p_\mu(I(\mathbf{h}_t^\alpha(s^{t-1}))) > 0$ . Hence,

$$\alpha(\mu(\cdot|\mathbf{h}_t^\alpha(s^{t-1}))) \in \arg \max_{a \in A} \hat{V}(\alpha/a, \mu(\cdot|\mathbf{h}_t^\alpha(s^{t-1}))).$$

Now, let  $a \notin A_\infty$ , and fix  $a^* \in A_\infty$ . We have that

$$\begin{aligned}
 \lim_{t \rightarrow \infty} \hat{V}(\alpha/a, \mu(\cdot|\mathbf{h}_t^\alpha(s^{t-1}))) &= \hat{V}(\alpha/a, \mu_{s^\infty}^\alpha) \\
 &< \max_{a' \in A} \hat{V}(\alpha/a', \mu_{s^\infty}^\alpha) = \hat{V}(\alpha/a^*, \mu_{s^\infty}^\alpha) \\
 &= \lim_{t \rightarrow \infty} \hat{V}(\alpha/a^*, \mu(\cdot|\mathbf{h}_t^\alpha(s^{t-1}))).
 \end{aligned}$$

Hence there exists  $T_{s^\infty}^a$  such that  $a \notin \alpha(\mu(\cdot|\mathbf{h}_t^\alpha(s^{t-1})))$  for every  $t \geq T_{s^\infty}^a$ . Let  $T_{s^\infty}^* = \max_{a \in A \setminus A_\infty} T_{s^\infty}^a$ . Then, from  $T_{s^\infty}^*$  onward, the only actions played are in  $A_\infty$ , that is, they satisfy the one-deviation property with respect to the limit beliefs  $\mu_{s^\infty}^\alpha$ . Let  $\hat{T}_{s^\infty} = \max \{T, T_{s^\infty}^*\}$ ; we have that from  $\hat{T}_{s^\infty}$  onward the action prescribed by strategy  $\alpha$ ,  $\mathbf{a}_t^\alpha(s^{t-1})$ , satisfies the one-deviation property with respect to beliefs  $\mu_{s^\infty}^\alpha$ , and  $\mu_{s^\infty}^\alpha$  is confirmed given such action. By Lemma 10, this implies that  $(\mathbf{a}_t^\alpha(s^{t-1}), \mu_{s^\infty}^\alpha, \bar{\theta})$  is an SCE for every  $t \geq \hat{T}_{s^\infty}$ .  $\square$

**Proof of Proposition 4.** By hypothesis, we know that  $(\mathbf{a}_t^\alpha(s^\infty), \mu(\cdot|\mathbf{h}_t^\alpha(s^\infty)))$  converges to an SCE. Therefore, there exists  $\hat{T}$  such that, for every  $t \geq \hat{T}$ , the triple  $(\mathbf{a}_t^\alpha(s^{t-1}), \mu_{s^\infty}^\alpha, \bar{\theta})$  is an SCE, and so:

$$\mathbf{a}_t^\alpha(s^{t-1}) \in \arg \max_{a \in A} \phi^{-1} \left( \int_{\Theta} \phi(R(a, \theta)) \mu_{s^\infty}^\alpha(d\theta) \right) = \{a^*\}.$$



It follows that  $(a^*, \mu_{s^\infty}^\alpha, \bar{\theta})$  is an SCE. Now, let  $a \neq a^*$ . By Lemmata 5 and 7,

$$\begin{aligned} & \lim_{t \rightarrow \infty} \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a^*, s) \theta(s) \right) \mu(d\theta | \mathbf{h}_t^\alpha(s^{t-1})) \right) \\ &= \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a^*, s) \theta(s) \right) \mu_{s^\infty}^\alpha(d\theta) \right) \\ &> \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a, s) \theta(s) \right) \mu_{s^\infty}^\alpha(d\theta) \right) \\ &= \lim_{t \rightarrow \infty} \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a, s) \theta(s) \right) \mu(d\theta | \mathbf{h}_t^\alpha(s^{t-1})) \right). \end{aligned}$$

Thus, there exists  $\bar{T}_{a,s^\infty} > \hat{T}$  such that  $t > \bar{T}_{a,s^\infty}$  implies:

$$\begin{aligned} & \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a^*, s) \theta(s) \right) \mu(d\theta | \mathbf{h}_t^\alpha(s^{t-1})) \right) \\ &> \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a, s) \theta(s) \right) \mu(d\theta | \mathbf{h}_t^\alpha(s^{t-1})) \right). \end{aligned}$$

Let  $\bar{T}_{s^\infty} = \max_{a \in A \setminus a^*} \bar{T}_{a,s^\infty}$ . We have that  $t > \bar{T}_{s^\infty}$  implies

$$\begin{aligned} & \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a^*, s) \theta(s) \right) \mu(d\theta | \mathbf{h}_t^\alpha(s^{t-1})) \right) \\ &= \max_{a \in A} \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a, s) \theta(s) \right) \mu(d\theta | \mathbf{h}_t^\alpha(s^{t-1})) \right). \end{aligned}$$

The thesis follows.  $\square$

**Proof of Corollary 3.** By Proposition 3, there exists a set  $E \subseteq S^\infty$  with  $p_{\bar{\theta}}(E) = 1$  such that convergence to an SCE happens on that set. By Proposition 1, there exists a set  $E^* \subseteq S^\infty$  with  $p_{\bar{\theta}}(E^*) = 1$  such that, if  $s^\infty \in E^*$ , then for every  $a \in \mathbf{a}_\infty^\alpha(s^\infty)$ :

$$\lim_{t \rightarrow \infty} \mu(\{\theta \in \Theta : F(a, \theta) = F(a, \bar{\theta})\} | \mathbf{h}_t^\alpha(s^\infty)) = 1. \tag{8}$$

Let  $G = E \cap E^*$ . Own-action independence and 8 imply that, if  $s^\infty \in G$ , then for every  $\hat{a} \in A$ :

$$\lim_{t \rightarrow \infty} \mu(\{\theta \in \Theta : F(\hat{a}, \theta) = F(\hat{a}, \bar{\theta})\} | \mathbf{h}_t^\alpha(s^\infty)) = 1.$$

Thus, for every  $s^\infty \in G \subseteq E^*$ ,

$$\phi^{-1} \left( \int_{\Theta} \phi(R(\hat{a}, \theta)) \mu_{s^\infty}^\alpha(d\theta) \right) = R(\hat{a}, \bar{\theta}); \tag{9}$$

that is, the value of each action under the limit belief is equal to the objective value. Since  $s^\infty \in G \subseteq E$ , there exists a finite time  $t$  such that  $(\mathbf{a}_t^\alpha(s^{\tau-1}), \mu_{s^\infty}^\alpha, \bar{\theta})$ ,  $\tau \geq t$ , forms an SCE. By definition of SCE and Equation (9), this means that for  $\tau \geq t$  only the objective myopic best reply is played.  $\square$

#### A.4. Comparative dynamics

Since this section deals with different levels of ambiguity attitude, we make the dependence of the value on the ambiguity attitude explicit by writing  $\hat{V}_\phi(\alpha, \mu)$  in place of  $\hat{V}(\alpha, \mu)$ .

**Proof of Proposition 5.** Since  $(\alpha, \mu)$  is rational under ambiguity neutrality, it satisfies the one-deviation property for every  $h_t$  such that  $p_\mu(I(h_t)) > 0$  and  $v = \mu(\cdot|h_t)$ :

$$V_{Id}(\alpha, \mu|h_t) \geq V_{Id}(\alpha/(h_t, a), \mu|h_t)$$

for every  $a \in A$ . Since the spaces of action and state are finite and  $\delta < 1$ , our problem is continuous at infinity. By the one-deviation principle for ambiguity neutral agents (see e.g., Theorem 4.2 in Fudenberg and Tirole, 1991), this implies that for every alternative strategy  $\gamma$  and every history  $h_t$  with  $p_\mu(I(h_t)) > 0$ :

$$V_{Id}(\alpha, \mu|h_t) \geq V_{Id}(\gamma, \mu|h_t). \tag{10}$$

Let  $v$  be in the belief range of  $\mu$ , say  $v = \mu(\cdot|h_t)$  and suppose, by way of contradiction, that  $\beta(v)$  prescribes an ambiguous action while  $\alpha(v)$  is unambiguous given  $v$ . Then,

$$\begin{aligned} V_\phi(\beta, \mu|h_t) &\geq V_\phi(\beta/(\alpha(v), h_t), \mu|h_t) \\ &= R(\alpha(v), p_v) + \delta \sum_{m \in M} F(\alpha(v), p_v)(m) V_\phi(\beta, \mu|(h_t, (\alpha(v), m))) \\ &= R(\alpha(v), p_v) + \delta V_\phi(\beta, \mu|h_t) \end{aligned}$$

where the inequality follows from the one-deviation property, the second equality follows from the assumption that  $\alpha(v)$  is unambiguous, and the third equality follows from  $\alpha(v)$  being  $v$ -unambiguous and Lemma 4. Therefore,

$$V_\phi(\beta, \mu|h_t) \geq \frac{R(\alpha(v), p_v)}{1 - \delta}.$$

Since  $\alpha(v)$  is  $v$ -unambiguous,

$$\begin{aligned} V_\phi(\alpha, \mu|h_t) &= R(\alpha(v), p_v) + \delta \sum_{m \in M} F(\alpha(v), p_v)(m) V_\phi(\alpha, \mu|(h_t, (\alpha(v), m))) \\ &= \sum_{k=0}^\infty \delta^k R(\alpha(v), p_v) = \frac{R(\alpha(v), p_v)}{1 - \delta}. \end{aligned}$$

Then by (10), and Jensen’s inequality (see page 294 in Billingsley, 2012 for the version used here),

$$\frac{R(\alpha(v), p_v)}{1 - \delta} = V_{\text{Id}}(\alpha, v|h_t) \geq V_{\text{Id}}(\beta, v|h_t) > V_\phi(\beta, v|h_t) \geq \frac{R(\alpha(v), p_v)}{1 - \delta},$$

a contradiction.  $\square$

**Proof of Proposition 6.** Let  $p_{\bar{\theta}}(\mathbf{h}_t^\alpha(s^{t-1})) > 0$ , and let  $(\alpha(\mu(\cdot|h_t^\alpha(s^{t-1}))), \mu(\cdot|h_t^\alpha(s^{t-1})), \bar{\theta})$  be an SCE under ambiguity neutrality. By Assumption 2, an action that is unambiguous given  $\mu$  will be unambiguous given  $\mu(\cdot|h_t^\alpha(s^{t-1}))$ . Therefore, when an action  $a \neq a^*$  is chosen, with probability 1 it does not induce any updating in beliefs. Then, stationarity of the strategy implies that the same action  $a$  will be chosen in the following period. Thus, the sequence  $(\alpha(\mu(\cdot|h_\tau^\alpha(s^{\tau-1}))))_{\tau=1}^t$  has almost surely the form  $(a^*, \dots, a^*, a, \dots, a)$  for some  $a \in A$ .<sup>35</sup>

The same reasoning guarantees that  $(\beta(\mu(\cdot|h_\tau^\beta(s^{\tau-1}))))_{\tau=1}^t$  also has almost surely the form  $(a^*, \dots, a^*, a, \dots, a)$  for some  $a \in A$ . We show that for almost every  $s^\infty$ , the sequence of  $a^*$  is (weakly) shorter under  $\beta$ . In particular, suppose that  $t \in \mathbb{N}$ ,  $\bar{\theta}(s^{t-1}) > 0$ , and let:

$$(\alpha(\mu(\cdot|h_\tau^\alpha(s^{\tau-1}))))_{\tau=1}^t = (a^*, \dots, a^*) = (\beta(\mu(\cdot|h_\tau^\beta(s^{\tau-1}))))_{\tau=1}^t.$$

Then,

$$\mathbf{h}_t^\alpha(s^{t-1}) = (a^*, f(a^*, s_t))_{\tau=1}^{t-1} = \mathbf{h}_t^\beta(s^{t-1}),$$

and by Assumption 2:

$$\mu(\cdot|h_\tau^\alpha(s^{\tau-1})) = \mu(\cdot|h_\tau^\beta(s^{\tau-1})).$$

By Proposition 5,  $\alpha(\mu(\cdot|h_t^\alpha(s^{t-1}))) \neq a^*$  implies  $\beta(\mu(\cdot|h_t^\beta(s^{t-1}))) \neq a^*$ ; with probability 1, if strategy  $\alpha$  ends experimentation after  $t - 1$  periods,  $\beta$  ends experimentation in at most  $t - 1$  periods. Therefore,  $(\beta(\mu(\cdot|h_\tau^\beta(s^{\tau-1}))))_{\tau=1}^t$  has the form  $(a^*, \dots, a^*, a, \dots, a)$  with a (weakly) shorter sequence of  $a^*$ . If  $\beta(\mu(\cdot|h_\tau^\beta(s^{t-1}))) \neq a^*$ ,  $(\beta(\mu(\cdot|h_\tau^\beta(s^t))), \mu(\cdot|h_\tau^\beta(s^t)), \bar{\theta})$  is an SCE by Lemma 10 and rationality of  $(\beta, \mu)$ . If  $(\beta(\mu(\cdot|h_\tau^\beta(s^{\tau-1}))))_{\tau=1}^t = (a^*)_{\tau=1}^t$ , then  $\alpha$  and  $\beta$  have prescribed the same action  $a^*$  at every node, and therefore:

$$\begin{aligned} & (\alpha(\mu(\cdot|h_t^\alpha(s^{t-1}))), \mu(\cdot|h_t^\alpha(s^{t-1})), \bar{\theta}) \\ &= (\beta(\mu(\cdot|h_t^\beta(s^{t-1}))), \mu(\cdot|h_t^\beta(s^{t-1})), \bar{\theta}). \end{aligned}$$

Again, by Lemma 10, the fact that the L.H.S. is an SCE with ambiguity neutrality, and the definition of SCE, the R.H.S. is an SCE under ambiguity aversion.  $\square$

**Proof of Proposition 7.** First, suppose that the objectively-optimal action played in the SCE is  $a^*$ . Then, given the original belief  $\mu$  and infinite history  $s^\infty$ ,  $(\mathbf{a}_t^\beta, \mu(\cdot|h_t^\beta))$  converges to an SCE  $(a^*, \mu^*, \bar{\theta})$  if and only if  $\beta(\mu(\cdot|h_\tau^\beta(s^\infty))) = a^*$  for every  $\tau$ . But Proposition 5 guarantees that this can happen only if  $\alpha(\mu(\cdot|h_\tau^\alpha(s^\infty))) = a^*$  for every  $\tau$ , and therefore  $(\mathbf{a}_t^\alpha, \mu(\cdot|h_t^\alpha))$  also converges to  $(a^*, \mu^*, \bar{\theta})$ .

<sup>35</sup> Possibly including the cases  $(a^*, \dots, a^*)$  and  $(a, \dots, a)$ .

Second, suppose that the objectively-optimal action is  $a \neq a^*$ . Consider  $s^\infty$  where convergence of  $(\mathbf{a}_t^\beta, \mu(\cdot|\mathbf{h}_t^\beta))$  to an SCE  $(a, \nu, \bar{\theta})$  happens. Then, by definition of SCE,  $R(a, \theta) = R(a, \bar{\theta})$  is constant on  $\text{supp } \nu$ . At the same time, Proposition 3 guarantees convergence of  $(\mathbf{a}_t^\alpha, \mu(\cdot|\mathbf{h}_t^\alpha))$  to an SCE  $(\hat{a}, \hat{\mu}, \bar{\theta})$ . Assumption 2 implies that  $\text{supp } \mu(\cdot|\mathbf{h}_t^\alpha) \subset \text{supp } \mu(\cdot)$  for every  $h_t$  and therefore  $R(a, \theta) = R(a, \bar{\theta})$  is constant on  $\text{supp } \hat{\mu}$ . Since, by definition of SCE,  $R(\hat{a}, \theta) = R(\hat{a}, \bar{\theta})$  is also constant on  $\text{supp } \hat{\mu}$  and:

$$\hat{a} \in \arg \max_{a' \in A} \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a', s) \theta(s) \right) \mu^*(d\theta) \right),$$

we must have  $R(\hat{a}, \bar{\theta}) \geq R(a, \bar{\theta}) = \max_{a' \in A} R(a', \bar{\theta})$ .  $\square$

**Proof of Proposition 8.** Let  $(\bar{a}, \mu, \bar{\theta})$  be an SCE under ambiguity attitude  $\phi$ . The optimality condition for an SCE gives:

$$\bar{a} \in \arg \max_{a' \in A} \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a', s) \theta(s) \right) \mu(d\theta) \right). \tag{11}$$

Therefore, there exists  $\theta^*$  in  $\text{supp } \mu$  with:

$$\sum_{s \in S} r(\bar{a}, s) \theta^*(s) \geq \sum_{s \in S} r(a^*, s) \theta^*(s).$$

Next, consider an arbitrary  $a \neq a^*$ ; such actions are unambiguous given  $\mu$ . By (11):

$$\begin{aligned} \sum_{s \in S} r(\bar{a}, s) \theta^*(s) &= \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(\bar{a}, s) \theta(s) \right) \mu(d\theta) \right) \\ &\geq \phi^{-1} \left( \int_{\Theta} \phi \left( \sum_{s \in S} r(a, s) \theta(s) \right) \mu(d\theta) \right) \\ &= \sum_{s \in S} r(a, s) \theta^*(s). \end{aligned}$$

We have established that, for every  $\sigma \in \Delta(A)$ ,  $\sum_{s \in S} r(\bar{a}, s) \theta^*(s) \geq \sum_{a' \in A} \sigma(a') \times \sum_{s \in S} r(a', s) \theta^*(s)$ . Now, define the function  $\bar{v} : \Delta(A) \times \Delta(\text{supp } \mu) \rightarrow \mathbb{R}$  as  $\bar{v}(\sigma, \nu) = \bar{V}_{\text{Id}}(\bar{a}, \nu) - \sum_{a' \in A} \sigma(a') \bar{V}_{\text{Id}}(a', \nu)$ . Then, for every  $\sigma \in \Delta(A)$ ,  $\bar{v}(\sigma, \delta_{\theta^*}) \geq 0$ , and so:

$$\min_{\sigma \in \Delta(A)} \max_{\nu \in \Delta(\text{supp } \mu)} \bar{v}(\sigma, \nu) \geq 0.$$

By Sion’s minimax theorem in Sion (1958),

$$\max_{\nu \in \Delta(\text{supp } \mu)} \min_{\sigma \in \Delta(A)} \bar{v}(\sigma, \nu) = \min_{\sigma \in \Delta(A)} \max_{\nu \in \Delta(\text{supp } \mu)} \bar{v}(\sigma, \nu) \geq 0.$$

The function  $\bar{v}$  is continuous and its domain is compact, so the set of maximizers is non-empty.<sup>36</sup> Pick any:

<sup>36</sup> Compactness of the domain of  $\bar{v}$  follows from the fact that  $\text{supp } \mu(\cdot|h_t)$  is a closed subset of  $\Delta$ .

$$\mu' \in \arg \max_{\nu \in \Delta(\text{supp } \mu)} \left( \min_{\sigma \in \Delta(A)} \bar{v}(\sigma, \nu) \right).$$

We have, for every  $a' \in A$ ,

$$\bar{V}_{\text{Id}}(\bar{a}, \mu') - \bar{V}_{\text{Id}}(a', \mu') \geq \min_{a' \in A} \bar{v}(\delta_{a'}, \mu') \geq 0.$$

It follows that  $(\bar{a}, \mu', \bar{\theta})$  is an SCE under ambiguity neutrality. By Theorem 1 in Battigalli et al. (2015),  $(\bar{a}, \mu', \bar{\theta})$  is an SCE under ambiguity attitude  $\phi'$ .  $\square$

**Proof of Proposition 9.** Since  $(\alpha, \mu)$  is rational under  $\phi$ , and the DM is myopic, the one-deviation property reads:

$$\forall a \in A, \phi^{-1} \left( \int_{\Theta} \phi \left( R(\alpha(\nu), \hat{\theta}) \right) \nu(d\hat{\theta}) \right) \geq \phi^{-1} \left( \int_{\Theta} \phi \left( R(a, \hat{\theta}) \right) \nu(d\hat{\theta}) \right),$$

where  $\nu$  is as in the proof of Proposition 5. Let  $(\beta, \mu)$  be rational under  $\phi'$ . Suppose, by way of contradiction, that  $\beta(\nu)$  prescribes an ambiguous action. Then,

$$\begin{aligned} \hat{V}_{\phi'}(\beta, \nu) &= (\phi')^{-1} \left( \int_{\Theta} \phi' \left( R(\beta(\nu), \hat{\theta}) \right) \nu(d\hat{\theta}) \right) \\ &\geq (\phi')^{-1} \left( \int_{\Theta} \phi' \left( R(\alpha(\nu), \hat{\theta}) \right) \nu(d\hat{\theta}) \right) \\ &= R(\alpha(\nu), p_{\nu}), \end{aligned}$$

where the inequality follows from the one-deviation property, while the second equality follows from the assumption that  $\alpha(\nu)$  is unambiguous. Therefore,

$$\hat{V}_{\phi'}(\beta, \nu) \geq R(\alpha(\nu), p_{\nu}).$$

By Jensen’s inequality, since  $\phi'$  is a strictly concave transformation of  $\phi$ ,

$$R(\alpha(\nu), p_{\nu}) = \hat{V}_{\phi}(\alpha, \nu) \geq \hat{V}_{\phi}(\beta, \nu) > \hat{V}_{\phi'}(\beta, \nu) \geq R(\alpha(\nu), p_{\nu}),$$

a contradiction.  $\square$

**Proof of Proposition 10.** The proof is very similar to that of Proposition 6, invoking Proposition 9 instead of Proposition 5. Further details are omitted.  $\square$

**Proof of Proposition 11.** As with the previous proposition, invoke Proposition 9 instead of Proposition 5.  $\square$

**References**

Anderson, C., 2012. Ambiguity aversion in multi-armed bandit problems. *Theory Decis.* 72, 15–33.  
 Arrow, K.J., 1971. *Essays in the Theory of Risk-Bearing*. Markham.  
 Arrow, K.J., Green, J.R., 1973. *Notes on Expectations Equilibria in Bayesian Settings*. Institute for Mathematical Studies in the Social Sciences. Working Paper 33.

- Battigalli, P., Catonini, E., Lanzani, G., Marinacci, M., 2019. Ambiguity attitudes and self-confirming equilibrium in sequential games. *Games Econ. Behav.* 115, 1–29.
- Battigalli, P., Cerreia-Vioglio, S., Maccheroni, F., Marinacci, M., 2015. Self-confirming equilibrium and model uncertainty. *Am. Econ. Rev.* 105, 646–677.
- Billingsley, P., 2012. *Probability and Measure*. John Wiley & Sons, Ltd.
- Blanchard, O., 1985. Debt, deficit, and finite horizon. *J. Polit. Econ.* 93, 223–247.
- Cerreia-Vioglio, S., Maccheroni, F., Marinacci, M., Montrucchio, L., 2013a. Ambiguity and robust statistics. *J. Econ. Theory* 148, 974–1049.
- Cerreia-Vioglio, S., Maccheroni, F., Marinacci, M., Montrucchio, L., 2013b. Classical subjective expected utility. *Proc. Natl. Acad. Sci. USA* 110, 6754–6759.
- Doob, J.L., 1949. Application of the theory of martingales. *Colloq. Int. Cent. Natl. Rech. Sci.*, 23–27.
- Dynkin, E., 1965. Controlled random sequences. *Theory Probab. Appl.* 10, 1–14.
- Easley, D., Kiefer, N.M., 1988. Controlling a stochastic process with unknown parameters. *Econometrica* 5, 1045–1064.
- Epstein, L.G., Schneider, M., 2007. Learning under ambiguity. *Rev. Econ. Stud.* 74, 1275–1303.
- Esponda, I., Pouzo, D., 2016. Berk–Nash equilibrium: a framework for modeling agents with misspecified models. *Econometrica* 84, 1093–1130.
- Folland, G., 2013. *Real Analysis: Modern Techniques and Their Applications*. John Wiley & Sons, Ltd.
- Fudenberg, D., He, S., 2018. Learning and type compatibility in signalling games. *Econometrica* 86, 1215–1255.
- Fudenberg, D., Kreps, D.M., 1995. Learning in extensive-form games, I: self-confirming equilibria. *Games Econ. Behav.* 8, 20–55.
- Fudenberg, D., Levine, D.K., 1993. Steady state learning and Nash equilibrium. *Econometrica* 61, 547–573.
- Fudenberg, D., Levine, D.K., 1998. *The Theory of Learning in Games*. MIT Press, Cambridge, MA.
- Fudenberg, D., Tirole, J., 1991. *Game Theory*. MIT Press, Cambridge, MA.
- Gilboa, I., Schmeidler, D., 1989. Maxmin expected utility with a non-unique prior. *J. Math. Econ.* 18, 141–153.
- Gittins, J.C., 1989. *Multi-Armed Bandit Allocation Indices*. Wiley-Interscience Series in Systems and Optimization. John Wiley & Sons, Ltd., Chichester.
- Hanany, E., Klibanoff, P., 2009. Updating ambiguity averse preferences. *B.E. J. Theor. Econ.* 9.
- Hanany, E., Klibanoff, P., Mukerji, S., 2019. Incomplete information games with ambiguity averse players. *Am. Econ. J. Microecon.* Forthcoming.
- Hellwig, M., Leininger, W., 1987. On the existence of subgame-perfect equilibrium in infinite-action games of perfect information. *J. Econ. Theory* 43, 55–75.
- Hinderer, K., 1970. *Foundations of Non-stationary Dynamic Programming with Discrete Time Parameter*. Springer Science & Business Media.
- Jewitt, I.E., Mukerji, S., 2017. Ordering ambiguous acts. *J. Econ. Theory* 171, 213–267.
- Kalai, E., Lehrer, E., 1993. Rational learning leads to Nash equilibrium. *Econometrica* 1, 1019–1045.
- Kalai, E., Lehrer, E., 1995. Subjective games and equilibria. *Games Econ. Behav.* 8, 123–163.
- Klibanoff, P., Marinacci, M., Mukerji, S., 2005. A smooth model of decision making under ambiguity. *Econometrica* 73, 1849–1892.
- Li, J., 2019. The k-armed bandit problem with multiple-priors. *J. Math. Econ.* 80, 22–38.
- Maccheroni, F., Marinacci, M., Rustichini, A., 2006. Dynamic variational preferences. *J. Econ. Theory* 128, 4–44.
- Marinacci, M., 2002. Learning from ambiguous urns. *Stat. Pap.* 43, 143–151.
- Marinacci, M., 2015. Model uncertainty. *J. Eur. Econ. Assoc.* 13, 998–1076.
- Rothschild, M., 1974. A two-armed bandit theory of market pricing. *J. Econ. Theory* 9, 185–202.
- Sion, M., 1958. On general minimax theorems. *Pac. J. Math.* 171 (176).
- Williams, D., 1991. *Probability with Martingales*. Cambridge University Press, Cambridge.