

Learning from ambiguous and misspecified models*

Massimo Marinacci

Department of Decision Sciences and IGIER, Bocconi University

Filippo Massari

School of Banking and Finance, UNSW

December 22, 2018

Abstract

We model inter-temporal ambiguity as the scenario in which a Bayesian learner holds more than one prior distribution over a set of models and provide sufficient conditions for ambiguity to fade away because of learning. Our conditions apply to most learning environments: iid and non-iid model-classes, well-specified and misspecified model-classes/prior support pairs. We show that ambiguity fades away if the empirical evidence supports a set of models with identical predictions, a condition much weaker than learning the truth.

Keywords: Ambiguity, Learning.

JEL Classification: D81, D83, C11

1 Introduction

Let M be a family of models and \mathcal{C} a set of prior distributions on it. If \mathcal{C} contains more than one prior distribution, its multiplicity represents the *a priori* ambiguity

*We thank Werner Ploberger for his comments. Massimo Marinacci acknowledges the financial support of the European Research Council (advanced grant INDIMACRO).

perceived by a Bayesian decision maker (DM). This setting has been used to highlight the interaction between learning and ambiguity.¹

Marinacci (2002) formalizes the intuition that if a DM observes repeated draws (with replacement) from the same ambiguous urn, ambiguity fades away over time because he eventually learns the true composition. If the learning problem is well-specified — in the sense that the true probability belongs to the model-class/prior support pair adopted by the DM — ambiguity fades away because all posterior distributions converge to a Dirac distribution on the true model.

Here, we generalize the result in Marinacci (2002) to the case in which the DM does not learn the true probability because his prior view of the world is incorrect — that is, when the learning problem is misspecified in the sense that the model-class/prior support pair does not contain the true model/parameter. We show that ambiguity fades away if the data clearly designates a unique most accurate model (or a set of models with equivalent predictions), a condition that is always satisfied in well-specified learning problems but in some misspecified learning problems. In a nutshell, ambiguity fades away in all cases in which the empirical evidence eventually dominates the effect of heterogeneity in the prior distributions. On the contrary, ambiguity persists in those sequences in which two or more models with different predictions have comparable likelihood infinitely often. When this happens, the posteriors are “split” between these models with weights that depend on the priors, and the DM perceives ambiguity.

Our key contribution is to formalize sufficient conditions for the posteriors obtained from all priors to concentrate on the same model. Our findings rely on and generalize standard results in statical learning theory. With a unique prior, a sufficient condition for the Bayesian posterior to concentrate on the true model (consistency) is that the prior μ attaches a positive mass to the true parameter(s) (Doob, 1949; Freedman, 1963). In a multiple priors setting, this result continues to hold: if all priors give positive mass to the true model, then all posteriors concentrate on it and ambiguity fades away (Marinacci, 2002). On the other hand, in an iid setting and if the true

¹Epstein and Schneider (2003) provides an axiomatization of prior-by-prior updating which requires the process of conditional preferences to be dynamically consistent. Because we are focusing on one-step-ahead decisions, the consistency issue has no bite in our setting.

parameter set does not belong to the prior support, the posterior concentrates on the model that is the closest in terms of K-L divergence to the truth if it is unique (Berk, 1966; White, 1982). In a multiple priors setting, this result suggests that if the minimizer of the K-L divergence, P^* , is unique and all priors give it a positive weight, then ambiguity fades away because all posteriors concentrate on P^* .

Theorem 2 proves this conjecture and generalizes it to the non-iid setting. Theorem 1 provides an empirical condition for the posteriors derived from all priors to concentrate on a unique model which does not depend on a-priori knowledge of the truth. Last, Theorem 3 drops the uniqueness requirement for the most accurate model providing a condition for ambiguity to fade away because the posteriors derived from all priors to concentrate on a set of models with identical predictions.

Discussion We prove that a Bayesian agent with multiple priors does not suffer from long-run ambiguity in all those cases in which the data support a unique model (or a set of models with identical predictions). How common are these situations? A precise answer to this question is hard to give because it depends on the true probability measure, which is typically unknown. If M counts finitely many iid probabilistic models, then the set of parameters characterizing an iid data generating process such that at least two models in M have identical average K-L divergence (a situation that violates all our sufficient conditions and may generate long-run ambiguity) is nongeneric, thus suggesting that ambiguity should be the exception, rather than the norm. However, we are cautious about concluding that ambiguity typically fades away in real world situations because models and parameters are hardly iid and chosen at random. For example, consider the standard problem of predicting stock market returns. Several models have been proposed and, to date, it is not clear which model is the closest to the truth — there is no definite statistical test that favors one unique model over another. Because the empirical evidence does not support a unique model, an investor with a set of priors on available models of stock market returns suffers ambiguity despite the large amount of available financial data.

Our, multiple-prior, model describes a DM who is uncertain about the a priori

probability of each model in the support but updates each model in a unique way. On the other hand, the multiple-likelihoods model (e.g., Epstein and Schneider, 2007; Epstein and Seo, 2015) describes a DM who believes that signals have multiple, hence uncertain, interpretations. Such signals can generate ambiguity even where none is present a priori. Learning models that accommodate such a possibility generate posterior sets different from those defined in this paper, and they lead to different results regarding if/when ambiguity fades away.

2 Probabilities

We consider a family of models $M = \{P_\theta : \theta \in \Theta\}$ with a finite parameter set $\Theta \subset \mathbb{R}^n$, defined on a σ -algebra Σ^∞ of subsets of X^∞ with representative element $x^\infty = x_1, x_2, \dots$, where $X^\infty := \times^\infty X$ is the infinite Cartesian product of a finite observation space X with representative element x and σ -algebra Σ . With a slight abuse of notation, we use $P_\theta(x^t)$ to denote the probability that model P_θ attaches to the cylinder with base x^t , as well as the likelihood that model P_θ attaches to the partial sequence (x_1, \dots, x_t) . The prior information about the parameters is summarized by prior distributions $\mu \in \Delta\Theta$. The set of prior distributions is \mathcal{C} . For any prior distribution $\mu \in \mathcal{C}$ the joint distribution of the parameters and the observations is $P_\mu \in \Delta(\Theta \times X^\infty)$. Defined by, for all sets $A \subseteq \Theta$ and all cylinders x^t ,

$$P_\mu(A \times x^t) := \int_A P_\theta(x^t) d\mu.$$

We denote by $\mu(\cdot|x^t) \in \Delta\Theta$ the usual posterior given the observations x^t ,² while $P_\mu(\cdot|x^t) \in \Delta(\Theta \times X)$ is the one-step-ahead predictive distribution of x_{t+1} , given observations x^t . By definition, for all $A \subseteq \Theta$ we have

$$P_\mu(A \times x_{t+1}|x^t) := \int_A P_\theta(x_{t+1}|x^t) d\mu(\cdot|x^t) := \int_A P_\theta(x_{t+1}|x^t) \frac{P_\theta(x^t) d\mu}{\int_\Theta P_\theta(x^t) d\mu}.$$

²We rule out the possibility of observing an event which is impossible according to all models in M .

3 Decisions

Let C be the space of consequences on which the DM has a bounded utility function $u : C \rightarrow \mathbb{R}$. We consider one-step-ahead acts, i.e., Σ -measurable maps $f : X \rightarrow C$ that associates a consequence to each observation in X . The decision criterion adopted by the DM depends on the quality of his prior information. For illustrative purposes, we briefly provide examples of the DM's decision criterion when facing *risk*, *unambiguous uncertainty*, and *ambiguity*.

Suppose there is an urn with 3 balls, each of which is either white, x_W , or red, x_R . Suppose the DM chooses a color and draws a ball from the urn. If this ball matches the DM's color, he wins \$100. Otherwise, he gets nothing. The consequence space is $C = \{\$0, \$100\}$, the observation space $X = \{x_R, x_W\}$, and the DM can choose between two acts: f_R , he bets on a red ball; and f_W , he bets on a white ball. The following table summarizes this decision problem:

	x_R	x_W	
f_R	\$100	0	(1)
f_W	0	\$100	

Finally, θ is the fraction of white balls in the urn, so that $\Theta = \{0, 1/3, 2/3, 1\}$. If draws are made with replacement from the same urn, M is the iid Bernoulli distribution family with parameter set Θ .

- **Scenario 1: Risky Urns.** The DM knows the true composition of the urn θ_0 (e.g., he knows that it contains exactly two white balls). In this case, the DM's choice criterion is, for every act f , given by:

$$\int_X u(f(x)) dP_{\theta_0}.$$

- **Scenario 2: Bayesian Urns.** The DM does not know the composition of the urn but has enough prior information to uniquely pin down a prior distribution μ on the set of possible compositions Θ . That is, \mathcal{C} is a singleton. For example, the DM might believe that all the compositions of the urn are equally likely. Unlike

the previous case, the DM's choice criterion now changes over time because of learning. In the first period, the DM's choice criterion is, for every act f , given by:

$$\int_{\Theta} \left[\int_X u(f(x)) dP_{\theta} \right] d\mu = \int_X u(f(x)) dP_{\mu}(x|\emptyset).$$

Subsequently, as the DM incorporates past realization, x^t , to his prior distribution using Bayes' rule, his choice criterion becomes:

$$\int_{\Theta} \left[\int_X u(f(x)) dP_{\theta} \right] d\mu(\cdot|x^t) = \int_X u(f(x)) dP_{\mu}(x|x^t).$$

- **Scenario 3: Ambiguous Urns.** The DM does not know the composition of the urn and does not have enough prior information to uniquely pin down a distribution on the set of possible compositions of the urn. That is, \mathcal{C} is not a singleton. For example, the DM might only know that every composition has at least a 1/10 probability of being the correct one: $\mathcal{C} := \{\mu \in \Delta : \forall \theta \in \Theta, \mu(\theta) \geq 1/10\}$. In evaluating an act in this scenario, the DM has to consider, for each act f , the set

$$\left\{ \int_X u(f(x)) dP_{\mu}(x|\emptyset) : \mu \in \mathcal{C} \right\}.$$

Subsequently, as the DM incorporates past realizations using Bayes' rule, the DM has to consider, for each act f , the set:

$$\left\{ \int_X u(f(x)) dP_{\mu}(x|x^t) : \mu \in \mathcal{C} \right\}.$$

Possible summaries of this set are the infimum and supremum:

$$\sup_{\mu \in \mathcal{C}} \int_X u(f(x)) dP_{\mu}(x|x^t) ; \quad \inf_{\mu \in \mathcal{C}} \int_X u(f(x)) dP_{\mu}(x|x^t).$$

4 Long-run ambiguity

As in Marinacci (2002), we consider the difference between the DM's expected utility under the most advantageous prior and under the least advantageous prior in \mathcal{C} to be a

measure of the ambiguity that a DM perceives in evaluating an act f . We are ultimately interested in verifying whether this quantity converges to 0 as the number of past observations goes to infinity and each prior gets independently updated using Bayes' rule. A tight sufficient condition for the most conservative and the least conservative expected utility to coincide is to require that the posteriors calculated from all priors in \mathcal{C} eventually coincide (see Lemma 1 in Appendix). We say that

Definition 1. *Ambiguity fades away at path $x^\infty \in X^\infty$ if,*

$$\lim_{t \rightarrow \infty} \left[\sup_{\mu', \mu'' \in \mathcal{C}} \int_X |dP_{\mu''}(x|x^t) - dP_{\mu'}(x|x^t)| \right] = 0 \quad (2)$$

where, for each $t > 0$, x^t indicates the first t realizations of path x^∞ .

Definition 1 requires that all posteriors concentrate on the same model (or on a set of models with identical predictions) on the realized path. Unlike the definition proposed by Marinacci (2002) — which requires all the posteriors to converge to a Dirac measure on the true model on a set of sequences of true measure 1 — ours does not assume an iid structure, and it does not depend on the true model. Thus, it can be used to discuss long-run ambiguity when the model class support contains models with learning, a dependence-structure, or is misspecified. In those cases in which all posteriors concentrate on the true model, our definition is equivalent to the notion of *weak merging* (Lehrer and Smorodinsky, 1996).

5 Main result

In this section, we present three sufficient conditions for ambiguity to fade away. The driving force of our results is the observation that the key component of Bayesian learning is the existence of a unique most accurate model, rather than the true model belonging to the prior support. For instance, Berk (1966) shows that if all models in the support and the truth are iid, then the posterior obtained from a unique prior eventually assigns probability 1 to the set of parameters that minimize the K-L divergence from the truth, if unique. Here, we generalize Berk (1966)'s result to the case of multi-prior,

non-iid setting. We provide two conditions (Theorems 1 and 2) that are sufficient for all posteriors to concentrate on the same model and one condition (Theorem 3) that is sufficient for all posteriors to concentrate on a set of models that deliver the same predictions. Let us start by formalizing an appropriate generalization of the notion of unique most accurate model.

Definition 2. *Given a path $x^\infty \in X^\infty$ and a family of models $M = \{P_\theta : \theta \in \Theta\}$, we say that $\hat{\theta} := \hat{\theta}(x^\infty, \Theta) \in \Theta$ is a strong maximum likelihood (SML) model if, for every $\theta \in \Theta$, the limit $\lim_{t \rightarrow \infty} P_\theta(x^t)/P_{\hat{\theta}}(x^t)$ exists and is finite.*

If $\hat{\theta}$ is unique, the SML model is the model whose likelihood converges to zero at the slowest rate on path x^∞ — i.e., $\forall \theta \neq \hat{\theta}, \lim_{t \rightarrow \infty} P_\theta(x^t)/P_{\hat{\theta}}(x^t) = 0$ on x^∞ . Otherwise, let $\Theta^* \subset \Theta$ be the subset of models of Θ whose likelihoods converge to zero at the slowest rate on path x^∞ — i.e., $\forall \theta \in \Theta \setminus \Theta^*, \forall \hat{\theta} \in \Theta^*, \lim_{t \rightarrow \infty} P_\theta(x^t)/P_{\hat{\theta}}(x^t) = 0$ on x^∞ . The set of SML models is non-empty if and only if all models in Θ^* eventually deliver the same next period prediction on x^∞ .

Our first result shows that the existence of a unique SML is a sufficient condition for ambiguity to fade away.

Theorem 1. *Let $M = \{P_\theta : \theta \in \Theta\}$ be a family of models and \mathcal{C} a compact set of strictly positive prior distributions on Θ . Ambiguity fades away at path x^∞ if a unique $\hat{\theta}(x^\infty, \Theta)$ exists.*

By definition, if the SML model is unique, its likelihood eventually dominates that of all other models. Consequently, the posteriors calculated from all priors eventually attach unitary weight to $\hat{\theta}$ and ambiguity fades away.

An alternative condition for ambiguity to fade away can be obtained by noticing that a sufficient condition for the existence of a unique SML model P -a.s. is the presence of a unique model with the lowest average K-L divergence.

Definition 3. *The average K-L divergence from P_θ to the true probability P_{θ_0} is*

$$\bar{D}(P_{\theta_0} || P_\theta) := \lim_{t \rightarrow \infty} \frac{1}{t} E_{P_{\theta_0}} \left[\ln \frac{P_{\theta_0}(x^t)}{P_\theta(x^t)} \right].$$

Next, we show that the existence of a unique model with minimal K-L divergence is a sufficient condition for ambiguity to fade away.

Theorem 2. *Let $M = \{P_\theta : \theta \in \Theta\}$ be a family of models and \mathcal{C} a compact set of strictly positive priors on Θ . Ambiguity fades away P_{θ_0} -a.s. if $\operatorname{argmin}_{\theta \in \Theta} \bar{D}(P_{\theta_0} || P_\theta)$ exists and is unique.*

This delivers a sufficient condition for ambiguity to fade away P_{θ_0} -a.s. which generalizes Berk (1966)'s results to the non-iid setting and includes Marinacci (2002)'s condition as a special case. When \mathcal{C} is a singleton, if all models in M and the true measure are iid, Berk (1966)'s result follows from Theorem 2 because the average K-L divergence coincides with the K-L divergence P_{θ_0} -a.s. by the Strong Law of Large Numbers. Whereas, if all models in M are iid and the truth belongs to M , Marinacci (2002)'s condition follows because the true model is the unique maximizer of the K-L divergence.

The condition of Theorem 1 depends only on the properties of the sequence of realizations, while Theorem 2's condition requires a priori knowledge of the true probability distribution—to calculate the K-L divergence. In a nutshell, the difference between the two conditions is as follows: Theorem 1 tells us that ambiguity persists if the data is inconclusive, while Theorem 2 tells us that ambiguity persists if the true probability generates inconclusive data P_{θ_0} -a.s.. These two conditions are not directly comparable, but there is a sense in which the former is more informative than the latter. There are cases in which two models have identical average K-L divergence but a diverging likelihood ratio P -a.s.³

Last, we present a sufficient condition for ambiguity to fade away that relaxes the uniqueness requirement for the SML model of Theorem 1 by asking for a (minimal) assumption about the true model. If all states have a positive probability to be visited after every history, then all SML models must deliver identical next period predictions, and the existence of at least one SML model is a sufficient condition for ambiguity to

³For example, in a market selection framework Massari (2017) analyzes models with diverging log-likelihood ratio and yet the same average K-L divergence because the divergence rate of the former $O(t^{-5})$ is dominated by the averaging factor $O(t^{-1})$.

fade away. Theorem 3 condition allows discussing long-run ambiguity when members of the prior support are learning models.

Theorem 3. *Let $M = \{P_\theta : \theta \in \Theta\}$ be a family of models and \mathcal{C} a compact set of strictly positive priors on Θ . Ambiguity fades away P_{θ_0} -a.s. if at least a SML exists and the true data generating process P_{θ_0} assigns strictly positive probability to every state after every history.⁴*

The additional assumption about the true model that we introduce in Theorem 3 is needed to guarantee that if there are multiples SML, all the SML deliver the same prediction. Scenario 9 below illustrates the role played by our conditions on the truth.⁵ If some states have zero probability to be visited, then it is possible to have multiple SML models and long-run ambiguity because these models make different predictions about states that are never empirically tested.

We conclude by presenting scenarios illustrating our conditions. Scenarios 4-6 show cases in which ambiguity fades away because the truth generates sequences that support only one model among those believed possible by the DM. Conversely, scenarios 6-9 show cases in which ambiguity persists in the long run because the truth generates sequences which equally endorse at least two models with different predictions in the prior support infinitely often. While we were unable to prove a necessary counterpart to our sufficient conditions, these scenarios suggest that our conditions are tight.

Suppose a DM confronted with decision problem 1 subjectively believes that he is facing iid realizations from an ambiguous urn with three balls, two of which have the same color. In our notation, he believes that M is the class of iid Bernoulli distributions with possible parameters $\Theta = \{1/3, 2/3\}$. His prior information is accurate enough to reduce ambiguity to only two possible priors on the composition of the urn: $\mathcal{C} = \{\mu'(\theta), \mu''(\theta)\}$. These are,

$$\mu'(\theta) = \{\mu'(\frac{1}{3}) = \frac{1}{2}, \mu'(\frac{2}{3}) = \frac{1}{2}\} \text{ and } \mu''(\theta) = \{\mu''(\frac{1}{3}) = \frac{1}{4}, \mu''(\frac{2}{3}) = \frac{3}{4}\}.$$

⁴In symbols, $\forall x^{t-1}, \forall x_t \in X, P_{\theta_0}(x_t|x^{t-1}) > 0$.

⁵We thank an anonymous referee for providing us with this example.

- **Scenario 4: Ambiguity fades away in well-specified learning settings.** Draws are indeed iid from an urn whose composition is $\theta_0 = \frac{2}{3}$. Because the learning problem is well-specified ($\theta_0 \in \Theta$), by the Strong Law of Large Numbers θ_0 is the SML model P_{θ_0} -a.s.. Thus, Theorem 1 (and Marinacci, 2002) implies that ambiguity fades away P_{θ_0} -a.s..
- **Scenario 5: Ambiguity fades away in a misspecified learning setting with incorrect dependence-structure.** Draws are not iid, as the DM incorrectly believes. Instead, the urn is secretly changed before every draw to deliver the deterministic sequence $x^\infty := \{W, W, R, W, W, R, \dots\}$. Because the frequency of W converges to $\frac{2}{3}$, then $\hat{\theta} = \frac{2}{3}$ is the SML parameter in Θ . By Theorem 1, all posteriors concentrate on $\hat{\theta}$ and ambiguity fades away. Although the DM fails to realize that draws are not iid, he successfully learns the best parameter in Θ and ambiguity fades away.
- **Scenario 6: Ambiguity fades away in a misspecified learning setting with correct dependence-structure.** Draws are iid from an urn whose composition is $\theta_0 = \frac{3}{5}$. In this case, the dependence-structure is correctly specified because draws are indeed iid, but the learning problem is misspecified because Θ does not contain the true parameter: $\theta_0 \notin \Theta$. It is easy to verify that $\hat{\theta} = \frac{2}{3}$ is the strong maximum likelihood P_{θ_0} -a.s..⁶ Thus, Theorem 1 implies that both posteriors concentrate on $\hat{\theta}$. Although the DM cannot learn the true model, ambiguity fades away because the data clearly indicates which model is the most accurate.
- **Scenario 7: Ambiguity persists in a misspecified learning setting with incorrect dependence-structure.** Draws are not iid, as the DM incorrectly believes. Instead, he is facing the deterministic sequence $x^\infty := \{W, R, W, R, \dots\}$. It is easy to verify that the conditions of Theorems 1-3 are not satisfied. The following argument shows that ambiguity does not fade away. In every even period $P_{\frac{1}{3}}$ and $P_{\frac{2}{3}}$ have identical likelihood. Therefore, each odd period prediction obtained from priors μ' and μ'' coincides with their first period prediction. Because

⁶By the Strong Law of Large Numbers, $\lim_{t \rightarrow \infty} \frac{P_{\frac{1}{3}}(x^t)}{P_{\frac{2}{3}}(x^t)} = P\text{-a.s.} \lim_{t \rightarrow \infty} \left(\frac{\left(\frac{1}{3}\right)^{\frac{3}{5t}} \left(\frac{2}{3}\right)^{\frac{2}{5t}}}{\left(\frac{2}{3}\right)^{\frac{3}{5t}} \left(\frac{1}{3}\right)^{\frac{2}{5t}}} \right)^t = 0$.

μ' 's and μ'' 's first period predictions differ, their predictions differ for every odd period, and ambiguity does not fade away.⁷

- **Scenario 8: Ambiguity persists in a misspecified learning setting with correct dependence-structure.** Draws are iid from an urn whose composition is $\theta_0 = \frac{1}{2}$. In this case, the dependence-structure is correctly specified because draws are iid, but Θ is not, since $\theta_0 \notin \Theta$. A symmetry argument can be used to show that the conditions of Theorems 1-3 are not satisfied.⁸ When this happens the predictions of μ' and μ'' differ and the DM suffers ambiguity, by the same argument used in Scenario 5.
- **Scenario 9: Ambiguity persists in a misspecified learning setting with multiple SML and degenerate truth.** Consider a DM facing a series of draws from three-color urns $X := \{a, b, c\}$ which he believe to be iid from two possible models $\Theta = \{\theta', \theta''\}$ with

$$P_{\theta'} = \left[\frac{1}{2}, \frac{3}{8}, \frac{1}{8} \right] \text{ and } P_{\theta''} = \left[\frac{1}{2}, \frac{1}{8}, \frac{3}{8} \right].$$

The DM has two priors $\mathcal{C} = \{\mu, \mu'\}$ on these models: $\mu(\theta') = .3 = 1 - \mu(\theta'')$. Draws are not iid, as the DM incorrectly believes. Instead, he is facing the deterministic sequence $x^\infty := \{a, a, a, \dots\}$. The condition of Theorem 1 is not satisfied because both models are SML — they attach an identical likelihood to x^∞ . The condition of Theorem 2 is not satisfied because both models have the same average K-L divergence.⁹ The condition of Theorem 3 is not satisfied because only state a has positive probability of occurring. Ambiguity does not fade away because the

⁷It is straightforward to verify that $P_{\mu'}(x_W|x^t) \neq P_{\mu''}(x_W|x^t)$ for every t even:

$$P_{\mu'}(x_W|x^t) = \frac{1}{2} \frac{\left(\frac{1}{3}\frac{2}{3}\right)^{\frac{t}{2}} \frac{1}{3}}{\left(\frac{1}{3}\frac{2}{3}\right)^{\frac{t}{2}} \frac{1}{2} + \left(\frac{1}{2}\frac{2}{3}\right)^{\frac{t}{2}} \frac{1}{2}} + \frac{1}{2} \frac{\left(\frac{1}{3}\frac{2}{3}\right)^{\frac{t}{2}} \frac{2}{3}}{\left(\frac{1}{3}\frac{2}{3}\right)^{\frac{t}{2}} \frac{1}{2} + \left(\frac{1}{3}\frac{2}{3}\right)^{\frac{t}{2}} \frac{1}{2}} = \frac{1}{2}$$

$$P_{\mu''}(x_W|x^t) = \frac{1}{4} \frac{\left(\frac{1}{3}\frac{2}{3}\right)^{\frac{t}{2}} \frac{1}{3}}{\left(\frac{1}{3}\frac{2}{3}\right)^{\frac{t}{2}} \frac{1}{4} + \left(\frac{3}{4}\frac{2}{3}\right)^{\frac{t}{2}} \frac{3}{4}} + \frac{3}{4} \frac{\left(\frac{1}{3}\frac{2}{3}\right)^{\frac{t}{2}} \frac{2}{3}}{\left(\frac{1}{3}\frac{2}{3}\right)^{\frac{t}{2}} \frac{1}{4} + \left(\frac{1}{3}\frac{2}{3}\right)^{\frac{t}{2}} \frac{3}{4}} = \frac{7}{12}$$

⁸Massari (2013) provides a formal argument based on an application of the Law of Iterated Logarithms showing that $P_{\frac{2}{3}}$ and $P_{\frac{1}{3}}$ have identical likelihood infinitely often.

⁹Under the standard convention $0 \ln 0 = 0$

predictive probabilities for state b and c calculated from the two time zero priors differ, and the two priors remain constant over time since $P_{\theta'}(a) = P_{\theta''}(a)$ and only state a occurs.

6 Appendix

In this appendix,

- Given two functions on the real line, f and g , $f = o(g)$, abbreviates $\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = 0$;
- $\hat{\theta}_t := \hat{\theta}(x^t)$ is the maximum likelihood model on the partial history x^t ;
- $D(P_{\hat{\theta}_t} || P_{\theta}) := E_{P_{\hat{\theta}_t}} \ln \frac{P_{\hat{\theta}_t}(x)}{P_{\theta}(x)}$ is the K-L divergence from P_{θ} to $P_{\hat{\theta}_t}$.

Proof of Theorems 1 and 3

Proof. Because \mathcal{C} is compact, then $\operatorname{argmax}_{\mu', \mu'' \in \mathcal{C}} \lim_{t \rightarrow \infty} \int_X |dP_{\mu''}(x|x^t) - dP_{\mu'}(x|x^t)|$ is non empty. Thus, it suffices to prove that

$$\exists SML \Leftrightarrow \forall \mu', \mu'' \in \mathcal{C}, \lim_{t \rightarrow \infty} \int_X |dP_{\mu''}(x|x^t) - dP_{\mu'}(x|x^t)| = 0.$$

- Theorem 1: by assumption, there is a unique SML, $\hat{\theta}$. Therefore,

$$\begin{aligned} & \lim_{t \rightarrow \infty} \int_X |dP_{\mu''}(x|x^t) - dP_{\mu'}(x|x^t)| \\ &:= \lim_{t \rightarrow \infty} \int_X \left| \sum_{\theta} P_{\theta}(x) \frac{P_{\theta}(x^t) \mu''(\theta)}{\sum_{\theta} P_{\theta}(x^t) \mu''(\theta)} - \sum_{\theta} P_{\theta}(x) \frac{P_{\theta}(x^t) \mu'(\theta)}{\sum_{\theta} P_{\theta}(x^t) \mu'(\theta)} \right| dx \\ &=^a \int_X \lim_{t \rightarrow \infty} \left| \frac{P_{\hat{\theta}}(x)}{1 + \sum_{\theta \neq \hat{\theta}} \frac{P_{\theta}(x^t) \mu''(\theta)}{P_{\hat{\theta}}(x^t) \mu''(\hat{\theta})}} + \frac{\sum_{\theta \neq \hat{\theta}} P_{\theta}(x) \frac{P_{\theta}(x^t) \mu''(\theta)}{P_{\hat{\theta}}(x^t) \mu''(\hat{\theta})}}{1 + \sum_{\theta \neq \hat{\theta}} \frac{P_{\theta}(x^t) \mu''(\theta)}{P_{\hat{\theta}}(x^t) \mu''(\hat{\theta})}} \right. \\ & \quad \left. - \frac{P_{\hat{\theta}}(x)}{1 + \sum_{\theta \neq \hat{\theta}} \frac{P_{\theta}(x^t) \mu'(\theta)}{P_{\hat{\theta}}(x^t) \mu'(\hat{\theta})}} - \frac{\sum_{\theta \neq \hat{\theta}} P_{\theta}(x) \frac{P_{\theta}(x^t) \mu'(\theta)}{P_{\hat{\theta}}(x^t) \mu'(\hat{\theta})}}{1 + \sum_{\theta \neq \hat{\theta}} \frac{P_{\theta}(x^t) \mu'(\theta)}{P_{\hat{\theta}}(x^t) \mu'(\hat{\theta})}} \right| dx \\ &= \int_X \left| \frac{P_{\hat{\theta}}(x)}{1 + o(1)} + o(1) - \frac{P_{\hat{\theta}}(x)}{1 + o(1)} - o(1) \right| dx \quad , \text{ by definition of SML;} \\ &= 0 \end{aligned}$$

Step a : The Lebesgue's Dominated Convergence Theorem allows exchanging integral and limit signs (Williams, 1991).¹⁰

¹⁰Let $\{r_t(x)\}_{t=1}^{\infty} := \{|P_{\mu''}(x|x^t) - P_{\mu'}(x|x^t)|\}_{t=1}^{\infty}$ and note that $|r_1|, |r_2|, \dots$ are bounded above.

- Theorem 3: If the SML is unique, the proof above holds. Otherwise, let $\hat{\Theta} \subseteq \Theta$ be the subset of models which satisfy the SML condition: $\forall \hat{\theta}, \bar{\theta} \in \hat{\Theta}, \lim_{t \rightarrow \infty} \frac{P_{\hat{\theta}}(x^t)}{P_{\bar{\theta}}(x^t)} = k > 0$. By Lemma 2, all models in $\hat{\Theta}$ eventually deliver the same predictions on all states. Call this prediction \bar{P} , then the result follows substituting \bar{P} and $\bar{\mu} = \sum_{\mu(\bar{\theta})} \mu(\theta)$ for $P_{\hat{\theta}}$ and $\hat{\mu}$ in Equation *a* above, respectively.

□

Proof of Theorem 2

Proof. Let $\hat{\theta}$ be the unique element of $\operatorname{argmin}_{\theta \in \Theta} \bar{D}(P_{\theta_0} || P_{\theta})$. Thus $\exists \epsilon > 0$:

$$\begin{aligned}
& \forall \theta \neq \hat{\theta}, \bar{D}(P_{\theta_0} || P_{\theta}) < \bar{D}(P_{\theta_0} || P_{\hat{\theta}}) - \epsilon \\
& \Rightarrow \forall \theta \neq \hat{\theta}, \lim_{t \rightarrow \infty} \frac{1}{t} E_{P_{\theta_0}} \left[\ln \frac{P_{\theta_0}(x^t)}{P_{\hat{\theta}}(x^t)} \right] - \lim_{t \rightarrow \infty} \frac{1}{t} E_{P_{\theta_0}} \left[\ln \frac{P_{\theta_0}(x^t)}{P_{\theta}(x^t)} \right] < -\epsilon \\
& \Rightarrow^a \forall \theta \neq \hat{\theta}, \lim_{t \rightarrow \infty} E_{P_{\theta_0}} \left[\frac{1}{t} \sum_{\tau=1}^t E_{P_{\theta_0}} \left[\ln \frac{P_{\theta_0}(x|x^{\tau-1})}{P_{\hat{\theta}}(x|x^{\tau-1})} \middle| x^{\tau-1} \right] - \frac{1}{t} \sum_{\tau=1}^t E_{P_{\theta_0}} \left[\ln \frac{P_{\theta_0}(x|x^{\tau-1})}{P_{\theta}(x|x^{\tau-1})} \middle| x^{\tau-1} \right] \right] < -\epsilon \\
& \Rightarrow^b \forall \theta \neq \hat{\theta}, \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau=1}^t \ln \frac{P_{\theta_0}(x_{\tau}|x^{\tau-1})}{P_{\hat{\theta}}(x_{\tau}|x^{\tau-1})} - \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau=1}^t \ln \frac{P_{\theta_0}(x_{\tau}|x^{\tau-1})}{P_{\theta}(x_{\tau}|x^{\tau-1})} < -\epsilon \quad P_{\theta_0}\text{-a.s.} \\
& \Rightarrow \forall \theta \neq \hat{\theta}, \lim_{t \rightarrow \infty} \sum_{\tau=1}^t \ln \frac{P_{\theta}(x_{\tau}|x^{\tau-1})}{P_{\hat{\theta}}(x_{\tau}|x^{\tau-1})} = -\infty \quad P_{\theta_0}\text{-a.s.} \\
& \Rightarrow \forall \theta \neq \hat{\theta}, \lim_{t \rightarrow \infty} \ln \frac{P_{\theta}(x^t)}{P_{\hat{\theta}}(x^t)} = -\infty \quad P_{\theta_0}\text{-a.s.} \\
& \Rightarrow \forall \theta \neq \hat{\theta}, \lim_{t \rightarrow \infty} \frac{P_{\theta}(x^t)}{P_{\hat{\theta}}(x^t)} = 0 \quad P_{\theta_0}\text{-a.s.} \\
& \Rightarrow \hat{\theta} \text{ is SML } P_{\theta_0}\text{-a.s.}
\end{aligned}$$

a) Telescoping the log and using the tower property of expectation.

b) The Strong Law of Large Numbers for Martingale Differences (Williams, 1991) allows substituting the limit average sum of conditional expected values with the limit average sum of realized values P -a.s.:

$$\forall \theta \in \Theta \quad \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau=1}^t \left(E_{P_{\theta_0}(x|x^{\tau-1})} \left[\ln \frac{P_{\theta_0}(x|x^{\tau-1})}{P_{\theta}(x|x^{\tau-1})} \right] - \ln \frac{P_{\theta_0}(x_{\tau}|x^{\tau-1})}{P_{\theta}(x_{\tau}|x^{\tau-1})} \right) = 0 \quad P_{\theta_0}\text{-a.s.}$$

□

Lemma 1. Let μ' and μ'' be two prior on Θ , if u is bounded, then

$$\lim_{t \rightarrow \infty} \int_X |dP_{\mu''}(x|x^t) - dP_{\mu'}(x|x^t)| = 0 \Rightarrow \lim_{t \rightarrow \infty} \left[\int_X u(f(x)) dP_{\mu''}(x|x^t) - \int_X u(f(x)) dP_{\mu'}(x|x^t) \right] = 0$$

Proof.

$$\begin{aligned}
& \lim_{t \rightarrow \infty} \int_X |dP_{\mu''}(x|x^t) - dP_{\mu'}(x|x^t)| = 0 \\
\Rightarrow \lim_{t \rightarrow \infty} A_t^1 &= \lim_{t \rightarrow \infty} \max_{x \in X} |u(f(x))| \int_X |dP_{\mu''}(x|x^t) - dP_{\mu'}(x|x^t)| = 0 && \text{because } u \text{ is bounded;} \\
\Rightarrow \lim_{t \rightarrow \infty} A_t^2 &= \lim_{t \rightarrow \infty} \int_X |u(f(x))| |dP_{\mu''}(x|x^t) - dP_{\mu'}(x|x^t)| = 0 && \text{because, } \forall t, A_t^1 \geq A_t^2 \geq 0; \\
\Rightarrow \lim_{t \rightarrow \infty} A_t^3 &= \lim_{t \rightarrow \infty} \left| \int_X u(f(x)) dP_{\mu''}(x|x^t) - \int_X u(f(x)) dP_{\mu'}(x|x^t) \right| = 0 && \text{because } \forall t, A_t^2 \geq A_t^3 \geq 0; \\
\Rightarrow \lim_{t \rightarrow \infty} & \left[\int_X u(f(x)) dP_{\mu''}(x|x^t) - \int_X u(f(x)) dP_{\mu'}(x|x^t) \right] = 0.
\end{aligned}$$

□

Lemma 2. *If the true data generating process P_{θ_0} assigns strictly positive probability to every state after every history, then all SML eventually attach identical probability to all states, i.e., if $\forall x^{t-1}, \forall x_t, P_{\theta_0}(x_t|x^{t-1}) > 0$,*

$$\lim_{t \rightarrow \infty} \frac{P_{\hat{\theta}}(x^t)}{P_{\bar{\theta}}(x^t)} = k \in (0, \infty) \Rightarrow \lim_{t \rightarrow \infty} \|P_{\hat{\theta}}(x|x^{t-1}) - P_{\bar{\theta}}(x|x^{t-1})\| = 0 \text{ } P_{\theta_0}\text{-a.s.}$$

Proof. We prove the contrapositive statement. If $\forall x^{t-1}, \forall x_t, P_{\theta_0}(x_t|x^{t-1}) > 0$,

$$\exists \epsilon > 0 : \|P_{\hat{\theta}}(x|x^{t-1}) - P_{\bar{\theta}}(x|x^{t-1})\| > \epsilon \text{ infinitely often} \Rightarrow \neg \lim_{t \rightarrow \infty} \frac{P_{\hat{\theta}}(x^t)}{P_{\bar{\theta}}(x^t)} = k \in (0, \infty) \text{ } P_{\theta_0}\text{-a.s.}$$

Our proof is an application of the Levy's extension of the Borel-Cantelli Lemma from which we borrow most notation (see, Williams, 1991, pg. 124.)

Consider two models $\hat{\theta}$ and $\bar{\theta}$ that deliver different predictions in at least a state infinitely often. Call the subsequence of such periods $(t_\tau)_{\tau=1}^\infty$, so that

$$\forall x_t, P_{\theta_0}(x_t) > 0 \Rightarrow \exists \epsilon_1 > 0 : \forall t_\tau, \xi_{t_\tau} := P_{\theta_0} \left(\frac{P_{\hat{\theta}}(x_{t_\tau}|x^{t_\tau-1})}{P_{\bar{\theta}}(x_{t_\tau}|x^{t_\tau-1})} \notin (1 \pm \epsilon_1) \right) > 0.$$

$$\text{Let } Z_n := \frac{1}{n} \sum_{\tau=1}^n I \left\{ \frac{P_{\hat{\theta}}(x_{t_\tau}|x^{t_\tau-1})}{P_{\bar{\theta}}(x_{t_\tau}|x^{t_\tau-1})} \notin (1 \pm \epsilon_1) \right\}, \text{ and } Y^n := \sum_{\tau=1}^n \xi_{t_\tau}.$$

Note that $\lim_{n \rightarrow \infty} Y^n = \infty$.

The Levy extension of the Borel-Cantelli Lemma guarantees that P_{θ_0} -a.s.,

$$\lim_{n \rightarrow \infty} Y^n = \infty \Rightarrow \lim_{n \rightarrow \infty} \frac{Z_n}{Y^n} = 1.$$

The result follows because

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{Z_n}{Y^n} = 1 &\Rightarrow \exists \epsilon_0 : P_{\theta_0} \left(\left\{ \frac{P_{\hat{\theta}}(x_{t_\tau} | x^{t_\tau-1})}{P_{\hat{\theta}}(x_{t_\tau} | x^{t_\tau-1})} \notin (1 \pm \epsilon_1) \text{ infinitely often} \right\} \right) = 1 \\ &\Rightarrow P_{\theta_0} \left(\left\{ \lim_{t \rightarrow \infty} \frac{P_{\hat{\theta}}(x^t)}{P_{\hat{\theta}}(x^t)} = \{0, \infty\} \text{ or } \lim_{t \rightarrow \infty} \frac{P_{\hat{\theta}}(x^t)}{P_{\hat{\theta}}(x^t)} \text{ does not exist} \right\} \right) = 1 \\ &\Leftrightarrow P_{\theta_0} \left(\neg \left\{ \lim_{t \rightarrow \infty} \frac{P_{\hat{\theta}}(x^t)}{P_{\hat{\theta}}(x^t)} = k \in (0, \infty) \right\} \right) = 1. \end{aligned}$$

□

References

- Berk, R. H. (1966). Limiting behavior of posterior distributions when the model is incorrect. *The Annals of Mathematical Statistics*, 37(1):51–58.
- Doob, J. L. (1949). *Application of the theory of martingales*. Colloques Internationaux du Centre National de la Recherche Scientifique Paris.
- Epstein, L. G. and Schneider, M. (2003). Recursive multiple-priors. *Journal of Economic Theory*, 113(1):1–31.
- Epstein, L. G. and Schneider, M. (2007). Learning under ambiguity. *The Review of Economic Studies*, 74(4):1275–1303.
- Epstein, L. G. and Seo, K. (2015). Exchangeable capacities, parameters and incomplete theories. *Journal of Economic Theory*, 157:879–917.
- Freedman, D. A. (1963). On the asymptotic behavior of bayes’ estimates in the discrete case. *The Annals of Mathematical Statistics*, pages 1386–1403.
- Lehrer, E. and Smorodinsky, R. (1996). Compatible measures and merging. *Mathematics of Operations Research*, 21(3):697–706.
- Marinacci, M. (2002). Learning from ambiguous urns. *Statistical Papers*, 43(1):143–151.
- Massari, F. (2013). Comment on ‘if you’re so smart, why aren’t you rich? belief selection in complete and incomplete markets’. *Econometrica*, 81(2):849–851.
- Massari, F. (2017). Markets with heterogeneous beliefs: A necessary and sufficient condition for a trader to vanish. *Journal of Economic Dynamics and Control*, 78:190–205.
- White, H. (1982). Maximum likelihood estimation of misspecified models. *Econometrica: Journal of the Econometric Society*, pages 1–25.
- Williams, D. (1991). *Probability with martingales*. Cambridge University Press.