

## CONSISTENCY OF MAXIMUM LIKELIHOOD ESTIMATION FOR SOME DYNAMICAL SYSTEMS

BY KEVIN MCGOFF<sup>1</sup>, SAYAN MUKHERJEE<sup>2</sup>,  
ANDREW NOBEL<sup>3</sup> AND NATESH PILLAI<sup>4</sup>

*Duke University, Duke University, University of North Carolina and  
Harvard University*

We consider the asymptotic consistency of maximum likelihood parameter estimation for dynamical systems observed with noise. Under suitable conditions on the dynamical systems and the observations, we show that maximum likelihood parameter estimation is consistent. Our proof involves ideas from both information theory and dynamical systems. Furthermore, we show how some well-studied properties of dynamical systems imply the general statistical properties related to maximum likelihood estimation. Finally, we exhibit classical families of dynamical systems for which maximum likelihood estimation is consistent. Examples include shifts of finite type with Gibbs measures and Axiom A attractors with SRB measures.

**1. Introduction.** Maximum likelihood estimation is a common, well-studied and powerful technique for statistical estimation. In the context of a statistical model with an unknown parameter, the maximum likelihood estimate of the unknown parameter is, by definition, any parameter value under which the observed data is most likely; such parameter values are said to maximize the likelihood function with respect to the observed data. In classical statistical models, one typically thinks of the unknown parameter as a real number or possibly a finite dimensional vector of real numbers. Here we consider maximum likelihood estimation for statistical models in

---

Received June 2013; revised July 2014.

<sup>1</sup>Supported by NSF Grant DMS-10-45153.

<sup>2</sup>Supported by NIH (Systems Biology): 5P50-GM081883, AFOSR: FA9550-10-1-0436, NSF CCF-1049290 and NSF DMS-12-09155.

<sup>3</sup>Supported in part by NSF Grants DMS-09-07177 and DMS-13-10002.

<sup>4</sup>Supported by NSF Grant DMS-11-07070.

*AMS 2000 subject classifications.* Primary 37A50, 37A25, 62B10, 62F12, 62M09; secondary 37D20, 60F10, 62M05, 62M10, 94A17.

*Key words and phrases.* Dynamical systems, hidden Markov models, maximum likelihood estimation, strong consistency.

<p>This is an electronic reprint of the original article published by the <a href="#">Institute of Mathematical Statistics</a> in <i>The Annals of Statistics</i>, 2015, Vol. 43, No. 1, 1–29. This reprint differs from the original in pagination and typographic detail.</p>
---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

which each parameter value corresponds to a stochastic system observed with noise.

Hidden Markov models (HMMs) provide a natural setting in which to study both stochastic systems with observational noise and maximum likelihood estimation. In this setting, one has a parametrized family of stochastic processes that are assumed to be Markov, and one attempts to perform inference about the underlying parameters from noisy observations of the process. There has been a substantial amount of work on statistical inference for HMMs, and we do not attempt a complete survey of that area here. In the 1960s, Baum and Petrie [5, 37] studied consistency of maximum likelihood estimation for finite state HMMs. Since that time, several other authors have shown that maximum likelihood estimation is consistent for HMMs under increasingly general conditions [13, 16, 18, 29–31], culminating with the work of Douc et al. [15], which currently provides the most general conditions on HMMs under which maximum likelihood estimation has been shown to be consistent.

We focus here on the consistency of maximum likelihood estimation for parametrized families of deterministic systems observed with noise. Inference methods for deterministic systems from noisy observations are of interest in a variety of scientific areas; for a few examples, see [19, 20, 28, 38–40, 46, 49].

For the purpose of this article, the terms deterministic system and dynamical system refer to a map  $T: X \rightarrow X$ . The set  $X$  is referred to as the state space, and the transformation  $T$  governs the evolution of states over one (discrete) time increment. Our main interest here lies in families of dynamical systems observed with noise. More precisely, we consider a state space  $X$  and a parameter space  $\Theta$ , and to each  $\theta$  in  $\Theta$ , we associate a dynamical system  $T_\theta: X \rightarrow X$ . Note that the state space  $X$  does not depend on  $\theta$ . For each  $\theta$  in  $\Theta$ , we assume that the system is started at equilibrium from a  $T_\theta$ -invariant measure  $\mu_\theta$ . See Section 2 for precise definitions. We are particularly interested in situations in which the family of dynamical systems is observed via noisy measurements (or observations). We consider a general observation model specified by a family of probability densities  $\{g_\theta(\cdot|x): \theta \in \Theta, x \in X\}$ , where  $g_\theta(\cdot|x)$  prescribes the distribution of an observation given that the state of the dynamical system is  $x$  and the state of nature is  $\theta$ . Under some additional conditions (see Section 3), our first main result states that maximum likelihood estimation is a consistent method of estimation of the parameter  $\theta$ .

We have chosen to state the conditions of our main consistency result in terms of statistical properties of the family of dynamical systems and the observations. However, these particular statistical properties have not been directly studied in the dynamical systems literature. In the interest of applying our general result to specific systems, we also establish several connections between well-studied properties of dynamical systems and the

statistical properties relevant to maximum likelihood estimation. Finally, we apply these results to some examples, including shifts of finite type with Gibbs measures and Axiom A attractors with SRB (Sinai–Ruelle–Bowen) measures. It is widely accepted in the field of ergodic theory and dynamical systems that these classes of systems have “good” statistical properties, and our results may be viewed as a precise confirmation of this view.

1.1. *Previous work.* There has been a substantial amount of work on statistical inference for HMMs, and a complete survey of that area is beyond the scope of this work. The asymptotic consistency of maximum likelihood estimation for HMMs has been studied at least since the work of Baum and Petrie [5, 37] under the assumption that both the hidden state space  $X$  and the observation space  $Y$  are finite sets. Leroux extended this result to the setting where  $Y$  is a general space and  $X$  is a finite set [29]. Several other authors have shown that maximum likelihood estimation is consistent for HMMs under increasingly general conditions [13, 16, 18, 30, 31], culminating with the work of Douc et al. [15], which currently provides the most general conditions for HMMs under which maximum likelihood estimation has been shown to be consistent.

Let us now discuss the results of Douc et al. [15] in greater detail. Consider parametrized families of HMMs in which both the hidden state space  $X$  and the observation space  $Y$  are complete, separable metric spaces. The main result of [15] shows that under several conditions, maximum likelihood estimation is a consistent method of estimation of the unknown parameter. These conditions involve some requirements on the transition kernel of the hidden Markov chain, as well as basic integrability conditions on the observations. The proof of that result relies on information-theoretic arguments, in combination with the application of some mixing conditions that follow from the assumptions on the transition kernel. To prove our consistency result, we take a similar information-theoretic approach, but instead of placing explicit restrictions on the transition kernel, we identify and study mixing conditions suitable for dynamical systems. See Remarks 2.4 and 3.3 for further discussion of our results in the context of HMMs.

Other directions of study regarding inference for HMMs include the behavior of MLE for misspecified HMMs [14], asymptotic normality for parameter estimates [8, 23], the dynamics of Bayesian updating [44] and starting the hidden process away from equilibrium [15]. Extending these results to dynamical systems is of potential interest.

The topic of statistical inference for dynamical systems has been widely studied in a variety of fields. Early interest from the statistical point of view is reflected in the following surveys [6, 12, 21, 22]. For a recent review of this area with many references, see [33]. There has been significant methodological work in the area of statistical inference for dynamical systems (for a

few recent examples, see [19, 20, 38, 46, 49]), but in this section we attempt to describe some of the more theoretical work in this area. The relevant theoretical work to date falls (very) roughly into three classes:

- state estimation (also known as denoising or filtering) for dynamical systems with observational noise;
- prediction for dynamical systems with observational noise;
- system reconstruction from dynamical systems without noise.

Let us now mention some representative works from these lines of research.

In the setting of dynamical systems with observational noise, Lalley introduced several ideas regarding state estimation in [25]. These ideas were subsequently generalized and developed in [26, 27]. Key results from this line of study include both positive and negative results on the consistency of denoising a dynamical system under additive observational noise. In short, the magnitude of the support of the noise seems to determine whether consistent denoising is possible. In related work, Judd [24] demonstrated that MLE can fail (in a particular sense) in state estimation when noise is large. It is perhaps interesting to note that there are examples of Axiom A systems with Gaussian observational noise for which state estimation cannot be consistent (by results of [26, 27]) and yet MLE provides consistent parameter estimation (by Theorem 5.7).

Steinwart and Anghel considered the problem of consistency in prediction accuracy for dynamical systems with observational noise [45]. They were able to show that support vector machines are consistent in terms of prediction accuracy under some conditions on the decay of correlations of the dynamical system.

The work of Adams and Nobel uses ideas from regression to study reconstruction of measure-preserving dynamical systems [1, 34, 35] without noise. These results show that certain types of inference are possible under fairly mild ergodicity assumptions. A sample result from this line of work is that a measure-preserving transformation may be consistently reconstructed from a typical trajectory observed without noise, assuming that the transformation preserves a measure that is absolutely continuous (with Radon–Nikodym derivative bounded away from 0 and infinity) with respect to a known reference measure.

*1.2. Organization.* In Section 2, we give some necessary background on dynamical systems observed with noise. Section 3 contains a statement and discussion of our main result (Theorem 3.1), which asserts that under some general statistical conditions, maximum likelihood parameter estimation is consistent for families of dynamical systems observed with noise. The purpose of Section 4 is to establish connections between well-studied properties

of dynamical systems and the (statistical) conditions appearing in Theorem 3.1. Section 5 gives several examples of widely studied families of dynamical systems to which we apply Theorem 3.1 and therefore establish consistency of maximum likelihood estimation. The proofs of our main results appear in Section 6, and we conclude with some final remarks in Section 7.

**2. Setting and notation.** Recall that our primary objects of study are parametrized families of dynamical systems. In this section we introduce these objects in some detail. First let us recall some terminology regarding dynamical systems and ergodic theory. We use  $\mathsf{X}$  to denote a state space, which we assume to be a complete separable metric space endowed with its Borel  $\sigma$ -algebra  $\mathcal{X}$ . Then a measurable dynamical system on  $\mathsf{X}$  is defined by a measurable map  $T: \mathsf{X} \rightarrow \mathsf{X}$ , which governs the evolution of states over one (discrete) time increment. For a probability measure  $\mu$  on the measurable space  $(\mathsf{X}, \mathcal{X})$ , we say that  $T$  preserves  $\mu$  (or  $\mu$  is  $T$ -invariant) if  $\mu(T^{-1}E) = \mu(E)$  for each set  $E$  in  $\mathcal{X}$ . We refer to the quadruple  $(\mathsf{X}, \mathcal{X}, T, \mu)$  as a measure-preserving system. To generate a trajectory  $(X_k)$  from such a measure-preserving system, one chooses  $X_0$  according to  $\mu$  and sets  $X_k = T^k(X_0)$  for  $k \geq 0$ . Note that  $(X_k)$  is then a stationary  $\mathsf{X}$ -valued stochastic process. Finally, the measure-preserving system  $(\mathsf{X}, \mathcal{X}, T, \mu)$  is said to be ergodic if  $T^{-1}E = E$  implies  $\mu(E) \in \{0, 1\}$ . See the books [36, 48] for an introduction to measure-preserving systems and ergodic theory.

Let us now introduce the setting of parametrized families of dynamical systems. We denote the parameter space by  $\Theta$ , which is assumed to be a compact metric space endowed with its Borel  $\sigma$ -algebra. Fix a state space  $\mathsf{X}$  and its Borel  $\sigma$ -algebra  $\mathcal{X}$  as above. To each parameter  $\theta$  in  $\Theta$ , we associate a measurable transformation  $T_\theta: \mathsf{X} \rightarrow \mathsf{X}$ , which prescribes the dynamics corresponding to the parameter  $\theta$ . Finally, we need to specify some initial conditions. In this article, we consider the case that the system is started from equilibrium. More precisely, we associate to each  $\theta$  in  $\Theta$  a  $T_\theta$ -invariant Borel probability measure  $\mu_\theta$  on  $(\mathsf{X}, \mathcal{X})$ . Thus, to each  $\theta$  in  $\Theta$ , we associate a measure-preserving system  $(\mathsf{X}, \mathcal{X}, T_\theta, \mu_\theta)$ , and we refer to the collection  $(\mathsf{X}, \mathcal{X}, T_\theta, \mu_\theta)_{\theta \in \Theta}$  as a parametrized family of dynamical systems. For ease of notation, we will refer to  $(T_\theta, \mu_\theta)_{\theta \in \Theta}$  as a family of dynamical systems on  $(\mathsf{X}, \mathcal{X})$ , instead of referring to the family of quadruples  $(\mathsf{X}, \mathcal{X}, T_\theta, \mu_\theta)_{\theta \in \Theta}$ .

We would like to study the situation that such a family of dynamical systems is observed via noisy measurements. Here we describe the specifics of our observation model. We suppose that we have a complete, separable metric space  $\mathsf{Y}$ , endowed with its  $\sigma$ -algebra  $\mathcal{Y}$ , which serves as our observation space. We also assume that we have a family of Borel probability densities  $\{g_\theta(\cdot|x) : \theta \in \Theta, x \in \mathsf{X}\}$  with respect to a fixed reference measure  $\nu$  on  $\mathsf{Y}$ . The density  $g_\theta(\cdot|x)$  prescribes the distribution of our observation given that the state of the dynamical system is  $x$  and the state of nature is  $\theta$ . Finally, we

assume that the noise involved in successive observations is conditionally independent given  $\theta$  and the underlying trajectory of the dynamical system. Thus our full model consists of a parametrized family of dynamical systems  $(T_\theta, \mu_\theta)_{\theta \in \Theta}$  on a measurable space  $(\mathsf{X}, \mathcal{X})$  with corresponding observation densities  $\{g_\theta(\cdot|x) : \theta \in \Theta, x \in \mathsf{X}\}$ .

In general, we would like to estimate the parameter  $\theta$  from our observations. Maximum likelihood estimation provides a basic method for performing such estimation. Our first main result states that maximum likelihood estimation is a consistent estimator of  $\theta$  under some general conditions on the family of systems and the noise. In order to state these results precisely, we now introduce the likelihood for our model. For the sake of notation, it will be convenient to denote finite sequences  $(x_i, \dots, x_j)$  with the notation  $x_i^j$ .

As we have assumed that our observations are conditionally independent given  $\theta$  and a trajectory  $(X_k)$ , we have that for  $\theta \in \Theta$  and  $y_0^n \in \mathsf{Y}^{n+1}$ , the likelihood of observing  $y_0^n$  given  $\theta$  and  $(X_k)$  is

$$p_\theta(y_0^n | X_0^n) = \prod_{j=0}^n g_\theta(y_j | X_j).$$

Since  $X_k = T_\theta^k(X_0)$  given  $\theta$  and  $X_0$ , the conditional likelihood of  $y_0^n$  given  $\theta$  and  $X_0 = x$  is

$$p_\theta(y_0^n | x) = \prod_{j=0}^n g_\theta(y_j | T_\theta^j(x)).$$

Since our model also assumes that  $X_0$  is distributed according to  $\mu_\theta$ , we have that for  $\theta \in \Theta$  and  $y_0^n \in \mathsf{Y}^{n+1}$ , the marginal likelihood of observing  $y_0^n$  given  $\theta$  is

$$(2.1) \quad p_\theta(y_0^n) = \int p_\theta(y_0^n | x) d\mu_\theta(x).$$

We denote by  $\nu^n$  the product measure on  $\mathsf{Y}^{n+1}$  with marginals equal to  $\nu$ . Let  $\mathbb{P}_\theta$  be the probability measure on  $\mathsf{X} \times \mathsf{Y}^{\mathbb{N}}$  such that for Borel sets  $A \subset \mathsf{X}$  and  $B \subset \mathsf{Y}^{n+1}$ , it holds that

$$\mathbb{P}_\theta(A \times B) = \int \int \mathbf{1}_A(x) \mathbf{1}_B(y_0^n) p_\theta(y_0^n | x) d\nu^n(y_0^n) d\mu_\theta(x),$$

which is well defined by Kolmogorov's consistency theorem. Let  $\mathbb{E}_\theta$  denote expectation with respect to  $\mathbb{P}_\theta$ , and let  $\mathbb{P}_\theta^{\mathsf{Y}}$  be the marginal of  $\mathbb{P}_\theta$  on  $\mathsf{Y}^{\mathbb{N}}$ .

Before we define consistency, let us first consider the issue of identifiability. Our notion of identifiability is captured by the following equivalence relation.

DEFINITION 2.1. Define an equivalence relation on  $\Theta$  as follows: let  $\theta \sim \theta'$  if  $\mathbb{P}_\theta^Y = \mathbb{P}_{\theta'}^Y$ . Denote by  $[\theta]$  the equivalence class of  $\theta$  with respect to this equivalence relation.

In a strong theoretical sense, if  $\theta'$  is in  $[\theta]$ , then the systems corresponding to the parameter values  $\theta'$  and  $\theta$  cannot be distinguished from each other based on observations of the system.

Now we fix a distinguished element  $\theta_0$  in  $\Theta$ . Here and in the rest of the paper, we assume that  $\theta_0$  is the “true” parameter; that is, the data are generated from the measure  $\mathbb{P}_{\theta_0}^Y$ . Hence, one may think of  $[\theta_0]$  as the set of parameters that cannot be distinguished from the true parameter.

DEFINITION 2.2. An approximate maximum likelihood estimator (MLE) is a sequence of measurable functions  $\hat{\theta}_n : (\mathcal{Y})^{n+1} \rightarrow \Theta$  such that

$$(2.2) \quad \frac{1}{n} \log p_{\hat{\theta}_n(Y_0^n)}(Y_0^n) \geq \sup_{\theta} \frac{1}{n} \log p_{\theta}(Y_0^n) - o_{\text{a.s.}}(1),$$

where  $o_{\text{a.s.}}(1)$  denotes a process that tends to zero  $\mathbb{P}_{\theta_0}$ -a.s. as  $n$  tends to infinity.

REMARK 2.1. Several notions in this article, including the definition of approximate MLE above, involve taking suprema over  $\theta$  in  $\Theta$ . In many situations of interest to us,  $\mathcal{X}$  and  $\Theta$  are compact, and all relevant functions are continuous in these arguments. In such cases, we have sufficient regularity to guarantee that suprema over  $\theta$  in  $\Theta$  are measurable. However, in the general situation, such suprema are not guaranteed to be measurable, and one must take some care. As all our measurable spaces are Polish (complete, separable metric spaces); such functions are always universally measurable [7], Proposition 7.47. Similarly, a Borel-measurable (approximate) maximum likelihood estimator need not exist, but the Polish assumption ensures the existence of universally measurable maximum likelihood estimators [7], Proposition 7.50. Thus all probabilities and expectations may be unambiguously extended to such quantities.

REMARK 2.2. In this work, we do not consider specific schemes for constructing an approximate MLE. Based on the existing results regarding denoising and system reconstruction (e.g., [1, 25–27, 34, 35], which are briefly discussed in Section 1.1), explicit construction of an approximate MLE may be possible under suitable conditions. Although the description and study of such constructive methods could be interesting, it is outside of the scope of this work.



REMARK 2.3. In principle, one could consider inference based on the conditional likelihood  $p_\theta(\cdot|x_0)$  in place of the marginal likelihood  $p_\theta(\cdot)$ . However, we do not pursue this direction in this work. For nonlinear dynamical systems, even the conditional likelihood  $p_\theta(\cdot|x_0)$  may depend very sensitively on  $x_0$ ; see [6], for example. Thus optimizing over  $x_0$  is essentially no more “tractable” than marginalizing the likelihood via an invariant measure.

REMARK 2.4. The framework of this paper may be translated into the language of Markov chains as follows. For each  $\theta \in \Theta$ , we define a (degenerate) Markov transition kernel  $Q_\theta$  as follows:

$$Q_\theta(x, y) = \delta_{T_\theta(x)}(y).$$

In other words, for each  $\theta \in \Theta$ ,  $x \in \mathsf{X}$ , and Borel set  $A \subset \mathsf{X}$ , the probability that  $X_1 \in A$  conditioned on  $X_0 = x$  is

$$Q_\theta(x, A) = \delta_{T_\theta(x)}(A),$$

where  $\delta_x$  is defined to be a point mass at  $x$ .

In all previous work on consistency of maximum likelihood estimation for HMMs (including [13, 15, 16, 18, 30, 31]), there have been significant assumptions placed on the Markovian structure of the hidden chain. For example, the central hypothesis appearing in [15] requires that there is a  $\sigma$ -finite measure  $\lambda$  on  $\mathsf{X}$  such that for some  $L \geq 0$ , the  $L$ -step transition kernel  $Q_\theta^L(x, \cdot)$  is absolutely continuous with respect to  $\lambda$  with bounded Radon–Nikodym derivative. If  $\mathsf{X}$  is uncountable, then the degeneracy of  $Q_\theta$ , which arises directly from the fact that we are considering deterministic systems, makes the existence of such a dominating measure impossible. In short, it is precisely the determinism in our hidden processes that prevents previous theorems for HMMs from applying to dynamical systems.

Nonetheless, there is a special case of systems that we consider in Section 5.1 that overlaps with the systems considered in the HMM literature. If  $\mathsf{X}$  is a shift of finite type,  $T_\theta$  is the shift map  $\sigma: \mathsf{X} \rightarrow \mathsf{X}$  for all  $\theta$ ,  $\mu_\theta$  is a (1-step) Markov measure for all  $\theta$ , and  $g_\theta(\cdot|x)$  depends only  $\theta$  and the zero coordinate  $x_0$ , then both the present work and the results in [15] apply to this setting and guarantee consistency of any approximate MLE under additional assumptions on the noise.

**3. Consistency of MLE.** In this section, we show that under suitable conditions, any approximate MLE is consistent for families of dynamical systems observed with noise. To make this statement precise, we make the following definition of consistency.

DEFINITION 3.1. An approximate MLE  $(\hat{\theta}_n)_n$  is consistent at  $\theta_0$  if  $\hat{\theta}_n(Y_0^n)$  converges to  $[\theta_0]$ ,  $\mathbb{P}_{\theta_0}$ -a.s. as  $n$  tends to infinity.



For the sake of notation, define the function  $\gamma: \Theta \times \mathbf{Y} \rightarrow \mathbb{R}_+$ , where

$$\gamma_\theta(y) = \sup_{x \in \mathbf{X}} g_\theta(y|x).$$

Also, for  $x > 0$ , let  $\log^+ x = \max(0, \log(x))$ .

Consider the following conditions on a family of dynamical systems observed with noise:

(S1) *Ergodicity.*

The system  $(T_{\theta_0}, \mu_{\theta_0})$  on  $(\mathbf{X}, \mathcal{X})$  is ergodic.

(S2) *Logarithmic integrability at  $\theta_0$ .*

It holds that

$$\mathbb{E}_{\theta_0}[\log^+ \gamma_{\theta_0}(Y_0)] < \infty$$

and

$$\mathbb{E}_{\theta_0} \left[ \left| \log \int g_{\theta_0}(Y_0|x) d\mu_{\theta_0}(x) \right| \right] < \infty.$$

(S3) *Logarithmic integrability away from  $\theta_0$ .*

For each  $\theta' \notin [\theta_0]$ , there exists a neighborhood  $U$  of  $\theta'$  such that

$$\mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log^+ \gamma_\theta(Y_0) \right] < \infty.$$

(S4) *Upper semi-continuity of the likelihood.*

For each  $\theta' \notin [\theta_0]$  and  $n \geq 0$ , the function  $\theta \mapsto p_\theta(Y_0^n)$  is upper semi-continuous at  $\theta'$ ,  $\mathbb{P}_{\theta_0}$ -a.s.

(S5) *Mixing condition.*

There exists  $\ell \geq 0$  such that for each  $m \geq 0$ , there exists a measurable function  $C_m: \Theta \times \mathbf{Y}^{m+1} \rightarrow \mathbb{R}_+$  such that if  $t \geq 1$  and  $w_0, \dots, w_t \in \mathbf{Y}^{m+1}$ , then

$$\int \prod_{j=0}^t p_\theta(w_j | T_\theta^{j(m+\ell)} x) d\mu_\theta(x) \leq \prod_{j=0}^t C_m(\theta, w_j) \prod_{j=0}^t p_\theta(w_j).$$

Furthermore, for each  $\theta' \notin [\theta_0]$ , there exists a neighborhood  $U$  of  $\theta'$  such that

$$\sup_m \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log C_m(\theta, Y_0^m) \right] < \infty.$$

(S6) *Exponential identifiability.*

For each  $\theta \notin [\theta_0]$ , there exists a sequence of measurable sets  $A_n \subset \mathbf{Y}^{n+1}$  such that

$$\liminf_n \mathbb{P}_{\theta_0}^Y(A_n) > 0 \quad \text{and} \quad \limsup_n \frac{1}{n} \log \mathbb{P}_\theta^Y(A_n) < 0.$$

The following theorem is our main general result.

**THEOREM 3.1.** *Suppose that  $(T_\theta, \mu_\theta)_{\theta \in \Theta}$  is a parametrized family of dynamical systems on  $(X, \mathcal{X})$  with corresponding observation densities  $(g_\theta)_{\theta \in \Theta}$ . If conditions (S1)–(S6) hold, then any approximate MLE is consistent at  $\theta_0$ .*

The proof of Theorem 3.1 is given in Section 6. In the following remark, we discuss conditions (S1)–(S6).

**REMARK 3.2.** Conditions (S1)–(S3) involve basic irreducibility and integrability conditions, and similar conditions have appeared in previous work on consistency of maximum likelihood estimation for HMMs; see, for example, [15, 29]. Taken together, conditions (S1) and (S2) ensure the almost sure existence and finiteness of the entropy rate for the process  $(Y_n)$ ,

$$h(\theta_0) = \lim_n \frac{1}{n} \log p_{\theta_0}(Y_0^n).$$

Condition (S3) serves as a basic integrability condition in the proof of Theorem 3.1, in which one must essentially show that for  $\theta \notin [\theta_0]$ ,

$$\limsup_n \frac{1}{n} \log p_\theta(Y_0^n) < h(\theta_0).$$

Conditions (S4)–(S6) are more interesting from the point of view of dynamical systems, and we discuss them in greater detail below.

The upper semi-continuity of the likelihood (S4) is closely related to the continuity of the map  $\theta \mapsto \mu_\theta$ . In general, the continuous dependence of  $\mu_\theta$  on  $\theta$  places nontrivial restrictions on a family of dynamical systems. This property (continuity of  $\theta \mapsto \mu_\theta$ ) is often called “statistical stability” in the dynamical systems and ergodic theory literature, and it has been studied for some families of systems; for example, see [2, 17, 42, 47] and references therein. In Section 4.1, we show how statistical stability of the family of dynamical systems may be used to establish the upper semi-continuity of the likelihood (S4).

The mixing condition (S5) involves control of the correlations of the observation densities along trajectories of the underlying dynamical system. Although the general topic of decay of correlations has been widely studied in dynamical systems (see [3] for an overview), condition (S5) is not implied by the particular decay of correlations properties that are typically studied for dynamical systems. Nonetheless, we show in Section 4.2 how some well-studied mixing properties of dynamical systems imply the mixing condition (S5).

Finally, condition (S6) involves the exponential identifiability of the true parameter  $\theta_0$ . We show in Section 4.3 how large deviations for a family

of dynamical systems may be used to establish exponential identifiability (S6). Large deviations estimates for dynamical systems have been studied in [41, 50], and our main goal in Section 4.3 is to connect such results to exponential identifiability (S6).

**REMARK 3.3.** Suppose one has a family of bi-variate stochastic processes  $\{(X_k^\theta, Y_k^\theta) : \theta \in \Theta\}$ , where  $(X_k^\theta)$  is interpreted as a hidden process and  $(Y_k^\theta)$  as an observation process. If the observations have conditional densities with respect to a common measure given  $(X_k^\theta)$  and  $\theta$ , then it makes sense to ask whether maximum likelihood estimation is a consistent method of inference for the parameter  $\theta$ .

It is well known that the setting of stationary stochastic processes may be translated into the deterministic setting of dynamical systems, which may be carried out as follows. Let  $\{(X_k^\theta) : \theta \in \Theta\}$  be a family of stationary stochastic processes on a measurable space  $(X, \mathcal{X})$ . Consider the product space  $\hat{X} = X^{\otimes \mathbb{Z}}$  with corresponding  $\sigma$ -algebra  $\hat{\mathcal{X}}$ . Each process  $(X_k^\theta)$  corresponds to a probability measure  $\mu_\theta$  on  $(\hat{X}, \hat{\mathcal{X}})$  with the property that  $\mu_\theta$  is invariant under the left-shift map  $T : \hat{X} \rightarrow \hat{X}$  given by  $\mathbf{x} = (x_i)_i \mapsto T(\mathbf{x}) = (x_{i+1})_i$ . With this translation, Theorem 3.1 shows that maximum likelihood estimation is consistent for families of hidden stochastic processes  $(X_k^\theta)$  observed with noise, whenever the corresponding family of dynamical systems  $(T, \mu_\theta)$  on  $(\hat{X}, \hat{\mathcal{X}})$  with observation densities satisfy conditions (S1)–(S6).

With the above translation, Theorem 3.1 applies to some families of processes allowing infinite-range dependence in both the hidden process  $(X_k^\theta)$  and the observation process  $(Y_k^\theta)$ . From this point of view, Theorem 3.1 highlights the fact that maximum likelihood estimation is consistent for dependent processes observed with noise as long as they satisfy some general conditions: ergodicity, logarithmic integrability of observations, continuous dependence on the parameters and some mixing of the observation process. It is interesting to note that the existing work on consistency of maximum likelihood estimation for HMMs [11, 13, 15, 16, 18, 29–31] makes assumptions of precisely this sort in the specific context of Markov chains.

**4. Statistical properties of dynamical systems.** In our main consistency result (Theorem 3.1), we establish the consistency of any approximate MLE under conditions (S1)–(S6). We have chosen to formulate our result in these terms because they reflect general statistical properties of dynamical systems observed with noise that are relevant to parameter inference. However, these conditions have not been explicitly studied in the dynamical systems literature, despite the fact that much effort has been devoted to understanding certain statistical aspects of dynamical systems. In this section, we make connections between the general statistical conditions appearing in Theorem 3.1 and some well-studied properties of dynamical systems. Section 4.1

shows how the notion of statistical stability may be used to verify the upper semi-continuity of the likelihood (S4). Section 4.2 connects well-known mixing properties of some measure-preserving dynamical systems to the mixing property (S5). In Section 4.3, we show how large deviations for dynamical systems may be used to deduce the exponential identifiability condition (S6). Proofs of statements in this section, as well as additional discussion, appear in Supplementary Appendix A [32].

4.1. *Statistical stability and continuity of  $p_\theta$ .* As discussed in Remark 3.2, the upper semi-continuity condition (S4) places nontrivial restrictions on the family of dynamical systems under consideration. In this section, we establish sufficient conditions for (S4) to hold. The continuous dependence of  $\mu_\theta$  on  $\theta$  is a property called statistical stability in the dynamical systems literature [2, 17, 42, 47]. Let us state this property precisely. Let  $M(\mathsf{X})$  denote the space of Borel probability measures on  $\mathsf{X}$ . Endow  $M(\mathsf{X})$  with the topology of weak convergence:  $\mu_n$  converges to  $\mu$  if  $\int f d\mu_n$  converges to  $\int f d\mu$  as  $n$  tends to infinity, for each continuous, bounded function  $f: \mathsf{X} \rightarrow \mathbb{R}$ . The family of dynamical systems  $(T_\theta, \mu_\theta)_{\theta \in \Theta}$  on  $(\mathsf{X}, \mathcal{X})$  is said to have statistical stability if the map  $\theta \mapsto \mu_\theta$  is continuous with respect to the weak topology on  $M(\mathsf{X})$ .

The following proposition shows that under some continuity and compactness assumptions, statistical stability of the family of dynamical systems implies upper semi-continuity of the likelihood (S4).

PROPOSITION 4.1. *Suppose that  $\mathsf{X}$  and  $\Theta$  are compact, and the maps  $T: \Theta \times \mathsf{X} \rightarrow \mathsf{X}$  and  $g: \Theta \times \mathsf{X} \times \mathsf{Y} \rightarrow \mathbb{R}_+$  are continuous. If the family  $(T_\theta, \mu_\theta)_{\theta \in \Theta}$  has statistical stability, then upper semi-continuity of the likelihood (S4) holds.*

The proof of Proposition 4.1 appears in Supplementary Appendix A.1 [32].

4.2. *Mixing.* In this section, we focus on mixing condition (S5). Recall that (S5) involves a nontrivial restriction on the correlations of the observation densities  $g_\theta$  along trajectories of the underlying dynamical system. Although mixing conditions have been widely studied in the dynamics literature, the particular type of condition appearing in (S5) appears not to have been investigated. Nonetheless, we show that a well-studied mixing property for dynamical systems implies the statistical mixing property (S5).

In order to study mixing for dynamical systems, one typically places restrictions on the type of events or observations that one considers (by considering certain functionals of the process). For example, in some situations a substantial amount work has been devoted to finding particular partitions of

state space with respect to which the system possess good mixing properties; an example of such partitions are the well-known Markov partitions [9]. If a system has good mixing properties with respect to a particular partition, and if that partition possesses certain (topological) regularity properties, then it is often possible to show that the system also has good mixing properties for related function classes, such as Lipschitz or Hölder continuous observables. For variations of this approach to mixing in dynamical systems, see the vast literature on decay of correlations; for an introduction, see the survey [3].

In this section, we follow the above approach to study mixing condition (S5) for dynamical systems observed with noise. First, we define a mixing property for families of dynamical systems with respect to a partition (M1). Second, we define a regularity property for partitions (M2). Third, we define a topological regularity property for a family of observation densities (M3). Finally, in the main result of this section (Proposition 4.2), we show how these three properties together imply the mixing condition (S5).

Here and in the rest of this section, we consider only invertible transformations. It is certainly possible to modify the definitions slightly to handle the noninvertible case, but we omit such modifications.

We will have need to consider finite partitions of  $X$ . The join of two partitions  $\mathcal{C}_0$  and  $\mathcal{C}_1$  is defined to be the common refinement of  $\mathcal{C}_0$  and  $\mathcal{C}_1$ , and it is denoted  $\mathcal{C}_0 \vee \mathcal{C}_1$ . Note that for any measurable transformation  $T: X \rightarrow X$ , if  $\mathcal{C}$  is a partition, then so is  $T^{-1}\mathcal{C} = \{T^{-1}A : A \in \mathcal{C}\}$ . For a fixed partition  $\mathcal{C}$  and  $i \leq j$ , let  $\mathcal{C}_i^j = \bigvee_{k=i}^j T_{\theta}^{-k}\mathcal{C}$ . Notice that  $\mathcal{C}_i^j$  depends on  $\theta$  through  $T_{\theta}$ , although we suppress this dependence in our notation. Now consider the following alternative conditions, which may be used in place of condition (S5):

(M1) *Mixing condition with respect to the partition  $\mathcal{C}$ .*

There exists  $L: \Theta \rightarrow \mathbb{R}_+$  and  $\ell \geq 0$  such that for all  $\theta \in \Theta$ ,  $m, n \geq 0$ ,  $A \in \mathcal{C}_0^m$  and  $B \in \mathcal{C}_0^n$ , it holds that

$$\mu_{\theta}(A \cap T_{\theta}^{-(m+\ell)}B) \leq L_{\theta}\mu_{\theta}(A)\mu_{\theta}(B).$$

Furthermore, for each  $\theta' \notin [\theta_0]$  there exists a neighborhood  $U$  of  $\theta'$  such that

$$\sup_{\theta \in U} L_{\theta} < \infty.$$

(M2) *Regularity of the partition  $\mathcal{C}$ .* There exists  $\beta \in (0, 1)$  such that for all  $\theta \in \Theta$  and  $m, n \geq 0$ , if  $A \in \mathcal{C}_{-m}^n$  and  $x, z \in A$ , then

$$d(x, z) \leq \beta^{\min(m, n)}.$$

(M3) *Regularity of observations.* There exists a function  $K: \Theta \times Y \rightarrow \mathbb{R}_+$  such that for  $y \in Y$  and  $x, z \in X$ ,

$$g_{\theta}(y|x) \leq g_{\theta}(y|z) \exp(K(\theta, y)d(x, z)).$$

Furthermore, for each  $\theta' \notin [\theta_0]$ , there exists a neighborhood  $U$  of  $\theta'$  such that

$$\mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} K(\theta, Y_0) \right] < \infty.$$

Let us now state the main proposition of this section, whose proof is deferred to Supplementary Appendix A.2 [32].

**PROPOSITION 4.2.** *Suppose  $(T_\theta, \mu_\theta)_{\theta \in \Theta}$  is a family of dynamical systems on  $(X, \mathcal{X})$  with corresponding observation densities  $(g_\theta)_{\theta \in \Theta}$ . If there exists a partition  $\mathcal{C}$  of  $X$  such that conditions (M1) and (M2) are satisfied, and if the observation regularity condition (M3) is satisfied, then mixing property (S5) holds.*

**4.3. Exponential identifiability.** In this section, we study exponential identifiability condition (S6). We show how large deviations for dynamical systems may be used in combination with some regularity of the observation densities to establish exponential identifiability (S6).

Let  $X_1$  and  $X_2$  be metric spaces with metrics  $d_1$  and  $d_2$ , respectively. Recall that a function  $f: X_1 \rightarrow X_2$  is said to be Hölder continuous if there exist  $\alpha > 0$  and  $C > 0$  such that for each  $x, z$  in  $X_1$ , it holds that

$$d_2(f(x), f(z)) \leq C d_1(x, z)^\alpha.$$

If  $(T, \mu)$  is a dynamical system on  $(X, \mathcal{X})$  such that  $T: X \rightarrow X$  is Hölder continuous, then we refer to  $(T, \mu)$  as a Hölder continuous dynamical system. For many dynamical systems, the class of Hölder continuous functions  $f: X \rightarrow \mathbb{R}$  provides a natural class of observables whose statistical properties are fairly well understood and satisfy some large deviations estimates [41, 50].

Consider the following conditions, which we later show are sufficient to guarantee exponential identifiability (S6):

(L1) *Large deviations.* For each  $\theta \notin [\theta_0]$ , for each Hölder continuous function  $f: X \rightarrow \mathbb{R}$ , and for each  $\delta > 0$ , it holds that

$$\limsup_n \frac{1}{n} \log \mu_\theta \left( \left| \frac{1}{n} \sum_{k=0}^{n-1} f(T_\theta^k(x)) - \int f d\mu_\theta \right| > \delta \right) < 0.$$

(L2) *Regularity of observations.* There exists  $\alpha > 0$  and  $K: \Theta \times Y \rightarrow \mathbb{R}_+$  such that for each  $x$  and  $z$  in  $X$ , it holds that

$$g_\theta(y|x) \leq g_\theta(y|z) \exp(K(\theta, y)d(x, z)^\alpha).$$

Furthermore, for  $\theta \in \Theta$  and  $C > 0$ , it holds that

$$\sup_x \int \exp(CK(\theta, y))g_\theta(y|x) d\nu(y) < \infty.$$

The following proposition relates large deviations for dynamical systems to the exponential identifiability condition (S6).

**PROPOSITION 4.3.** *Suppose that  $(T_\theta, \mu_\theta)_{\theta \in \Theta}$  is a family of Hölder continuous dynamical systems on the  $(X, \mathcal{X})$  with corresponding observation densities  $(g_\theta)_{\theta \in \Theta}$ . Further suppose that the large deviations property (L1) and the observation regularity property (L2) are satisfied. Then the exponential identifiability condition (S6) holds.*

The proof of Proposition 4.3 appears in Supplementary Appendix A.3 [32].

**5. Examples.** In this section we present some classical families of dynamical systems for which maximum likelihood estimation is consistent. We begin in Section 5.1 by considering symbolic dynamical systems called shifts of finite type. The state space for such systems consists of (bi-)infinite sequences of symbols from a finite set, and the transformation on the state space is always given by the “left-shift” map, which just shifts each point one coordinate to the left. Such systems are considered models of “chaotic” dynamical systems that may be defined by a finite amount of combinatorial information. In this setting Gibbs measures form a natural class of invariant measures, which have been studied due to their connections to statistical physics. These measures play a central role in a topic called the *thermodynamic formalism*, which is well described in the books [10, 43]. Note that  $k$ th order finite state Markov chains form a special case of Gibbs measures. The main result of this section is Theorem 5.1, which states that under sufficient regularity conditions, any approximate maximum likelihood estimator is consistent for families of Gibbs measures on a shift of finite type. The crucial assumptions for this theorem involve continuous dependence of the Gibbs measures on  $\theta$  and sufficiently regular dependence of  $g_\theta(y|x)$  on  $x$ . Additional proofs and discussion for this section appear in the Supplementary Appendix B [32].

Having established consistency of maximum likelihood estimation for families of Gibbs measures on a shift of finite type, we deduce in Section 5.2 that maximum likelihood estimation is consistent for families of Axiom A attractors observed with noise. Axiom A systems are well studied differentiable dynamical systems on manifolds that, like shifts of finite type, exhibit “chaotic” behavior; for a thorough treatment of Axiom A systems, see the book [10]. In related statistical work, Lalley [25] considered the problem of denoising the trajectories of Axiom A systems. For these systems, there is a natural class of measures, known as SRB (Sinai–Ruelle–Bowen) measures. See the article [52] for an introduction to these measures with discussion of their interpretation and importance. With the construction of Markov partitions [9, 10], one may view an Axiom A attractor with its SRB measure



as a factor of a shift of finite type with a Gibbs measure. Using this natural factor structure, we establish the consistency of any approximate maximum likelihood estimator for Axiom A systems. Proofs and discussion of these topics appear in the Supplementary Appendix C [32].

5.1. *Gibbs measures.* In this section, we consider the setting of symbolic dynamics, shifts of finite type and Gibbs measures. We prove that any approximate maximum likelihood estimator is consistent for these systems (Theorem 5.1) under some general assumptions on the observations. Finally, we consider two examples of observations in greater detail. In the first example, we consider “discrete” observations, corresponding to a “noisy channel.” In the second example, we consider making real-valued observations with Gaussian observational noise. For a brief introduction to shifts of finite type and Gibbs measures that contains everything needed in this work, see the Supplementary Appendix B [32]. For a complete introduction to shifts of finite type and Gibbs measures, see [10].

Let us now consider some families of measure-preserving systems on SFTs. Let  $A$  be an alphabet, and let  $M$  be a binary matrix with dimensions  $|A| \times |A|$ . Let  $\mathbf{X} = X_M$  be the associated SFT, and let  $\mathcal{X}$  be the Borel  $\sigma$ -algebra on  $\mathbf{X}$ . For  $\alpha > 0$ , let  $f : \Theta \rightarrow C^\alpha(\mathbf{X})$  be a continuous map, and let  $\mu_\theta$  be the Gibbs measure associated to the potential function  $f_\theta$ . In this setting, we refer to  $(\mu_\theta)_{\theta \in \Theta}$  as a continuously parametrized family of Gibbs measures on  $(\mathbf{X}, \mathcal{X})$ .

**THEOREM 5.1.** *Suppose  $\mathbf{X} = X_M$  is a mixing shift of finite type and  $(\mu_\theta)_{\theta \in \Theta}$  is a continuously parametrized family of Gibbs measures on  $(\mathbf{X}, \mathcal{X})$ . If the family of observation densities  $(g_\theta)_{\theta \in \Theta}$  satisfies the integrability conditions (S2) and (S3) and the regularity conditions (M3) and (L2), then any approximate maximum likelihood estimator is consistent.*

The proof of Theorem 5.1 is based on an appeal to Theorem 3.1. However, in order to verify the hypotheses of Theorem 3.1, we combine the results of Section 4 with some well-known properties of Gibbs measures. This proof appears in the Supplementary Appendix B [32].

**REMARK 5.2.** There is an analogous theory of “one-sided” symbolic dynamics and Gibbs measures, in which  $A^{\mathbb{Z}}$  is replaced by  $A^{\mathbb{N}}$  and appropriate modifications are made in the definitions. The two-sided case deals with invertible dynamical systems, whereas the one-sided case handles noninvertible systems. We have stated Theorem 5.1 in the invertible setting, although it applies as well in the noninvertible setting, with the obvious modifications.

EXAMPLE 5.3. In this example, we consider families of dynamical systems  $(T_\theta, \mu_\theta)$  on  $(X, \mathcal{X})$ , where  $X$  is a mixing shift of finite type,  $T_\theta = \sigma|_X$ , and  $\mu_\theta$  is a continuous family of Gibbs measures on  $X$  (as in Theorem 5.1). Here we consider the particular observation model in which our observations of  $X$  are passed through a discrete, memoryless, noisy channel. Suppose that  $Y$  is a finite set,  $\nu$  is counting measure on  $Y$  and for each symbol  $a$  in  $A$  and parameter  $\theta$  in  $\Theta$ , we have a probability distribution  $\pi_\theta(\cdot|a)$  on  $Y$ . We consider the case that our observation densities  $g_\theta$  satisfy  $g_\theta(\cdot|x) = \pi_\theta(\cdot|x_0)$ . This situation is covered by Theorem 5.1, since the following conditions may be easily verified: observation integrability (S2) and (S3) and observation regularity (M3) and (L2).

EXAMPLE 5.4. In this example, we once again consider families of dynamical systems  $(T_\theta, \mu_\theta)$  on  $(X, \mathcal{X})$ , such that  $X$  is a mixing shift of finite type,  $T_\theta = \sigma|_X$  and  $\mu_\theta$  is a continuous family of Gibbs measures on  $X$  (as in Theorem 5.1). Here we consider the particular observation model in which we make real-valued, parameter-dependent measurements of the system, which are corrupted by Gaussian noise with parameter-dependent variance. More precisely, let us assume that  $Y = \mathbb{R}$ , and there exists a Lipschitz continuous  $\varphi: \Theta \times X \rightarrow \mathbb{R}$  and continuous  $s: \Theta \rightarrow (0, \infty)$  such that

$$g_\theta(y|x) = \frac{1}{s(\theta)\sqrt{2\pi}} \exp\left(-\frac{1}{2s(\theta)^2}(\varphi_\theta(x) - y)^2\right).$$

We now proceed to verify conditions (S2), (S3), (M3) and (L2). First, by compactness and continuity, there exist  $C_1, C_2, C_3 > 0$  such that for  $\theta$  in  $\Theta$ ,  $y$  in  $Y$  and  $x$  in  $X$ , it holds that

$$(5.1) \quad C_1^{-1} \exp(-C_2 y^2) \leq g_\theta(y|x) \leq C_1 \exp(-C_3 y^2).$$

From (5.1), one easily obtains the observation integrability conditions (S2) and (S3). Furthermore, there exists  $C_4, C_5 > 0$  such that for  $x, z \in X$ , it holds that

$$\begin{aligned} & \frac{g_\theta(y|x)}{g_\theta(y|z)} \\ &= \exp\left(-\frac{1}{2s(\theta)^2}[(\varphi_\theta(x) - y)^2 - (\varphi_\theta(z) - y)^2]\right) \\ (5.2) \quad &= \exp\left(-\frac{1}{2s(\theta)^2}[(\varphi_\theta(x) - \varphi_\theta(z))(\varphi_\theta(x) + \varphi_\theta(z)) + 2y(\varphi_\theta(z) - \varphi_\theta(x))]\right) \\ &\leq \exp((C_4 + C_5|y|)|\varphi_\theta(x) - \varphi_\theta(z)|). \end{aligned}$$

Let  $\varphi$  be Lipschitz continuous with constant  $C_6$ , and let  $K(\theta, y) = C_6(C_4 + C_5|y|)$ . With this choice of  $K$  and (5.2), one may easily verify the observation regularity conditions (M3) and (L2).

REMARK 5.5. Similar calculations to those in Example 5.4 imply that any approximate maximum likelihood estimator is also consistent if the observational noise is “double-exponential” [i.e.,  $g_\theta(y|x) \propto e^{-|y-x|}$ ]. Indeed, these calculations should hold for most members of the exponential family, although we do not pursue them here.

5.2. *Axiom A systems.* In this section, we show how the previous results may be applied to some smooth (differentiable) families of dynamical systems. These results follow easily from the results in Section 5.1, using the work of Bowen and others (see [9, 10] and references therein) in constructing Markov partitions for these systems. With Markov partitions, Axiom A systems may be viewed as factors of the shifts of finite type with Gibbs measures. For a brief introduction of Axiom A systems that contains the details necessary for this work, see the Supplementary Appendix C [32].

The basic fact that allows us to transfer our results from shifts of finite type to Axiom A systems is that consistency of maximum likelihood estimation is preserved under taking appropriate factors. Let us now make this statement precisely. Suppose that  $(T_\theta, \mu_\theta)_{\theta \in \Theta}$  is a family of dynamical systems on  $(X, \mathcal{X})$  with observation densities  $(g_\theta)_{\theta \in \Theta}$ . Further, suppose that there are continuous maps  $\pi: \Theta \times \tilde{X} \rightarrow X$  and  $\tilde{T}: \Theta \times \tilde{X} \rightarrow \tilde{X}$  such that:

- (i) for each  $\theta$ , we have that  $\pi_\theta \circ \tilde{T}_\theta = T_\theta \circ \pi_\theta$ ;
- (ii) for each  $\theta$ , there is a unique probability measure  $\tilde{\mu}_\theta$  on  $\tilde{X}$  such that  $\tilde{\mu}_\theta \circ \pi_\theta^{-1} = \mu_\theta$ ;
- (iii) for each  $\theta$ , the map  $\pi_\theta$  is injective  $\tilde{\mu}_\theta$ -a.s.

For  $x$  in  $\tilde{X}$  and  $\theta$  in  $\Theta$ , define  $\tilde{g}_\theta(\cdot|x) = g_\theta(\cdot|\pi_\theta(x))$ . Then  $(\tilde{T}_\theta, \tilde{\mu}_\theta)_{\theta \in \Theta}$  is a family of dynamical systems on  $(\tilde{X}, \tilde{\mathcal{X}})$  with observation densities  $(\tilde{g}_\theta)_{\theta \in \Theta}$ . In this situation, we say that  $(T_\theta, \mu_\theta, g_\theta)_{\theta \in \Theta}$  is an isomorphic factor of  $(\tilde{T}_\theta, \tilde{\mu}_\theta, \tilde{g}_\theta)_{\theta \in \Theta}$ , and  $\pi$  is the factor map. The following proposition addresses the consistency of maximum likelihood estimation for isomorphic factors. Its proof is straightforward and omitted.

PROPOSITION 5.6. *Suppose that  $(T_\theta, \mu_\theta, g_\theta)_{\theta \in \Theta}$  is an isomorphic factor of  $(\tilde{T}_\theta, \tilde{\mu}_\theta, \tilde{g}_\theta)_{\theta \in \Theta}$ . Then maximum likelihood estimation is consistent for  $(T_\theta, \mu_\theta, g_\theta)_{\theta \in \Theta}$  if and only if maximum likelihood estimation is consistent for  $(\tilde{T}_\theta, \tilde{\mu}_\theta, \tilde{g}_\theta)_{\theta \in \Theta}$ .*

For the sake of brevity, we defer precise definitions for Axiom A systems to Supplementary Appendix C [32].

We consider families of Axiom A systems as follows. Suppose that  $f: \Theta \times X \rightarrow X$  is a parametrized family of diffeomorphisms such that:

- (i)  $\theta \mapsto f_\theta$  is Hölder continuous;

- (ii) there exists  $\alpha > 0$  such that for each  $\theta$ , the map  $f_\theta$  is  $C^{1+\alpha}$ ;
- (iii) for each  $\theta$ ,  $\Omega(f_\theta)$  is an Axiom A attractor and the restriction  $f_\theta|_{\Omega(f_\theta)}$  is topologically mixing;
- (iv) for each  $\theta$ , the measure  $\mu_\theta$  is the unique SRB measure corresponding to  $f_\theta$  [10], Theorem 4.1.

If these conditions are satisfied, then we say that  $(f_\theta, \mu_\theta)_{\theta \in \Theta}$  is a parametrized family of Axiom A systems on  $(\mathbf{X}, \mathcal{X})$ .

**THEOREM 5.7.** *Suppose that  $(f_\theta, \mu_\theta)_{\theta \in \Theta}$  is a parametrized family of Axiom A systems on  $(\mathbf{X}, \mathcal{X})$ . Further, suppose that  $(g_\theta)_{\theta \in \Theta}$  is a family of observations densities satisfying the following conditions: observation integrability (S2) and (S3) and observation regularity (M3) and (L2). Then maximum likelihood estimation is consistent.*

The proof of Theorem 5.7 appears in the Supplementary Appendix C [32].

**6. Proof of the main result.** Propositions 6.1–6.5 are used in the proof of Theorem 3.1, which is given at the end of the present section.

**PROPOSITION 6.1.** *Suppose that condition (S1) (ergodicity) holds. Then the process  $(Y_k)$  is ergodic under  $\mathbb{P}_{\theta_0}^Y$ .*

**PROOF.** Let  $m > 0$  be arbitrary, and let  $A$  and  $B$  be Borel subsets of  $Y^{m+1}$ . To obtain the ergodicity of  $\{Y_k\}_k$ , it suffices to show that (see [36])

$$(6.1) \quad \lim_n \frac{1}{n} \sum_{k=0}^n \mathbb{P}_{\theta_0}^Y(Y_0^m \in A, Y_k^{k+m} \in B) = \mathbb{P}_{\theta_0}^Y(Y_0^m \in A) \mathbb{P}_{\theta_0}^Y(Y_0^m \in B).$$

For  $x \in \mathbf{X}$ , define

$$\eta_A(x) = \int \mathbf{1}_A(y_0^m) p_{\theta_0}(y_0^m | x) d\nu^m(y_0^m),$$

and define  $\eta_B(x)$  similarly. For  $k > m$ , by the conditional independence of  $Y_0^m$  and  $Y_k^{k+m}$  given  $\theta_0$  and  $X_0 = x$ , we have that

$$\begin{aligned} & \mathbb{P}_{\theta_0}^Y(Y_0^m \in A, Y_k^{k+m} \in B) \\ &= \int \int \mathbf{1}_A(y_0^m) \mathbf{1}_B(y_k^{k+m}) p_{\theta_0}(y_0^{n+m} | x) d\nu^{n+m}(y_0^{n+m}) d\mu_{\theta_0}(x) \\ &= \int \left( \int \mathbf{1}_A(y_0^m) p_{\theta_0}(y_0^m | x) d\nu^m(y_0^m) \right. \\ & \quad \left. \times \int \mathbf{1}_B(y_k^{k+m}) p_{\theta_0}(y_k^{k+m} | T_{\theta_0}^k(x)) d\nu^m(y_k^{k+m}) \right) d\mu_{\theta_0}(x) \\ &= \int \eta_A(x) \eta_B(T_{\theta_0}^k(x)) d\mu_{\theta_0}(x), \end{aligned}$$

where we have used Fubini's theorem. Since  $m$  is fixed, we have that

$$\begin{aligned} & \lim_n \frac{1}{n} \sum_{k=0}^n \mathbb{P}_{\theta_0}^Y(Y_0^m \in A, Y_k^{k+m} \in B) \\ &= \lim_n \left( \frac{1}{n} \sum_{k=0}^m \mathbb{P}_{\theta_0}^Y(Y_0^m \in A, Y_k^{k+m} \in B) \right. \\ & \quad \left. + \frac{1}{n} \sum_{k=m+1}^n \int \eta_A(x) \eta_B(T_{\theta_0}^k(x)) d\mu_{\theta_0}(x) \right) \\ &= \lim_n \frac{1}{n} \sum_{k=m+1}^n \int \eta_A(x) \eta_B(T_{\theta_0}^k(x)) d\mu_{\theta_0}(x). \end{aligned}$$

Since  $(T_{\theta_0}, \mu_{\theta_0})$  is ergodic, an alternative characterization of ergodicity (see [36]) gives that

$$\begin{aligned} \lim_n \frac{1}{n} \sum_{k=0}^n \mathbb{P}_{\theta_0}^Y(Y_0^m \in A, Y_k^{k+m} \in B) &= \lim_n \frac{1}{n} \sum_{k=m+1}^n \int \eta_A(x) \eta_B(T_{\theta_0}^k(x)) d\mu_{\theta_0}(x) \\ &= \int \eta_A(x) d\mu_{\theta_0}(x) \int \eta_B(x) d\mu_{\theta_0}(x) \\ &= \mathbb{P}_{\theta_0}^Y(Y_0^m \in A) \mathbb{P}_{\theta_0}^Y(Y_0^m \in B). \end{aligned}$$

Thus we have verified equation (6.1), and the proof is complete.  $\square$

For the following propositions, recall our notation that

$$\gamma_{\theta}(y) = \sup_x g_{\theta}(y|x).$$

PROPOSITION 6.2. *Suppose that conditions (S1) and (S2) hold. Then there exists  $h(\theta_0) \in (-\infty, \infty)$  such that*

$$h(\theta_0) = \lim_n \mathbb{E}_{\theta_0} \left( \frac{1}{n} \log p_{\theta_0}(Y_0^n) \right).$$

Moreover, the following equality holds  $\mathbb{P}_{\theta_0}$ -a.s.:

$$h(\theta_0) = \lim_n \frac{1}{n} \log p_{\theta_0}(Y_0^n).$$

PROOF. The proposition is a direct application of Barron's generalized Shannon–McMillan–Breiman theorem [4]. Here we simply check that the hypotheses of that theorem hold in our setting. Since condition (S1) (ergodicity) holds, Proposition 6.1 gives  $(Y_k)$  is stationary and ergodic under  $\mathbb{P}_{\theta_0}$ . By

definition,  $Y_0^n$  has density  $p_{\theta_0}(Y_0^n)$  with respect to the  $\sigma$ -finite measure  $\nu^n$ . The measure  $\nu^n$  is a product of the measure  $\nu$  taken  $n + 1$  times. As such, the sequence  $\{\nu^n\}$  clearly satisfies Barron's condition that this sequence is "Markov with stationary transitions." Define  $D_n = \mathbb{E}_{\theta_0}(\log p_{\theta_0}(Y_0^{n+1})) - \mathbb{E}_{\theta_0}(\log p_{\theta_0}(Y_0^n))$ . Let us show that for  $n > 0$ , we have that

$$(6.2) \quad \mathbb{E}_{\theta_0}(|\log p_{\theta_0}(Y_0^n)|) < \infty,$$

which clearly implies that  $-\infty < D_n < \infty$ . Once (6.2) is established, we will have verified all of the hypotheses of Barron's generalized Shannon–McMillan–Breiman theorem, and the proof of the proposition will be complete.

Observe that the first part of the integrability condition (S2) gives that

$$(6.3) \quad \mathbb{E}_{\theta_0}[\log^+ p_{\theta_0}(Y_0^n)] \leq (n + 1)\mathbb{E}_{\theta_0}[\log^+ \gamma_{\theta_0}(Y_0)] < \infty.$$

Then the second part of the integrability condition (S2) implies that

$$(6.4) \quad \begin{aligned} \mathbb{E}_{\theta_0}[\log p_{\theta_0}(Y_0^n)] &= \mathbb{E}_{\theta_0} \left[ \log \frac{p_{\theta_0}(Y_0^n)}{\prod_{k=0}^n \int g_{\theta_0}(Y_k|x) d\mu_{\theta_0}(x)} \right] \\ &\quad + \mathbb{E}_{\theta_0} \left[ \sum_{k=0}^n \log \int g_{\theta_0}(Y_k|x) d\mu_{\theta_0}(x) \right] \\ &\geq -(n + 1)\mathbb{E}_{\theta_0} \left[ \left| \log \int g_{\theta_0}(Y_0|x) d\mu_{\theta_0}(x) \right| \right] \\ &> -\infty, \end{aligned}$$

where we have used that relative entropy is nonnegative. By (6.3) and (6.4), we conclude that (6.2) holds, which completes the proof.  $\square$

The following proposition is used in the proof of Theorem 3.1 to given an almost sure bound for the normalized log-likelihoods in terms of quantities involving only expectations.

**PROPOSITION 6.3.** *Suppose that conditions (S1), (S3) and (S5) hold. Let  $\ell$  be as in condition (S5). Then for  $\theta' \notin [\theta_0]$ , there exists a neighborhood  $U$  of  $\theta'$  such that for each  $m > 0$ , the following inequality holds  $\mathbb{P}_{\theta_0}$ -a.s.:*

$$\begin{aligned} \limsup_{n \rightarrow \infty} \sup_{\theta \in U} \frac{1}{n} \log p_{\theta}(Y_0^n) &\leq \frac{1}{m + \ell} \mathbb{E}_{\theta_0} \left( \sup_{\theta \in U} \log p_{\theta}(Y_0^m) \right) \\ &\quad + \frac{\ell}{m + \ell} \mathbb{E}_{\theta_0} \left( \sup_{\theta \in U} \log^+ \gamma_{\theta}(Y_0) \right) \\ &\quad + \frac{1}{m + \ell} \mathbb{E}_{\theta_0} \left( \sup_{\theta \in U} \log C_m(\theta, Y_0^m) \right). \end{aligned}$$

Informally, in the proof of Proposition 6.3, we use the mixing property from condition (S5) to parse a sequence of observations into alternating sequences of “large blocks” and “small blocks,” and then the ergodicity and integrability conditions finish the proof. More specifically, we break up the sequence of observations  $Y_0^n$  into alternating blocks of length  $m$  and  $\ell$ , where  $\ell$  is given by condition (S5).

PROOF. Let  $\theta' \notin [\theta_0]$ . Fix a neighborhood  $U$  of  $\theta'$  so that the conclusions of both condition (S3) and condition (S5) hold. Let  $m > 0$  be arbitrary, and let  $\ell$  be as in condition (S5). We consider sequences of observations of length  $n$ , where  $n$  is a large integer. These sequences of observations will be parsed into alternating blocks of lengths  $m$  and  $\ell$ , respectively, starting from an offset of size  $s$  and possibly ending with a remainder sequence. For the sake of notation, we use interval notation to denote intervals of integers. For  $n > 2(m + \ell)$  and  $s$  in  $[0, m + \ell)$ , let  $R = R(s, m, \ell, n) \in [0, m + \ell)$  and  $k = k(s, m, \ell, n) \geq 0$  be defined by the condition  $n = s + k(m + \ell) + R$ . Then we partition  $[0, n]$  as follows:

$$\begin{aligned} B_s &= [0, s), \\ I_s(j) &= [s + (m + \ell)(j - 1), s + (m + \ell)(j - 1) + m) \quad \text{for } 1 \leq j \leq k, \\ J_s(j) &= [s + (m + \ell)(j - 1) + m, s + (m + \ell)j) \quad \text{for } 1 \leq j \leq k, \\ E_s &= [s + t(m + \ell), n]. \end{aligned}$$

Given a sequence  $Y_0^n$  of observations, we define the following subsequences of  $Y_0^n$  according to the above partitions of  $[0, n]$ :

$$\begin{aligned} b_s &= Y|_{B_s}, \\ w_s(j) &= Y|_{I_s(j)} \quad \text{for } 1 \leq j \leq k, \\ v_s(j) &= Y|_{J_s(j)} \quad \text{for } 1 \leq j \leq k, \\ e_s &= Y|_{E_s}. \end{aligned}$$

For a sequence  $y_0^t$  in  $\mathcal{Y}^{t+1}$ , define

$$\gamma_\theta(y_0^t) = \prod_{j=0}^t \gamma_\theta(y_j) = \prod_{j=0}^t \sup_x g_\theta(y_j|x).$$

Then for  $\theta$  in  $U$ , it follows from condition (S5) that

$$p_\theta(Y_0^n) \leq \gamma_\theta(b_s) \gamma_\theta(e_s) \cdot \prod_{j=1}^k \gamma_\theta(v_s(j)) \cdot \prod_{j=1}^k C_m(\theta, w_s(j)) \cdot \prod_{j=1}^k p_\theta(w_s(j)).$$



Taking the logarithm of both sides and averaging over  $s$  in  $[0, m + \ell)$ , we obtain

$$\begin{aligned}
 \log p_\theta(Y_0^n) &\leq \frac{1}{m + \ell} \sum_{s=0}^{m+\ell-1} \sum_{j=1}^k [\log p_\theta(w_s(j)) + \log C_m(\theta, w_s(j))] \\
 (6.5) \quad &+ \frac{1}{m + \ell} \sum_{s=0}^{m+\ell-1} \sum_{j=1}^k \log \gamma_\theta(v_s(j)) \\
 &+ \frac{1}{m + \ell} \sum_{s=0}^{m+\ell-1} [\log \gamma_\theta(b_s) + \log \gamma_\theta(e_s)].
 \end{aligned}$$

Let us now take the supremum over  $\theta$  in  $U$  in (6.5) and evaluate the limits of the three terms on the right-hand side as  $n$  tends to infinity.

Let  $\xi_1: \mathbf{Y}^{m+1} \rightarrow \mathbb{R}$  and  $\xi_2: \mathbf{Y}^{m+1} \rightarrow \mathbb{R}$  be defined by

$$\begin{aligned}
 \xi_1(y_0^m) &= \sup_{\theta \in U} \log p_\theta(y_0^m), \\
 \xi_2(y_0^m) &= \sup_{\theta \in U} \log C_m(\theta, y_0^m).
 \end{aligned}$$

With this notation, we have that

$$\begin{aligned}
 &\frac{1}{n} \sum_{s=0}^{m+\ell-1} \sum_{j=1}^k \left[ \sup_{\theta \in U} \log p_\theta(w_s(j)) + \sup_{\theta \in U} \log C_m(\theta, w_s(j)) \right] \\
 &= \frac{1}{n} \sum_{i=0}^n [\xi_1(Y_i^{i+m}) + \xi_2(Y_i^{i+m})].
 \end{aligned}$$

Since  $(Y_k)$  is ergodic (by Proposition 6.1), it follows from Birkhoff's ergodic theorem and conditions (S3) and (S5) that the following limit exists  $\mathbb{P}_{\theta_0}$ -a.s.:

$$\begin{aligned}
 &\lim_n \frac{1}{n} \sum_{s=0}^{m+\ell-1} \sum_{j=1}^k \left[ \sup_{\theta \in U} \log p_\theta(w_s(j)) + \sup_{\theta \in U} \log C_m(\theta, w_s(j)) \right] \\
 &= \lim_n \frac{1}{n} \sum_{i=0}^n [\xi_1(Y_i^{i+m}) + \xi_2(Y_i^{i+m})] \\
 (6.6) \quad &= \mathbb{E}_{\theta_0}[\xi_1(Y_0^m)] + \mathbb{E}_{\theta_0}[\xi_2(Y_0^m)] \\
 &= \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log p_\theta(Y_0^m) \right] \\
 &\quad + \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log C_m(\theta, Y_0^m) \right].
 \end{aligned}$$

Similarly, using Birkhoff's ergodic theorem and condition (S3), we have that the following holds  $\mathbb{P}_{\theta_0}$ -a.s.:

$$\begin{aligned}
(6.7) \quad \limsup_n \frac{1}{n} \sum_{s=0}^{m+\ell-1} \sum_{j=1}^k \sup_{\theta \in U} \log \gamma_\theta(v_s(j)) &\leq \limsup_n \frac{1}{n} \sum_{i=0}^n \sup_{\theta \in U} \log^+ \gamma_\theta(Y_i^{i+\ell-1}) \\
&\leq \ell \limsup_n \frac{1}{n} \sum_{i=0}^n \sup_{\theta \in U} \log^+ \gamma_\theta(Y_i) \\
&= \ell \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log^+ \gamma_\theta(Y_0) \right].
\end{aligned}$$

Finally, Birkhoff's ergodic theorem and condition (S3) again imply that the following limit holds  $\mathbb{P}_{\theta_0}$ -a.s.:

$$(6.8) \quad \lim_n \frac{1}{n} \sum_{s=0}^{m+\ell-1} \left[ \sup_{\theta \in U} \log^+ \gamma_\theta(b_s) + \sup_{\theta \in U} \log^+ \gamma_\theta(e_s) \right] = 0,$$

where we have used that  $\max(|B_s|, |E_s|) \leq m + \ell$ .

Combining the inequalities in (6.5)–(6.8), we obtain that

$$\begin{aligned}
\limsup_{n \rightarrow \infty} \sup_{\theta \in U} \frac{1}{n} \log p_\theta(Y_0^n) &\leq \frac{1}{m + \ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log p_\theta(Y_0^m) \right] \\
&\quad + \frac{1}{m + \ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log C_m(\theta, Y_0^m) \right] \\
&\quad + \frac{\ell}{m + \ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log^+ \gamma_\theta(Y_0) \right],
\end{aligned}$$

as desired.  $\square$

The following proposition is a direct application of Lemma 10 in [15] to the present setting, and we omit the proof.

**PROPOSITION 6.4.** *Suppose that the following conditions hold: ergodicity (S1), logarithmic integrability at  $\theta_0$  (S2) and exponential identifiability (S6). Then for  $\theta \notin [\theta_0]$ , it holds that*

$$\limsup_n \frac{1}{n} \mathbb{E}_{\theta_0} [\log p_\theta(Y_0^n)] < h(\theta_0).$$

The following proposition provides an essential estimate in the proof of Theorem 3.1.

PROPOSITION 6.5. *Suppose that conditions (S1)–(S6) hold, and let  $\ell$  be as in (S5). Then for  $\theta' \notin [\theta_0]$ , there exists  $m > 0$  and a neighborhood  $U$  of  $\theta'$  such that*

$$\begin{aligned} h(\theta_0) &> \frac{1}{m+\ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log p_{\theta}(Y_0^m) \right] \\ &+ \frac{\ell}{m+\ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log^+ \gamma_{\theta}(Y_0) \right] \\ &+ \frac{1}{m+\ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log C_m(\theta', Y_0^m) \right]. \end{aligned}$$

PROOF. Suppose  $\theta' \notin [\theta_0]$ . By Proposition 6.4, there exists  $\varepsilon > 0$  such that

$$(6.9) \quad \limsup_n \frac{1}{n} \mathbb{E}_{\theta_0} [\log p_{\theta'}(Y_0^n)] < h(\theta_0) - \varepsilon.$$

By conditions (S3) (logarithmic integrability away from  $\theta_0$ ) and (S5) (mixing), there exists a neighborhood  $U'$  of  $\theta'$  and  $m_0 > 0$  such that for  $m \geq m_0$ , we have that

$$(6.10) \quad \begin{aligned} \varepsilon/2 &> \frac{\ell}{m+\ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U'} \log^+ \gamma_{\theta}(Y_0) \right] \\ &+ \frac{1}{m+\ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U'} \log C_m(\theta, Y_0^m) \right]. \end{aligned}$$

Fix  $m \geq m_0$  such that

$$(6.11) \quad \frac{1}{m+\ell} \mathbb{E}_{\theta_0} [\log p_{\theta'}(Y_0^m)] < \limsup_n \frac{1}{n} \mathbb{E}_{\theta_0} [\log p_{\theta'}(Y_0^n)] + \varepsilon/4.$$

For  $\eta > 0$ , let  $B(\theta', \eta)$  denote the ball of radius  $\eta$  about  $\theta'$  in  $\Theta$ . For  $\eta$  such that  $B(\theta', \eta) \subset U'$ , we have that

$$\sup_{\theta \in B(\theta', \eta)} \log p_{\theta}(Y_0^m) \leq \sum_{k=0}^m \sup_{\theta \in U'} \log^+ \gamma_{\theta}(Y_k).$$

The sum above is integrable with respect to  $\mathbb{P}_{\theta_0}$  and does not depend on  $\eta$ . Then (the reverse) Fatou's Lemma implies that

$$\limsup_{\eta \rightarrow 0} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in B(\theta', \eta)} \log p_{\theta}(Y_0^m) \right] \leq \mathbb{E}_{\theta_0} \left[ \limsup_{\eta \rightarrow 0} \sup_{\theta \in B(\theta', \eta)} \log p_{\theta}(Y_0^m) \right].$$

By condition (S4) [upper semi-continuity of  $\theta \mapsto p_{\theta}(Y_0^m)$ ], we see that

$$\mathbb{E}_{\theta_0} \left[ \limsup_{\eta \rightarrow 0} \sup_{\theta \in B(\theta', \eta)} \log p_{\theta}(Y_0^m) \right] \leq \mathbb{E}_{\theta_0} [\log p_{\theta'}(Y_0^m)].$$

Now by an appropriate choice of  $\eta > 0$ , we have shown that there exists a neighborhood  $U \subset U'$  of  $\theta'$  such that

$$(6.12) \quad \frac{1}{m+\ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log p_{\theta}(Y_0^m) \right] < \frac{1}{m+\ell} \mathbb{E}_{\theta_0} [\log p_{\theta'}(Y_0^m)] + \varepsilon/4.$$

Combining estimates (6.9)–(6.12), we obtain the desired inequality.  $\square$

**PROOF OF THEOREM 3.1.** Let  $h(\theta_0)$  be defined as in Proposition 6.2. We prove the theorem by showing the following statement: for each closed set  $C$  in  $\Theta$  such that  $C \cap [\theta_0] = \emptyset$ , it holds that

$$(6.13) \quad \limsup_n \sup_{\theta \in C} \frac{1}{n} \log p_{\theta}(Y_0^n) < h(\theta_0).$$

Let  $C$  be a closed subset of  $\Theta$  such that  $C \cap [\theta_0] = \emptyset$ . Since  $\Theta$  is compact,  $C$  is compact. Suppose that for each  $\theta' \in C$ , there exists a neighborhood  $U$  of  $\theta'$  such that

$$(6.14) \quad \limsup_n \sup_{\theta \in U \cap C} \frac{1}{n} \log p_{\theta}(Y_0^n) < h(\theta_0).$$

Then by compactness, we would conclude that (6.13) holds and thus complete the proof of the theorem.

Let  $\theta'$  be in  $C$ . Let us now show that there exists a neighborhood  $U$  of  $\theta'$  such that (6.14) holds. Since  $\theta'$  is in  $C$ , we have that  $\theta' \notin [\theta_0]$ . Let  $\ell$  be as in (S5). By Proposition 6.5, there exists  $m > 0$  and a neighborhood  $U'$  of  $\theta'$  such that

$$(6.15) \quad \begin{aligned} h(\theta_0) &> \frac{1}{m+\ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U'} \log p_{\theta}(Y_0^m) \right] \\ &+ \frac{\ell}{m+\ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U'} \log^+ \gamma_{\theta}(Y_0) \right] \\ &+ \frac{1}{m+\ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U'} \log C_m(\theta, Y_0^m) \right]. \end{aligned}$$

By Proposition 6.3, there exists a neighborhood  $U \subset U'$  of  $\theta'$  such that

$$(6.16) \quad \begin{aligned} \limsup_{n \rightarrow \infty} \sup_{\theta \in U} \frac{1}{n} \log p_{\theta}(Y_0^n) &\leq \frac{1}{m+\ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log p_{\theta}(Y_0^m) \right] \\ &+ \frac{\ell}{m+\ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log^+ \gamma_{\theta}(Y_0) \right] \\ &+ \frac{1}{m+\ell} \mathbb{E}_{\theta_0} \left[ \sup_{\theta \in U} \log C_m(\theta', Y_0^m) \right]. \end{aligned}$$

Combining (6.15) and (6.16), we obtain (6.14), which completes the proof of the theorem.  $\square$

**7. Concluding remarks.** In this paper, we demonstrate how the properties of a family of dynamical systems affect the asymptotic consistency of maximum likelihood parameter estimation. We have exhibited a collection of general statistical conditions on families of dynamical systems observed with noise, and we have shown that under these general conditions, maximum likelihood estimation is a consistent method of parameter estimation. Furthermore, we have shown that these general conditions are indeed satisfied by some classes of well-studied families of dynamical systems. As mentioned in the [Introduction](#), our results can be considered as a theoretical validation of the notion from dynamical systems that these classes of systems have “good” statistical properties.

However, there remain interesting families of systems to which our results do not apply, including some classes of systems that are also believed to have “good” statistical properties. In particular, the class of systems modeled by Young towers with exponential tail [51] has exponential decay of correlations and certain large deviations estimates [41]. These families include a positive measure set of maps from the quadratic family  $[x \mapsto ax(1-x)]$  and the Hénon family, as well as certain billiards and many other systems of physical and mathematical interest [51]. In short, the setting of systems modeled by Young towers with exponential tail provides a very attractive setting in which to consider consistency of maximum likelihood estimation. Unfortunately, our proof does not apply to systems in this setting in general, mainly due to the presence of the mixing condition (S5), which is not satisfied by these systems in general.

A natural next step might be to obtain rates of convergence and derive central limit theorems for maximum likelihood estimation. To this end, it might be possible to build off of analogous results for HMMs [8, 23]. We leave these questions for future work.

## SUPPLEMENTARY MATERIAL

**Supplement to “Consistency of maximum likelihood estimation for some dynamical systems”** (DOI: [10.1214/14-AOS1259SUPP](https://doi.org/10.1214/14-AOS1259SUPP); .pdf). We provide three technical appendices. In Appendix A, we present proofs of Propositions 4.1, 4.2 and 4.3. In Appendix B, we discuss shifts of finite type and Gibbs measures and prove Theorem 5.1. Finally, Appendix C contains definitions for Axiom A systems, as well as a proof of Theorem 5.7.

## REFERENCES

- [1] ADAMS, T. M. and NOBEL, A. B. (2001). Finitary reconstruction of a measure preserving transformation. *Israel J. Math.* **126** 309–326. [MR1882042](#)
- [2] ALVES, J. F., CARVALHO, M. and FREITAS, J. M. (2010). Statistical stability and continuity of SRB entropy for systems with Gibbs–Markov structures. *Comm. Math. Phys.* **296** 739–767. [MR2628821](#)

- [3] BALADI, V. (2001). Decay of correlations. In *Smooth Ergodic Theory and Its Applications (Seattle, WA, 1999)*. *Proc. Sympos. Pure Math.* **69** 297–325. Amer. Math. Soc., Providence, RI. [MR1858537](#)
- [4] BARRON, A. R. (1985). The strong ergodic theorem for densities: Generalized Shannon–McMillan–Breiman theorem. *Ann. Probab.* **13** 1292–1303. [MR0806226](#)
- [5] BAUM, L. E. and PETRIE, T. (1966). Statistical inference for probabilistic functions of finite state Markov chains. *Ann. Math. Statist.* **37** 1554–1563. [MR0202264](#)
- [6] BERLINER, L. M. (1992). Statistics, probability and chaos. *Statist. Sci.* **7** 69–122. [MR1173418](#)
- [7] BERTSEKAS, D. P. and SHREVE, S. E. (1978). *Stochastic Optimal Control: The Discrete Time Case*. *Mathematics in Science and Engineering* **139**. Academic Press, New York. [MR0511544](#)
- [8] BICKEL, P. J., RITOV, Y. and RYDÉN, T. (1998). Asymptotic normality of the maximum-likelihood estimator for general hidden Markov models. *Ann. Statist.* **26** 1614–1635. [MR1647705](#)
- [9] BOWEN, R. (1970). Markov partitions for Axiom A diffeomorphisms. *Amer. J. Math.* **92** 725–747. [MR0277003](#)
- [10] BOWEN, R. (2008). *Equilibrium States and the Ergodic Theory of Anosov Diffeomorphisms*, revised ed. *Lecture Notes in Math.* **470**. Springer, Berlin. [MR2423393](#)
- [11] CAPPÉ, O., MOULINES, E. and RYDÉN, T. (2005). *Inference in Hidden Markov Models*. Springer, New York. [MR2159833](#)
- [12] CHATTERJEE, S. and YILMAZ, M. R. (1992). Chaos, fractals and statistics. *Statist. Sci.* **7** 49–68. [MR1173417](#)
- [13] DOUC, R. and MATIAS, C. (2001). Asymptotics of the maximum likelihood estimator for general hidden Markov models. *Bernoulli* **7** 381–420. [MR1836737](#)
- [14] DOUC, R. and MOULINES, E. (2012). Asymptotic properties of the maximum likelihood estimation in misspecified hidden Markov models. *Ann. Statist.* **40** 2697–2732. [MR3097617](#)
- [15] DOUC, R., MOULINES, E., OLSSON, J. and VAN HANDEL, R. (2011). Consistency of the maximum likelihood estimator for general hidden Markov models. *Ann. Statist.* **39** 474–513. [MR2797854](#)
- [16] DOUC, R., MOULINES, É. and RYDÉN, T. (2004). Asymptotic properties of the maximum likelihood estimator in autoregressive models with Markov regime. *Ann. Statist.* **32** 2254–2304. [MR2102510](#)
- [17] FREITAS, J. M. and TODD, M. (2009). The statistical stability of equilibrium states for interval maps. *Nonlinearity* **22** 259–281. [MR2475546](#)
- [18] GENON-CATALOT, V. and LAREDO, C. (2006). Leroux’s method for general hidden Markov models. *Stochastic Process. Appl.* **116** 222–243. [MR2197975](#)
- [19] IONIDES, E., BRETÓ, C. and KING, A. (2006). Inference for nonlinear dynamical systems. *Proc. Natl. Acad. Sci. USA* **103** 18438–18443.
- [20] IONIDES, E. L., BHADRA, A., ATCHADÉ, Y. and KING, A. (2011). Iterated filtering. *Ann. Statist.* **39** 1776–1802. [MR2850220](#)
- [21] ISHAM, V. (1993). Statistical aspects of chaos: A review. In *Networks and Chaos—Statistical and Probabilistic Aspects*. *Monogr. Statist. Appl. Probab.* **50** 124–200. Chapman & Hall, London. [MR1314654](#)
- [22] JENSEN, J. L. (1993). Chaotic dynamical systems with a view towards statistics: A review. In *Networks and Chaos—Statistical and Probabilistic Aspects*. *Monogr. Statist. Appl. Probab.* **50** 201–250. Chapman & Hall, London. [MR1314655](#)
- [23] JENSEN, J. L. and PETERSEN, N. V. (1999). Asymptotic normality of the maximum likelihood estimator in state space models. *Ann. Statist.* **27** 514–535. [MR1714719](#)

- [24] JUDD, K. (2007). Failure of maximum likelihood methods for chaotic dynamical systems. *Phys. Rev. E (3)* **75** 036210.
- [25] LALLEY, S. P. (1999). Beneath the noise, chaos. *Ann. Statist.* **27** 461–479. [MR1714721](#)
- [26] LALLEY, S. P. (2001). Removing the noise from chaos plus noise. In *Nonlinear Dynamics and Statistics (Cambridge, 1998)* 233–244. Birkhäuser, Boston, MA. [MR1937487](#)
- [27] LALLEY, S. P. and NOBEL, A. B. (2006). Denoising deterministic time series. *Dyn. Partial Differ. Equ.* **3** 259–279. [MR2271730](#)
- [28] LAW, K. J. H. and STUART, A. M. (2012). Evaluating data assimilation algorithms. *Monthly Weather Review* **140** 3757–3782.
- [29] LEROUX, B. G. (1992). Maximum-likelihood estimation for hidden Markov models. *Stochastic Process. Appl.* **40** 127–143. [MR1145463](#)
- [30] LE GLAND, F. and MEVEL, L. (2000). Basic properties of the projective product with application to products of column-allowable nonnegative matrices. *Math. Control Signals Systems* **13** 41–62. [MR1742139](#)
- [31] LE GLAND, F. and MEVEL, L. (2000). Exponential forgetting and geometric ergodicity in hidden Markov models. *Math. Control Signals Systems* **13** 63–93. [MR1742140](#)
- [32] MCGOFF, K., MUKHERJEE, S., NOBEL, A. and PILLAI, N. (2014). Supplement to “Consistency of maximum likelihood estimation for some dynamical systems.” DOI:[10.1214/14-AOS1259SUPP](#).
- [33] MCGOFF, K., MUKHERJEE, S. and PILLAI, N. (2013). Statistical inference for dynamical systems: A review. Available at [arXiv:1204.6265](#).
- [34] NOBEL, A. (2001). Consistent estimation of a dynamical map. In *Nonlinear Dynamics and Statistics (Cambridge, 1998)* 267–280. Birkhäuser, Boston, MA. [MR1937489](#)
- [35] NOBEL, A. B. and ADAMS, T. M. (2001). Estimating a function from ergodic samples with additive noise. *IEEE Trans. Inform. Theory* **47** 2895–2902. [MR1872848](#)
- [36] PETERSEN, K. (1989). *Ergodic Theory. Cambridge Studies in Advanced Mathematics* **2**. Cambridge Univ. Press, Cambridge. [MR1073173](#)
- [37] PETRIE, T. (1969). Probabilistic functions of finite state Markov chains. *Ann. Math. Statist.* **40** 97–115. [MR0239662](#)
- [38] PISARENKO, V. F. and SORNETTE, D. (2004). Statistical methods of parameter estimation for deterministically chaotic time series. *Phys. Rev. E (3)* **69** 036122, 12. [MR2096393](#)
- [39] POOLE, D. and RAFTERY, A. E. (2000). Inference for deterministic simulation models: The Bayesian melding approach. *J. Amer. Statist. Assoc.* **95** 1244–1255. [MR1804247](#)
- [40] RAMSAY, J. O., HOOKER, G., CAMPBELL, D. and CAO, J. (2007). Parameter estimation for differential equations: A generalized smoothing approach. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **69** 741–796. [MR2368570](#)
- [41] REY-BELLET, L. and YOUNG, L.-S. (2008). Large deviations in non-uniformly hyperbolic dynamical systems. *Ergodic Theory Dynam. Systems* **28** 587–612. [MR2408394](#)
- [42] RUELLE, D. (1997). Differentiation of SRB states. *Comm. Math. Phys.* **187** 227–241. [MR1463827](#)
- [43] RUELLE, D. (2004). *Thermodynamic Formalism: The Mathematical Structures of Equilibrium Statistical Mechanics*, 2nd ed. Cambridge Univ. Press, Cambridge. [MR2129258](#)
- [44] SHALIZI, C. R. (2009). Dynamics of Bayesian updating with dependent data and misspecified models. *Electron. J. Stat.* **3** 1039–1074. [MR2557128](#)



- [45] STEINWART, I. and ANGHEL, M. (2009). Consistency of support vector machines for forecasting the evolution of an unknown ergodic dynamical system from observations with unknown noise. *Ann. Statist.* **37** 841–875. [MR2502653](#)
- [46] TONI, T., WELCH, D., STRELKOWA, N., IPSEN, A. and STUMPF, M. (2009). Approximate Bayesian computation scheme for parameter inference and model selection in dynamical systems. *J. R. Soc. Interface* **6** 187–202.
- [47] VÁSQUEZ, C. H. (2007). Statistical stability for diffeomorphisms with dominated splitting. *Ergodic Theory Dynam. Systems* **27** 253–283. [MR2297096](#)
- [48] WALTERS, P. (1982). *An Introduction to Ergodic Theory. Graduate Texts in Mathematics* **79**. Springer, New York. [MR0648108](#)
- [49] WOOD, S. N. (2010). Statistical inference for noisy nonlinear ecological dynamic systems. *Nature* **466** 1102–1104.
- [50] YOUNG, L.-S. (1990). Large deviations in dynamical systems. *Trans. Amer. Math. Soc.* **318** 525–543. [MR0975689](#)
- [51] YOUNG, L.-S. (1998). Statistical properties of dynamical systems with some hyperbolicity. *Ann. of Math. (2)* **147** 585–650. [MR1637655](#)
- [52] YOUNG, L.-S. (2002). What are SRB measures, and which dynamical systems have them? *J. Stat. Phys.* **108** 733–754. [MR1933431](#)

K. MCGOFF  
 DEPARTMENT OF MATHEMATICS  
 DUKE UNIVERSITY  
 DURHAM, NORTH CAROLINA 27708  
 USA  
 E-MAIL: [mcgoff@math.duke.edu](mailto:mcgoff@math.duke.edu)

A. NOBEL  
 DEPARTMENT OF STATISTICS  
 AND OPERATIONS RESEARCH  
 UNIVERSITY OF NORTH CAROLINA  
 CHAPEL HILL, NORTH CAROLINA 27599-3260  
 USA  
 E-MAIL: [nobel@email.unc.edu](mailto:nobel@email.unc.edu)

S. MUKHERJEE  
 DEPARTMENTS OF STATISTICAL SCIENCE,  
 COMPUTER SCIENCE, AND MATHEMATICS  
 INSTITUTE FOR GENOME SCIENCES & POLICY  
 DUKE UNIVERSITY  
 DURHAM, NORTH CAROLINA 27708  
 USA  
 E-MAIL: [sayan@stat.duke.edu](mailto:sayan@stat.duke.edu)

N. PILLAI  
 DEPARTMENT OF STATISTICS  
 HARVARD UNIVERSITY  
 CAMBRIDGE, MASSACHUSETTS 02138  
 USA  
 E-MAIL: [pillai@fas.harvard.edu](mailto:pillai@fas.harvard.edu)