

Consistent Order Estimation and Minimal Penalties

Elisabeth Gassiat and Ramon van Handel

Abstract—Consider an i.i.d. sequence of random variables whose distribution f^* lies in one of the nested families of models \mathcal{M}_q , $q \geq 1$. The smallest index q^* such that \mathcal{M}_{q^*} contains f^* is called the model order. The aim of this paper is to explore the consistency properties of penalized likelihood model order estimators such as Bayesian information criterion. We show in a general setting that the minimal strongly consistent penalty is of order $\eta(q) \log \log n$, where $\eta(q)$ is a dimensional quantity. In contrast to previous work, an *a priori* upper bound on the model order is not assumed. The results rely on a sharp characterization of the pathwise fluctuations of the generalized likelihood ratio statistic under entropy assumptions on the model classes. Our results are applied to the geometrically complex problem of location mixture order estimation, which is widely used but poorly understood.

Index Terms—Consistent order estimation, location mixtures, penalized likelihood, uniform law of iterated logarithm.

I. INTRODUCTION

LET $(X_k)_{k \geq 1}$ be a sequence of random variables whose distribution f^* lies in one of the nested families of models $(\mathcal{M}_q)_{q \geq 1}$, indexed (and ordered) by the integers. We define the model order as the smallest index q^* such that the true distribution f^* lies in the corresponding model class. The model order typically determines the most parsimonious representation of the true distribution of the underlying model (for example, it might determine the parametrization of the model which has the smallest possible dimension). On the other hand, the model order often has a concrete interpretation in terms of the modeling of the underlying phenomenon (for example, the estimation of the number of clusters in a dataset, or the number of regimes in an economic time series). Therefore, the problem of estimating the model order from observed data is of significant practical, as well as theoretical, interest.

Of course, a satisfactory solution to this problem must provide an estimation method that does not assume prior knowledge on the unknown distribution f^* . In particular, prior bounds on model order and on parameter sets should be avoided. Yet, in this light, even one of the most widely used model selection criteria—the Bayesian information criterion (BIC) of Schwarz—is poorly understood. The chief motivation for the use of BIC (as opposed to other model selection criteria, such as Akaike’s information criterion) is that it is expected to yield a strongly con-

Manuscript received February 21, 2012; revised August 06, 2012; accepted September 23, 2012. Date of publication October 02, 2012; date of current version January 16, 2013.

E. Gassiat is with the Laboratoire de Mathématiques d’Orsay, Université Paris-Sud, 91405 Orsay Cedex, France (e-mail: elisabeth.gassiat@math.u-psud.fr).

R. van Handel is with the Department of Operations Research and Financial Engineering, Princeton University, Princeton, NJ 08544 USA (e-mail: rvan@princeton.edu).

Communicated by G. Moustakides, Associate Editor for Detection and Estimation.

Digital Object Identifier 10.1109/TIT.2012.2221122

sistent estimator of the model order. However, almost all existing consistency proofs assume a prior upper bound on the order as well as compactness of the parameter set. As is emphasized by Csiszár and Shields [1], this is hardly satisfactory from the theoretical point of view and provides little confidence in the basic motivation for this method. More delicate questions, such as the minimal penalty that yields a consistent order estimator in absence of a prior bound on the order, remain open (the problem of identifying the minimal penalty, which minimizes the probability of underestimating the order, is also raised in [1]).

In this paper, we consider a general class of penalized likelihood order estimators of the form

$$\hat{q}_n = \operatorname{argmax}_{q \geq 1} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \operatorname{pen}(n, q) \right\}$$

where $\operatorname{pen}(n, q)$ is a penalty function and $\ell_n(f)$ is the likelihood of $(X_k)_{1 \leq k \leq n}$ under the distribution f . Our aim is to understand what penalties yield strong consistency of the order estimator, i.e., $\hat{q}_n \rightarrow q^*$ as $n \rightarrow \infty$ a.s. Characterizing strong consistency hinges on a precise understanding of the pathwise fluctuations of the likelihood ratio statistic

$$\sup_{f \in \mathcal{M}_q} \ell_n(f) - \sup_{f \in \mathcal{M}_{q^*}} \ell_n(f)$$

as $n \rightarrow \infty$, uniformly in the model order $q > q^*$. When there is a known upper bound on the order $q^* \leq q_{\max} < \infty$ and the model classes \mathcal{M}_q are regularly parametrized by a compact subset of Euclidean space, an upper bound on the pathwise fluctuations can be obtained by classical parametric methods: Taylor expansion of the likelihood and an application of a law of iterated logarithm. This approach forms the basis for most consistency proofs for penalized likelihood order estimators in the literature (see, for example, [2]–[6]). However, such techniques fail in the absence of a prior upper bound: even though each model class \mathcal{M}_q is finite dimensional, the full model $\mathcal{M} = \bigcup_q \mathcal{M}_q$ is infinite dimensional and, as such, the problem in the absence of a prior upper bound is inherently nonparametric. When the classes \mathcal{M}_q are noncompact, one must introduce sieves $\mathcal{M}_q^n \subset \mathcal{M}_q^{n+1} \subset \dots \subset \mathcal{M}_q$, complicating the problem further (in this case, even the parametric theory remains poorly understood [7]–[9]). An entirely different approach based on universal coding theory [4], [10]–[14] yields bounds on the pathwise fluctuations that do not require prior bounds on the order or compactness of the models. However, these bounds are far from tight and cannot even establish consistency of BIC, let alone smaller penalties (this appears to be a fundamental limitation of this approach due to Rissanen’s theorem; see [15] and [16]).

The problem area that is investigated in this paper was initiated in the work of Csiszár and Shields [1], [17], who proved consistency of BIC for Markov chain order estimation in absence of a prior bound on the order (see also [18]). To our

knowledge, little progress has been made on this subject beyond their work. The proofs in [1] and [17] rely heavily on the availability of an explicit expression for the maximum likelihood for Markov chains and employ delicate estimates specific to that setting. Their techniques are, therefore, not well suited to investigating such problems in other settings. Moreover, the methods in [1] and [17] do not yield minimal penalties. However, the Markov chain case was recently reconsidered in [19] using very different techniques based on empirical process theory, which are potentially much more generally applicable and which shed light on minimal penalties.

The main results of this paper provide generally applicable upper and lower bounds on the pathwise fluctuations of the likelihood ratio statistic uniformly in the model order $q > q^*$, for the case of i.i.d. observations $(X_k)_{k \geq 1}$, without a prior bound on the model order and in possibly noncompact parameter spaces. These results are then used to investigate strong consistency of penalized likelihood order estimators. We use empirical process methods as in [19], but the difficulties to be surmounted in the present setting are of a different nature. The main difficulty for Markov chain models in [1], [17], and [19] is their dependence structure; in this paper, we assume i.i.d. models. On the other hand, the geometric structure of Markov chains is exceedingly simple: the family of q th-order Markov chains in the Hellinger distance is simply a Euclidean ball when viewed in the appropriate parametrization. In contrast, in general order estimation problems, one is often faced with model classes that are geometrically very complex. An important case study that will be considered in this paper is location mixture models (widely used in practice for clustering), which possess a notoriously complicated nonregular geometry. We will be able, for example, to establish strong consistency of BIC for mixture order estimation in absence of a prior bound on the order or on the parameter set, providing a counterpart to the results of Csiszár and Shields [1] in a setting very different than that of Markov chains.

The techniques developed here originate in our attempts to understand the order estimation problem for hidden Markov models (HMM) [12]. In that setting, consistency of BIC (even with a prior bound on the order) remains unknown. The two cases considered here and in [19]—Markov chains and i.i.d. mixtures—can be viewed as two extreme cases of HMM. While our approach provides a substantial step toward understanding the HMM setting, a striking and as of yet poorly understood breakdown in the ergodicity of HMM [20] has so far prohibited further progress in this direction.

The remainder of this paper is organized as follows. Section II introduces the general model under consideration and states our results on the pathwise fluctuations of the likelihood ratio statistic. Section III derives the consequences for order estimation and considers also the special case of location mixture models. Proofs are given in the appendixes.

II. PATHWISE FLUCTUATIONS OF THE LIKELIHOOD

A. Basic Setting and Notation

Let (E, \mathcal{E}, μ) be a measure space. For each $q, n \geq 1$, let \mathcal{M}_q^n be a given family of strictly positive probability densities with respect to μ (that is, we assume that $\int f d\mu = 1$ and that

$f > 0$ μ -a.e. for every $f \in \mathcal{M}_q^n$). Moreover, we assume that $(\mathcal{M}_q^n)_{q,n \geq 1}$ is a nested family of models in the sense that $\mathcal{M}_q^n \subseteq \mathcal{M}_{q+1}^n$ and $\mathcal{M}_q^n \subseteq \mathcal{M}_q^{n+1}$ for all $q, n \geq 1$. Let $\mathcal{M}_q = \bigcup_n \mathcal{M}_q^n$, $\mathcal{M}^n = \bigcup_q \mathcal{M}_q^n$, and $\mathcal{M} = \bigcup_{q,n} \mathcal{M}_q^n$.

Consider an i.i.d. sequence of E -valued random variables $(X_k)_{k \geq 1}$ whose common distribution under the measure \mathbf{P}^* is $f^* d\mu$, where $f^* \in \mathcal{M}_{q^*} \setminus \text{cl } \mathcal{M}_{q^*-1}$ for some $q^* \geq 1$ (here $\text{cl } \mathcal{M}_q$ denotes the $L^1(d\mu)$ -closure of \mathcal{M}_q). The index q^* is called the *model order*. Let us define

$$\ell_n(f) = \sum_{i=1}^n \log f(X_i), \quad f \in \mathcal{M}.$$

Evidently, $\ell_n(f)$ is the log-likelihood of the i.i.d. sequence $(X_k)_{k \leq n}$ when $X_k \sim f d\mu$. Our aim is to study the pathwise fluctuations of the likelihood ratio statistic

$$\sup_{f \in \mathcal{M}_q^n} \ell_n(f) - \sup_{f \in \mathcal{M}_{q^*}^n} \ell_n(f)$$

as $n \rightarrow \infty$, uniformly over the order parameter $q \geq q^*$. Pathwise upper and lower bounds on the likelihood ratio statistic are the key ingredient in the study of strong consistency of penalized likelihood order estimators (see Section III).

Example II.1 (Location Mixtures): The guiding example for our theory, the case of location mixtures, will be studied in detail in Section III-B. We presently introduce this example in order to clarify our basic setup.

Let $E = \mathbb{R}^d$ (with its Borel σ -field \mathcal{E}) and let μ be the Lebesgue measure on \mathbb{R}^d . We fix a strictly positive probability density f_0 with respect to μ , and define $f_\theta(x) = f_0(x - \theta)$ for $x, \theta \in \mathbb{R}^d$. Fix a sequence $T(n) \uparrow \infty$ and define

$$\mathcal{M}_q^n = \left\{ \sum_{i=1}^q \pi_i f_{\theta_i} : \pi_i \geq 0, \sum_{i=1}^q \pi_i = 1, \|\theta_i\| \leq T(n) \right\}.$$

Then, \mathcal{M}_q is the family of all q -component mixtures of translates of the density f_0 , while \mathcal{M}_q^n is the subset of the mixtures \mathcal{M}_q whose translation parameters $(\theta_i)_{i=1,\dots,q}$ are restricted to a ball of radius $T(n)$. The number of components q^* of the true mixture $f^* \in \mathcal{M}$ can be estimated from observations using the order estimator

$$\hat{q}_n = \underset{q \geq 1}{\operatorname{argmax}} \left\{ \sup_{f \in \mathcal{M}_q^n} \ell_n(f) - \text{pen}(n, q) \right\}.$$

Pathwise control of the likelihood ratio statistic allows us to identify what penalties $\text{pen}(n, q)$ and cutoff sequences $T(n)$ yield strong consistency of \hat{q}_n (cf., Section III-B).

Remark II.2: To avoid measurability problems and other technical complications, we employ throughout this paper the simplifying convention that all uncountable suprema (such as $\sup_{f \in \mathcal{M}_q^n} \ell_n(f)$) are interpreted as essential suprema with respect to the measure \mathbf{P}^* . In the majority of applications, the model classes \mathcal{M}_q^n will be separable, in which case the supremum and essential supremum coincide.

In the sequel, we will denote by $\|\cdot\|_p$ the $L^p(f^* d\mu)$ -norm, i.e., $\|g\|_p^p = \int |g(x)|^p f^*(x) \mu(dx)$, and we denote by

$\langle f, g \rangle = \int f(x)g(x)f^*(x)\mu(dx)$ the Hilbert space inner product in $L^2(f^*d\mu)$. Define the Hellinger distance

$$h(f, g)^2 = \int (\sqrt{f} - \sqrt{g})^2 d\mu, \quad f, g \in \mathcal{M}.$$

It is easily seen that $h(f, f^*) = \|\sqrt{f/f^*} - 1\|_2$. Finally, we will denote by $\mathcal{N}(\mathcal{Q}, \delta)$ for any class of functions \mathcal{Q} and $\delta > 0$ the minimal number of brackets of $L^2(f^*d\mu)$ -width δ needed to cover \mathcal{Q} : that is, $\mathcal{N}(\mathcal{Q}, \delta)$ is the smallest cardinality N of a collection of pairs of functions $\{g_i^L, g_i^U\}_{i=1,\dots,N}$ such that $\max_{i \leq N} \|g_i^U - g_i^L\|_2 \leq \delta$ and for every $g \in \mathcal{Q}$, we have $g_i^L \leq g \leq g_i^U$ pointwise for some $i \leq N$.

B. Upper Bound

We aim to obtain a pathwise upper bound on the likelihood ratio statistic that holds *uniformly* in $q > q^*$. To this end, define for $q, n \geq 1$ and $\varepsilon > 0$ the Hellinger ball

$$\mathcal{H}_q^n(\varepsilon) = \{\sqrt{f/f^*} : f \in \mathcal{M}_q^n, h(f, f^*) \leq \varepsilon\}.$$

Note that the definition of $\mathcal{H}_q^n(\varepsilon)$ depends on f^* (which is fixed throughout this paper). The following result shows that the geometry of the Hellinger balls $\mathcal{H}_q^n(\varepsilon)$ controls the pathwise fluctuations of the likelihood ratio statistic.

Theorem II.3: Suppose that for all n sufficiently large

$$\mathcal{N}(\mathcal{H}_q^n(\varepsilon), \delta) \leq \left(\frac{K(n)\varepsilon}{\delta} \right)^{\eta(q)}$$

for all $q \geq q^*$ and $\delta \leq \varepsilon$, where $K(n) \geq 1$ and $\eta(q) \geq q$ are increasing functions. Then

$$\begin{aligned} \overline{\lim}_{n \rightarrow \infty} \frac{1}{\log K(2n) \vee \log \log n} \times \\ \sup_{q \geq q^*} \frac{1}{\eta(q)} \left\{ \sup_{f \in \mathcal{M}_q^n} \ell_n(f) - \sup_{f \in \mathcal{M}_{q^*}^n} \ell_n(f) \right\} \leq C \end{aligned}$$

\mathbf{P}^* -a.s., where $C > 0$ is a universal constant.

The proof of Theorem II.3 is given in Appendix A.

The assumption of Theorem II.3 on the entropy of the Hellinger balls $\mathcal{H}_q^n(\varepsilon)$ states, roughly speaking, that the class of densities \mathcal{M}_q^n endowed with the Hellinger distance has the same metric structure as a Euclidean ball of dimension $\eta(q)$ and radius of order $K(n)$, at least locally in a neighborhood of the true density f^* . The effective dimension $\eta(q)$ controls the fluctuations of the likelihood ratio statistic as a function of the model order, while the effective radius $K(n)$ controls the fluctuations as a function of time up to a minimal rate of order $\log \log n$. In the following section, we will see that the minimal $\log \log n$ rate is indeed optimal.

Remark II.4: A bound on $\mathcal{N}(\mathcal{H}_q^n(\varepsilon), \delta)$ of the form required by Theorem II.3 is easily obtained if \mathcal{M}_q^n are regularly parametrized classes. That is, suppose that we can write

$$\mathcal{M}_q^n = \{f_\theta : \theta \in \Theta_q^n\}, \quad \Theta_q^n \subset \mathbb{R}^{\eta(q)}$$

where we have a pointwise Lipschitz estimate of the form

$$|\sqrt{f_\theta(x)/f^*(x)} - \sqrt{f_{\theta'}(x)/f^*(x)}| \leq F(x) \|\theta - \theta'\|$$

for some function F in L^2 and norm $\|\cdot\|$ on $\mathbb{R}^{\eta(q)}$, and

$$h(f_\theta, f^*) \geq c \|\theta - \theta^*\|$$

with $c > 0$. Then, the requisite bound on $\mathcal{N}(\mathcal{H}_q^n(\varepsilon), \delta)$ follows easily (cf., [21, Example 19.7]). This covers many cases of practical interest. However, geometrically complex models such as finite mixtures do not admit a regular parametrization, while our results are nonetheless sufficiently general to apply to such models (see Section III-B). In nonregular models, the entropy bound required by Theorem II.3 is far from obvious, and the requisite geometric analysis is of independent interest. Such problems are investigated by the authors in [22] and form the basis for the results in Section III-B.

C. Lower Bound

Throughout this section, we specialize to the case that $\mathcal{M}_q^n = \mathcal{M}_q$ does not depend on n (this implies essentially that \mathcal{M}_q is compact). In this setting, Theorem II.3 yields an upper bound of order $\log \log n$ on the pathwise fluctuations of the likelihood ratio statistic. The aim of this section is to obtain a matching lower bound of order $\log \log n$, which shows that the minimal rate in Theorem II.3 is essentially optimal. For the purposes of a lower bound, uniformity in q is irrelevant, so it suffices to restrict attention to some fixed $q > q^*$. We will in fact obtain a much stronger result in this case that completely characterizes the pathwise asymptotics of the likelihood ratio statistic for fixed q in sufficiently smooth families.

The geometric structure required in this section is somewhat different than that of Theorem II.3. Instead of Hellinger balls, we consider the classes of weighted densities $\mathcal{D}_q = \{d_f : f \in \mathcal{M}_q, f \neq f^*\}$ and $\mathcal{D} = \bigcup_q \mathcal{D}_q$, where

$$d_f = \frac{\sqrt{f/f^*} - 1}{h(f, f^*)}, \quad f \in \mathcal{M}, \quad f \neq f^*.$$

Define for $\varepsilon > 0$ and $q \geq 1$ the local weighted classes

$$\begin{aligned} \mathcal{D}_q(\varepsilon) &= \{d_f : f \in \mathcal{M}_q, 0 < h(f, f^*) \leq \varepsilon\} \\ \bar{\mathcal{D}}_q &= \bigcap_{\varepsilon > 0} \text{cl } \mathcal{D}_q(\varepsilon) \end{aligned}$$

where the closure $\text{cl } \mathcal{D}_q(\varepsilon)$ is in $L^2(f^*d\mu)$. Clearly, $\bar{\mathcal{D}}_q$ is the set of all possible limit points of d_f as $h(f, f^*) \rightarrow 0$ in \mathcal{M}_q . If the neighborhoods of $\bar{\mathcal{D}}_q$ are sufficiently rich, such limits can be taken along a continuous path in the following sense.

Definition II.5: A point $d \in \bar{\mathcal{D}}_q$ is called continuously accessible if there is a path $(f_t)_{t \in [0,1]} \subset \mathcal{M}_q \setminus \{f^*\}$ such that the map $t \mapsto h(f_t, f^*)$ is continuous, $h(f_t, f^*) \rightarrow 0$ as $t \rightarrow 0$, and $d_{f_t} \rightarrow d$ in $L^2(f^*d\mu)$ as $t \rightarrow 0$. The subset of all continuously accessible points in $\bar{\mathcal{D}}_q$ is denoted as $\bar{\mathcal{D}}_q^c$.

We can now formulate the main result of this section.

Theorem II.6: Let $q^* \leq p < q$. Assume that

$$\int_0^1 \sqrt{\log \mathcal{N}(\mathcal{D}_q, u)} du < \infty$$

and that $|d| \leq D$ for all $d \in \mathcal{D}_q$ with $D \in L^{2+\alpha}(f^* d\mu)$ for some $\alpha > 0$. Then

$$\overline{\lim}_{n \rightarrow \infty} \frac{1}{\log \log n} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \sup_{f \in \mathcal{M}_p} \ell_n(f) \right\} \geq \sup_{g \in L_0^2(f^* d\mu)} \left\{ \sup_{f \in \bar{\mathcal{D}}_q^c} (\langle f, g \rangle)_+^2 - \sup_{f \in \bar{\mathcal{D}}_p} (\langle f, g \rangle)_+^2 \right\} \quad \mathbf{P}^*\text{-a.s.}$$

as well as

$$\overline{\lim}_{n \rightarrow \infty} \frac{1}{\log \log n} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \sup_{f \in \mathcal{M}_p} \ell_n(f) \right\} \leq \sup_{g \in L_0^2(f^* d\mu)} \left\{ \sup_{f \in \bar{\mathcal{D}}_q} (\langle f, g \rangle)_+^2 - \sup_{f \in \bar{\mathcal{D}}_p^c} (\langle f, g \rangle)_+^2 \right\} \quad \mathbf{P}^*\text{-a.s.}$$

where $L_0^2(f^* d\mu) = \{g \in L^2(f^* d\mu) : \|g\|_2 \leq 1, \langle 1, g \rangle = 0\}$.

Only the first (lower bound) part of the theorem is needed to conclude optimality of the minimal $\log \log n$ rate in Theorem II.3. Indeed, we will obtain as a corollary the following lower bound counterpart to Theorem II.3.

Corollary II.7: Suppose there exists $q > q^*$ such that the following hold.

- 1) There is an envelope function $D : E \rightarrow \mathbb{R}$ such that $|d| \leq D$ for all $d \in \mathcal{D}_q$ and $D \in L^{2+\alpha}(f^* d\mu)$ for some $\alpha > 0$. Moreover, $\int_0^1 \sqrt{\log \mathcal{N}(\mathcal{D}_q, u)} du < \infty$.
- 2) $\bar{\mathcal{D}}_q^c \setminus \bar{\mathcal{D}}_{q^*}$ is nonempty.

Let $\eta(q) > 0$ be an arbitrary positive function. Then

$$\overline{\lim}_{n \rightarrow \infty} \frac{1}{\log \log n} \sup_{q \geq q^*} \frac{1}{\eta(q)} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \sup_{f \in \mathcal{M}_{q^*}} \ell_n(f) \right\} \geq C$$

$\mathbf{P}^*\text{-a.s.}$, where $C > 0$ is nonrandom but may depend on f^*, η .

The proofs of Theorem II.6 and Corollary II.7 are given in Appendix B.

The fact that the geometric assumptions in Theorem II.6 and Corollary II.7 are expressed in terms of weighted classes is not surprising, as the sharp asymptotic expression provided by Theorem II.6 for the pathwise fluctuations of the likelihood are expressed in terms of a variational problem on the weighted classes. Nonetheless, we are naturally led to ask whether there is any relation between the geometric assumptions imposed in the upper bound Theorem II.3 and the lower bound Theorem II.6, which appear to be quite different at first sight. In [22], we show that the global entropy of the weighted class is closely related to local entropy, so that the geometric assumptions for the upper and lower bounds are not too far apart.

Remark II.8: When $\bar{\mathcal{D}}_q$ and $\bar{\mathcal{D}}_p$ each contain an $L^2(f^* d\mu)$ -dense subset of continuously accessible points

(which is typically the case in sufficiently smooth models), Theorem II.6 provides the exact characterization

$$\overline{\lim}_{n \rightarrow \infty} \frac{1}{\log \log n} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \sup_{f \in \mathcal{M}_p} \ell_n(f) \right\} = \sup_{g \in L_0^2(f^* d\mu)} \left\{ \sup_{f \in \bar{\mathcal{D}}_q} (\langle f, g \rangle)_+^2 - \sup_{f \in \bar{\mathcal{D}}_p} (\langle f, g \rangle)_+^2 \right\} \quad \mathbf{P}^*\text{-a.s..}$$

Besides its intrinsic interest, this result has a surprising consequence. In the case that \mathcal{M}_q and \mathcal{M}_p are regular parametric models with $\dim(\mathcal{M}_q) > \dim(\mathcal{M}_p)$, one can choose $g \in \bar{\mathcal{D}}_q$ which is orthogonal to $\bar{\mathcal{D}}_p$. As $\bar{\mathcal{D}}_q, \bar{\mathcal{D}}_p \subseteq L_0^2(f^* d\mu)$ (see the proof of Corollary II.7), it follows easily that in this case, the right-hand side of the previous equation display is precisely equal to 1. In particular, we obtain the curious conclusion that in regular parametric models, the magnitude of the fluctuations of the likelihood ratio statistic does not depend on the dimensions $\dim(\mathcal{M}_q)$ and $\dim(\mathcal{M}_p)$. In contrast, it is well known that in regular parametric models, the likelihood ratio statistic itself converges weakly to a chi-square distribution with $\dim(\mathcal{M}_q) - \dim(\mathcal{M}_p)$ degrees of freedom, so the tails of the distribution of the likelihood ratio statistic do in fact depend strongly on the dimensions $\dim(\mathcal{M}_q)$ and $\dim(\mathcal{M}_p)$. Of course, the dimension independence of the pathwise fluctuations will also cease to hold if we are interested in a result that is uniform in the order q , as in Theorem II.3. This highlights the fact that the problems investigated in this paper are fundamentally different depending on whether or not one assumes a prior upper bound on the model order.

III. STRONGLY CONSISTENT ORDER ESTIMATION

The goal of this section is to apply the results of Section II to identify what penalties and cutoffs yield strongly consistent order estimators. We first develop some general consistency and inconsistency results and then consider specifically the challenging problem of mixture order estimation.

A. Consistency and Minimal Penalties

In this section, we consider the general setting introduced in Section II-A. We now suppose, however, that the true model order q^* (as well as the true density f^*) is not known, so that we must estimate q^* from an observation sequence $(X_k)_{k \geq 1}$. To this end, define the penalized likelihood order estimator

$$\hat{q}_n = \operatorname{argmax}_{q \geq 1} \left\{ \sup_{f \in \mathcal{M}_q^n} \ell_n(f) - \operatorname{pen}(n, q) \right\}$$

where $\operatorname{pen}(n, q)$ is a penalty function. Our goal is to show that the penalized likelihood order estimator is strongly consistent, i.e., $\hat{q}_n \rightarrow q^*$ as $n \rightarrow \infty$ $\mathbf{P}^*\text{-a.s.}$, for a suitable choice of the penalty (that does not depend on q^* or f^*). Let us emphasize that the maximum in the definition of \hat{q}_n is taken over *all* model orders $q \geq 1$, that is, we do not assume that an *a priori* upper bound on the order is available, in contrast to most previous work on this topic.

We obtain the following general result.

Theorem III.1: Suppose that for all n sufficiently large

$$\mathcal{N}(\mathcal{H}_q^n(\varepsilon), \delta) \leq \left(\frac{K(n)\varepsilon}{\delta} \right)^{\eta(q)}$$

for all $q \geq q^*$ and $\delta \leq \varepsilon$, where $K(n) \geq 1$ and $\eta(q) \geq q$ are increasing functions and we assume that $\log K(n) = o(n)$. Let $\text{pen}(n, q)$ be a penalty that is increasing in q and

$$\begin{aligned} \lim_{n \rightarrow \infty} \sup_{q > q^*} \frac{\eta(q)\{\log K(2n) \vee \log \log n\}}{\text{pen}(n, q) - \text{pen}(n, q^*)} &= 0 \\ \lim_{n \rightarrow \infty} \max_{q < q^*} \frac{\text{pen}(n, q)}{n} &= 0. \end{aligned}$$

Then, $\hat{q}_n \rightarrow q^*$ as $n \rightarrow \infty$ \mathbf{P}^* -a.s.

Theorem III.1 is proved in Appendix C.

Let us now specialize to the case that $\mathcal{M}_q^n = \mathcal{M}_q$ does not depend on n , as in Section II-C. In this case, Theorem III.1 immediately yields the following corollary.

Corollary III.2: Suppose that for all $q \geq q^*$ and $\delta \leq \varepsilon$

$$\mathcal{N}(\mathcal{H}_q(\varepsilon), \delta) \leq \left(\frac{K\varepsilon}{\delta} \right)^{\eta(q)}$$

where $K \geq 1$ and $\eta(q) \geq q$ is a strictly increasing function. Define the penalty

$$\text{pen}(n, q) = \eta(q)\varpi(n)$$

where $\varpi(n)$ is any function such that

$$\lim_{n \rightarrow \infty} \frac{\log \log n}{\varpi(n)} = 0, \quad \lim_{n \rightarrow \infty} \frac{\varpi(n)}{n} = 0.$$

Then, $\hat{q}_n \rightarrow q^*$ as $n \rightarrow \infty$ \mathbf{P}^* -a.s.

Corollary III.2 states that, when $\mathcal{M}_q^n = \mathcal{M}_q$ does not depend on n , the penalized likelihood order estimator is strongly consistent provided the penalty grows faster than $\log \log n$ and slower than n . Clearly, the $\log \log n$ rate is the minimal one attainable by applying Theorem III.1. This raises the question whether the $\log \log n$ rate is indeed minimal, in the sense that smaller penalties yield inconsistent estimators. The following result shows that this is indeed the case, so that the result of Corollary III.2 is essentially optimal.

Corollary III.3: Suppose there exists $q > q^*$ such that

- 1) there is an envelope function $D : E \rightarrow \mathbb{R}$ such that $|d| \leq D$ for all $d \in \mathcal{D}_q$, $D \in L^{2+\alpha}(f^* d\mu)$ for some $\alpha > 0$, and $\int_0^1 \sqrt{\log \mathcal{N}(\mathcal{D}_q, u)} du < \infty$;
- 2) $\overline{\mathcal{D}}_q^c \setminus \overline{\mathcal{D}}_{q^*}$ is nonempty.

Let $\eta(q) > 0$ be any strictly increasing function, and let

$$\text{pen}(n, q) = C \eta(q) \log \log n.$$

If $C > 0$ is sufficiently small, $\hat{q}_n \neq q^*$ infinitely often \mathbf{P}^* -a.s.

The proof of Corollary III.3 is given in Appendix C. Let us note that the proof of Corollary III.3 actually shows that $\sup_{f \in \mathcal{M}_q} \ell_n(f) - \text{pen}(n, q) > \sup_{f \in \mathcal{M}_{q^*}} \ell_n(f) - \text{pen}(n, q^*)$ infinitely often \mathbf{P}^* -a.s., so Corollary III.3 is not altered even if we were to impose a prior upper bound on the order.

In conclusion, we have shown that when $\mathcal{M}_q^n = \mathcal{M}_q$ does not depend on n , penalties growing faster than $\log \log n$ are consistent, while the penalty $C \eta(q) \log \log n$ is inconsistent when the constant C is sufficiently small. From the proof of Theorem III.1, we can also see that the penalty $C \eta(q) \log \log n$ is consistent when C is sufficiently large. However, the critical value of C may depend on the unknown parameter f^* , so that this *minimal* penalty may not be implementable. On the other hand, assuming that $\eta(q)$ does not depend on f^* (as is typically the case), penalties satisfying the assumptions of Theorem III.1 obviously do not depend on the unknown parameter f^* and therefore define admissible estimators. When \mathcal{M}_q^n depends on n , larger penalties may be required to ensure consistency, depending on the growth rate of $K(n)$.

B. Location Mixture Order Estimation

We finally apply our general results to location mixture order estimation. Throughout this section, let $E = \mathbb{R}^d$ and let μ be the Lebesgue measure on \mathbb{R}^d . Fix a strictly positive probability density f_0 with respect to μ , and define

$$\mathcal{M}_q^n = \left\{ \sum_{i=1}^q \pi_i f_{\theta_i} : \pi_i \geq 0, \sum_{i=1}^q \pi_i = 1, \theta_i \in \Theta(n) \right\}$$

where $f_\theta(x) = f_0(x - \theta)$ and $\dots \subseteq \Theta(n) \subseteq \Theta(n+1) \subseteq \dots \subset \mathbb{R}^d$ is an increasing family of bounded subsets of \mathbb{R}^d . We fix $f^* \in \mathcal{M}$ throughout this section. Let

$$\begin{aligned} H_0(x) &= \sup_{\theta \in \Theta} f_\theta(x)/f^*(x) \\ H_1(x) &= \sup_{\theta \in \Theta} \max_{i=1,\dots,d} |\partial f_\theta(x)/\partial \theta^i|/f^*(x) \\ H_2(x) &= \sup_{\theta \in \Theta} \max_{i,j=1,\dots,d} |\partial^2 f_\theta(x)/\partial \theta^i \partial \theta^j|/f^*(x) \\ H_3(x) &= \sup_{\theta \in \Theta} \max_{i,j,k=1,\dots,d} |\partial^3 f_\theta(x)/\partial \theta^i \partial \theta^j \partial \theta^k|/f^*(x) \end{aligned}$$

when f_0 is sufficiently differentiable, and let

Assumption A: The following hold:

- 1) $f_0 \in C^3$ and $f_0(x), (\partial f_0/\partial \theta^i)(x)$ vanish as $\|x\| \rightarrow \infty$.
- 2) $H_k \in L^4(f^* d\mu)$ for $k = 0, 1, 2$ and $H_3 \in L^2(f^* d\mu)$.

In the following, we consider two separate cases. The first case is that of a compact parameter set, where $\Theta(n) = \Theta$ does not depend on n . In this setting, we obtain a general result. Then, we consider the noncompact case in the setting of Gaussian mixtures and illustrate how Theorem III.1 can be used to obtain consistency results in this case. To be able to use Theorem III.1, we need suitable estimates on the local entropy of mixtures. The following result is given in [22].

Theorem III.4: Suppose that Assumption A holds. Then, if $\Theta(n) = \Theta$ is a bounded subset of \mathbb{R}^d with diameter $2T$

$$\mathcal{N}(\mathcal{H}_q^n(\varepsilon), \delta) \leq \left(\frac{C_\Theta \varepsilon}{\delta} \right)^{10(d+1)q+1}$$

for all $q \geq q^*$ and $\delta/\varepsilon \leq 1$, where

$$C_\Theta = L^*(T \vee 1)^{1/3} (\|H_0\|_4^4 \vee \|H_1\|_4^4 \vee \|H_2\|_4^4 \vee \|H_3\|_2^2)^{5/4}$$

and L^* is a constant that depends only on d, q^* , and f^* .

Example III.5 (Gaussian Mixtures): Consider mixtures of Gaussian densities $f_0(x) = (2\pi)^{-d/2} |\Sigma|^{-1/2} e^{-x^T \Sigma^{-1} x/2}$, where Σ is a positive-definite $d \times d$ matrix, $|\Sigma|$ denotes its determinant, and x^T the transpose vector of x , and let $\Theta(T) = \{\theta \in \mathbb{R}^d : \|\theta\| \leq T\}$. Fix a nondegenerate mixture f^* , and define $T^* = \max_{i=1,\dots,q^*} \|\theta_i^*\|$. Denote by $\mathcal{H}_q(\varepsilon, T)$ the Hellinger ball associated with the parameter set $\Theta(T)$. Then

$$\mathcal{N}(\mathcal{H}_q(\varepsilon, T), \delta) \leq \left(\frac{C_1^* e^{C_2^* T^2} \varepsilon}{\delta} \right)^{10(d+1)q+1}$$

for all $q \geq q^*$, $T \geq T^*$, and $\delta/\varepsilon \leq 1$, where C_1^* and C_2^* are constants that depend on d , q^* , and f^* only. To prove this, it suffices to show that Assumption A holds and that $\|H_k\|_4$ for $k = 0, 1, 2$ and $\|H_3\|_2$ are of order e^{CT^2} . These facts are readily verified by a straightforward computation.

Let us first consider the case of a compact parameter set. We obtain a general consistency result under Assumption A.

Proposition III.6: Suppose that the parameter set $\Theta(n) = \Theta$ is a bounded subset of \mathbb{R}^d independent of n , and that Assumption A holds. If we choose a penalty of the form

$$\text{pen}(n, q) = q \omega(n), \quad \lim_{n \rightarrow \infty} \frac{\log \log n}{\omega(n)} = \lim_{n \rightarrow \infty} \frac{\omega(n)}{n} = 0,$$

then $\hat{q}_n \rightarrow q^*$ as $n \rightarrow \infty$ \mathbf{P}^* -a.s. On the other hand, if

$$\text{pen}(n, q) = C q \log \log n$$

where $C > 0$ is a sufficiently small constant, then we have $\hat{q}_n \neq q^*$ infinitely often \mathbf{P}^* -a.s.

We, therefore, find that in the setting of location mixtures with a compact parameter set, the minimal penalty is of order $\log \log n$. Moreover, the popular BIC penalty

$$\text{pen}(n, q) = \frac{dq + q - 1}{2} \log n \quad (\text{III.1})$$

yields a strongly consistent mixture order estimator in this setting, without a prior upper bound on the order. The requisite Assumption A is mild, which highlights the broad applicability of this result. However, the assumption of a compact parameter space can be quite restrictive in practice.

Let us, therefore, consider a case where the parameter space is noncompact. For simplicity, we restrict our attention to Gaussian mixtures, that is, we choose $f_0(x) = (2\pi)^{-d/2} e^{-\|x\|^2/2}$, and we choose the restricted parameter sets $\Theta(n) = \{\theta \in \mathbb{R}^d : \|\theta\| \leq T(n)\}$ for some sequence $T(n) \uparrow \infty$. Our aim is to choose the penalty $\text{pen}(n, q)$ and cutoff $T(n)$ so that the penalized likelihood order estimator is strongly consistent.

In this setting, we obtain the following result.

Proposition III.7: For the case $f_0(x) = (2\pi)^{-d/2} |\Sigma|^{-1/2} e^{-x^T \Sigma^{-1} x/2}$ and $\Theta(n) = \{\theta \in \mathbb{R}^d : \|\theta\| \leq T(n)\}$, consider a penalty of the form $\text{pen}(n, q) = q \omega(n)$. If

$$\lim_{n \rightarrow \infty} \frac{\log \log n}{\omega(n)} = \lim_{n \rightarrow \infty} \frac{\omega(n)}{n} = 0, \quad T(n) = O(\sqrt{\log \log n})$$

then $\hat{q}_n \rightarrow q^*$ as $n \rightarrow \infty$ \mathbf{P}^* -a.s.

On the other hand, the BIC penalty (III.1) yields a strongly consistent order estimator if $T(n) = o(\sqrt{\log n})$.

This result illustrates that our theory can establish consistency of the penalized likelihood mixture order estimator without any prior upper bounds on the model order or the magnitude of the

true parameters. Let us note that there is nothing particularly special about the Gaussian case: a similar result can be obtained, in principle, for any mixture distribution f_0 , as long as one can obtain suitable estimates on the quantities $\|H_i\|_4$ that appear in Theorem III.4 (see Example III.5 for the Gaussian case).

The proofs of Propositions III.6 and III.7 are given in Appendix D.

C. Remarks

The above results provide a fairly sharp theoretical understanding of the consistency properties of penalized likelihood order estimators. We have seen that penalties of order $\log \log n$ are minimal and that penalties of order n are maximal (the latter is easily seen in the proof of Theorem III.1) without a prior upper bound on the order. This rather broad range leaves significant freedom in the choice of a penalty in practice.

There can be no optimal penalty for all situations. The choice of penalty allows us to tradeoff between underestimation and overestimation of the order: a larger penalty increases the probability of underestimation but decreases the probability of overestimation. Which effect is more undesirable will depend on the application at hand. Moreover, we have only addressed the consistency properties of order estimators: other considerations, such as finite-sample performance, may be of significant importance. Finite-sample bounds that appear in our proofs (as in the proof of Theorem II.3) suggest that the rate of convergence of $\log \log n$ penalties can be very slow.

Many practitioners use the BIC penalty, whose performance is supported by numerous empirical studies in the literature. The aim of this paper is not to suggest an alternative methodology, but to explore the fundamental limitations to a broad class of penalized likelihood order estimators in the absence of a prior upper bound on the order. When the parameter space is not compact, our results can also provide some guidance to the practical implementation of order estimators such as BIC.

APPENDIX A PROOF OF THEOREM II.3

The proof of Theorem II.3 is based on the following deviation bound for the log-likelihood ratio. This bound is essentially from [23, Corollary 7.5], but the additional maximum inside the probability is essential for our purposes.

Theorem A.1: Let \mathcal{M} be a family of strictly positive probability densities with respect to a reference measure μ , fix some $f^* \in \mathcal{M}$, and define the Hellinger ball $\mathcal{H}(\varepsilon) = \{\sqrt{f/f^*} : f \in \mathcal{M}, h(f, f^*) \leq \varepsilon\}$ where $h(f, g)^2 = \int (\sqrt{f} - \sqrt{g})^2 d\mu$. Suppose that for some constants $K \geq 1$, $p \geq 1$ and all $\delta \leq \varepsilon$

$$\mathcal{N}(\mathcal{H}(\varepsilon), \delta) \leq \left(\frac{K\varepsilon}{\delta} \right)^p$$

where $\mathcal{N}(\mathcal{H}(\varepsilon), \delta)$ is the minimal number of brackets of $L^2(f^* d\mu)$ -width δ needed to cover $\mathcal{H}(\varepsilon)$. Let $(X_i)_{i \in \mathbb{N}}$ be i.i.d. with distribution $f^* d\mu$. Then

$$\mathbf{P} \left[\max_{n \leq k \leq 2n} \sup_{f \in \mathcal{M}} \sum_{j=1}^k \log \left(\frac{f(X_j)}{f^*(X_j)} \right) \geq \alpha \right] \leq C e^{-\alpha/C}$$

for all $\alpha \geq Cp(1 + \log K)$, $n \geq 1$ [C is a universal constant].

Proof: Define $\bar{f} = (f + f^*)/2$ for any $f \in \mathcal{M}$, and define the empirical process $\nu_n(g) = n^{-1/2} \sum_{k=1}^n \{g(X_k) - \mathbb{E}[g(X_k)]\}$. Using concavity of $\log x$, we have

$$\sum_{j=1}^k \log \left(\frac{f(X_j)}{f^*(X_j)} \right) \leq 2k^{1/2} \nu_k(\log(\bar{f}/f^*)) - 2kD(f^*\|\bar{f})$$

where $D(f^*\|f) = \int \log(f^*/f) f^* d\mu$ is relative entropy. As $D(f^*\|f) \geq h(f, f^*)^2$, we can estimate

$$\begin{aligned} & \mathbf{P} \left[\max_{n \leq k \leq 2n} \sup_{f \in \mathcal{M}} \sum_{j=1}^k \log \left(\frac{f(X_j)}{f^*(X_j)} \right) \geq \alpha \right] \\ & \leq \mathbf{P} \left[\max_{n \leq k \leq 2n} \sup_{f \in \mathcal{M}} \{\sqrt{k} \nu_k(\log(\bar{f}/f^*)) - kh(\bar{f}, f^*)^2\} \geq \frac{\alpha}{2} \right] \\ & \leq \sum_{s=0}^S \mathbf{P} \left[\max_{n \leq k \leq 2n} \sup_{f \in \mathcal{M}^s} |\sqrt{k} \nu_k(\log(\bar{f}/f^*))| \geq \alpha 2^{s-1} \right] \\ & \leq 3 \sum_{s=0}^S \max_{n \leq k \leq 2n} \mathbf{P} \left[\sup_{f \in \mathcal{M}^s} |\nu_k(\log(\sqrt{\bar{f}/f^*}))| \geq \alpha 2^{s-5}/\sqrt{n} \right] \end{aligned}$$

where $\mathcal{M}^s = \{f \in \mathcal{M} : \alpha 2^{s-1} < nh(\bar{f}, f^*)^2 \leq \alpha 2^s\}$, $1 \leq s \leq S$, $S = \min\{s : \alpha 2^s n^{-1} > 2\}$, $\mathcal{M}^0 = \{f \in \mathcal{M} : nh(\bar{f}, f^*)^2 \leq \alpha\}$ and we have used Lemma A.2 below for the last inequality. The remainder of the proof is identical to that of [23, Th. 7.4], provided we show that for $\bar{\mathcal{H}}(\varepsilon) = \{\sqrt{\bar{f}/f^*} : f \in \mathcal{M}, h(\bar{f}, f^*) \leq \varepsilon\}$

$$\mathcal{N}(\bar{\mathcal{H}}(\varepsilon), \delta) \leq \left(\frac{2\sqrt{2}K\varepsilon}{\delta} \right)^p.$$

To this end, fix $\delta \leq \varepsilon$, and note that $h(f, f^*) \leq 4h(\bar{f}, f^*)$ by [23, Lemma 4.2], so that $\{f \in \mathcal{M} : h(\bar{f}, f^*) \leq \varepsilon\} \subseteq \{f \in \mathcal{M} : h(f, f^*) \leq 4\varepsilon\}$. By assumption, there exist $N \leq (2\sqrt{2}K\varepsilon/\delta)^p$ and functions $g_1, \dots, g_N, h_1, \dots, h_N$ such that $\|h_i - g_i\|_2 \leq \delta\sqrt{2}$ for every i , and for every $u \in \mathcal{H}(4\varepsilon)$ there is an i such that $g_i \leq u \leq h_i$. But for every $f \in \mathcal{M}$ such that $h(\bar{f}, f^*) \leq \varepsilon$, we then have for some i

$$2^{-1/2} \sqrt{g_i^2 + 1} \leq \sqrt{\bar{f}/f^*} \leq 2^{-1/2} \sqrt{h_i^2 + 1}.$$

Using $|\sqrt{a+c} - \sqrt{b+c}| \leq |\sqrt{a} - \sqrt{b}|$ for $a, b, c \geq 0$, we have

$$\left\| 2^{-1/2} \sqrt{h_i^2 + 1} - 2^{-1/2} \sqrt{g_i^2 + 1} \right\|_2 \leq 2^{-1/2} \|h_i - g_i\|_2 \leq \delta.$$

The result now follows directly. \blacksquare

The following variant of Etemadi's inequality was used in the proof. The proof follows closely that of the classical Etemadi inequality (see [24, Appendix M19]).

Lemma A.2: Let \mathcal{Q} be a family of measurable functions $f : E \rightarrow \mathbb{R}$. Then, we have for every $\alpha > 0$ and $m, n \in \mathbb{N}$, $m \leq n$

$$\begin{aligned} & \mathbf{P}^* \left[\max_{k=m, \dots, n} \sup_{f \in \mathcal{Q}} |S_k(f)| \geq 3\alpha \right] \\ & \leq 3 \max_{k=m, \dots, n} \mathbf{P}^* \left[\sup_{f \in \mathcal{Q}} |S_k(f)| \geq \alpha \right] \end{aligned}$$

where $S_n(f) = n^{1/2} \nu_n(f)$.

Proof: Define the stopping time

$$\tau = \inf \left\{ k \geq m : \sup_{f \in \mathcal{Q}} |S_k(f)| \geq 3\alpha \right\}.$$

Then

$$\begin{aligned} & \mathbf{P}^* \left[\max_{k=m, \dots, n} \sup_{f \in \mathcal{Q}} |S_k(f)| \geq 3\alpha \right] = \mathbf{P}^*[\tau \leq n] \\ & \leq \mathbf{P}^* \left[\sup_{f \in \mathcal{Q}} |S_n(f)| \geq \alpha \right] + \\ & \quad \sum_{k=m}^n \mathbf{P}^* \left[\tau = k \text{ and } \sup_{f \in \mathcal{Q}} |S_n(f)| < \alpha \right]. \end{aligned}$$

But on the event $\{\tau = k \text{ and } \sup_{f \in \mathcal{Q}} |S_n(f)| < \alpha\}$, we have

$$2\alpha \leq \sup_{f \in \mathcal{Q}} |S_k(f)| - \sup_{f \in \mathcal{Q}} |S_n(f)| \leq \sup_{f \in \mathcal{Q}} |S_k(f) - S_n(f)|.$$

Therefore, we can estimate

$$\begin{aligned} & \mathbf{P}^* \left[\max_{k=m, \dots, n} \sup_{f \in \mathcal{Q}} |S_k(f)| \geq 3\alpha \right] \\ & \leq \mathbf{P}^* \left[\sup_{f \in \mathcal{Q}} |S_n(f)| \geq \alpha \right] + \\ & \quad \sum_{k=m}^n \mathbf{P}^* \left[\tau = k \text{ and } \sup_{f \in \mathcal{Q}} |S_n(f) - S_k(f)| \geq 2\alpha \right] \\ & \leq \mathbf{P}^* \left[\sup_{f \in \mathcal{Q}} |S_n(f)| \geq \alpha \right] + \\ & \quad \max_{k=m, \dots, n} \mathbf{P}^* \left[\sup_{f \in \mathcal{Q}} |S_n(f) - S_k(f)| \geq 2\alpha \right] \end{aligned}$$

where we have used that $\sup_{f \in \mathcal{Q}} |S_n(f) - S_k(f)|$ and $\{\tau = k\}$ are independent to obtain the last inequality. The remainder of the proof is now easily completed. \blacksquare

We can now complete the proof of Theorem II.3.

Proof of Theorem II.3: By assumption, we have $f^* \in \mathcal{M}_q^n$ for all $q \geq q^*$ when n is sufficiently large. Then, by Theorem A.1, we have for n sufficiently large

$$\mathbf{P}^* \left[\max_{n \leq k \leq 2n} \sup_{f \in \mathcal{M}_q^{2n}} \{\ell_k(f) - \ell_k(f^*)\} \geq \alpha \right] \leq C e^{-\alpha/C}$$

for all $\alpha \geq C\eta(q)(1 + \log K(2n))$ and $q \geq q^*$. Define

$$\Delta_k(q, q^*) = \sup_{f \in \mathcal{M}_q^k} \ell_k(f) - \sup_{f \in \mathcal{M}_{q^*}^k} \ell_k(f).$$

Using that $\mathcal{M}_q^k \subseteq \mathcal{M}_q^{2n}$ for $n \leq k \leq 2n$ and $\ell_k(f^*) \leq \sup_{f \in \mathcal{M}_{q^*}^k} \ell_k(f)$, we have for n sufficiently large

$$\mathbf{P}^* \left[\max_{n \leq k \leq 2n} \sup_{q \geq q^*} \frac{1}{\eta(q)} \Delta_k(q, q^*) \geq \alpha \right] \leq \sum_{q=q^*}^{\infty} C e^{-\alpha\eta(q)/C}$$

for all $\alpha \geq C(1 + \log K(2n))$. Let $\beta(n)$ be an increasing function. Then, for all n sufficiently large

$$\mathbf{P}^* \left[\max_{2^n \leq k \leq 2^{n+1}} \frac{1}{\beta(k)} \sup_{q \geq q^*} \frac{1}{\eta(q)} \Delta_k(q, q^*) \geq 2C \right] \leq \frac{2C}{n^2}$$

provided that $\beta(2^n) \geq \log K(2^{n+1}) \vee \log \log 2^n$. The proof is now easily completed using the Borel–Cantelli lemma. ■

APPENDIX B PROOF OF THEOREM II.6

The proof of Theorem II.6 is based on a sequence of auxiliary results. First, we will need a compact law of iterated logarithm for the Strassen functional

$$I_n(g) = \frac{1}{\sqrt{2n \log \log n}} \sum_{i=1}^n \{g(X_i) - \mathbf{E}^*(g(X_1))\}.$$

We state the requisite result for future reference.

Theorem B.1: Let \mathcal{Q} be a family of measurable functions from E to \mathbb{R} such that

$$\int_0^1 \sqrt{\log \mathcal{N}(\mathcal{Q}, u)} du < \infty.$$

Then, \mathbf{P}^* -a.s., the sequence $(I_n)_{n \geq 0}$ is relatively compact in $\ell_\infty(\mathcal{Q})$, and its set of cluster points coincides precisely with the set $\mathcal{K} = \{f \mapsto \langle f, g \rangle : g \in L_0^2(f^* d\mu)\}$.

Proofs of this result can be found in [25, Th. 4.2] or in [26, Th. 9]. We will also need the following simple well-known fact, whose proof is omitted.

Lemma B.2: Let $(X_i)_{i \geq 1}$ be an i.i.d. sequence such that $\mathbf{E}[|X_1|^p] < \infty$. Then, $n^{-1/p} \max_{i=1, \dots, n} |X_i| \rightarrow 0$ a.s.

Finally, we will need the following likelihood inequality that relates the log-likelihood ratio $\ell_n(f) - \ell_n(f^*)$ to the empirical process. Related inequalities appear in [6], [27], and [28], but the following form is perhaps the most natural.

Lemma B.3: For any probability density $f \neq f^*$

$$\ell_n(f) - \ell_n(f^*) \leq |\nu_n(d_f)|^2$$

where $\nu_n(g) = n^{-1/2} \sum_{k=1}^n \{g(X_k) - \mathbf{E}^*[g(X_k)]\}$.

Proof: Note that

$$h(f, f^*)^2 = 2 - \int 2\sqrt{ff^*} d\mu = -2 h(f, f^*) \mathbf{E}^*(d_f(X_1)).$$

Using $\log(1+x) \leq x$, we can estimate

$$\begin{aligned} \ell_n(f) - \ell_n(f^*) &= \sum_{i=1}^n 2 \log(1 + h(f, f^*) d_f(X_i)) \\ &\leq \sum_{i=1}^n 2 h(f, f^*) d_f(X_i) \\ &= 2 \nu_n(d_f) h(f, f^*) \sqrt{n} - h(f, f^*)^2 n \\ &\leq \sup_{p \in \mathbb{R}} \{2 \nu_n(d_f) p - p^2\}. \end{aligned}$$

The proof is easily completed. ■

We can now obtain the following asymptotic expansion of the log-likelihood, which provides a pathwise counterpart to the weak convergence theory in [27] and [28].

Proposition B.4: Let $q \geq q^*$. Assume that

$$\int_0^1 \sqrt{\log \mathcal{N}(\mathcal{D}_q, u)} du < \infty.$$

Moreover, suppose that $|d| \leq D$ for all $d \in \mathcal{D}_q$ with $D \in L^{2+\alpha}(f^* d\mu)$ for some $\alpha > 0$. Then

$$\begin{aligned} \sup_{f \in \mathcal{M}_q(4\sqrt{\log \log n/n})} &\left\{ 2 I_n(d_f) h(f, f^*) \sqrt{\frac{2n}{\log \log n}} \right. \\ &\left. - h(f, f^*)^2 \frac{2n}{\log \log n} \right\} \\ &- \frac{1}{\log \log n} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \ell_n(f^*) \right\} \xrightarrow{n \rightarrow \infty} 0 \quad \mathbf{P}^*\text{-a.s.} \end{aligned}$$

where we have defined $\mathcal{M}_q(\varepsilon) = \{f \in \mathcal{M}_q : h(f, f^*) \leq \varepsilon\}$.

Proof: We proceed in several steps.

Step 1 (Localization): As $q \geq q^*$ (hence $f^* \in \mathcal{M}_q$), clearly

$$\sup_{f \in \mathcal{M}_q} \ell_n(f) - \ell_n(f^*) = \sup_{f \in \mathcal{M}_q : \ell_n(f) - \ell_n(f^*) \geq 0} \{\ell_n(f) - \ell_n(f^*)\}.$$

Now note that, as in the proof of Lemma B.3,

$$\ell_n(f) - \ell_n(f^*) \leq 2 \nu_n(d_f) h(f, f^*) \sqrt{n} - h(f, f^*)^2 n.$$

Therefore, we can estimate

$$\begin{aligned} \sup_{f \in \mathcal{M}_q : \ell_n(f) - \ell_n(f^*) \geq 0} &h(f, f^*) \\ &\leq \sup_{f \in \mathcal{M}_q : \ell_n(f) - \ell_n(f^*) \geq 0} \left\{ h(f, f^*) + \frac{\ell_n(f) - \ell_n(f^*)}{n h(f, f^*)} \right\} \\ &\leq \frac{2}{\sqrt{n}} \sup_{f \in \mathcal{M}_q : \ell_n(f) - \ell_n(f^*) \geq 0} \nu_n(d_f) \\ &\leq \sqrt{\frac{8 \log \log n}{n}} \sup_{d \in \mathcal{D}_q} I_n(d). \end{aligned}$$

Now note that we can estimate

$$\begin{aligned} \sup_{d \in \mathcal{D}_q} I_n(d) &\leq \inf_{g \in L_0^2(f^* d\mu)} \sup_{d \in \mathcal{D}_q} |I_n(d) - \langle d, g \rangle| \\ &\quad + \sup_{d \in \mathcal{D}_q} \sup_{g \in L_0^2(f^* d\mu)} \langle d, g \rangle. \end{aligned}$$

The first term on the right converges to zero \mathbf{P}^* -a.s. as $n \rightarrow \infty$ by Theorem B.1, while the second term is easily seen to equal $\sup_{d \in \mathcal{D}_q} \|d - \langle 1, d \rangle\|_2 \leq 1$. Therefore

$$\sup_{f \in \mathcal{M}_q : \ell_n(f) - \ell_n(f^*) \geq 0} h(f, f^*) \leq (1 + \varepsilon) \sqrt{\frac{8 \log \log n}{n}}$$

eventually as $n \rightarrow \infty$ \mathbf{P}^* -a.s. for any $\varepsilon > 0$. In particular

$$\{f \in \mathcal{M}_q : \ell_n(f) - \ell_n(f^*) \geq 0\} \subseteq$$

$$\left\{ f \in \mathcal{M}_q : h(f, f^*) \leq 4\sqrt{\log \log n/n} \right\}$$

eventually as $n \rightarrow \infty$ \mathbf{P}^* -a.s. This implies that

$$\begin{aligned} & \sup_{f \in \mathcal{M}_q} \ell_n(f) - \ell_n(f^*) \\ & \leq \sup_{f \in \mathcal{M}_q : h(f, f^*) \leq 4\sqrt{\log \log n/n}} \{\ell_n(f) - \ell_n(f^*)\} \end{aligned}$$

eventually as $n \rightarrow \infty$ \mathbf{P}^* -a.s. But the reverse inequality clearly holds for all $n \geq 0$, so that in fact

$$\sup_{f \in \mathcal{M}_q} \ell_n(f) - \ell_n(f^*) = \sup_{f \in \mathcal{M}_q (4\sqrt{\log \log n/n})} \{\ell_n(f) - \ell_n(f^*)\}$$

eventually as $n \rightarrow \infty$ \mathbf{P}^* -a.s.

Step 2 (Taylor Expansion): Taylor expansion gives $2\log(1+x) = 2x - x^2 + x^2 R(x)$, where $R(x) \rightarrow 0$ as $x \rightarrow 0$. Thus, we can write, for any $f \in \mathcal{M}_q$

$$\begin{aligned} \ell_n(f) - \ell_n(f^*) &= \sum_{i=1}^n 2 \log(1 + h(f, f^*) d_f(X_i)) = \\ &= 2 h(f, f^*) \sum_{i=1}^n \left\{ d_f(X_i) + \frac{1}{2} h(f, f^*) \right\} \\ &\quad - h(f, f^*)^2 \sum_{i=1}^n (d_f(X_i))^2 - n h(f, f^*)^2 \\ &\quad + h(f, f^*)^2 \sum_{i=1}^n (d_f(X_i))^2 R(h(f, f^*) d_f(X_i)). \end{aligned}$$

Using that $\mathbf{E}^*(d_f(X_1)) = -h(f, f^*)/2$, we therefore have

$$\begin{aligned} \frac{1}{\log \log n} \{\ell_n(f) - \ell_n(f^*)\} &= R_{f,n} \frac{n h(f, f^*)^2}{\log \log n} \\ &\quad + 2 I_n(d_f) h(f, f^*) \sqrt{\frac{2n}{\log \log n}} - h(f, f^*)^2 \frac{2n}{\log \log n} \end{aligned}$$

where we have defined

$$\begin{aligned} R_{f,n} &= \frac{1}{n} \sum_{i=1}^n \{1 - (d_f(X_i))^2\} \\ &\quad + \frac{1}{n} \sum_{i=1}^n (d_f(X_i))^2 R(h(f, f^*) d_f(X_i)). \end{aligned}$$

It follows easily that

$$\begin{aligned} & \left| \sup_{f \in \mathcal{M}_q (4\sqrt{\log \log n/n})} \left\{ 2 I_n(d_f) h(f, f^*) \sqrt{\frac{2n}{\log \log n}} \right. \right. \\ & \quad \left. \left. - h(f, f^*)^2 \frac{2n}{\log \log n} \right\} \right| \\ & \quad - \frac{1}{\log \log n} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \ell_n(f^*) \right\} \Bigg| \\ & \leq \sup_{f \in \mathcal{M}_q (4\sqrt{\log \log n/n})} |R_{f,n}| \frac{n h(f, f^*)^2}{\log \log n} \\ & \leq 16 \sup_{f \in \mathcal{M}_q (4\sqrt{\log \log n/n})} |R_{f,n}| \end{aligned}$$

eventually as $n \rightarrow \infty$ \mathbf{P}^* -a.s.

Step 3 (End of Proof): We can easily estimate

$$\begin{aligned} & \sup_{f \in \mathcal{M}_q (4\sqrt{\log \log n/n})} |R_{f,n}| \leq \sup_{f \in \mathcal{M}_q} \left| \frac{1}{n} \sum_{i=1}^n \{d_f(X_i)^2 - 1\} \right| \\ & \quad + \left(\sup_{|x| \leq 4\sqrt{\log \log n/n} \max_{i \leq n} D(X_i)} |R(x)| \right) \frac{1}{n} \sum_{i=1}^n (D(X_i))^2. \end{aligned}$$

As $\mathcal{N}(\mathcal{D}_q, \delta) < \infty$ for every $\delta > 0$, the class $\{d^2 : d \in \mathcal{D}_q\}$ can be covered by a finite number of brackets with arbitrary small $L^1(f^* d\mu)$ -norm and is therefore \mathbf{P}^* -Glivenko–Cantelli. Moreover, by construction, $\mathbf{E}^*[(d_f(X_i))^2] = 1$ for all $f \in \mathcal{M}_q$. Therefore, the first term in this expression converges to zero as $n \rightarrow \infty$ \mathbf{P}^* -a.s. On the other hand, by Lemma B.2 and the fact that $D \in L^{2+\alpha}(f^* d\mu)$, we have \mathbf{P}^* -a.s.

$$\begin{aligned} & \sqrt{\log \log n/n} \max_{i=1, \dots, n} D(X_i) = \\ & \quad \frac{\sqrt{\log \log n}}{n^{\alpha/2(2+\alpha)}} n^{-1/(2+\alpha)} \max_{i=1, \dots, n} D(X_i) \xrightarrow{n \rightarrow \infty} 0. \end{aligned}$$

Therefore, the second term converges to zero also, and the proof is evidently complete. \blacksquare

Proposition B.5: Let $q \geq q^*$. Assume that

$$\int_0^1 \sqrt{\log \mathcal{N}(\mathcal{D}_q, u)} du < \infty.$$

Moreover, suppose that $|d| \leq D$ for all $d \in \mathcal{D}_q$ with $D \in L^{2+\alpha}(f^* d\mu)$ for some $\alpha > 0$. Then

$$\lim_{n \rightarrow \infty} \left\{ \sup_{d \in \bar{\mathcal{D}}_q} (I_n(d))_+^2 - \frac{1}{\log \log n} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \ell_n(f^*) \right\} \right\}$$

is nonnegative \mathbf{P}^* -a.s.

Proof: By Proposition B.4, we have

$$\begin{aligned} & \varliminf_{n \rightarrow \infty} \left\{ \sup_{d \in \bar{\mathcal{D}}_q} (I_n(d))_+^2 - \frac{1}{\log \log n} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \ell_n(f^*) \right\} \right\} \\ & \geq \varliminf_{n \rightarrow \infty} \left\{ \sup_{d \in \bar{\mathcal{D}}_q} (I_n(d))_+^2 \right. \\ & \quad \left. - \sup_{f \in \mathcal{M}_q (4\sqrt{\log \log n/n})} \sup_{p \geq 0} \{2 I_n(d_f) p - p^2\} \right\} \\ & = \varliminf_{n \rightarrow \infty} \left\{ \sup_{d \in \bar{\mathcal{D}}_q} (I_n(d))_+^2 - \sup_{f \in \mathcal{M}_q (4\sqrt{\log \log n/n})} (I_n(d_f))_+^2 \right\}. \end{aligned}$$

Suppose that the right-hand side is negative with positive probability. Then, there exists $\varepsilon > 0$ and a sequence $\tau_n \uparrow \infty$ of random times such that

$$\sup_{d \in \bar{\mathcal{D}}_q} (I_{\tau_n}(d))_+^2 - \sup_{f \in \mathcal{M}_q (4\sqrt{\log \log \tau_n/\tau_n})} (I_{\tau_n}(d_f))_+^2 \leq -\varepsilon \tag{B1}$$

for all n with positive probability. We will show that this entails a contradiction.

By Theorem B.1 (which can be applied here as $\mathcal{N}(\mathcal{D}_q, \delta) = \mathcal{N}(\text{cl } \mathcal{D}_q, \delta)$ for all $\delta > 0$), the process $(I_{\tau_n})_{n \geq 0}$ is \mathbf{P}^* -a.s. relatively compact in $\ell_\infty(\text{cl } \mathcal{D}_q)$ with

$$\inf_{g \in L_0^2(f^* d\mu)} \sup_{d \in \text{cl } \mathcal{D}_q} |I_{\tau_n}(d) - \langle d, g \rangle| \xrightarrow{n \rightarrow \infty} 0 \quad \mathbf{P}^*\text{-a.s.} \quad (\text{B2})$$

Then, there is a set of positive probability on which (B1) and (B2) hold simultaneously. We now concentrate our attention on a single sample path in this set. For any such path, we can clearly find a further subsequence $\sigma_n \uparrow \infty$ such that $\sup_{d \in \text{cl } \mathcal{D}_q} |I_{\sigma_n}(d) - \langle d, g \rangle| \rightarrow 0$ as $n \rightarrow \infty$ for some element $g \in L_0^2(f^* d\mu)$. Therefore, we obtain

$$\begin{aligned} \sup_{d \in \text{cl } \mathcal{D}_q} |(I_{\sigma_n}(d))_+^2 - (\langle d, g \rangle)_+^2| &\leq \sup_{d \in \text{cl } \mathcal{D}_q} |I_{\sigma_n}(d) - \langle d, g \rangle| \xrightarrow{n \rightarrow \infty} 0 \\ &+ 2 \sup_{d \in \text{cl } \mathcal{D}_q} |I_{\sigma_n}(d) - \langle d, g \rangle| \sup_{d \in \text{cl } \mathcal{D}_q} |\langle d, g \rangle| \xrightarrow{n \rightarrow \infty} 0 \end{aligned}$$

where we have used the elementary estimate

$$\begin{aligned} |a_+^2 - b_+^2| &= |a_+ - b_+|(a_+ + b_+) \\ &\leq |a_+ - b_+|(|a_+ - b_+| + 2b_+) \leq |a - b|(|a - b| + 2|b|) \end{aligned}$$

for any $a, b \in \mathbb{R}$, and the fact that $\sup_{d \in \text{cl } \mathcal{D}_q} |\langle d, g \rangle| \leq \sup_{d \in \text{cl } \mathcal{D}_q} \|d\|_2 \|g\|_2 \leq 1$. Thus, (B1) gives

$$\begin{aligned} \lim_{n \rightarrow \infty} \left\{ \sup_{d \in \bar{\mathcal{D}}_q} (\langle d, g \rangle)_+^2 - \sup_{f \in \mathcal{M}_q(4\sqrt{\log \log \sigma_n / \sigma_n})} (\langle d_f, g \rangle)_+^2 \right\} &= \\ \lim_{n \rightarrow \infty} \left\{ \sup_{d \in \bar{\mathcal{D}}_q} (I_{\sigma_n}(d))_+^2 - \sup_{f \in \mathcal{M}_q(4\sqrt{\log \log \sigma_n / \sigma_n})} (I_{\sigma_n}(d_f))_+^2 \right\} & \\ \leq -\varepsilon. \end{aligned}$$

But as the map $d \mapsto \langle d, g \rangle$ is continuous in $L^2(f^* d\mu)$ and $\text{cl } \mathcal{D}_q(4\sqrt{\log \log \sigma_n / \sigma_n})$ is compact in $L^2(f^* d\mu)$ (this follows from $\mathcal{N}(\mathcal{D}_q, \delta) < \infty$ for all $\delta > 0$), we have

$$\begin{aligned} \sup_{f \in \mathcal{M}_q(4\sqrt{\log \log \sigma_n / \sigma_n})} (\langle d_f, g \rangle)_+^2 &= \\ \sup_{d \in \text{cl } \mathcal{D}_q(4\sqrt{\log \log \sigma_n / \sigma_n})} (\langle d, g \rangle)_+^2 &\xrightarrow{n \rightarrow \infty} \\ \sup_{d \in \bigcap_{n \geq 0} \text{cl } \mathcal{D}_q(4\sqrt{\log \log \sigma_n / \sigma_n})} (\langle d, g \rangle)_+^2 &= \sup_{d \in \bar{\mathcal{D}}_q} (\langle d, g \rangle)_+^2. \end{aligned}$$

Thus, we have a contradiction, completing the proof. \blacksquare

We now obtain a converse to the previous result.

Proposition B.6: Let $q \geq q^*$. Assume that

$$\int_0^1 \sqrt{\log \mathcal{N}(\mathcal{D}_q, u)} du < \infty.$$

Moreover, suppose that $|d| \leq D$ for all $d \in \mathcal{D}_q$ with $D \in L^{2+\alpha}(f^* d\mu)$ for some $\alpha > 0$. Then

$$\overline{\lim}_{n \rightarrow \infty} \left\{ \sup_{d \in \bar{\mathcal{D}}_q^c} (I_n(d))_+^2 - \frac{1}{\log \log n} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \ell_n(f^*) \right\} \right\}$$

is nonpositive \mathbf{P}^* -a.s.

Proof: Suppose the result is false. By Proposition B.4, there is $\varepsilon > 0$ and a sequence $\tau_n \uparrow \infty$ of random times so that

$$\begin{aligned} &\sup_{d \in \bar{\mathcal{D}}_q^c} (I_{\tau_n}(d))_+^2 \\ &- \sup_{f \in \mathcal{M}_q(4\sqrt{\log \log \tau_n / \tau_n})} \left\{ -h(f, f^*)^2 \frac{2\tau_n}{\log \log \tau_n} \right. \\ &\quad \left. + 2 I_{\tau_n}(d_f) h(f, f^*) \sqrt{\frac{2\tau_n}{\log \log \tau_n}} \right\} \geq \varepsilon \end{aligned}$$

for all n with positive probability. Proceeding as in the proof of Proposition B.5, we can then show that there is a sequence of times $\sigma_n \uparrow \infty$ and some $g \in L_0^2(f^* d\mu)$ such that

$$\begin{aligned} &\overline{\lim}_{n \rightarrow \infty} \left\{ \sup_{d \in \bar{\mathcal{D}}_q^c} (\langle d, g \rangle)_+^2 \right. \\ &- \sup_{f \in \mathcal{M}_q(4\sqrt{\log \log \sigma_n / \sigma_n})} \left\{ -h(f, f^*)^2 \frac{2\sigma_n}{\log \log \sigma_n} \right. \\ &\quad \left. \left. + 2 \langle d_f, g \rangle h(f, f^*) \sqrt{\frac{2\sigma_n}{\log \log \sigma_n}} \right\} \right\} \geq \varepsilon. \end{aligned}$$

We will show that this entails a contradiction.

Let $d_0 \in \bar{\mathcal{D}}_q$ be a continuously accessible point. Then, there exists an $\alpha_0 > 0$ (depending on d_0) and a path $(f_\alpha)_{\alpha \in [0, \alpha_0]}$ such that $h(f_\alpha, f^*) = \alpha$ for all $\alpha \in [0, \alpha_0]$ and $d_{f_\alpha} \rightarrow d_0$ in $L^2(f^* d\mu)$ as $\alpha \rightarrow 0$. Now choose the sequence

$$\alpha_n = \{(\langle d_0, g \rangle)_+ + \sigma_n^{-1}\} \sqrt{\frac{\log \log \sigma_n}{2\sigma_n}}.$$

As $(\langle d_0, g \rangle)_+ \leq \|d_0\|_2 \|g\|_2 \leq 1$, we clearly have

$$0 < \alpha_n < \alpha_0 \wedge 4\sqrt{\log \log \sigma_n / \sigma_n}$$

for all n sufficiently large. In particular, it follows that $f_{\alpha_n} \in \mathcal{M}_q(4\sqrt{\log \log \sigma_n / \sigma_n})$, so that

$$\begin{aligned} &\sup_{f \in \mathcal{M}_q(4\sqrt{\log \log \sigma_n / \sigma_n})} \left\{ 2 \langle d_f, g \rangle h(f, f^*) \sqrt{\frac{2\sigma_n}{\log \log \sigma_n}} \right. \\ &\quad \left. - h(f, f^*)^2 \frac{2\sigma_n}{\log \log \sigma_n} \right\} \\ &\geq 2 \langle d_{f_{\alpha_n}}, g \rangle \{(\langle d_0, g \rangle)_+ + \sigma_n^{-1}\} - \{(\langle d_0, g \rangle)_+ + \sigma_n^{-1}\}^2. \end{aligned}$$

Therefore, we have

$$\begin{aligned} &\overline{\lim}_{n \rightarrow \infty} \left\{ \sup_{d \in \bar{\mathcal{D}}_q^c} (\langle d, g \rangle)_+^2 \right. \\ &- \sup_{f \in \mathcal{M}_q(4\sqrt{\log \log \sigma_n / \sigma_n})} \left\{ -h(f, f^*)^2 \frac{2\sigma_n}{\log \log \sigma_n} \right. \\ &\quad \left. \left. + 2 \langle d_f, g \rangle h(f, f^*) \sqrt{\frac{2\sigma_n}{\log \log \sigma_n}} \right\} \right\} \\ &\leq \sup_{d \in \bar{\mathcal{D}}_q^c} (\langle d, g \rangle)_+^2 - (\langle d_0, g \rangle)_+^2 \end{aligned}$$

for any continuously accessible element $d_0 \in \bar{\mathcal{D}}_q$. But clearly, we can choose d_0 to make the right-hand side of this expression arbitrarily small. Thus, we have the desired contradiction. ■

We can now complete the proof of Theorem II.6.

Proof of Theorem II.6: We obtain separately the lower and upper bounds.

Lower Bound: By Propositions B.5 and B.6, we have

$$\begin{aligned} \overline{\lim}_{n \rightarrow \infty} \frac{1}{\log \log n} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \sup_{f \in \mathcal{M}_p} \ell_n(f) \right\} &\geq \\ \overline{\lim}_{n \rightarrow \infty} \left\{ \sup_{d \in \bar{\mathcal{D}}_q^c} (I_n(d))_+^2 - \sup_{d \in \bar{\mathcal{D}}_p} (I_n(d))_+^2 \right\} &\quad \mathbf{P}^* \text{-a.s.} \end{aligned}$$

Now fix any $g \in L_0^2(f^* d\mu)$. By Theorem B.1 (which applies here as $\mathcal{N}(\mathcal{D}_q, \delta) = \mathcal{N}(\text{cl } \mathcal{D}_q, \delta) \geq \mathcal{N}(\bar{\mathcal{D}}_q, \delta)$ for all $\delta > 0$), there is a sequence $\tau_n \uparrow \infty$ of random times such that $I_{\tau_n} \rightarrow \langle \cdot, g \rangle$ in $\ell_\infty(\bar{\mathcal{D}}_q)$ \mathbf{P}^* -a.s. Therefore

$$\begin{aligned} \sup_{d \in \bar{\mathcal{D}}_q^c} (I_{\tau_n}(d))_+^2 - \sup_{d \in \bar{\mathcal{D}}_p} (I_{\tau_n}(d))_+^2 &\xrightarrow{n \rightarrow \infty} \\ \sup_{d \in \bar{\mathcal{D}}_q^c} (\langle d, g \rangle)_+^2 - \sup_{d \in \bar{\mathcal{D}}_p} (\langle d, g \rangle)_+^2 &\quad \mathbf{P}^* \text{-a.s.} \end{aligned}$$

so that certainly

$$\begin{aligned} \overline{\lim}_{n \rightarrow \infty} \frac{1}{\log \log n} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \sup_{f \in \mathcal{M}_p} \ell_n(f) \right\} &\geq \\ \sup_{d \in \bar{\mathcal{D}}_q^c} (\langle d, g \rangle)_+^2 - \sup_{d \in \bar{\mathcal{D}}_p} (\langle d, g \rangle)_+^2 & \end{aligned}$$

\mathbf{P}^* -a.s. But as this inequality holds for every $g \in L_0^2(f^* d\mu)$, taking the supremum over g gives the requisite lower bound.

Upper Bound: By Propositions B.5 and B.6, we have

$$\begin{aligned} \overline{\lim}_{n \rightarrow \infty} \frac{1}{\log \log n} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \sup_{f \in \mathcal{M}_p} \ell_n(f) \right\} &\leq \\ \overline{\lim}_{n \rightarrow \infty} \left\{ \sup_{d \in \bar{\mathcal{D}}_q} (I_n(d))_+^2 - \sup_{d \in \bar{\mathcal{D}}_p^c} (I_n(d))_+^2 \right\} &\quad \mathbf{P}^* \text{-a.s.} \end{aligned}$$

It is elementary that for any $d, d' \in \bar{\mathcal{D}}_q$ and $g \in L_0^2(f^* d\mu)$

$$\begin{aligned} (I_n(d))_+^2 - (I_n(d'))_+^2 & \\ \leq |(I_n(d))_+^2 - (\langle d, g \rangle)_+^2| + |(I_n(d'))_+^2 - (\langle d', g \rangle)_+^2| & \\ + (\langle d, g \rangle)_+^2 - (\langle d', g \rangle)_+^2 & \\ \leq 2 \sup_{d \in \bar{\mathcal{D}}_q} |(I_n(d))_+^2 - (\langle d, g \rangle)_+^2| + (\langle d, g \rangle)_+^2 - (\langle d', g \rangle)_+^2. & \end{aligned}$$

Taking the supremum over $d \in \bar{\mathcal{D}}_q$ and the infimum over $d' \in \bar{\mathcal{D}}_p^c$, we find that

$$\begin{aligned} \sup_{d \in \bar{\mathcal{D}}_q} (I_n(d))_+^2 - \sup_{d \in \bar{\mathcal{D}}_p^c} (I_n(d))_+^2 & \\ \leq 2 \sup_{d \in \bar{\mathcal{D}}_q} |(I_n(d))_+^2 - (\langle d, g \rangle)_+^2| & \\ + \sup_{d \in \bar{\mathcal{D}}_q} (\langle d, g \rangle)_+^2 - \sup_{d \in \bar{\mathcal{D}}_p^c} (\langle d, g \rangle)_+^2 & \end{aligned}$$

$$\begin{aligned} &\leq 2 \sup_{d \in \bar{\mathcal{D}}_q} |(I_n(d))_+^2 - (\langle d, g \rangle)_+^2| \\ &+ \sup_{g \in L_0^2(f^* d\mu)} \left\{ \sup_{d \in \bar{\mathcal{D}}_q} (\langle d, g \rangle)_+^2 - \sup_{d \in \bar{\mathcal{D}}_p^c} (\langle d, g \rangle)_+^2 \right\}. \end{aligned}$$

But as this holds for any $g \in L_0^2(f^* d\mu)$, we finally obtain

$$\begin{aligned} &\sup_{d \in \bar{\mathcal{D}}_q} (I_n(d))_+^2 - \sup_{d \in \bar{\mathcal{D}}_p^c} (I_n(d))_+^2 \\ &\leq 2 \inf_{g \in L_0^2(f^* d\mu)} \sup_{d \in \bar{\mathcal{D}}_q} |(I_n(d))_+^2 - (\langle d, g \rangle)_+^2| \\ &+ \sup_{g \in L_0^2(f^* d\mu)} \left\{ \sup_{d \in \bar{\mathcal{D}}_q} (\langle d, g \rangle)_+^2 - \sup_{d \in \bar{\mathcal{D}}_p^c} (\langle d, g \rangle)_+^2 \right\}. \end{aligned}$$

It follows as in the proof of Proposition B.5 that the first term in this expression converges to zero \mathbf{P}^* -a.s. The requisite upper bound follows immediately. ■

Finally, we now complete the proof of Corollary II.7.

Proof of Corollary II.7: It evidently suffices to prove that

$$\Gamma := \sup_{g \in L_0^2(f^* d\mu)} \left\{ \sup_{d \in \bar{\mathcal{D}}_q^c} (\langle d, g \rangle)_+^2 - \sup_{d \in \bar{\mathcal{D}}_q^{*\star}} (\langle d, g \rangle)_+^2 \right\} > 0. \quad (\text{B3})$$

To this end, note that by direct computation

$$\langle 1, d_f \rangle = \frac{\int \sqrt{ff^*} d\mu - 1}{h(f, f^*)} = -\frac{h(f, f^*)}{2}.$$

Choose $(f_n)_{n \geq 0} \subset \mathcal{M}_q \setminus \{f^*\}$ such that $h(f_n, f^*) \rightarrow 0$ and $d_{f_n} \rightarrow d_0 \in \bar{\mathcal{D}}_q$; then

$$\langle 1, d_0 \rangle = \lim_{n \rightarrow \infty} \langle 1, d_{f_n} \rangle = -\lim_{n \rightarrow \infty} \frac{h(f_n, f^*)}{2} = 0.$$

Moreover, it is immediate that $\|d_0\|_2 \leq 1$. We have, therefore, shown that $\bar{\mathcal{D}}_q \subset L_0^2(f^* d\mu)$. Now choose $g \in \bar{\mathcal{D}}_q^c \setminus \bar{\mathcal{D}}_q^{*\star}$. As $\bar{\mathcal{D}}_q^{*\star}$ is closed, it follows directly that

$$\sup_{d \in \bar{\mathcal{D}}_q^c} (\langle d, g \rangle)_+^2 = 1, \quad \sup_{d \in \bar{\mathcal{D}}_q^{*\star}} (\langle d, g \rangle)_+^2 < 1.$$

Therefore, (B3) holds, and the proof is complete. ■

APPENDIX C PROOF OF THEOREM III.1

The proof of Theorem III.1 is based on Theorem II.3 and the following result.

Proposition C.1: Let \mathcal{M}^n for $n \geq 1$ be a family of strictly positive probability densities with respect to a reference measure μ such that $\mathcal{M}^n \subseteq \mathcal{M}^{n+1}$ for all n . Define $\mathcal{M} = \bigcup_n \mathcal{M}^n$, and let f^* be another probability density with respect to μ such that $f^* \notin \text{cl } \mathcal{M}$, where $\text{cl } \mathcal{M}$ denotes the $L^1(d\mu)$ -closure of \mathcal{M} . Let $\mathcal{H}^n = \{\sqrt{f/f^*} : f \in \mathcal{M}^n\}$, and suppose there exist $K(n) \geq 1$ and $p \geq 1$ so that

$$\mathcal{N}(\mathcal{H}^n, \delta) \leq \left(\frac{K(n)}{\delta} \right)^p$$

for all $\delta \leq 1$ and $n \geq 1$, where $\mathcal{N}(\mathcal{H}^n, \delta)$ is the minimal number of brackets of $L^2(f^* d\mu)$ -width δ needed to cover \mathcal{H}^n . Let $(X_i)_{i \in \mathbb{N}}$ be i.i.d. with distribution $f^* d\mu$. If in addition $\log K(n) = o(n)$, then we have

$$\overline{\lim}_{n \rightarrow \infty} \sup_{f \in \mathcal{M}^n} \frac{1}{n} \sum_{j=1}^n \log \left(\frac{f(X_j)}{f^*(X_j)} \right) < 0 \quad \text{a.s.}$$

Proof: As in the proof of Theorem A.1, we have

$$\begin{aligned} \frac{1}{n} \sum_{j=1}^n \log \left(\frac{f(X_j)}{f^*(X_j)} \right) \\ \leq 4n^{-1/2} \nu_n(\log(\{\bar{f}/f^*\}^{1/2})) - 2D(f^* \|\bar{f}). \end{aligned}$$

The following claim will be proved below:

$$\lim_{n \rightarrow \infty} \sup_{f \in \mathcal{M}^n} n^{-1/2} \nu_n(\log(\{\bar{f}/f^*\}^{1/2})) = 0 \quad \text{a.s.}$$

Using the claim, the proof is easily completed: indeed, if the claim holds, then we have a.s.

$$\overline{\lim}_{n \rightarrow \infty} \sup_{f \in \mathcal{M}^n} \frac{1}{n} \sum_{j=1}^n \log \left(\frac{f(X_j)}{f^*(X_j)} \right) \leq -2 \inf_{f \in \mathcal{M}} D(f^* \|\bar{f}) < 0$$

where the last inequality follows from Pinsker's inequality and $f^* \notin \text{cl } \mathcal{M}$.

It, therefore, remains to prove the claim. To this end, we apply [23, Th. 5.11] as in the proof of [23, Th. 7.4] (cf., Theorem A.1), which yields

$$\mathbf{P} \left[\sup_{f \in \mathcal{M}^n} |n^{-1/2} \nu_n(\log(\{\bar{f}/f^*\}^{1/2}))| \geq \alpha \right] \leq C e^{-n\alpha^2/C}$$

for every $\alpha > 0$ such that $C\sqrt{p}(1 + \sqrt{\log K(n)}) \leq \alpha\sqrt{n} \leq 32\sqrt{n}$ and $n \geq 1$, where C is a universal constant. As $\log K(n) = o(n)$, we have

$$\sum_{n \geq 1} \mathbf{P} \left[\sup_{f \in \mathcal{M}^n} |n^{-1/2} \nu_n(\log(\{\bar{f}/f^*\}^{1/2}))| \geq \alpha \right] < \infty$$

for $0 < \alpha \leq 32$, so the claim follows from Borel–Cantelli. ■

We can now complete the proof of Theorem III.1.

Proof of Theorem III.1: Define

$$\Delta_n(q, q^*) = \sup_{f \in \mathcal{M}_q^n} \ell_n(f) - \sup_{f \in \mathcal{M}_{q^*}^n} \ell_n(f).$$

By Theorem II.3 and easy manipulations, \mathbf{P}^* -a.s.

$$\begin{aligned} & \overline{\lim}_{n \rightarrow \infty} \sup_{q > q^*} \frac{1}{\text{pen}(n, q) - \text{pen}(n, q^*)} \Delta_n(q, q^*) \\ & \leq \lim_{n \rightarrow \infty} \sup_{q > q^*} \frac{\eta(q)\{\log K(2n) \vee \log \log n\}}{\text{pen}(n, q) - \text{pen}(n, q^*)} \times \\ & \quad \overline{\lim}_{n \rightarrow \infty} \frac{1}{\log K(2n) \vee \log \log n} \sup_{q > q^*} \frac{1}{\eta(q)} \Delta_n(q, q^*) = 0. \end{aligned}$$

Therefore, \mathbf{P}^* -a.s. eventually as $n \rightarrow \infty$

$$\sup_{f \in \mathcal{M}_q^n} \ell_n(f) - \text{pen}(n, q) < \sup_{f \in \mathcal{M}_{q^*}^n} \ell_n(f) - \text{pen}(n, q^*)$$

for all $q > q^*$. It follows that $\overline{\lim}_{n \rightarrow \infty} \hat{q}_n \leq q^*$ \mathbf{P}^* -a.s., that is, the penalized likelihood order estimator does not asymptotically overestimate the order.

On the other hand, we note that for every $q < q^*$

$$\begin{aligned} & \overline{\lim}_{n \rightarrow \infty} \frac{1}{n} \left\{ \sup_{f \in \mathcal{M}_q^n} \ell_n(f) - \sup_{f \in \mathcal{M}_{q^*}^n} \ell_n(f) \right\} \\ & \leq \overline{\lim}_{n \rightarrow \infty} \sup_{f \in \mathcal{M}_q^n} \frac{1}{n} \sum_{j=1}^n \log \left(\frac{f(X_j)}{f^*(X_j)} \right) \end{aligned}$$

which is strictly negative \mathbf{P}^* -a.s. by Proposition C.1, where we have used that $\log K(n) = o(n)$ and that $\mathcal{N}(\mathcal{H}_q^n(2), \delta) \leq \mathcal{N}(\mathcal{H}_{q^*}^n(2), \delta) \leq (2K(n)/\delta)^{\eta(q^*)}$ for all $\delta \leq 2$ and n sufficiently large. As $\text{pen}(n, q)/n \rightarrow 0$ as $n \rightarrow \infty$ for $q < q^*$

$$\overline{\lim}_{n \rightarrow \infty} \max_{q < q^*} \frac{1}{n} \{ \Delta_n(q, q^*) - \text{pen}(n, q) + \text{pen}(n, q^*) \} < 0$$

\mathbf{P}^* -a.s. In particular, we find that \mathbf{P}^* -a.s. eventually as $n \rightarrow \infty$

$$\sup_{f \in \mathcal{M}_q^n} \ell_n(f) - \text{pen}(n, q) < \sup_{f \in \mathcal{M}_{q^*}^n} \ell_n(f) - \text{pen}(n, q^*)$$

for all $q < q^*$. It follows that $\underline{\lim}_{n \rightarrow \infty} \hat{q}_n \geq q^*$ \mathbf{P}^* -a.s., that is, the penalized likelihood order estimator does not asymptotically underestimate the order. ■

Finally, let us prove Corollary III.3.

Proof of Corollary III.3: It is shown in the proof of Corollary II.7 that

$$\Gamma := \sup_{g \in L_0^2(f^* d\mu)} \left\{ \sup_{d \in \bar{\mathcal{D}}_q^c} (\langle d, g \rangle)_+^2 - \sup_{d \in \bar{\mathcal{D}}_{q^*}} (\langle d, g \rangle)_+^2 \right\} > 0.$$

By Theorem II.6, we have \mathbf{P}^* -a.s.

$$\begin{aligned} & \overline{\lim}_{n \rightarrow \infty} \frac{1}{\text{pen}(n, q) - \text{pen}(n, q^*)} \left\{ \sup_{f \in \mathcal{M}_q} \ell_n(f) - \sup_{f \in \mathcal{M}_{q^*}} \ell_n(f) \right\} \\ & \geq \frac{1}{C\{\eta(q) - \eta(q^*)\}} \times \\ & \quad \sup_{g \in L_0^2(f^* d\mu)} \left\{ \sup_{d \in \bar{\mathcal{D}}_q^c} (\langle d, g \rangle)_+^2 - \sup_{d \in \bar{\mathcal{D}}_{q^*}} (\langle d, g \rangle)_+^2 \right\}. \end{aligned}$$

Therefore, choosing $C < \Gamma/\{\eta(q) - \eta(q^*)\}$, we find that

$$\sup_{f \in \mathcal{M}_q} \ell_n(f) - \text{pen}(n, q) > \sup_{f \in \mathcal{M}_{q^*}} \ell_n(f) - \text{pen}(n, q^*)$$

infinitely often \mathbf{P}^* -a.s., so $\hat{q}_n \neq q^*$ infinitely often \mathbf{P}^* -a.s. ■

APPENDIX D

PROOF OF PROPOSITION III.6

The proofs of consistency in Propositions III.6 and III.7 follow almost immediately from Theorem III.1, Theorem III.4, and Example III.5. Let us begin with Proposition III.7.

Proof of Proposition III.7: By Example III.5, the assumption of Theorem III.1 holds with $\eta(q) = 10(d+1)q + 1$ and $\log K(n) = \log C_1^* + C_2^* T(n)^2$. The desired consistency results now follow immediately from Theorem III.1. ■

The consistency part of Proposition III.6 follows similarly. The main difficulty here is to establish the condition $\bar{\mathcal{D}}_q^c \setminus \bar{\mathcal{D}}_{q^*} \neq \emptyset$

\emptyset of Corollary III.3, which is needed to prove the inconsistency part of Proposition III.6. In the proof of the latter condition, we rely on the geometric results on mixtures established in [22]. In the remainder of this section, we always assume that we are in the setting of Proposition III.6.

Lemma D.1: Suppose that Assumption A holds. Then, we have

$$\bar{\mathcal{D}}_{q^*} = \left\{ \frac{L}{\|L\|_2} : L = \sum_{i=1}^{q^*} \left\{ \eta_i \frac{f_{\theta_i^*}}{f^*} + \beta_i^* \frac{D_1 f_{\theta_i^*}}{f^*} \right\}, \right. \\ \left. \eta_i \in \mathbb{R}, \beta_i \in \mathbb{R}^d, \sum_{i=1}^{q^*} \eta_i = 0 \right\}.$$

Proof: Let $(f_n)_{n \geq 1} \subset \mathcal{M}_{q^*}$ be such that $h(f_n, f^*) \rightarrow 0$ and $d_{f_n} \rightarrow d_0 \in \bar{\mathcal{D}}_{q^*}$. By [22, Th. 3.10], we may assume without loss of generality that $f_n = \sum_{i=1}^{q^*} \pi_i^n f_{\theta_i^n}$ with $\theta_i^n \rightarrow \theta_i^*$ and $\pi_i^n \rightarrow \pi_i^*$ for $i = 1, \dots, q^*$. Taylor expansion gives

$$\frac{f_n - f^*}{f^*} = L_n + R_n, \quad |R_n| \leq \frac{d}{2} H_2 \sum_{i=1}^{q^*} \pi_i^n \|\theta_i^n - \theta_i^*\|^2$$

where

$$L_n = \sum_{i=1}^{q^*} \left\{ (\pi_i^n - \pi_i^*) \frac{f_{\theta_i^*}}{f^*} + \pi_i^n (\theta_i^n - \theta_i^*)^* \frac{D_1 f_{\theta_i^*}}{f^*} \right\}.$$

Proceeding as in [22, Lemmas 3.15–3.16], we can estimate

$$\left\| d_{f_n} - \frac{L_n}{\|L_n\|_2} \right\|_2 \\ \leq 2\|S\|_4^2 \{2\|S\|_2 + 1\} h(f_n, f^*) + \{\|S\|_2 + 1\} \frac{\|R_n\|_2}{\|L_n\|_2}.$$

But using [22, Th. 3.10], we have for n sufficiently large

$$\|L_n\|_2 \geq \|L_n\|_1 \geq c^* \sum_{i=1}^{q^*} \pi_i^n \|\theta_i^n - \theta_i^*\|.$$

Thus, we have

$$\frac{\|R_n\|_2}{\|L_n\|_2} \leq \frac{d\|H_2\|_2}{2c^*} \frac{\sum_{i=1}^{q^*} \pi_i^n \|\theta_i^n - \theta_i^*\|^2}{\sum_{i=1}^{q^*} \pi_i^n \|\theta_i^n - \theta_i^*\|} \\ \leq \frac{d\|H_2\|_2}{2c^*} \max_{i=1, \dots, q^*} \|\theta_i^n - \theta_i^*\| \xrightarrow{n \rightarrow \infty} 0.$$

Therefore, $L_n/\|L_n\|_2 \rightarrow d_0$ in $L^2(f^* d\mu)$. Now define

$$\eta_i^n = \frac{\pi_i^n - \pi_i^*}{Z_n}, \quad \beta_i^n = \frac{\pi_i^n (\theta_i^n - \theta_i^*)}{Z_n} \\ Z_n = \sum_{i=1}^{q^*} \{|\pi_i^n - \pi_i^*| + \|\pi_i^n (\theta_i^n - \theta_i^*)\|\}.$$

As $\sum_{i=1}^{q^*} \{|\eta_i^n| + \|\beta_i^n\|\} = 1$ for all n , we may extract a subsequence such that $\eta_i^n \rightarrow \eta_i$, $\beta_i^n \rightarrow \beta_i$, and $\sum_{i=1}^{q^*} \{|\eta_i| + \|\beta_i\|\} = 1$. We obtain immediately

$$d_0 = \frac{L}{\|L\|_2}, \quad L = \sum_{i=1}^{q^*} \left\{ \eta_i \frac{f_{\theta_i^*}}{f^*} + \beta_i^* \frac{D_1 f_{\theta_i^*}}{f^*} \right\}.$$

Clearly, $\sum_{i=1}^{q^*} \eta_i = 0$. Thus, we have shown that any $d_0 \in \bar{\mathcal{D}}_{q^*}$ has the desired form.

It remains to show that any function of the desired form is in fact an element of $\bar{\mathcal{D}}_{q^*}$. To this end, fix $\eta_i \in \mathbb{R}$, $\beta_i \in \mathbb{R}^d$ with $\sum_{i=1}^{q^*} \eta_i = 0$, and define f_t for $t > 0$ as

$$f_t = \sum_{i=1}^{q^*} (\pi_i^* + t\eta_i) f_{\theta_i^* + \beta_i t / \pi_i^*}.$$

Clearly, $f_t \in \mathcal{M}_{q^*}$ for all t sufficiently small, and $f_t \rightarrow f^*$ as $t \rightarrow 0$. But

$$\frac{f_t - f^*}{t} = \sum_{i=1}^{q^*} \pi_i^* \frac{f_{\theta_i^* + \beta_i t / \pi_i^*} - f_{\theta_i^*}}{t} + \sum_{i=1}^{q^*} \eta_i f_{\theta_i^* + \beta_i t / \pi_i^*}.$$

Therefore, clearly

$$\frac{1}{t} \frac{f_t - f^*}{f^*} \xrightarrow{t \rightarrow 0} \sum_{i=1}^{q^*} \left\{ \eta_i \frac{f_{\theta_i^*}}{f^*} + \beta_i^* \frac{D_1 f_{\theta_i^*}}{f^*} \right\} = L.$$

Using [22, Lemma 3.15], we obtain

$$\lim_{t \rightarrow 0} d_{f_t} = \lim_{t \rightarrow 0} \frac{(f_t - f^*)/tf^*}{\|(f_t - f^*)/tf^*\|_2} = \frac{L}{\|L\|_2}.$$

Thus, any function of the desired form is in $\bar{\mathcal{D}}_{q^*}$. ■

Remark D.2: The above proof in fact shows that $\bar{\mathcal{D}}_{q^*} = \bar{\mathcal{D}}_{q^*}^c$.

We can now complete the proof of Proposition III.6.

Proof of Proposition III.6: We begin by proving consistency of the penalty $\text{pen}(n, q) = q\omega(n)$. Note that by Theorem III.4, the assumption of Corollary III.2 holds with $\eta(q) = 10(d+1)q + 1 \leq 11(d+1)q$. Thus, consistency of $\text{pen}(n, q) = q\omega(n)$ follows directly from Corollary III.2 using $\varpi(n) = \omega(n)/11(d+1)$.

To prove that the penalty $\text{pen}(n, q) = Cq \log \log n$ is inconsistent for $C > 0$ sufficiently small, it suffices to show that $\bar{\mathcal{D}}_{q^*+1} \setminus \bar{\mathcal{D}}_{q^*}$ is nonempty. Indeed, if this is the case, then we can apply Corollary III.3 with $q = q^* + 1$, where the requisite entropy assumption follows from Theorem III.1.

Fix $v \in \mathbb{R}^d$, and consider f_t defined for $t > 0$ as follows:

$$f_t = \frac{\pi_1^*}{2} (f_{\theta_1^* + vt} + f_{\theta_1^* - vt}) + \sum_{i=2}^{q^*} \pi_i^* f_{\theta_i^*}.$$

Clearly, $f_t \in \mathcal{M}_{q^*+1}$ for all t sufficiently small, $f_t \rightarrow f^*$ as $t \rightarrow 0$, and

$$\frac{f_t - f^*}{t^2} = \frac{\pi_1^*}{2} \frac{f_{\theta_1^* + vt} - 2f_{\theta_1^*} + f_{\theta_1^* - vt}}{t^2} \xrightarrow{t \rightarrow 0} \frac{\pi_1^*}{2} v^* D_2 f_{\theta_1^*} v.$$

As in the proof of Lemma D.1, we find that

$$\lim_{t \rightarrow 0} d_{f_t} = \lim_{t \rightarrow 0} \frac{(f_t - f^*)/t^2 f^*}{\|(f_t - f^*)/t^2 f^*\|_2} = \frac{v^* D_2 f_{\theta_1^*} v}{\|v^* D_2 f_{\theta_1^*} v\|_2} = d_0.$$

By construction, $d_0 \in \bar{\mathcal{D}}_{q^*+1}^c$. But by [22, Th. 3.10], the functions $f_{\theta_i^*}$, $D_1 f_{\theta_i^*}$, and $v^* D_2 f_{\theta_i^*} v$ ($i = 1, \dots, q^*$) are all linearly independent. Together with Lemma D.1, this shows that $d_0 \notin \bar{\mathcal{D}}_{q^*}$. Thus, $d_0 \in \bar{\mathcal{D}}_{q^*+1}^c \setminus \bar{\mathcal{D}}_{q^*}$. ■

ACKNOWLEDGMENT

The authors would like to thank Michel Ledoux for suggesting some helpful references.

REFERENCES

- [1] I. Csiszár and P. C. Shields, “The consistency of BIC Markov order estimator,” *Ann. Stat.*, vol. 28, pp. 1601–1619, 2000.
- [2] E. J. Hannan and B. G. Quinn, “The determination of the order of an autoregression,” *J. Roy. Statist. Soc. Ser. B*, vol. 41, pp. 190–195, 1979.
- [3] R. Nishii, “Maximum likelihood principle and model selection when the true model is unspecified,” *J. Multivariate Anal.*, vol. 27, pp. 392–403, 1988.
- [4] L. Finesso, “Consistent estimation of the order for Markov and hidden Markov chains,” Ph.D. dissertation, Univ. Maryland, College Park, 1990.
- [5] C. Keribin, “Consistent estimation of the order of mixture models,” *Sankhya Ser. A*, vol. 62, pp. 49–66, 2000.
- [6] A. Chambaz, “Testing the order of a model,” *Ann. Statist.*, vol. 34, pp. 1166–1203, 2006.
- [7] J. A. Hartigan, “A failure of likelihood asymptotics for normal mixtures,” in *Proceedings of the Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer, Vol. II (Berkeley, Calif., 1983)*. Belmont, CA: Wadsworth, 1985, pp. 807–810.
- [8] P. Bickel and H. Chernoff, “Asymptotic Distribution of the Likelihood Ratio Statistic in a Prototypical Non Regular Problem,” in *Statistics and Probability: A Raghu Raj Bahadur Festschrift*. New Delhi, India: Wiley, 1993, pp. 83–96.
- [9] X. Liu and Y. Shao, “Asymptotics for the likelihood ratio test in a two-component normal mixture model,” *J. Statist. Plann. Infer.*, vol. 123, no. 1, pp. 61–81, 2004.
- [10] J. C. Kieffer, “Strongly consistent code-based identification and order estimation for constrained finite-state model classes,” *IEEE Trans. Inf. Theory*, vol. 39, no. 3, pp. 893–902, May 1993.
- [11] C.-C. Liu and P. Narayan, “Order estimation and sequential universal data compression of a hidden Markov source by the method of mixtures,” *IEEE Trans. Inf. Theory*, vol. 40, no. 4, pp. 1167–1180, Jul. 1994.
- [12] E. Gassiat and S. Boucheron, “Optimal error exponents in hidden Markov model order estimation,” *IEEE Trans. Inf. Theory*, vol. 48, no. 4, pp. 964–980, Apr. 2003.
- [13] O. Cappé, E. Moulines, and T. Rydén, *Inference in Hidden Markov Models*. New York: Springer, 2005.
- [14] A. Chambaz, A. Garivier, and E. Gassiat, “A MDL approach to HMM with Poisson and Gaussian emissions. Application to order identification,” *J. Stat. Plann. Inf.*, vol. 139, pp. 962–977, 2009.
- [15] J. Rissanen, “Stochastic complexity and modeling,” *Ann. Statist.*, vol. 14, pp. 1080–1100, 1986.
- [16] A. Barron, J. Rissanen, and B. Yu, “The minimum description length principle in coding and modeling,” *IEEE Trans. Inf. Theory*, vol. 44, no. 6, pp. 2743–2760, Oct. 1998.
- [17] I. Csiszar, “Large-scale typicality of Markov sample paths and consistency of MDL order estimators,” *IEEE Trans. Inf. Theory*, vol. 48, no. 6, pp. 1616–1628, Jun. 2002.
- [18] I. Csiszar and Z. Talata, “Context tree estimation for not necessarily finite memory processes, via BIC and MDL,” *IEEE Trans. Inf. Theory*, vol. 52, no. 3, pp. 1007–1616, Mar. 2006.
- [19] R. van Handel, “On the minimal penalty for Markov order estimation,” *Probab. Th. Rel. Fields*, vol. 150, pp. 709–738, 2011.
- [20] E. Gassiat and C. Keribin, “The likelihood ratio test for the number of components in a mixture with Markov regime,” *ESAIM Probab. Statist.*, vol. 4, pp. 25–52, 2000.
- [21] A. W. van der Vaart, *Asymptotic Statistics*. Cambridge, U.K.: Cambridge Univ. Press, 1998.
- [22] E. Gassiat and R. van Handel, Preprint “The local geometry of finite mixtures,” 2012.
- [23] S. A. van de Geer, *Applications of Empirical Process Theory*. Cambridge, U.K.: Cambridge Univ. Press, 2000.
- [24] P. Billingsley, *Convergence of Probability Measures*, 2nd ed. New York: Wiley, 1999.
- [25] M. Ossiander, “A central limit theorem under metric entropy with L_2 bracketing,” *Ann. Probab.*, vol. 15, pp. 897–919, 1987.
- [26] M. Ledoux and M. Talagrand, “Comparison theorems, random geometry and some limit theorems for empirical processes,” *Ann. Probab.*, vol. 17, pp. 596–631, 1989.
- [27] E. Gassiat, “Likelihood ratio inequalities with applications to various mixtures,” *Ann. Inst. H. Poincaré Probab. Statist.*, vol. 38, pp. 897–906, 2002.
- [28] X. Liu and Y. Shao, “Asymptotics for likelihood ratio tests under loss of identifiability,” *Ann. Statist.*, vol. 31, pp. 807–832, 2003.

Elisabeth Gassiat received the Ph.D. degree from Université Paris-Sud in 1988, and is currently Professor in the Mathematic Department at Université Paris-Sud. Her main research interest is in statistical theory, including semi-parametric statistics, mixture and hidden Markov modeling, Bayesian inference.

Ramon van Handel received the Ph.D. degree from the California Institute of Technology in 2007, and is currently on the faculty of the School of Engineering and Applied Science at Princeton University. His main research interests are in probability theory and related fields.