

Computable Bayesian Compression for Uniformly Discretizable Statistical Models

Łukasz Dębowski

Centrum Wiskunde & Informatica, 1098 XG Amsterdam, The Netherlands

Abstract. Supplementing Vovk and V’yugin’s ‘if’ statement, we show that Bayesian compression provides the best enumerable compression for parameter-typical data if and only if the parameter is Martin-Löf random with respect to the prior. The result is derived for uniformly discretizable statistical models, introduced here. They feature the crucial property that given a discretized parameter, we can compute how much data is needed to learn its value with little uncertainty. Exponential families and certain nonparametric models are shown to be uniformly discretizable.

1 Introduction

Algorithmic information theory inspires an appealing interpretation of Bayesian inference [1–4]. Literally, a fixed individual parameter cannot have the property of being distributed according to a distribution but, when it is represented as a sequence of digits, the parameter is almost surely algorithmically random. Thus, if you believe that a parameter obeys a prior, it may rather mean that you suppose that the parameter is algorithmically random with respect to the prior. We want to argue that this interpretation is valid.

We will assume that the parameter θ is, in some sense, effectively identifiable. Then one can disprove that a finite prefix of a fixed, not fully known θ is algorithmically random by estimating the prefix and showing that there exists a shorter description of that prefix. Hence, Bayesian beliefs seem admissible scientific hypotheses according to the Popperian philosophy, cf. [1].

Secondly, it follows that the Bayesian measure $\int P_\theta dQ(\theta)$ gives the best enumerable compression of P_θ -typical data *if and only if* parameter θ is algorithmically random with respect to prior Q . This statement is useful when P_θ is not computable for a fixed θ . Moreover, once we know where Bayesian compression fails, we should systematically adjust the prior to our hypotheses about the algorithmic complexity of θ in an application.

As we will show, this ‘*if and only if*’ result can be foreseen using the chain rule for prefix Kolmogorov complexity of finite objects [5], [6, Theorem 3.9.1]. The chain rule allows to relate randomness deficiencies for finite prefixes of the data and of the parameter in some specific statistical models, which we call uniformly discretizable. That yields a somewhat weaker ‘*if and only if*’ statement. Subsequently, the statement can be strengthened using the dual chain rule for impossibility levels of infinite sequences [1, Theorem 1] and extensions of Lamalgen’s theorem for conditionally random sequences [7], [4, Theorem 4.2 and

5.3]. The condition of uniform discretization can be completely removed from the ‘*if*’ part and relaxed to an effective identifiability of the parameter in the ‘*only if*’ part. Namely, given a prefix of the parameter, we must be able to compute how much data is needed to learn its value with a fixed uncertainty.¹

The organization of this paper is as follows. In Section 2, we discuss quality of Bayesian compression for individual parameters and we derive the randomness deficiency bounds for prefixes of the parameter and the parameter-typical data. These bounds hold for the newly introduced class of uniformly discretizable statistical models. In Section 3, we show that exponential families are uniformly discretizable. The assumptions on the prior and the proof look familiar to statisticians working in minimum description length (MDL) inference [8, 9]. An example of a ‘nonparametric’ uniformly discretizable model appears in Section 4. In the final Section 5, we prove that countable mixtures of uniformly discretizable models are uniformly discretizable if the Bayesian estimator consistently chooses the right submodel for the data.

The definition of uniformly discretizable models is given below. Condition (3) says that the parameter may be discretized to $m \geq \mu(n)$ digits for the sake of approximating the ‘true’ probability of data x^n . Condition (4) asserts that the parameter, discretized to m digits, can be predicted for all but finitely many m given data x^n of length $n \geq \nu(m)$. Functions μ and ν depend on a model.

To fix our notation in advance, we use a countable alphabet \mathbb{X} and a finite $\mathbb{Y} = \{0, 1, \dots, D - 1\}$, $D > 1$. The logarithm to base D is written as \log . An italic $x \in \mathbb{X}^+$ is a string, a boldface $\mathbf{x} \in \mathbb{X}^{\mathbb{N}}$ is an infinite sequence. The n -th symbol of \mathbf{x} is written as $x_n \in \mathbb{X}$ and x^n is the prefix of \mathbf{x} of length n : $\mathbf{x} = x_1x_2x_3\dots$ and $x^n = x_1x_2\dots x_n$. Capital boldface $\mathbf{Y} : \mathbb{X}^* \rightarrow \mathbb{R}$ denotes a distribution of strings normalized lengthwise, i.e., $0 \leq \mathbf{Y}(x)$, $\sum_a \mathbf{Y}(xa)\mathbf{1}_{\{|a|=n\}} = \mathbf{Y}(x)$, and $\mathbf{Y}(\lambda) = 1$ for the empty string λ . There is a unique measure on measurable sets of infinite sequences $\mathbf{x} \in \mathbb{X}^{\mathbb{N}}$, also denoted as \mathbf{Y} , such that $\mathbf{Y}(\{\mathbf{x} : x^n = x \text{ for } n = |x|\}) = \mathbf{Y}(x)$. Quantifier ‘ n -eventually’ means ‘for all but finitely many $n \in \mathbb{N}$ ’.

Definition 1. Fix a measurable subset $\Theta \subset \mathbb{Y}^{\mathbb{N}}$. Let $\mathbf{P} : \mathbb{X}^* \times \Theta \ni (x, \theta) \mapsto \mathbf{P}_{\theta}(x) \in \mathbb{R}$ be a probability kernel, i.e., $\mathbf{P}_{\theta} : \mathbb{X}^* \rightarrow \mathbb{R}$ is a probability measure for each $\theta \in \Theta$ and the mapping $\theta \mapsto \mathbf{P}_{\theta}$ is measurable. Let also $\mathbf{Q} : \mathbb{Y}^* \rightarrow \mathbb{R}$ be a probability measure on Θ , i.e., $\mathbf{Q}(\Theta) = 1$. A Bayesian statistical model (\mathbf{P}, \mathbf{Q}) is called (μ, ν) -uniformly discretizable if it satisfies the following.

(i) Define the measure $\mathbf{T} : \mathbb{X}^* \times \mathbb{Y}^* \rightarrow \mathbb{R}$ as

$$\mathbf{T}(x, \theta) := \int_{A(\theta)} \mathbf{P}_{\theta}(x) d\mathbf{Q}(\theta), \quad (1)$$

¹ Note added post scriptum: By the Proposition 1 of V.V. V’yugin, *On empirical meaning of randomness with respect to a real parameter*, <http://arxiv.org/abs/0806.4484>, 2008, the *only if* part and the disjointness of sets $\mathcal{L}_{\mathbf{P}|\theta}$ conjectured in Section 2 follow if there exist effectively strongly consistent estimators. Although such estimators exists for the examples of models discussed further, cf. G. Davie, *Ann. Probab.* 29:1426–1434, 2001, it does not seem that recursive models from our Definition 1 admit them in general.

where $A(\theta) := \{\theta \in \Theta : \theta \text{ is the prefix of } \theta\}$, and denote its other marginal

$$\mathbf{Y}(x) := \mathbf{T}(x, \lambda) = \int \mathbf{P}_\theta(x) d\mathbf{Q}(\theta). \quad (2)$$

(ii) Function $\mu : \mathbb{N} \rightarrow \mathbb{R}$ is nondecreasing and we require that for all $\theta \in \Theta$, \mathbf{P}_θ -almost all \mathbf{x} , and $m \geq \mu(n)$,

$$\lim_{n \rightarrow \infty} \frac{\log [\mathbf{Q}(\theta^m) \mathbf{P}_\theta(x^n) / \mathbf{T}(x^n, \theta^m)]}{\log m} = 0. \quad (3)$$

(iii) Function $\nu : \mathbb{N} \rightarrow \mathbb{R}$ is nondecreasing and we require that for all $\theta \in \Theta$, \mathbf{P}_θ -almost all \mathbf{x} , and $n \geq \nu(m)$,

$$\lim_{m \rightarrow \infty} \mathbf{T}(x^n, \theta^m) / \mathbf{Y}(x^n) = 1. \quad (4)$$

Remark: A Bayesian model $(\tilde{\mathbf{P}}, \tilde{\mathbf{Q}})$ with a kernel $\tilde{\mathbf{P}} : \mathbb{X}^* \times \tilde{\Theta} \rightarrow \mathbb{R}$ and a measure $\tilde{\mathbf{Q}}$ on $\tilde{\Theta}$ will be called (ρ, μ, ν) -uniformly discretizable if $(\tilde{\mathbf{P}}, \tilde{\mathbf{Q}})$ is (μ, ν) -uniformly discretizable for a bijection $\rho : \tilde{\Theta} \rightarrow \Theta$, $\mathbf{P}_\theta(x) := \tilde{\mathbf{P}}_{\rho^{-1}(\theta)}(x)$, and $\mathbf{Q} := \tilde{\mathbf{Q}} \circ \rho^{-1}$. We will write ‘ $(\rho, \mu(n), \nu(m))$ -uniformly discretizable’ when there are no convenient symbols for functions μ and ν .

A few words of comment to this definition are due. By condition (3), the support of prior \mathbf{Q} equals Θ , i.e., $\mathbf{Q}(\theta^m) > 0$ for all m and $\theta \in \Theta$. Condition (4) admits a consistent estimator if there is a function $\sigma : \mathbb{X}^* \rightarrow \mathbb{N}$, where $\nu(\sigma(x^n)) \leq n$, $\sigma(x^{n+1}) \geq \sigma(x^n)$, and $\lim_n \sigma(x^n) = \infty$. Define a discrete maximum a posteriori estimator $\text{MAP}(x; \sigma) \in \arg\max_{\theta \in \mathbb{Y}^m} \mathbf{T}(x, \theta)$ with $m = \sigma(x)$.² The estimator is called consistent if $\text{MAP}(x^n; \sigma) = \theta^{\sigma(x^n)}$ n -eventually for all $\theta \in \Theta$ and \mathbf{P}_θ -almost all \mathbf{x} . This property is indeed satisfied.

Four models presented in Sections 3 and 4 satisfy a stronger condition.

Definition 2. A (μ, ν) -uniformly discretizable model is called μ -uniformly discretizable if ν is recursive and $\mu(\nu(m)) \leq m^\alpha$ for an $\alpha > 0$.

These models feature $\log \mu(n)$ close to the logarithm of Shannon redundancy $-\log \mathbf{Y}(x^n) + \log \mathbf{P}_\theta(x^n)$. A heuristic rationale is as follows. If we had $\mu \circ \nu = \text{id}$, $-\log \mathbf{Q}(\theta^m) = \Omega(m)$, and we put $n = \nu(m)$ then

$$|-\log \mathbf{Y}(x^n) + \log \mathbf{P}_\theta(x^n) + \log \mathbf{Q}(\theta^m)| = o(\log m)$$

and hence $\mu(n) = m = O(-\log \mathbf{Y}(x^n) + \log \mathbf{P}_\theta(x^n))$. Whereas $-\log \mathbf{Q}(\theta^m) = \Omega(m)$ is a reasonable assumption, we rather observe $\mu(\nu(m)) > m$.

The present approach allows only discrete data. We hope, however, that uniformly discretizable models can be generalized to nondiscrete data so that consistency and algorithmic optimality of Bayesian procedures in density estimation could be characterized in a similar fashion, cf. [10]. Another interesting

² Note added post scriptum: Wrongly called a discrete maximum likelihood estimator in the version published in the proceedings.

path of development is to integrate the algorithmic perspective on Bayesianism with the present MDL framework [8, 9], where normalized maximum likelihood codes are discussed. By the algorithmic optimality of Bayesian compression, the normalized maximum likelihood measure, if it can be defined properly, should converge to the Bayesian measure $\int \mathbf{P}_\theta d\mathbf{Q}(\theta)$ in log-loss. We also suppose that reasonable luckiness functions, introduced to guarantee existence of modified normalized maximum likelihood codes [9, Section 11.3], may be close to algorithmic information between the parameter and the data.

2 Bounds for the data and parameter complexity

We will use a universal computer with an oracle, which can compute certain functions $\mathbb{R} \rightarrow \mathbb{R}$. To make it clear which these are, we adopt the following definitions, cf. [11], [6, Sections 1.7 and 3.1], [1, Section 2], [4, Section 5]:

- (i) A universal computer is an appropriate finite state machine that interacts with infinite tapes. The machine can move along the tapes in discrete steps, read and write on them single symbols from the finite set \mathbb{Y} , and announce the end of computation. We fix three one-sided tapes. At the beginning of computation, tape α contains a program, i.e., a string from a prefix-free subset of \mathbb{Y}^+ , and tape β contains an oracle, i.e., an element of $(0\mathbb{Y}^*) \cup (1\mathbb{Y}^{\mathbb{N}})$. At the end of computation, tape γ contains an output, i.e., a string from \mathbb{Y}^* .
- (ii) The prefix Kolmogorov complexity $K(y)$ of a string $y \in \mathbb{Y}^*$ is the minimal length of such a program on tape α that y is output on tape γ provided no symbol is read from tape β .
- (iii) The conditional complexity $K(y|\delta)$ for $y \in \mathbb{Y}^*$ and $\delta \in \mathbb{Y}^* \cup \mathbb{Y}^{\mathbb{N}}$ is the minimal length of such a program on tape α that y is output on tape γ given 0δ or 1δ , respectively, as an oracle on tape β .
- (iv) A function $f : \mathbb{Y}^* \cup \mathbb{Y}^{\mathbb{N}} \rightarrow \mathbb{Y}^*$ is *recursive* if there is such a program $z \in \mathbb{Y}^+$ that string $f(y)$ is output for all oracles $y \in \mathbb{Y}^* \cup \mathbb{Y}^{\mathbb{N}}$.
- (v) Function ϕ is a *prefix code* if it is an injection and its image is prefix-free.
- (vi) For certain prefix codes $\phi_{\mathbb{W}} : \mathbb{W} \rightarrow \mathbb{Y}^*$ and $\phi_{\mathbb{U}} : \mathbb{U} \rightarrow \mathbb{Y}^* \cup \mathbb{Y}^{\mathbb{N}}$ and arbitrary $w \in \mathbb{W}$ and $u \in \mathbb{U}$, we put $K(w) := K(\phi_{\mathbb{W}}(w))$ and $K(w|u) := K(\phi_{\mathbb{W}}(w)|\phi_{\mathbb{W}}(u))$. Fixing $\phi_{\mathbb{Y}^*}$ and $\phi_{\mathbb{Y}^* \cup \mathbb{Y}^{\mathbb{N}}}$ as identity functions, $f : \mathbb{U} \rightarrow \mathbb{W}$ is called *recursive* if so is $\phi_{\mathbb{W}} \circ f \circ \phi_{\mathbb{U}}^{-1}$.
- (vii) Numbers obey special conventions. Integers are Elias-coded [12], whereas $\phi_{\mathbb{Q}}(p/q) := \phi_{\mathbb{Z}}(p)_{\mathbb{N}}(q)$ for every irreducible fraction p/q . To convert a real number from $(-\infty, \infty)$ into a one-sided sequence, we assume that $\phi_{\mathbb{R}}(\mathbf{r}) = \theta$ satisfies $[1 + \exp(-\mathbf{r})] = \sum_{i=1}^{\infty} \theta_i D^{-i}$. This solves the problem of real arguments. A real-valued function $f : \mathbb{W} \rightarrow \mathbb{R}$ is called *enumerable* if there is a recursive function $g : \mathbb{W} \times \mathbb{N} \rightarrow \mathbb{Q}$ nondecreasing in k such that $\lim_k g(w, k) = f(w)$. A stronger condition, the f is called *recursive* if there is a recursive function $h : \mathbb{W} \times \mathbb{N} \rightarrow \mathbb{Q}$ such that $|f(w) - h(w, k)| < 1/k$.
- (viii) Pairs (w, u) enjoy the code $\phi_{\mathbb{W} \times \mathbb{U}}(w, u) := \phi_{\mathbb{W}}(w)\phi_{\mathbb{U}}(u)$. This code cannot be used if w is real. In the Proposition 2 of Section 3, where we need to string real vectors, Cantor's code is used instead.

- (ix) The concepts mentioned above are analogously extended to partial functions. Special care must be taken to assume computability of their domains, which is important to guarantee that the inverse of the Shannon-Fano-Elias code, used in Theorem 1, is recursive.

Last but not least, a *semimeasure* \mathbf{U} is a function $\mathbb{X}^* \rightarrow \mathbb{R}$ that satisfies $0 \leq \mathbf{U}(x)$, $\sum_a \mathbf{U}(xa) \mathbf{1}_{\{|a|=n\}} \leq \mathbf{U}(x)$, and $\mathbf{U}(\lambda) \leq 1$. Symbol $\overset{*}{<}$ denotes inequality up to a multiplicative constant.

Impossibility level

$$\mathcal{I}(\mathbf{x}; \mathbf{Y}) := \sup_{n \in \mathbb{N}} \frac{D^{-K(x^n)}}{\mathbf{Y}(x^n)} \quad (5)$$

is a natural measure of randomness deficiency for a sequence $\mathbf{x} \in \mathbb{X}^{\mathbb{N}}$ with respect to a recursive measure \mathbf{Y} , cf. [1], [6, Def. 4.5.10 and Thm. 4.5.5].³ The respective set of \mathbf{Y} -Martin-Löf random sequences

$$\mathcal{L}_{\mathbf{Y}} := \{\mathbf{x} : \mathcal{I}(\mathbf{x}; \mathbf{Y}) < \infty\} \quad (6)$$

has two important properties.

Firstly, $\mathcal{L}_{\mathbf{Y}}$ is the maximal set of sequences on which no enumerable semimeasure outperforms a recursive measure \mathbf{Y} more than by a multiplicative constant. Let \mathbf{M} be the universal enumerable semimeasure [6, Section 4.5.1]. By [2, Theorem 1 and Lemma 3], we have

$$\mathcal{I}(\mathbf{x}; \mathbf{Y}) \overset{*}{<} \liminf_{n \rightarrow \infty} \frac{\mathbf{M}(x^n)}{\mathbf{Y}(x^n)} \overset{*}{<} \sup_{n \in \mathbb{N}} \frac{\mathbf{M}(x^n)}{\mathbf{Y}(x^n)} \overset{*}{<} [\mathcal{I}(\mathbf{x}; \mathbf{Y})]^{1+\epsilon} \quad (7)$$

for a fixed $\epsilon > 0$ and recursive \mathbf{Y} . By the definition of \mathbf{M} , $\mathbf{U}(x^n) \overset{*}{<} \mathbf{M}(x^n)$ for any enumerable (semi)measure \mathbf{U} . Hence $\sup_{n \in \mathbb{N}} \mathbf{U}(x^n)/\mathbf{Y}(x^n) < \infty$ if $\mathbf{x} \in \mathcal{L}_{\mathbf{Y}}$. Moreover, $\mathcal{L}_{\mathbf{Y}} = \mathcal{L}_{\mathbf{U}}$ if \mathbf{Y} and \mathbf{U} are *mutually equivalent* recursive measures, i.e., $\sup_{n \in \mathbb{N}} \mathbf{U}(x^n)/\mathbf{Y}(x^n) < \infty \iff \sup_{n \in \mathbb{N}} \mathbf{Y}(x^n)/\mathbf{U}(x^n) < \infty$ for all $\mathbf{x} \in \mathbb{X}^{\mathbb{N}}$.

Secondly, the set $\mathcal{L}_{\mathbf{Y}}$ has full measure \mathbf{Y} . The fact is well-known, cf. e.g. [1, Remark 2], and it can be seen easily using the auxiliary statement below, which strengthens Barron's result [13, Theorem 3.1]. Whereas $\mathbf{Y}(\mathcal{L}_{\mathbf{Y}}) = 1$ follows for $|B(\cdot)| = K(\cdot)$, we shall use this lemma later also for $|B(\cdot)| = K(\cdot|\theta)$.

Lemma 1 (no hypercompression). *Let $B : \mathbb{X}^* \rightarrow \mathbb{Y}^+$ be a prefix code. Then*

$$|B(x^n)| + \log \mathbf{Y}(x^n) > 0 \quad (8)$$

n -eventually for \mathbf{Y} -almost all sequences \mathbf{x} .

Proof. Consider the function $W(x) := D^{-|B(x)|}$. By the Markov inequality,

$$\mathbf{Y}((8) \text{ is false}) = \mathbf{Y}\left(\frac{W(x^n)}{\mathbf{Y}(x^n)} \geq 1\right) \leq \mathbf{E}_{\mathbf{x} \sim \mathbf{Y}} \left[\frac{W(x^n)}{\mathbf{Y}(x^n)} \right] = \sum_x \mathbf{1}_{\{|x|=n\}} W(x).$$

³ Note added post scriptum: The version published in the proceedings contains, mistakenly, inf instead of sup in the formulae (5) and (17).

Hence $\sum_n \mathbf{Y}((8) \text{ is false}) \leq \sum_x D^{-|B(x)|} \leq 1 < \infty$ by the Kraft inequality. The claim now follows by the Borel-Cantelli lemma. \square

Now let \mathbf{Y} be a recursive Bayesian measure (2). In a prototypical case, measures \mathbf{P}_θ are not enumerable \mathbf{Q} -almost surely. But the data that are almost surely typical for these measures can be optimally compressed with the effectively computable measure \mathbf{Y} . That is, $\mathbf{P}_\theta(\mathcal{L}_\mathbf{Y}) = 1$ holds \mathbf{Q} -almost everywhere, as implied by the following statement.

Lemma 2 (cf. [14, Section 9]). *Equality $\mathbf{Y}(\mathcal{X}) = 1$ for $\mathbf{Y} = \int \mathbf{P}_\theta d\mathbf{Q}(\theta)$ implies $\mathbf{P}_\theta(\mathcal{X}) = 1$ for \mathbf{Q} -almost all θ .*

Proof. Let $\mathcal{G}_n := \{\theta \in \Theta : \mathbf{P}_\theta(\mathcal{X}) \geq 1 - 1/n\}$. We have $1 = \mathbf{Y}(\mathcal{X}) \leq \mathbf{Q}(\mathcal{G}_n) + \mathbf{Q}(\Theta \setminus \mathcal{G}_n)(1 - 1/n) = 1 - n^{-1}\mathbf{Q}(\Theta \setminus \mathcal{G}_n)$. Thus $\mathbf{Q}(\mathcal{G}_n) = 1$. By σ -additivity, $\mathbf{Q}(\mathcal{G}) = \inf_n \mathbf{Q}(\mathcal{G}_n) = 1$ follows for $\mathcal{G} := \{\theta \in \Theta : \mathbf{P}_\theta(\mathcal{X}) = 1\} = \bigcap_n \mathcal{G}_n$. \square

Notably, the Bayesian compressor can be shown optimal exactly when the parameter is incompressible. Strictly speaking, we will obtain $\mathbf{P}_\theta(\mathcal{L}_\mathbf{Y}) = 1$ if and only if θ is Martin-Löf random with respect to \mathbf{Q} . This holds, of course, under some tacit assumptions. For instance, if we take $\mathbf{P}_\theta \equiv \mathbf{Y}$ then $\mathbf{P}_\theta(\mathcal{L}_\mathbf{Y}) = 1$ for all $\theta \in \Theta$. We may thus suppose that the ‘if and only if’ statement holds provided the parameter can be effectively identified. The following two propositions form the first step to see what assumptions are needed exactly.

Lemma 3. *For a computer-dependent constant A , we have*

$$K(x|\theta) \leq A + K(x|\theta^m, K(\theta^m)) + K(K(\theta^m)) + K(m). \quad (9)$$

Proof. A certain program for computing x given θ operates as follows. It first calls a subroutine of length $K(m)$ to compute m and a subroutine of length $K(K(\theta^m))$ to compute $K(\theta^m)$. Then it reads the prefix θ^m of θ and passes θ^m and $K(\theta^m)$ to a subroutine of length $K(x|\theta^m, K(\theta^m))$ which returns x . \square

Theorem 1. *Let (\mathbf{P}, \mathbf{Q}) be a Bayesian statistical model with a recursive prior $\mathbf{Q} : \mathbb{Y}^* \rightarrow \mathbb{R}$ and a recursive kernel $\mathbf{P} : \mathbb{X}^* \times \Theta \rightarrow \mathbb{R}$.*

(i) *If (3) holds for \mathbf{P}_θ -almost all \mathbf{x} then*

$$K(x^n) + \log \mathbf{Y}(x^n) \geq K(\theta^m) + \log \mathbf{Q}(\theta^m) - 3 \log m + o(\log m) \quad (10)$$

is also true for \mathbf{P}_θ -almost all \mathbf{x} .

(ii) *If (4) holds for a recursive $\tau : \mathbb{Y}^* \rightarrow \mathbb{N}$ and $n = \tau(\theta^m)$ then*

$$K(x^n) + \log \mathbf{Y}(x^n) \leq K(\theta^m) + \log \mathbf{Q}(\theta^m) + O(1). \quad (11)$$

Proof. (i) For \mathbf{P}_θ -almost all \mathbf{x} we have both (3) and

$$K(x^n|\theta) + \log \mathbf{P}_\theta(x^n) \geq 0 \quad (12)$$

n -eventually, by Lemma 1 for $|B(\cdot)| = K(\cdot|\theta)$. Applying Lemma 3 to these sequences yields

$$\begin{aligned} K(x^n|\theta^m, K(\theta^m)) + \log \mathbf{T}(x^n, \theta^m) - \log \mathbf{Q}(\theta^m) \\ \geq -K(K(\theta^m)) - K(m) + o(\log m) = -2 \log m + o(\log m) \end{aligned}$$

because $K(\theta^m) \leq m + \log m + o(\log m)$ and $K(m) \leq \log m + o(\log m)$. Since

$$K(x^n|\theta^m, K(\theta^m)) + K(\theta^m) = K(x^n, \theta^m) + O(1) \quad (13)$$

by the chain rule for prefix complexity [6, Theorem 3.9.1], we obtain

$$K(x^n, \theta^m) + \log \mathbf{T}(x^n, \theta^m) \geq K(\theta^m) + \log \mathbf{Q}(\theta^m) - 2 \log m + o(\log m).$$

In the following, we apply (13) with x^n and θ^m switched, and observe that

$$K(\theta^m|x^n, K(x^n)) \leq A + K(m) - \log \frac{\mathbf{T}(x^n, \theta^m)}{\mathbf{Y}(x^n)}$$

follows by conditional Shannon-Fano-Elias coding of θ^m of an arbitrary length given x^n , cf. [15, Section 5.9]. Hence (10) holds for \mathbf{P}_θ -almost all \mathbf{x} .

(ii) By conditional Shannon-Fano-Elias coding of x^n given θ^m we obtain

$$K(x^n, \theta^m) \leq A' + K(\theta^m) - \log \frac{\mathbf{T}(x^n, \theta^m)}{\mathbf{Q}(\theta^m)}. \quad (14)$$

(This time, we need not specify the length of x^n separately since it can be computed from θ^m .) Substituting (4) into (14) and chaining the result with $K(x^n) \leq A'' + K(x^n, \theta^m)$ yields (11). \square

Theorem 1 applies to uniformly discretizable models if we plug in $m \geq \mu(n)$ and $\tau(\theta^m) \geq \nu(m)$. Hence we obtain the first, less elegant dichotomy.

Proposition 1. *Let (\mathbf{P}, \mathbf{Q}) be a μ -uniformly discretizable model with a recursive prior $\mathbf{Q} : \mathbb{Y}^* \rightarrow \mathbb{R}$ and a recursive kernel $\mathbf{P} : \mathbb{X}^* \times \Theta \rightarrow \mathbb{R}$. We have*

$$\mathbf{P}_\theta(\mathcal{L}_{\mathbf{Y}, \log \mu(n)}) = \begin{cases} 1 & \text{if } \theta \in \mathcal{L}_{\mathbf{Q}, \log n}, \\ 0 & \text{if } \theta \notin \mathcal{L}_{\mathbf{Q}, \log n}, \end{cases} \quad (15)$$

where the sets of $(\mathbf{Y}, g(n))$ -random sequences are defined as

$$\mathcal{L}_{\mathbf{Y}, g(n)} := \left\{ \mathbf{x} : \inf_{n \in \mathbb{N}} \frac{K(x^n) + \log \mathbf{Y}(x^n)}{g(n)} > -\infty \right\}. \quad (16)$$

In particular, $\mathcal{L}_{\mathbf{Y}, 1} = \mathcal{L}_{\mathbf{Y}}$.

Theorem 1(ii) suffices to prove $\mathbf{P}_\theta(\mathcal{L}_{\mathbf{Y}}) = 0$ for $\theta \notin \mathcal{L}_{\mathbf{Q}}$ but to show $\mathbf{P}_\theta(\mathcal{L}_{\mathbf{Y}}) = 1$ in the other case we need a stronger statement than Theorem 1(i). Here we can rely on the chain rule for conditional impossibility levels by Vovk and V'yugin [1, Theorem 1] and extensions of Lambalgen's theorem for

conditionally random sequences by Takahashi [4]. For a recursive kernel \mathbf{P} , let us define by analogy the conditional impossibility level

$$\mathcal{I}(\mathbf{x}; \mathbf{P}|\boldsymbol{\theta}) := \sup_{n \in \mathbb{N}} \frac{D^{-K(x^n|\boldsymbol{\theta})}}{\mathbf{P}_{\boldsymbol{\theta}}(x^n)} \quad (17)$$

and the set of conditionally random sequences

$$\mathcal{L}_{\mathbf{P}|\boldsymbol{\theta}} := \{\mathbf{x} \in \mathbb{X}^{\mathbb{N}} : \mathcal{I}(\mathbf{x}; \mathbf{P}|\boldsymbol{\theta}) < \infty\}. \quad (18)$$

We have $\mathbf{P}_{\boldsymbol{\theta}}(\mathcal{L}_{\mathbf{P}|\boldsymbol{\theta}}) = 1$ for all $\boldsymbol{\theta}$ by Lemma 1, as used in (12). Adjusting the proof of [6, Theorem 4.5.5] to computation with an oracle, we can show that the definition of $\mathcal{I}(\mathbf{x}; \mathbf{P}|\boldsymbol{\theta})$ given here is equivalent to the one given by [1], cf. [6, Def. 4.5.10]. Hence

$$\inf_{\boldsymbol{\theta} \in \Theta} [\mathcal{I}(\mathbf{x}; \mathbf{P}|\boldsymbol{\theta}) \mathcal{I}(\boldsymbol{\theta}; \mathbf{Q})] \stackrel{*}{<} \mathcal{I}(\mathbf{x}; \mathbf{Y}) \stackrel{*}{<} \inf_{\boldsymbol{\theta} \in \Theta} [\mathcal{I}(\mathbf{x}; \mathbf{P}|\boldsymbol{\theta}) [\mathcal{I}(\boldsymbol{\theta}; \mathbf{Q})]^{1+\epsilon}] \quad (19)$$

holds for $\mathbf{Y} = \int \mathbf{P}_{\boldsymbol{\theta}} d\mathbf{Q}(\boldsymbol{\theta})$ and $\epsilon > 0$ by [1, Corollary 4].

Inequality (19) and Theorem 1(ii) imply the main claim of this article.

Theorem 2. *Let (\mathbf{P}, \mathbf{Q}) be a Bayesian statistical model with a recursive prior $\mathbf{Q} : \mathbb{Y}^* \rightarrow \mathbb{R}$ and a recursive kernel $\mathbf{P} : \mathbb{X}^* \times \Theta \rightarrow \mathbb{R}$. Suppose that (4) holds for all $\boldsymbol{\theta} \in \Theta$, $\mathbf{P}_{\boldsymbol{\theta}}$ -almost all \mathbf{x} , and $n = \tau(\theta^m)$, where $\tau : \mathbb{Y}^* \rightarrow \mathbb{N}$ is recursive. Then we have*

$$\mathbf{P}_{\boldsymbol{\theta}}(\mathcal{L}_{\mathbf{Y}}) = \begin{cases} 1 & \text{if } \boldsymbol{\theta} \in \mathcal{L}_{\mathbf{Q}}, \\ 0 & \text{if } \boldsymbol{\theta} \notin \mathcal{L}_{\mathbf{Q}}. \end{cases} \quad (20)$$

The upper part of (20) can be strengthened as decomposition $\mathcal{L}_{\mathbf{Y}} = \bigcup_{\boldsymbol{\theta} \in \mathcal{L}_{\mathbf{Q}}} \mathcal{L}_{\mathbf{P}|\boldsymbol{\theta}}$, which holds for all recursive \mathbf{P} and \mathbf{Q} [4, Cor. 4.3 & Thm. 5.3]. (Our definition of a recursive \mathbf{P} corresponds to ‘uniformly computable’ in [4].) We suppose that, under the assumption of Theorem 2, sets $\mathcal{L}_{\mathbf{P}|\boldsymbol{\theta}}$ are disjoint for $\boldsymbol{\theta} \in \Theta$. This would strengthen the lower part of (20).

3 The case of exponential families

As shown in [16], k -parameter exponential families exhibit Shannon redundancy $-\log \mathbf{Y}(x^n) + \log \mathbf{P}_{\boldsymbol{\theta}}(x^n) = \frac{k}{2} \log n + \Theta(\log \log n)$. Here we shall prove that these models are uniformly discretizable with $\mu(n) = (\frac{k}{2} + \epsilon) \log n$ respectively. The result is established under a familiar condition. Namely, a prior $\tilde{\mathbf{Q}}$ on $\tilde{\Theta} \subset \mathbb{R}^k$ is *universally lower-bounded* by the Lebesgue measure $\boldsymbol{\lambda}$ if for each $\boldsymbol{\vartheta} \in \tilde{\Theta}$ there exists an open set $C \ni \boldsymbol{\vartheta}$ and a $w > 0$ such that $\tilde{\mathbf{Q}}(E) \geq w\boldsymbol{\lambda}(E)$ for every measurable $E \subset C$. This condition implies that $\tilde{\Theta}$ is the support of $\tilde{\mathbf{Q}}$ and is satisfied, in particular, if $\tilde{\mathbf{Q}}$ and $\boldsymbol{\lambda}$ restricted to $\tilde{\Theta}$ are mutually equivalent.

Let us write the components of vectors $\boldsymbol{\vartheta}, \boldsymbol{\vartheta}' \in \mathbb{R}^k$ as $\boldsymbol{\vartheta} = (\vartheta_1, \vartheta_2, \dots, \vartheta_k)$ and their Euclidean distance as $|\boldsymbol{\vartheta}' - \boldsymbol{\vartheta}| := \sqrt{\sum_{l=1}^k (\vartheta'_l - \vartheta_l)^2}$.

Example 1 (an exponential family). Let the kernel $\tilde{\mathbf{P}} : \mathbb{X}^* \times \tilde{\Theta} \ni (x, \vartheta) \mapsto \tilde{\mathbf{P}}_{\vartheta}(x) \in \mathbb{R}$ represent a regular k -parameter exponential family. That is:

- (i) Certain functions $p : \mathbb{X} \rightarrow (0, \infty)$ and $T : \mathbb{X} \rightarrow \mathbb{R}^k$ satisfy $\sum_{x \in \mathbb{X}} p(x) < \infty$ and $\forall \beta \in \mathbb{R}^k \setminus \mathbf{0} \forall c \in \mathbb{R} \exists x \in \mathbb{X} \sum_{l=1}^k \beta_l T_l(x) \neq c$ (i.e., T has affinely independent components).
- (ii) Let $Z(\beta) := \sum_{x \in \mathbb{X}} p(x) \exp\left(\sum_{l=1}^k \beta_l T_l(x)\right)$ and define measures

$$\tilde{\mathbf{P}}_{\beta}(x^n) := \prod_{i=1}^n p(x_i) \exp\left(\sum_{l=1}^k \beta_l T_l(x) - \ln Z(\beta)\right)$$

for $\beta \in \mathbf{B} := \{\beta' \in \mathbb{R}^k : Z(\beta') < \infty\}$.

- (iii) We require that \mathbf{B} is open. (It is not empty since $\mathbf{0} \in \mathbf{B}$.) Under this condition, $\vartheta(\cdot) : \mathbf{B} \ni \beta \mapsto \vartheta(\beta) := \mathbf{E}_{x \sim \tilde{\mathbf{P}}_{\beta}} T(x_i) \in \mathbb{R}^k$ is a twice differentiable injection [17], [9]. Thus assume $\tilde{\Theta} := \vartheta(\mathbf{B})$ and put $\tilde{\mathbf{P}}_{\vartheta} := \tilde{\mathbf{P}}_{\beta(\vartheta)}$ for $\beta(\cdot) := \vartheta^{-1}(\cdot)$.

Additionally, let the prior $\tilde{\mathbf{Q}}$ be universally lower-bounded by the Lebesgue measure on \mathbb{R}^k and let it satisfy $\tilde{\mathbf{Q}}(\tilde{\Theta}) = 1$.

Proposition 2. Use Cantor's code $\rho := \rho_s \circ \rho_n$, where $\rho_n : \tilde{\Theta} \rightarrow (0, 1)^k$ is a differentiable injection and $\rho_s : (0, 1)^k \rightarrow \mathbb{Y}^{\mathbb{N}}$ satisfies $\rho_s(\mathbf{y}) = \theta_1 \theta_2 \theta_3 \dots$ for any vector $\mathbf{y} \in (0, 1)^k$ with components $\mathbf{y}_l = \sum_{i=1}^{\infty} \theta_{(i-1)k+l} D^{-i}$. Then the model $(\tilde{\mathbf{P}}, \tilde{\mathbf{Q}})$ is $(\rho, (\frac{k}{2} + \epsilon) \log n, D^{(2/k+\epsilon)m})$ -uniformly discretizable for $\epsilon > 0$.

Proof. Let $\Theta := \rho(\tilde{\Theta})$, $\mathbf{P}_{\theta}(x) := \tilde{\mathbf{P}}_{\rho^{-1}(\theta)}(x)$, $\mathbf{Q} := \tilde{\mathbf{Q}} \circ \rho^{-1}$, and $A(\theta) := \{\theta' \in \Theta : \theta \text{ is the prefix of } \theta'\}$. Consider a $\theta \in \Theta$. Firstly, let $m \geq (\frac{k}{2} + \epsilon) \log n$. We have (21) for $\vartheta = \rho^{-1}(\theta)$ and $A_n = \rho^{-1}(A(\theta^m))$. Hence (3) holds by the Theorem 3(i) below. Secondly, let $n \geq D^{(2/k+\epsilon)m}$. We have (23) for $\vartheta = \rho^{-1}(\theta)$ and $B_n = \rho^{-1}(A(\theta^m))$. Hence (4) follows by Theorem 3(ii). \square

The statement below may look more familiar for statisticians.

Theorem 3. Fix a $\vartheta \in \tilde{\Theta}$ for the model specified in Example 1.

- (i) If we take sufficiently small measurable sets $A_n \subset \tilde{\Theta}$ which satisfy

$$\limsup_{n \rightarrow \infty} \frac{\sup_{\vartheta' \in A_n} |\vartheta' - \vartheta|}{\sqrt{n^{-1} \ln \ln n}} = 0 \quad (21)$$

and put $\tilde{\mathbf{P}}_n(x) := \int_{A_n} \tilde{\mathbf{P}}_{\vartheta'}(x) d\tilde{\mathbf{Q}}(\vartheta') / \int_{A_n} d\tilde{\mathbf{Q}}(\vartheta')$ then

$$\lim_{n \rightarrow \infty} \frac{\log \tilde{\mathbf{P}}_n(x^n) - \log \tilde{\mathbf{P}}_{\vartheta}(x^n)}{\ln \ln n} = 0 \quad (22)$$

for $\tilde{\mathbf{P}}_{\vartheta}$ -almost all x .

(ii) On the other hand, if we take sufficiently large measurable sets

$$B_n \supset \left\{ \boldsymbol{\vartheta}' \in \tilde{\Theta} : |\boldsymbol{\vartheta}' - \boldsymbol{\vartheta}| \geq n^{-1/2+\alpha} \right\} \quad (23)$$

for an arbitrary $\alpha \in (0, 1/2)$ then

$$\lim_{n \rightarrow \infty} \left(\log \int \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) d\tilde{\mathbf{Q}}(\boldsymbol{\vartheta}') - \log \int_{B_n} \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) d\tilde{\mathbf{Q}}(\boldsymbol{\vartheta}') \right) = 0 \quad (24)$$

for $\tilde{\mathbf{P}}_{\boldsymbol{\vartheta}}$ -almost all \boldsymbol{x} .

Proof. (i) Function $\hat{\boldsymbol{\vartheta}}(x^n) := n^{-1} \sum_{i=1}^n T(x_i)$ is the maximum likelihood estimator of $\boldsymbol{\vartheta}$, in the usual sense. Thus the Taylor expansion for any $\boldsymbol{\vartheta} \in \tilde{\Theta}$ yields

$$\log \tilde{\mathbf{P}}_{\hat{\boldsymbol{\vartheta}}(x^n)}(x^n) - \log \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}}(x^n) = n \sum_{l,m=1}^k R_{lm}(\boldsymbol{\vartheta}) S_{lm}(\boldsymbol{\vartheta}), \quad (25)$$

where $S_{lm}(\boldsymbol{\vartheta}) := (\boldsymbol{\vartheta}_l - \hat{\boldsymbol{\vartheta}}_l(x^n))(\boldsymbol{\vartheta}_m - \hat{\boldsymbol{\vartheta}}_m(x^n))$ and $R_{lm}(\boldsymbol{\vartheta}) := \int_0^1 (1-t) I_{lm}(t\boldsymbol{\vartheta} + (1-t)\hat{\boldsymbol{\vartheta}}(x^n)) dt$, whereas the observed Fisher information matrix $I_{lm}(\boldsymbol{\vartheta}) := -n^{-1} \partial_{\boldsymbol{\vartheta}_l} \partial_{\boldsymbol{\vartheta}_m} \log \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}}(x^n)$ does not depend on n and x^n . Consequently,

$$\begin{aligned} \log \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}}(x^n) - \log \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) = \\ n \sum_{l,m=1}^k [R_{lm}(\boldsymbol{\vartheta}') [S_{lm}(\boldsymbol{\vartheta}') - S_{lm}(\boldsymbol{\vartheta})] + [R_{lm}(\boldsymbol{\vartheta}') - R_{lm}(\boldsymbol{\vartheta})] S_{lm}(\boldsymbol{\vartheta})]. \end{aligned}$$

With C_n denote the intersection of $\tilde{\Theta}$ and the smallest ball containing A_n and $\hat{\boldsymbol{\vartheta}}(x^n)$. Let $d_n := |\boldsymbol{\vartheta} - \hat{\boldsymbol{\vartheta}}(x^n)|$ and $a_n := \sup_{\boldsymbol{\vartheta}' \in A_n} |\boldsymbol{\vartheta}' - \boldsymbol{\vartheta}|$. Hence we bound

$$\left| \log \tilde{\mathbf{P}}_n(x^n) - \log \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}}(x^n) \right| \leq n \sum_{l,m=1}^k [|R_{lm}^+| a_n (2d_n + a_n) + |R_{lm}^+ - R_{lm}^-| d_n^2],$$

where $R_{lm}^+ := \sup_{\boldsymbol{\vartheta}' \in C_n} R_{lm}(\boldsymbol{\vartheta}')$ and $R_{lm}^- := \inf_{\boldsymbol{\vartheta}' \in C_n} R_{lm}(\boldsymbol{\vartheta}')$. By continuity of Fisher information $I_{lm}(\boldsymbol{\vartheta})$ as a function of $\boldsymbol{\vartheta}$, R_{lm}^+ and R_{lm}^- tend to $I_{lm}(\boldsymbol{\vartheta})$ for $n \rightarrow \infty$. On the other hand, the law of iterated logarithm

$$\limsup_{n \rightarrow \infty} \frac{\hat{\boldsymbol{\vartheta}}_l(x^n) - \boldsymbol{\vartheta}_l}{\sigma_l \sqrt{2n^{-1} \ln \ln n}} = 1 \quad (26)$$

is satisfied for $\tilde{\mathbf{P}}_{\boldsymbol{\vartheta}}$ -almost all \boldsymbol{x} with variance $\sigma_l^2 := \text{Var}_{\boldsymbol{x} \sim \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}}} T_l(x_i)$ since the maximum likelihood estimator is unbiased, i.e., $\mathbf{E}_{\boldsymbol{x} \sim \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}}} \hat{\boldsymbol{\vartheta}}(x^n) = \boldsymbol{\vartheta}$. Consequently, we obtain (22) for (21).

(ii) The proof applies Laplace approximation as in [18] or in the proof of Theorem 8.1 of [9, pages 248–251]. First of all, we have

$$\log \int \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) d\tilde{\mathbf{Q}}(\boldsymbol{\vartheta}') - \log \int_{B_n} \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) d\tilde{\mathbf{Q}}(\boldsymbol{\vartheta}') \leq \frac{\int_{\tilde{\Theta} \setminus B_n} \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) d\tilde{\mathbf{Q}}(\boldsymbol{\vartheta}')}{\int_{B_n} \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) d\tilde{\mathbf{Q}}(\boldsymbol{\vartheta}')}.$$

In the following, we consider a sufficiently large n . Because of the law of iterated logarithm (26), $\hat{\boldsymbol{\vartheta}}(x^n)$ belongs to B_n for $\tilde{\mathbf{P}}_{\boldsymbol{\vartheta}}$ -almost all \mathbf{x} . Hence the robustness property and the convexity of Kullback-Leibler divergence for exponential families [9, Eq. (19.12) and Proposition 19.2] imply a bound for the numerator

$$\int_{\tilde{\Theta} \setminus B_n} \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) d\tilde{\mathbf{Q}}(\boldsymbol{\vartheta}') \leq \sup_{\boldsymbol{\vartheta}' \in \tilde{\Theta} \setminus B_n} \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) \leq \sup_{\boldsymbol{\vartheta}' \in \partial B_n} \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n),$$

where ∂B_n is the boundary of B_n . Using (25) gives further

$$\sup_{\boldsymbol{\vartheta}' \in \partial B_n} \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) \leq \tilde{\mathbf{P}}_{\hat{\boldsymbol{\vartheta}}(x^n)}(x^n) \exp[-nR^- \delta^2]$$

with $R^- := \inf_{\boldsymbol{\vartheta}' \in B_n} \left[\sum_{l=1}^k R_{lm}(\boldsymbol{\vartheta}') S_{lm}(\boldsymbol{\vartheta}') \right] / |\boldsymbol{\vartheta}' - \hat{\boldsymbol{\vartheta}}(x^n)|^2$ and $\delta := \inf_{\boldsymbol{\vartheta}' \in \partial B_n} |\boldsymbol{\vartheta}' - \hat{\boldsymbol{\vartheta}}(x^n)|$. Since the prior is universally lower-bounded by the Lebesgue measure, then (25) implies a bound for the denominator

$$\int_{B_n} \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) d\tilde{\mathbf{Q}}(\boldsymbol{\vartheta}') \geq w \tilde{\mathbf{P}}_{\hat{\boldsymbol{\vartheta}}(x^n)}(x^n) \int_{|t| < \delta} \exp[-nR^+ |t|^2] dt,$$

where $w > 0$ and $R^+ := \sup_{\boldsymbol{\vartheta}' \in B_n} \left[\sum_{l=1}^k R_{lm}(\boldsymbol{\vartheta}') S_{lm}(\boldsymbol{\vartheta}') \right] / |\boldsymbol{\vartheta}' - \hat{\boldsymbol{\vartheta}}(x^n)|^2$. Hence we obtain an inequality for the ratio

$$\frac{\int_{\tilde{\Theta} \setminus B_n} \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) d\tilde{\mathbf{Q}}(\boldsymbol{\vartheta}')}{\int_{B_n} \tilde{\mathbf{P}}_{\boldsymbol{\vartheta}'}(x^n) d\tilde{\mathbf{Q}}(\boldsymbol{\vartheta}')} \leq \frac{\sqrt{nR^+} \exp[-nR^- \delta^2/2]}{w \int_{|t| < \delta \sqrt{nR^+}} \exp[-|t|^2] dt}.$$

The right-hand side tends to zero with $n \rightarrow \infty$ since $\delta = \Omega(n^{-1/2+\alpha})$ whereas R^+ and R^- tend to strictly positive constants by continuity and strictly positive definiteness of the Fisher information matrix. \square

4 Less standard examples

In this section we shall present less standard examples of statistical models. We begin with two very simple models.

Example 2 (the data are the parameter). Put $\mathbf{P}_{\boldsymbol{\theta}}(x^n) := \mathbf{1}_{\{x^n = \boldsymbol{\theta}^n\}}$ for $\mathbb{X} = \mathbb{Y}$ and let $\mathbf{Q}(\boldsymbol{\theta}) > 0$ for $\boldsymbol{\theta} \in \mathbb{Y}^*$. This model is (n, m) -uniformly discretizable.

Example 3 (a singleton model). Each parameter $\boldsymbol{\theta}$ is random with respect to the prior \mathbf{Q} concentrated on this parameter, $\Theta = \{\boldsymbol{\theta}\}$. The respective singleton model (\mathbf{P}, \mathbf{Q}) is $(0, 0)$ -uniformly discretizable.

Now, a slightly more complex instance. Consider a class of stationary processes $(X_i)_{i \in \mathbb{Z}}$ of form $X_i := (K_i, \theta_{K_i})$, where the variables K_i are independent and distributed according to the hyperbolic distribution

$$P(K_i = k) = p(k) := \frac{k^{-1/\beta}}{\zeta(1/\beta)}, \quad k \in \mathbb{N}, \quad (27)$$

with a fixed $\beta \in (0, 1)$. This family of processes was introduced to model logical consistency of texts in natural language [19]. The distribution of variables X_i is equal to the measure $P(X_i \in \cdot) = \mathbf{P}_{\boldsymbol{\theta}}$ for the following Bayesian model.

Example 4 (an accessible description model). Put

$$\mathbf{P}_\theta(x^n) := \prod_{i=1}^n p(k_i) \mathbf{1}_{\{z_i = \theta_{k_i}\}} \quad (28)$$

for $x_i = (k_i, z_i) \in \mathbb{N} \times \mathbb{Y}$ and let $\mathbf{Q}(\theta) > 0$ for $\theta \in \mathbb{Y}^*$.

For this model, Shannon information between the data and the parameter equals $\mathbf{E}_{(\mathbf{x}, \theta) \sim \mathbf{T}} [-\log \mathbf{Y}(x^n) + \log \mathbf{P}_\theta(x^n)] = \Theta(n^\beta)$ asymptotically if $\mathbf{Q}(\theta) = D^{-|\theta|}$, cf. [19, Theorem 10]. As a consequence of the next statement, the accessible description model (28) is $(n^\nu, m^{1/\lambda})$ -uniformly discretizable for

$$\nu > 2\beta/(1 - \beta) \text{ and } \lambda < \beta.$$

Proposition 3. *For independent variables $(K_i)_{i \in \mathbb{Z}}$ with the distribution (27),*

$$\{K_1, K_2, \dots, K_n\} \setminus \{1, 2, \dots, \lceil n^\nu \rceil\} = \emptyset, \quad (29)$$

$$\{1, 2, \dots, \lfloor n^\lambda \rfloor\} \setminus \{K_1, K_2, \dots, K_n\} = \emptyset, \quad (30)$$

n-eventually almost surely.

Proof. To establish the first claim, put $U_n := \lceil n^\nu \rceil$ and observe

$$\begin{aligned} P(\{K_1, K_2, \dots, K_n\} \setminus \{1, 2, \dots, U_n\} \neq \emptyset) &\leq \sum_{j=U_n+1}^{\infty} P(j \in \{K_1, K_2, \dots, K_n\}) \\ &= \sum_{j=U_n+1}^{\infty} 1 - (1 - p(j))^n \leq \sum_{j=U_n+1}^{\infty} np(j) \\ &\leq \frac{n}{\zeta(1/\beta)} \int_{U_n}^{\infty} k^{-1/\beta} dk = \frac{n}{\zeta(1/\beta)} \frac{U_n^{1-1/\beta}}{1/\beta - 1} \leq \frac{n^{-1-\epsilon}}{\zeta(1/\beta)(1/\beta - 1)} \text{ for an } \epsilon > 0. \end{aligned}$$

Hence $\sum_{n=1}^{\infty} P(\{K_1, K_2, \dots, K_n\} \setminus \{1, 2, \dots, U_n\} \neq \emptyset) < \infty$ so (29) holds by the Borel-Cantelli lemma. As for the second claim, put $L_n := \lfloor n^\lambda \rfloor$ and observe

$$\begin{aligned} P(\{1, 2, \dots, L_n\} \setminus \{K_1, K_2, \dots, K_n\} \neq \emptyset) &\leq \sum_{j=1}^{L_n} P(j \notin \{K_1, K_2, \dots, K_n\}) \\ &= \sum_{j=1}^{L_n} (1 - p(j))^n \leq L_n (1 - p(L_n))^n = L_n \exp[n \log(1 - p(L_n))] \\ &\leq L_n \exp[-np(L_n)] \leq n^\beta \exp[-n^\epsilon] \text{ for an } \epsilon > 0. \end{aligned}$$

Hence $\sum_{n=1}^{\infty} P(\{1, 2, \dots, L_n\} \setminus \{K_1, K_2, \dots, K_n\} \neq \emptyset) < \infty$ so (30) is also satisfied by the Borel-Cantelli lemma. \square

To use the above statement for the Bayesian model, notice first that $\mathbf{P}_\theta(x^n) > 0$ for \mathbf{P}_θ -almost all \mathbf{x} . Hence equalities $z_i = \theta_{k_i}$ and

$$\begin{aligned} \mathbf{T}(x^n, \theta^m) &= \sum_{y^M \in \mathbb{Y}^M} \left(\prod_{i=1}^n p(k_i) \mathbf{1}_{\{z_i = y_{k_i}\}} \right) \left(\prod_{k=1}^m \mathbf{1}_{\{\theta_k = y_k\}} \right) \mathbf{Q}(y^M) \\ &= \mathbf{P}_\theta(x^n) \sum_{y^M \in \mathbb{Y}^M} \left(\prod_{k \in \{k_1, k_2, \dots, k_n\} \cup \{1, 2, \dots, m\}} \mathbf{1}_{\{\theta_k = y_k\}} \right) \mathbf{Q}(y^M) \end{aligned}$$

hold for \mathbf{P}_θ -almost all \mathbf{x} with $M := \max\{m, k_1, k_2, \dots, k_n\}$. Consequently,

$$\mathbf{Q}(\theta^m)\mathbf{P}_\theta(x^n) = \mathbf{T}(x^n, \theta^m) \quad \text{if} \quad \{k_1, k_2, \dots, k_n\} \setminus \{1, 2, \dots, m\} = \emptyset, \quad (31)$$

$$\mathbf{T}(x^n, \theta^m) = \mathbf{Y}(x^n) \quad \text{if} \quad \{1, 2, \dots, m\} \setminus \{k_1, k_2, \dots, k_n\} = \emptyset. \quad (32)$$

Thus the model given in Example 4 is $(n^\nu, m^{1/\lambda})$ -uniformly discretizable.

The last example is not uniformly discretizable. It stems from the observation that any probability measure on \mathbb{X}^∞ can be encoded with a single sequence from \mathbb{Y}^∞ . Such parameter is not identifiable, however.

Example 5 (a model that contains all distributions). For simplicity let $\mathbb{X} = \mathbb{N}$ and $\mathbb{Y} = \{0, 1\}$. The link between θ and \mathbf{P}_θ will be established by imposing equalities $\mathbf{P}_\theta(\lambda) = 1$ and

$$\mathbf{P}_\theta(x^n) = \left(\mathbf{P}_\theta(x^{n-1}) - \sum_{y < x_n} \mathbf{P}_\theta(x^{n-1}y) \right) \cdot \sum_{k=1}^{\infty} \theta_{\phi(x^n, k)} 2^{-k}, \quad (33)$$

where a recursive bijection $\phi : \mathbb{N}^+ \times \mathbb{N} \rightarrow \mathbb{N}$ is used. It is easy to see that \mathbf{P}_θ is a probability measure on \mathbb{X}^∞ for each θ . Conversely, each probability measure on \mathbb{X}^∞ equals \mathbf{P}_θ for at least one θ .

Let the prior be the uniform measure $\mathbf{Q}(\theta) := 2^{-|\theta|}$. Then the Bayesian measure $\mathbf{Y} = \int \mathbf{P}_\theta(x) d\mathbf{Q}(\theta)$ is recursive and equals

$$\mathbf{Y}(x^n) = \frac{1}{2} \left(\mathbf{Y}(x^{n-1}) - \sum_{y < x_n} \mathbf{Y}(x^{n-1}y) \right) \implies \log_2 \mathbf{Y}(x^n) = -\sum_{i=1}^n x_i.$$

Measure \mathbf{Y} is not only optimal for all \mathbf{Q} -random θ , in the sense of $\mathbf{P}_\theta(\mathcal{L}_\mathbf{Y}) = 1$, but it is also optimal for a certain $\theta \notin \mathcal{L}_\mathbf{Q}$ that satisfies $\mathbf{P}_\theta = \mathbf{Y}$. On the other hand, by the asymptotic equipartition property, $\mathbf{P}_\theta(\mathcal{L}_\mathbf{Y}) = 0$ for stationary measures \mathbf{P}_θ that have a different entropy rate than \mathbf{Y} [15, Section 15.7].

5 Countable unions of models

Bayesian mixtures of uniformly discretizable models are uniformly discretizable under the additional condition (34), which says that Bayesian model selection is consistent for each $\theta \in \Theta$. Let us write $\theta_k^m := \theta_k \theta_{k+1} \dots \theta_m$. Moreover, define \mathbf{T}^i and \mathbf{Y}^i via (1)–(2) for models $(\mathbf{P}^i, \mathbf{Q}^i)$ substituted for (\mathbf{P}, \mathbf{Q}) respectively.

Theorem 4. *Let models $(\mathbf{P}^i, \mathbf{Q}^i)$ be (μ_i, ν_i) -uniformly discretizable with kernels $\mathbf{P}_\theta^i(x)$ for $\theta \in \Theta^i$ and $i \in A$, a countable set. For a prefix code $c : A \rightarrow \mathbb{Y}^+$, put $\Theta := \bigcup_{i \in A} c(i)\Theta^i$. Consecutively, denote $\text{idx}(\theta) := i$ and $\text{trn}(\theta) := \vartheta$ for $\theta = c(i)\vartheta \in \Theta$. Define the kernel $\mathbf{P}_\theta(x) := \mathbf{P}_{\text{trn}(\theta)}^{\text{idx}(\theta)}(x)$ for $\theta \in \Theta$ and the prior $\mathbf{Q} := \sum_{i \in A} w(i)(\mathbf{Q}^i \circ \text{trn})$ for $\sum_{i \in A} w(i) = 1$ and $w(i) > 0$. The model (\mathbf{P}, \mathbf{Q}) is (μ, ν) -uniformly discretizable provided*

$$\begin{aligned} \mu(n) &:= \sup_{i \in A} (|c(i)| + \mu_i(n)) < \infty, \\ \nu(m) &:= \sup_{i \in A} \nu_i(m - |c(i)|) < \infty, \end{aligned}$$

and

$$\lim_{n \rightarrow \infty} \mathbf{Y}(x^n) / \mathbf{Y}^i(x^n) = w(i) \quad (34)$$

for $i = \text{idx}(\boldsymbol{\theta})$, $\mathbf{P}_{\boldsymbol{\theta}}$ -almost all \mathbf{x} , and all $\boldsymbol{\theta} \in \Theta$.

Remark: Assuming recursive models and mutually singular $\mathbf{P}_{\boldsymbol{\vartheta}}^i$, convergence (34) may fail only for $\boldsymbol{\theta}$ that are not \mathbf{Q} -random, cf. [20]. Put $\mathcal{X} := \{\mathbf{x} : \lim_n \mathbf{Y}(x^n) / \mathbf{Y}^i(x^n) = w(i)\}$. By the ordinary martingale convergence, $\mathbf{Y}^i(\mathcal{X}) = 1$, whereas by convergence of recursive martingales [4, Theorem 3.1], $\mathcal{X} \supset \mathcal{L}_{\mathbf{Y}^i}$. Next, by [4, Cor. 4.3 & Thm. 5.3], we obtain $\mathcal{L}_{\mathbf{Y}^i} \supset \mathcal{L}_{\mathbf{P}^i|\boldsymbol{\vartheta}}$ for \mathbf{Q}^i -random $\boldsymbol{\vartheta}$. Hence $\mathbf{P}_{\boldsymbol{\theta}}(\mathcal{X}) = 1$ if $\boldsymbol{\theta} \in \mathcal{L}_{\mathbf{Q}}$ and (35) holds true, in view of the Theorem 5 below.

Proof. Let $i = \text{idx}(\boldsymbol{\theta})$. Observe that $\mathbf{T}(x^n, \theta^m) = w(i) \mathbf{T}^i(x^n, \theta_{|c(i)|+1}^m)$ and $\mathbf{Q}(\theta^m) = w(i) \mathbf{Q}^i(\theta_{|c(i)|+1}^m)$ if $m \geq |c(i)|$. Hence for $\mathbf{P}_{\boldsymbol{\theta}}$ -almost all \mathbf{x} and $m \geq \mu(n)$, we have

$$\left| \log \frac{\mathbf{Q}(\theta^m) \mathbf{P}_{\boldsymbol{\theta}}(x^n)}{\mathbf{T}(x^n, \theta^m)} \right| = \left| \log \frac{\mathbf{Q}^i(\theta_{|c(i)|+1}^m) \mathbf{P}_{\text{trn}(\boldsymbol{\theta})}^i(x^n)}{\mathbf{T}^i(x^n, \theta_{|c(i)|+1}^m)} \right| = o(\log m).$$

On the other hand, for $\mathbf{P}_{\boldsymbol{\theta}}$ -almost all \mathbf{x} and $n \geq \nu(m)$,

$$\lim_{m \rightarrow \infty} \frac{\mathbf{T}(x^n, \theta^m)}{\mathbf{Y}(x^n)} = \lim_{m \rightarrow \infty} \left[\frac{w^i \mathbf{Y}^i(x^n)}{\mathbf{Y}(x^n)} \cdot \frac{\mathbf{T}^i(x^n, \theta_{|c(i)|+1}^m)}{\mathbf{Y}^i(x^n)} \right] = 1.$$

□

A complementary result says that the set of random parameters with respect to the mixture is the union of the respective sets for the combined models.

Theorem 5. *Consider the models from Theorem 4 and suppose that \mathbf{Q}^i satisfy*

$$\mathbf{Q}^i(\theta^k) / \mathbf{Q}^i(\theta^m) \geq a c^{k-m} \quad (35)$$

for all $k \geq m \geq 0$ and certain constants $c < 1$ and $a > 0$. Then for $g(n) = \Omega(1)$ we have $\boldsymbol{\theta} \in \mathcal{L}_{\mathbf{Q}, g(n)}$ if and only if $\text{trn}(\boldsymbol{\theta}) \in \mathcal{L}_{\mathbf{Q}^{\text{idx}(\boldsymbol{\theta}), g(n)}}$.

Proof. Let $i = \text{idx}(\boldsymbol{\theta})$. The claim is true if

$$\left| K(\theta^m) + \log \mathbf{Q}(\theta^m) - K(\theta_{|c(i)|+1}^{|c(i)|+m}) - \log \mathbf{Q}^i(\theta_{|c(i)|+1}^{|c(i)|+m}) \right| = O(1)$$

for $m \geq |c(i)|$. The latter condition is satisfied since $\left| K(\theta^m) - K(\theta_{|c(i)|+1}^{|c(i)|+m}) \right| \leq |c(i)| + O(1)$, whereas $\left| \log \mathbf{Q}(\theta^m) - \log \mathbf{Q}^i(\theta_{|c(i)|+1}^{|c(i)|+m}) \right| \leq |\log w(i)| + O(|c(i)|)$ by $\mathbf{Q}(\theta^m) = w(i) \mathbf{Q}^i(\theta_{|c(i)|+1}^m)$ and (35). □

These propositions may be useful when we seek a compressor that is optimal for all random and certain nonrandom parameters with respect to a given prior. A possible solution is to find priors against which the originally considered nonrandom parameters are random. Suppose that these priors and the original prior yield uniformly discretizable models and consistent Bayesian selection among these models is feasible. Then Theorems 2, 4, and 5 guarantee that the Bayesian mixture of all considered models achieves the best enumerable compression for all requested parameters and no so many others!

6 Acknowledgements

I would like to thank P. Grünwald, P. Harremoës, and J. Mielniczuk for discussions. Cordial acknowledgements are due to an anonymous referee for suggesting relevant references. They helped to improve this paper considerably. The research, supported under the PASCAL II Network of Excellence, IST-2002-506778, was done during the author's leave from the Institute of Computer Science, Polish Academy of Sciences.

References

1. Vovk, V.G., V'yugin, V.V.: On the empirical validity of the Bayesian method. *J. Roy. Statist. Soc. B* **55** (1993) 253–266
2. Vovk, V.G., V'yugin, V.V.: Prequential level of impossibility with some applications. *J. Roy. Statist. Soc. B* **56** (1994) 115–123
3. Vitányi, P., Li, M.: Minimum description length induction, Bayesianism and Kolmogorov complexity. *IEEE Trans. Inform. Theor.* **46** (2000) 446–464
4. Takahashi, H.: On a definition of random sequences with respect to conditional probability. *Inform. Comput.* **206** (2008) 1375–1382
5. Gács, P.: On the symmetry of algorithmic information. *Dokl. Akad. Nauk SSSR* **15** (1974) 1477–1480
6. Li, M., Vitányi, P.M.B.: *An Introduction to Kolmogorov Complexity and Its Applications*, 2nd ed. Springer (1997)
7. van Lambalgen, M.: *Random Sequences*. PhD thesis, Universiteit van Amsterdam (1987)
8. Barron, A., Rissanen, J., Yu, B.: The minimum description length principle in coding and modeling. *IEEE Trans. Inform. Theor.* **44** (1998) 2743–2760
9. Grünwald, P.D.: *The Minimum Description Length Principle*. The MIT Press (2007)
10. Yu, B., Speed, T.P.: Data compression and histograms. *Probab. Theor. Rel. Fields* **92** (1992) 195–229
11. Hopcroft, J.E., Ullman, J.D.: *Introduction to Automata Theory, Languages and Computation*. Addison-Wesley (1979)
12. Elias, P.: Universal codeword sets and representations for the integers. *IEEE Trans. Inform. Theor.* **21** (1975) 194–203
13. Barron, A.R.: *Logically Smooth Density Estimation*. PhD thesis, Stanford University (1985)

14. Dawid, A.: Statistical theory: The prequential approach. *J. Roy. Statist. Soc. A* **147** (1984) 278–292
15. Cover, T.M., Thomas, J.A.: *Elements of Information Theory*. Wiley (1991)
16. Li, L., Yu, B.: Iterated logarithmic expansions of the pathwise code lengths for exponential families. *IEEE Trans. Inform. Theor.* **46** (2000) 2683–2689
17. Barndorff-Nielsen, O.E.: *Information and Exponential Families*. Wiley (1978)
18. Jeffreys, H.: *Theory of Probability*. 3rd Edition. Oxford University Press (1961)
19. Dębowski, Ł.: On the vocabulary of grammar-based codes and the logical consistency of texts. E-print: <http://arxiv.org/abs/0810.3125> (2008)
20. Csiszar, I., Shields, P.C.: The consistency of the BIC Markov order estimator. *Ann. Statist.* **28** (2000) 1601–1619