



Coding on countably infinite alphabets

Stéphane Boucheron, Aurélien Garivier, Elisabeth Gassiat

► To cite this version:

Stéphane Boucheron, Aurélien Garivier, Elisabeth Gassiat. Coding on countably infinite alphabets. 2006. hal-00121892v1

HAL Id: hal-00121892

<https://hal.science/hal-00121892v1>

Preprint submitted on 22 Dec 2006 (v1), last revised 16 Jan 2008 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Coding on countably infinite alphabets

Stéphane Boucheron, Aurélien Garivier and Elisabeth Gassiat

Abstract

This paper describes universal lossless coding strategies for compressing sources on countably infinite alphabets. Classes of memoryless sources defined by an envelope condition on the marginal distribution provide benchmarks for coding techniques originating from the theory of universal coding over finite alphabets. We prove general upper-bounds on minimax regret and lower-bounds on minimax redundancy for such source classes. The general upper bounds emphasize the role of the Normalized Maximum Likelihood codes with respect to minimax regret in the infinite alphabet context. Lower bounds are derived by tailoring sharp bounds on the redundancy of Krichevsky-Trofimov coders for sources over finite alphabets. Up to logarithmic (resp. constant) factors the bounds are matching for source classes defined by algebraically declining (resp. exponentially vanishing) envelopes. Effective and (almost) adaptive coding techniques are described for the collection of source classes defined by algebraically vanishing envelopes. Those results extend our knowledge concerning universal coding to contexts where the key tools from parametric inference (Bernstein-Von Mises theorem, Wilks theorem) are known to fail.

keywords: NML; countable alphabets; redundancy; adaptive compression; minimax;

I. INTRODUCTION

This paper is concerned with the problem of universal coding on a countably infinite alphabet \mathcal{X} (say \mathbb{N}_+) as described for example by Orlitsky and Santhanam (2004). Throughout this paper, a source on the countable alphabet \mathcal{X} is a probability distribution on the set of infinite sequences of symbols from \mathcal{X} . The symbol Λ will be used to denote various classes of sources on the countably infinite alphabet \mathcal{X} . The sequence of symbols emitted by a source is denoted by the $\mathcal{X}^{\mathbb{N}}$ -valued random variable $\mathbf{X} = (X_n)_{n \in \mathbb{N}}$. If P denotes the distribution of \mathbf{X} , P^n denotes the distribution of $X_{1:n} = X_1, \dots, X_n$, and we let $\Lambda^n = \{P^n : P \in \Lambda\}$. For any countable set \mathcal{X} , let $\mathfrak{M}_1(\mathcal{X})$ be the set of all probability measures on \mathcal{X} .

From Shannon noiseless coding Theorem (see Cover and Thomas, 1991), the binary entropy of P^n , $H(X_{1:n}) = \mathbb{E}_{P^n} [-\log_2 P(X_{1:n})]$ provides a tight lower bound on the expected number of binary symbols needed to encode outcomes of P^n . On the other hand, thanks to arithmetic coding (see for example Cover and Thomas, 1991), any distribution $Q^n \in \mathfrak{M}_1(\mathcal{X}^n)$ defines a prefix code, that encodes string \mathbf{x} using $\lceil -\log_2 Q^n(\mathbf{x}) \rceil + 1$ bits. If the arithmetic code derived from distribution Q^n is used to encode outcomes from P^n , the *expected redundancy* of Q^n (with respect to P^n) is defined as the expected difference between the expected code length $\mathbb{E}_P [-\log_2 Q^n(X_{1:n})]$ and $H(X_{1:n})$. Up to a factor $\log 2$, it is equal to the Kullback-Leibler divergence (or relative entropy) $D(P^n, Q^n) = \sum_{\mathbf{x} \in \mathcal{X}^n} P^n\{\mathbf{x}\} \log \frac{P^n(\mathbf{x})}{Q^n(\mathbf{x})} = \mathbb{E}_{P^n} \left[\log \frac{P^n(X_{1:n})}{Q^n(X_{1:n})} \right]$. From now on, unless it is necessary, we will not specify the base of the logarithm.

At large, universal coding attempts to develop sequences of coding probabilities $(Q^n)_n$ so as to minimize expected redundancy over a whole class of sources. Technically speaking, several distinct notions of universality have been considered in the literature. A function $\rho(n)$ is said to be a strong (respectively weak) *universal redundancy rate* for a class of sources Λ if there exists a sequence of coding probabilities $(Q^n)_n$ such that for all n , $R^+(Q^n, \Lambda^n) = \sup_{P \in \Lambda} D(P^n, Q^n) \leq \rho(n)$ (respectively for all $P \in \Lambda$, there exists a constant $C(P)$ such that for all n , $D(P^n, Q^n) \leq C(P)\rho(n)$). A redundancy rate $\rho(n)$ is said to be non-trivial if $\lim_n \rho(n) = 0$. Finally a class Λ of sources will be said to be *feebly universal* if there exists a single sequence of coding probabilities $(Q^n)_n$ such that $\sup_{P \in \Lambda} \lim_n \frac{1}{n} D(P^n, Q^n) = 0$ (Note that this notion of feeble universality is usually called weak universality, (see Kieffer, 1978, Györfi et al., 1994), we deviate from the tradition, in order to avoid confusion with the notion of weak universal redundancy rate).

As far as finite alphabets are concerned, it is well-known that the class of stationary ergodic sources is feebly universal. This is witnessed by the performance of Lempel-Ziv codes (see Cover and Thomas, 1991). It is also known that the class of stationary ergodic sources over a finite alphabet does not admit any non-trivial weak universal redundancy rate (Shields, 1993). On the other hand, fairly large classes of sources admitting strong universal redundancy rates and non-trivial weak universal redundancy rates have been exhibited (see Barron et al., 1998, Catoni, 2004, and references therein).

In this paper, we will mostly focus on strong universal redundancy rates for classes of sources over infinite alphabets. Note that in the latter setting, even feeble universality should not be taken for granted: the class of memoryless processes on \mathbb{N}_+ is not feebly universal.

Kieffer (1978) characterized feebly universal classes, and the argument was simplified by Györfi et al. (1994). Recall that the entropy rate $H(P)$ of a stationary source is defined as $\lim_n H(P^n)/n$.

Proposition 1: A class Λ of stationary sources over a countable alphabet \mathcal{X} is feebly universal if and only if there exists a probability distribution $Q \in \mathfrak{M}_1(\mathcal{X})$ such that for every $P \in \Lambda$ with finite entropy rate, Q satisfies $\mathbb{E}_P \log \frac{1}{Q(X_1)} < \infty$ or equivalently $D(P^1, Q) < \infty$.

For sources over countably infinite alphabets, the characterization of strong universality has not yet reached such a degree of maturity.

The *maximal redundancy* of Q^n with respect to Λ is defined by:

$$R^+(Q^n, \Lambda^n) = \sup_{P \in \Lambda} D(P^n, Q^n).$$

The infimum of $R^+(Q^n, \Lambda^n)$ is called the *minimax redundancy* with respect to Λ :

$$R^+(\Lambda^n) = \inf_{Q^n \in \mathfrak{M}_1(\mathcal{X}^n)} R^+(Q^n, \Lambda^n).$$

It is the smallest strong universal redundancy rate for Λ . When finite, it is often called the information radius of Λ^n .

Assume that a subset of Λ is parametrized by Θ and that Θ can be equipped with (prior) probability distributions W in such a way that $\theta \mapsto P_\theta^n\{A\}$ is a random variable for every $A \subseteq \mathcal{X}^n$. A convenient way to derive lower

bounds on $R^+(\Lambda^n)$ consists in using the relation $\mathbb{E}_W[D(P_\theta^n, Q^n)] \leq R^+(Q^n, \Lambda^n)$.

The sharpest lower bound is obtained by optimizing the prior probability distributions (assuming $(P_\theta^n)_{\theta \in \Theta} = \Lambda^n$ taking the so-called least favorable prior in order to avoid confusion with the notion of weak universal redundancy rate), it is called the maximin bound

$$\sup_{W \in \mathfrak{M}_1(\Theta)} \inf_{Q^n \in \mathfrak{M}_1(\mathcal{X}^n)} \mathbb{E}_W[D(P_\theta^n, Q^n)].$$

It has been proved in a series of papers (Gallager, 1968, Davisson, 1973, Haussler, 1997) (and could also have been derived from a general minmax theorem by Sion, 1958) that such a lower bound is tight.

Theorem 1: Let Λ denote a class of sources over some finite or countably infinite alphabet. For each n , the minimax redundancy over Λ coincides with

$$R^+(\Lambda^n) = \sup_{\Theta, W \in \mathfrak{M}_1(\Theta)} \inf_{Q^n \in \mathfrak{M}_1(\mathcal{X}^n)} \mathbb{E}_W[D(P_\theta^n, Q^n)],$$

where Θ runs over all parametrizations of countable subsets of Λ .

If the set $\Lambda^n = \{P^n : P \in \Lambda\}$ is not pre-compact with respect to the topology of weak convergence, then both sides are infinite. Otherwise the maximin and minimax average redundancies are finite and coincide; moreover, the minimax redundancy is achieved by the mixture coding distribution $Q^n(\cdot) = \int_{\Theta} P_\theta^n(\cdot) W(d\theta)$ where W is the least favorable prior and Θ may be uncountable.

Another approach to universal coding considers *individual sequences* (see Feder et al., 1992, Cesa-Bianchi and Lugosi, 2006, and references therein). Let the *regret* of a coding distribution Q^n on string $\mathbf{x} \in \mathbb{N}_+^n$ with respect to Λ be $\sup_{P^n \in \Lambda} \log P^n(\mathbf{x})/Q^n(\mathbf{x})$. Taking the maximum with respect to $x \in \mathbb{N}_+^n$, and then optimizing over the choice of Q^n , we get the *minimax regret*:

$$R^*(\Lambda^n) = \inf_{Q^n \in \mathfrak{M}_1(\mathcal{X}^n)} \max_{x \in \mathbb{N}_+^n} \sup_{P \in \Lambda} \log \frac{P^n(x)}{Q^n(x)}.$$

In order to provide proper insight, let us recall the precise asymptotic bounds on minimax redundancy and regret for memoryless sources over finite alphabets (see Clarke and Barron, 1990; 1994, Barron et al., 1998, Xie and Barron, 1997; 2000, Orlitsky and Santhanam, 2004, Catoni, 2004, Szpankowski, 1998, and references therein).

Theorem 2: Let \mathcal{X} be an alphabet of m symbols, and Λ denote the class of memoryless processes on \mathcal{X} then

$$\begin{aligned} \lim_n R^+(\Lambda^n) - \frac{m-1}{2} \log \frac{n}{2\pi e} &= \log \left(\frac{\Gamma(1/2)^m}{\Gamma(m/2)} \right) \\ \lim_n R^*(\Lambda^n) - \frac{m-1}{2} \log \frac{n}{2\pi} &= \log \left(\frac{\Gamma(1/2)^m}{\Gamma(m/2)} \right). \end{aligned}$$

For all n , if $m < n$:

$$R^*(\Lambda^n) \leq \frac{m-1}{2} \log n + 2.$$

The last inequality is checked in the Appendix .

Remark 1: The set of memoryless sources over alphabet $\mathcal{X} = \{1, \dots, m\}$ is conveniently parametrized by $\Theta = \{\theta : \theta \in \mathbb{R}_+^{m-1}, \sum_{i=1}^{m-1} \theta[i] \leq 1\}$. To alleviate notations, we agree on $\theta[m] = 1 - \sum_{j=1}^{m-1} \theta[j]$. For any string \mathbf{x}

from \mathcal{X} , let $n_j = \sum_{i=1}^n 1_{\mathbf{x}_i=j}$ then $P_{\boldsymbol{\theta}}^n$ is defined by the probability mass function

$$P_{\boldsymbol{\theta}}^n(\mathbf{x}) = \prod_{j=1}^m (\theta[j])^{n_j}.$$

Jeffrey's prior has a density $w_J(\boldsymbol{\theta})$ proportional to the square root of the Fisher Information $J(\boldsymbol{\theta})$ where:

$$\sqrt{J(\boldsymbol{\theta})} = \prod_{j=1}^m (\theta[j])^{-1/2}.$$

Clarke and Barron (1994) have proved that Jeffrey's prior is asymptotically least favorable:

$$\text{letting } Q_{w_J}^n(\mathbf{x}) = \int_{\Theta} P_{\boldsymbol{\theta}}^n(\mathbf{x}) dw_J(\boldsymbol{\theta})$$

then

$$\lim_n \left(\mathbb{E}_{w_J} [D(P_{\boldsymbol{\theta}}^n, Q_{w_J}^n)] - \frac{m-1}{2} \log \frac{n}{2\pi e} \right) = \log \left(\frac{\Gamma(1/2)^m}{\Gamma(m/2)} \right).$$

Moreover, a sequence of modifications of Jeffrey's prior asymptotically achieves minimax redundancy (See Xie and Barron, 1997).

Remark 2: The phenomenon pointed out in Theorem 2 holds not only for the class of memoryless sources over a finite alphabet but also for classes of sources that are smoothly parametrized by finite dimensional sets (see again Clarke and Barron, 1990; 1994, Barron et al., 1998, Xie and Barron, 1997; 2000, Orlitsky and Santhanam, 2004, Catoni, 2004).

Theorem 2 can be considered as an information-theoretical refinement of a classical result in parametric statistics : the asymptotics of the maximin redundancy reflects the asymptotic normality of the rescaled posterior measure as asserted by the (Laplace-) Bernstein-Von Mises Theorem and the connexion between the entropy and the variance of Gaussian measures (see Clarke and Barron, 1990; 1994, Barron et al., 1998, van der Vaart, 1998).

Our interest in the coding problem for infinite alphabets stems partly from the fact that in non-parametric settings, the Bernstein-von Mises Theorem does not hold in full generality (See Cox, 1993, Freedman, 1999, Ghosh and Ramamoorthi, 2003, and references therein.).

Let us pay further attention to the minimax regret. For a source class Λ , for every $\mathbf{x} \in \mathcal{X}^n$, let the maximum likelihood $\hat{p}(\mathbf{x})$ be defined as $\sup_{P \in \Lambda} P^n(\mathbf{x})$. If $\sum_{\mathbf{x} \in \mathbb{N}_+^n} \hat{p}(\mathbf{x}) < \infty$, the *Normalized Maximum Likelihood* coding probability is well-defined and given by

$$Q_{\text{NML}}^n(\mathbf{x}) = \frac{\hat{p}(\mathbf{x})}{\sum_{\mathbf{x} \in \mathbb{N}_+^n} \hat{p}(\mathbf{x})}.$$

Shtarkov (1987) showed that the *Normalized Maximum Likelihood* coding probability achieves the same regret over all strings of length n and that this regret coincides with the *minimax regret*:

$$R^*(\Lambda^n) = \log \sum_{\mathbf{x} \in \mathbb{N}_+^n} \hat{p}(\mathbf{x}).$$

The maximum regret achieved by the mixture defined by Jeffrey's prior is within a non-null constant from the

minimax regret. Moreover, Xie and Barron (1997) have shown that:

$$D(Q_{w_J}^n, Q_{\text{NML}}^n) \rightarrow 0.$$

This holds for a variety of classes of sources smoothly parametrized by finite-dimensional sets (Barron et al., 1998).

Again, the relation between minimax regret and minimax redundancy for the set of memoryless sources over a finite alphabet can be linked to another classical result from asymptotic statistics. Let $\hat{\theta} \in \Theta$ denote the parameter that achieves maximum likelihood on \mathbf{x} . Simple algebra shows that

$$D(Q_{w_J}^n, Q_{\text{NML}}^n) = R^*(\Lambda^n) - R^-(\Lambda^n) - \int_{\Theta} w_J(\theta) \mathbb{E}_{P_{\theta}^n} \left[\log \frac{P_{\hat{\theta}}^n(X_{1:n})}{P_{\theta}^n(X_{1:n})} \right] d\theta.$$

Then, by a theorem due to Wilks (see van der Vaart, 1998), the last summand converges toward one half the expectation of a χ_{m-1}^2 distributed random variable (using natural logarithms on both sides).

Memoryless sources over finite alphabets are special cases of envelope classes. The latter will be of primary interest.

Definition 1: Let f be a mapping from \mathbb{N}_+ to $[0, 1]$. The envelope class Λ_f defined by function f is the collection of stationary memoryless sources with first marginal distribution dominated by f :

$$\Lambda_f = \{P : \forall x \in \mathbb{N}, P^1\{x\} \leq f(x), \text{ and } P \text{ is stationary and memoryless.}\}.$$

We will be concerned with the following topics.

- 1) Understanding general structural properties of minimax redundancy and minimax regret.
- 2) Characterizing those source classes that have finite minimax regret.
- 3) Quantitative relations between minimax redundancy or regret and integrability of the envelope function.
- 4) Developing effective coding techniques for source classes with known non-trivial minimax redundancy rate.
- 5) Developing adaptive coding schemes for collections of source classes that are too large to enjoy even a weak redundancy rate.

The paper is organized as follows. Section II describes some structural properties of minimax redundancies and regrets for classes of stationary memoryless sources. Those properties include monotonicity and sub-additivity. Proposition 5 characterizes those source classes that admit finite regret. This characterization emphasizes the role of Shtarkov Normalized Maximum Likelihood coding probability. Proposition 6 describes a simple source class for which the minimax regret is infinite, while the minimax redundancy is finite. Finally Proposition 3 asserts that such a contrast is not possible for the so-called envelope classes.

In Section III, Theorems 4 and 5 provide quantitative relations between the summability properties of the envelope function and minimax regrets and redundancies. Those results build on the non-asymptotic bounds on minimax redundancy derived by Xie and Barron (1997).

Section IV focuses on two kinds of envelope classes. This section serves as a benchmark for the two main

results from the preceding section. In Subsection IV-A, lower-bounds on minimax redundancy and upper-bounds on minimax regret for classes defined by envelope function $k \mapsto 1 \wedge Ck^{-\alpha}$ are described. Up to a factor $\log n$ those bounds are matching. In Subsection IV-B, lower-bounds on minimax redundancy and upper-bounds on minimax regret for classes defined by envelope function $k \mapsto 1 \wedge C \exp^{-\alpha k}$ are described. Up to a multiplicative constant, those bounds coincide and grow like $\log^2 n$.

In Sections V and VI, we turn to effective coding techniques geared toward source classes defined by power-law envelopes. In Section V, we elaborate on the ideas embodied in Proposition 4 from Section II, and combine mixture coding and Elias penultimate code (Elias, 1975) to match the upper-bounds on minimax redundancy described in Section IV. One of the messages from Section IV is that the union of envelope classes defined by power laws, does not admit a weak redundancy rate that grows at a rate slower than $n^{1/\beta}$ for any $\beta > 1$. In Section VI, we finally develop an adaptive coding scheme for the union of envelope classes defined by power laws. This adaptive coding scheme combines the censoring coding technique developed in the preceding subsection and an estimation of tail-heaviness.

II. STRUCTURAL PROPERTIES OF THE MINIMAX REDUNDANCY AND MINIMAX REGRET

Propositions 3 and 4 are sanity-check statements. In order to prove them, we will use the following proposition which emphasizes the role of the NML coder with respect to the minimax regret. At best, it is a comment on Shtarkov's original work (Shtarkov, 1987, Haussler and Oppen, 1997).

Proposition 2: Let Λ be a class of stationary memoryless sources over a countably infinite alphabet, the minimax regret with respect to Λ^n , $R^*(\Lambda^n)$ is finite if and only if the normalized maximum likelihood (Shtarkov) coding probability Q_{NML}^n is well-defined and given by

$$Q_{\text{NML}}^n(\mathbf{x}) = \frac{\hat{p}(\mathbf{x})}{\sum_{\mathbf{y} \in \mathcal{X}^n} \hat{p}(\mathbf{y})}$$

where $\hat{p}(\mathbf{x}) = \sup_{P \in \Lambda} P^n(\mathbf{x})$.

Note that the definition of Q_{NML}^n does not assume either that the maximum likelihood is achieved or that it is uniquely defined.

Proof: The fact that if Q_{NML}^n is well-defined, the minimax regret is finite and equal to

$$\log \left(\sum_{\mathbf{y} \in \mathcal{X}^n} \hat{p}(\mathbf{y}) \right)$$

is the fundamental observation of Shtarkov (1987).

On the other hand, if $R^*(\Lambda^n) < \infty$, there exists a probability distribution Q^n on \mathcal{X}^n and a finite number r such that for all $\mathbf{x} \in \mathcal{X}^n$,

$$\hat{p}(\mathbf{x}) \leq r \times Q^n(\mathbf{x}),$$

summing gives

$$\sum_{\mathbf{x} \in \mathcal{X}^n} \hat{p}(\mathbf{x}) \leq r < \infty.$$

■

Proposition 3: Let Λ denote a class of sources, then the minimax redundancy $R^+(\Lambda^n)$ and the minimax regret $R^*(\Lambda^n)$ are non-decreasing functions of n .

Proof: As far as R^+ is concerned, by Theorem 1, it is enough to check that the maximin (mutual information) lower bound is non-decreasing.

For any prior distribution W on a parameter set Θ (recall that $\{P_\theta : \theta \in \Theta\} \subseteq \Lambda$, and that the mixture coding probability Q^n is defined by $Q^n(A) = \mathbb{E}_W[P_\theta^n(A)]$)

$$\mathbb{E}_W [D(P_\theta^{n+1}, Q^{n+1})] = I(\theta; X_1^{n+1}) = I(\theta; (X_1^n, X_{n+1})) \geq I(\theta; X_1^n) = \mathbb{E}_W [D(P_\theta^n, Q^n)].$$

Let us now consider the minimax regret. It is enough to consider the case where $R^*(\Lambda^n)$ is finite. Thus we may rely on Proposition 2. Let n and m be two positive integers. Let ϵ be a small positive real. For any string $\mathbf{x} \in \mathcal{X}^n$, let $Q \in \Lambda$, be such that $Q\{\mathbf{x}\} \geq \hat{p}(\mathbf{x})(1 - \epsilon)$. Then

$$\begin{aligned} \hat{p}(\mathbf{x}x') &\geq Q(\mathbf{x}) \times Q(x' | \mathbf{x}) \\ &\geq \hat{p}(\mathbf{x})(1 - \epsilon) \times Q(x' | \mathbf{x}). \end{aligned}$$

Summing over all possible $x' \in \mathcal{X}$ we get

$$\sum_{x'} \hat{p}(\mathbf{x}x') \geq \hat{p}(\mathbf{x})(1 - \epsilon).$$

Summing now over all $\mathbf{x} \in \mathcal{X}^n$ and $x' \in \mathcal{X}$,

$$\sum_{\mathbf{x} \in \mathcal{X}^n, x' \in \mathcal{X}} \hat{p}(\mathbf{x}x') \geq \sum_{\mathbf{x} \in \mathcal{X}^n} \hat{p}(\mathbf{x})(1 - \epsilon).$$

So that by letting ϵ tend to 0,

$$\sum_{\mathbf{x} \in \mathcal{X}^{n+1}} \hat{p}(\mathbf{x}) \geq \sum_{\mathbf{x} \in \mathcal{X}^n} \hat{p}(\mathbf{x}).$$

■

Note that the proposition holds even though Λ is not a collection of memoryless sources. This Proposition can be easily completed when dealing with memoryless sources.

Proposition 4: If Λ is a class of stationary memoryless sources, then the functions $n \mapsto R^+(\Lambda^n)$ and $n \mapsto R^*(\Lambda^n)$ are sub-additive.

Proof: Here again, given Theorem 1, in order to establish sub-additivity for R^+ , it is enough to check the property for the maximin lower bound. Let n, m be two positive integers, and W be any prior on Θ (with $\{P_\theta : \theta \in \Theta\} \subseteq \Lambda$). As sources from Λ are memoryless, $X_{1:n}$ and X_{n+1}^{n+m} are independent conditionally on θ

and thus

$$\begin{aligned}
& I(X_{n+1}^{n+m}; \theta | X_1^n) \\
&= H(X_{n+1}^{n+m} | X_1^n) - H(X_{n+1}^{n+m} | X_1^n, \theta) \\
&= H(X_{n+1}^{n+m} | X_1^n) - H(X_{n+1}^{n+m} | \theta) \\
&\leq H(X_{n+1}^{n+m}) - H(X_{n+1}^{n+m} | \theta) \\
&= I(X_{n+1}^{n+m}; \theta) .
\end{aligned}$$

Hence, using the fact that under each P_θ , the process $(X_n)_{n \in \mathbb{N}_+}$ is stationary:

$$\begin{aligned}
I(X_1^{n+m}; \theta) &= I(X_1^n; \theta) + I(X_{n+1}^{n+m}; \theta | X_1^n) \\
&\leq I(X_1^n; \theta) + I(X_{n+1}^{n+m}; \theta) \\
&= I(X_1^n; \theta) + I(X_1^m; \theta) .
\end{aligned}$$

Let us now check the sub-additivity of minimax regret. For any $\epsilon > 0$, for $\mathbf{x} \in \mathcal{X}^{n+m}$, let $P \in \Lambda$ be such that $(1 - \epsilon)\hat{p}(\mathbf{x}) \leq P^{n+m}(\mathbf{x})$. As for $\mathbf{x} \in \mathcal{X}^n$ and $\mathbf{x}' \in \mathcal{X}^m$, $P^{n+m}(\mathbf{x}\mathbf{x}') = P^n(\mathbf{x}) \times P^m(\mathbf{x}')$, we have for any $\epsilon > 0$, and any $\mathbf{x} \in \mathcal{X}^n, \mathbf{x}' \in \mathcal{X}^m$

$$(1 - \epsilon)\hat{p}(\mathbf{x}\mathbf{x}') \leq \hat{p}(\mathbf{x}) \times \hat{p}(\mathbf{x}') .$$

Hence, letting ϵ tend to 0, and summing over all $\mathbf{x} \in \mathcal{X}^{n+m}$:

$$\begin{aligned}
& R^*(\Lambda^{n+m}) \\
&= \log \sum_{\mathbf{x} \in \mathcal{X}^n, \mathbf{x}' \in \mathcal{X}^m} \hat{p}(\mathbf{x}\mathbf{x}') \\
&\leq \log \sum_{\mathbf{x} \in \mathcal{X}^n} \hat{p}(\mathbf{x}) + \log \sum_{\mathbf{x}' \in \mathcal{X}^m} \hat{p}(\mathbf{x}') \\
&= R^*(\Lambda^n) + R^*(\Lambda^m) .
\end{aligned}$$

■

Remark 3: Counter-examples witness the fact that subadditivity of redundancies does not hold in full generality. The Fekete Lemma (see Dembo and Zeitouni, 1998) leads to:

Corollary 1: Let Λ denote a class of stationary memoryless sources over a countable alphabet. For both minimax redundancy R^+ and minimax regret R^* ,

$$\lim_{n \rightarrow \infty} \frac{R^+(\Lambda^n)}{n} = \inf_{n \in \mathbb{N}_+} \frac{R^+(\Lambda^n)}{n} \leq R^+(\Lambda^1) ,$$

and

$$\lim_{n \rightarrow \infty} \frac{R^*(\Lambda^n)}{n} = \inf_{n \in \mathbb{N}_+} \frac{R^*(\Lambda^n)}{n} \leq R^*(\Lambda^1) .$$

Hence, in order to prove that $R^+(\Lambda^n) < \infty$ (respectively $R^*(\Lambda^n) < \infty$), it is enough to check that $R^+(\Lambda^1) < \infty$ (respectively $R^*(\Lambda^1) < \infty$).

The following Proposition combines Propositions 2, 3 and 4. It can be rephrased as follows: a class of memoryless sources admits a non-trivial strong minimax regret if and only if Shtarkov NML coding probability is well-defined for $n = 1$.

Proposition 5: Let Λ be a class of stationary memoryless sources over a countably infinite alphabet. Let \hat{p} be defined by $\hat{p}(x) = \sup_{P \in \Lambda} P\{x\}$. The minimax regret with respect to Λ^n is finite if and only if the normalized maximum likelihood (Shtarkov) coding probability is well-defined and :

$$R^*(\Lambda^n) < \infty \Leftrightarrow \sum_{x \in \mathbb{N}_+} \hat{p}(x) < \infty.$$

Proof: The direct part follows from Proposition 2.

For the converse part, if $\sum_{x \in \mathbb{N}_+} \hat{p}(x) = \infty$, then $R^*(\Lambda^1) = \infty$ and from Proposition 3, $R^*(\Lambda^n) = \infty$ for every positive integer n . ■

When dealing with smoothly parametrized classes of sources over finite alphabets (see Barron et al., 1998, Xie and Barron, 2000) or even with the massive classes defined by renewal sources (Csiszár and Shields, 1996), the minimax regret and minimax redundancy are usually of the same order of magnitude. This can not be taken for granted when dealing with classes of stationary memoryless sources over a countable alphabet.

Proposition 6: Let f be a positive, strictly decreasing function defined on \mathbb{N} such that $f(1) < 1$. For $k \in \mathbb{N}$, let p_k be the probability mass function on \mathbb{N} defined by:

$$p_k(l) = \begin{cases} 1 - f(k) & \text{if } l = 0; \\ f(k) & \text{if } l = k; \\ 0 & \text{otherwise.} \end{cases}$$

Let $\Lambda^1 = \{p_1, p_2, \dots\}$, let Λ be the class of stationary memoryless sources with first marginal Λ^1 . The finiteness of the minimax redundancy with respect to Λ_f^n depends on the limiting behavior of $f(k) \log k$: for every positive integer n :

$$f(k) \log k \rightarrow_{k \rightarrow \infty} \infty \Leftrightarrow R^+(\Lambda^n) = \infty.$$

Remark 4: When $f(k) = \frac{1}{\log k}$, the minimax redundancy $R^+(\Lambda_f^n)$ is finite for all n . Note, however that this does not warrant the existence of a non-trivial strong universal redundancy rate. However, as $\sum_k f(k) = \infty$, minimax regret is infinite by Proposition 5.

A similar result appears in the discussion of Theorem 3 in (Haussler and Oppel, 1997) where classes with finite

minimax redundancy and infinite minimax regret are called irregular.

We will be able to refine those observations after the statement of Corollary 2.

Proof: Let us first prove the direct part. Assume that $f(k) \log k \rightarrow_{k \rightarrow \infty} \infty$. In order to check that $R^+(\Lambda^1) = \infty$, we resort to the mutual information lower bound (Theorem 1) and describe an appropriate collection of Bayesian games.

Let m be a positive integer and let θ be uniformly distributed over $\{1, 2, \dots, m\}$. Let X be distributed according to p_k conditionally on $\theta = k$. Let Z be the random variable equal to 1 if $X = \theta$ and equal to 0 otherwise. Obviously, $H(\theta|X, Z = 1) = 0$; moreover, as f is assumed to be non-increasing, $P(Z = 0|\theta = k) = 1 - f(k) \leq 1 - f(m)$ and thus:

$$\begin{aligned} H(\theta|X) &= H(Z|X) + H(\theta|Z, X) \\ &\leq 1 + P(Z = 0)H(\theta|X, Z = 0) \\ &\quad + P(Z = 1)H(\theta|X, Z = 1) \\ &\leq 1 + (1 - f(m)) \log m. \end{aligned}$$

Hence,

$$\begin{aligned} R^+(\Lambda^1) &\geq I(\theta, X) \\ &\geq \log m - (1 - f(m)) \log m \\ &= f(m) \log m \end{aligned}$$

which grows to infinity with m , so that as announced $R^+(\Lambda^1) = \infty$.

Let us now prove the converse part. Assume that the sequence $(f(k) \log k)_{k \in \mathbb{N}_+}$ is upper-bounded by some constant C . In order to check that $R^+(\Lambda^n) < \infty$, for all n , by Proposition 4, it is enough to check that $R^+(\Lambda_f^1) < \infty$, and thus, it is enough to exhibit a probability distribution Q over $\mathcal{X} = \mathbb{N}$ such that $\sup_{P \in \Lambda^1} D(P, Q) < \infty$.

Let Q be defined by $Q(k) = A / ((1 \vee (k(\log k)^2)))$ for $k \geq 2$, $Q(0), Q(1) > 0$ where A is a normalizing constant that ensures that Q is a probability distribution over \mathcal{X} .

Then for any $k > 3$ (which warrants $k(\log k)^2 > 1$), letting P_k be the probability defined by the probability mass function p_k :

$$\begin{aligned} D(P_k, Q) &= (1 - f(k)) \log \frac{(1 - f(k))}{Q(0)} + f(k) \log \left(\frac{f(k)k(\log k)^2}{A} \right) \\ &\leq -\log Q(0) + C + f(k) \left(2 \log^{(2)}(k) - \log(A) \right) \\ &\leq C + \log \frac{C^2}{A Q(0)}. \end{aligned}$$

This is enough to conclude that

$$R^+(\Lambda^1) \leq \left(C + \log \frac{C^2}{A Q(0)} \right) \vee D(P_1, Q) < \infty.$$

■

Remark 5: Note that the coding probability used in the proof of the converse part of the proposition corresponds to one of the simplest prefix codes for integers proposed by Elias (1975).

The following theorem shows that, as far as envelope classes are concerned, minimax redundancy and minimax regret are either both finite or both infinite.

Theorem 3: Let f be a non-negative function from \mathbb{N}_+ to $[0, 1]$, let Λ_f be the class of stationary memoryless sources defined by envelope f . Then

$$R^+(\Lambda_f^n) < \infty \Leftrightarrow R^*(\Lambda_f^n) < \infty.$$

Remark 6: We will refine this result after the statement of Corollary 2.

Recall from Proposition 5 that $R^*(\Lambda_f^n) < \infty \Leftrightarrow \sum_{k \in \mathbb{N}_+} f(k) < \infty$.

Proof:

In order to check that

$$\sum_{k \in \mathbb{N}_+} f(k) = \infty \Rightarrow R^+(\Lambda_f^n) = \infty,$$

it is enough to check that if $\sum_{k \in \mathbb{N}_+} f(k) = \infty$, the envelope class contains an infinite collection of mutually singular sources.

Let the infinite sequence of integers $(h_i)_{i \in \mathbb{N}}$ be defined recursively by $h_0 = 0$ and

$$h_{i+1} = \min \left\{ h : \sum_{k=h_i+1}^h f(k) > 1 \right\}.$$

The memoryless source P_i is defined by its first marginal P_i^1 which is given by

$$P_i^1(m) = \frac{f(m)}{\sum_{k=h_i+1}^{h_{i+1}} f(k)} \text{ for } m \in \{p_i + 1, \dots, p_{i+1}\}.$$

Taking any prior with infinite Shannon entropy over the $\{P_i^1 ; i \in \mathbb{N}_+\}$ shows that

$$R^+(\{P_i^1 ; i \in \mathbb{N}_+\}) = \infty.$$

■

III. ENVELOPE CLASSES

The next two theorems establish quantitative relations between minimax redundancy and regrets and the shape of the envelope function. The first one holds for any class of memoryless sources.

Theorem 4: If Λ is a class of memoryless sources, let the tail function \bar{F}_{Λ^1} be defined by $\bar{F}_{\Lambda^1}(u) = \sum_{k>u} \hat{p}(k)$, then:

$$R^*(\Lambda^n) \leq \inf_{u: u \leq n} \left[n \bar{F}_{\Lambda^1}(u) \log e + \frac{u-1}{2} \log n + \right] + 2.$$

Choosing a sequence $(u_n)_n$ of positive integers in such a way that $u_n \rightarrow \infty$ while $u_n/n \rightarrow 0$, this theorem allows to complete Proposition 5.

Corollary 2: Let Λ denote a class of memoryless sources, then the following holds:

$$R^*(\Lambda^n) < \infty \Leftrightarrow R^*(\Lambda^n) = o(n) \text{ and } R^+(\Lambda^n) = o(n),$$

Remark 7: We may now have a second look at Proposition 6 and Theorem 3. In the setting of Proposition 6, this Corollary asserts that if $\sum_k f(k) < \infty$, a non-trivial strong redundancy rate exists.

This corollary complements Theorem 3 by asserting that envelope classes have either non-trivial strong redundancy rates or infinite minimax redundancies.

Remark 8: Again, this statement has to be connected with related propositions from Haussler and Opper (1997). The last paper establishes bounds on minimax redundancy using geometric properties of the source class under Hellinger metric. For example, Theorem 4 in (Haussler and Opper, 1997) relates minimax redundancy and the metric dimension of the set Λ^n with respect to the Hellinger metric (which coincides with L_2 metric between the square roots of densities) under the implicit assumption that sources lying in small Hellinger balls have finite relative entropy (so that upper bounds in Lemma 7 there are finite). Envelope classes may not satisfy this assumption. Hence, there is no easy way to connect Theorem 4 and results from (Haussler and Opper, 1997).

Proof: (Theorem 4.) Any integer u defines a decomposition of a string $\mathbf{x} \in \mathbb{N}_+^n$ into two non-contiguous substrings: a substring \mathbf{z} made of the m symbols from \mathbf{x} that are larger than u , and one substring \mathbf{y} made of the

$n - m$ symbols that are smaller than u .

$$\begin{aligned}
& \sum_{\mathbf{x} \in \mathbb{N}_+^n} \hat{p}(\mathbf{x}) \\
& \stackrel{(a)}{=} \sum_{m=0}^n \binom{n}{m} \sum_{\mathbf{z} \in \{u, \dots\}^m} \sum_{\mathbf{y} \in \{1, 2, \dots, u\}^{n-m}} \hat{p}(\mathbf{zy}) \\
& \stackrel{(b)}{\leq} \sum_{m=0}^n \binom{n}{m} \sum_{\mathbf{z} \in \{u, \dots\}^m} \prod_{i=1}^m \hat{p}(\mathbf{z}_i) \sum_{\mathbf{y} \in \{1, 2, \dots, u\}^{n-m}} \hat{p}(\mathbf{y}) \\
& \stackrel{(c)}{\leq} \left(\sum_{m=0}^n \binom{n}{m} \bar{F}_{\Lambda^1}(u)^m \right) \left(\sum_{\mathbf{y} \in \{1, 2, \dots, u\}^n} \hat{p}(\mathbf{y}) \right) \\
& \stackrel{(d)}{\leq} (1 + \bar{F}_{\Lambda^1}(u))^n 2^{\frac{u-1}{2} \log n + 2}.
\end{aligned}$$

Equation (a) is obtained by reordering the symbols in the strings, Inequalities (b) and (c) follow respectively from Proposition 4 and Proposition 3. Inequality (d) is a direct consequence of the last inequality in Theorem 2.

Hence,

$$\begin{aligned}
R^*(\Lambda^n) & \leq n \log(1 + \bar{F}_{\Lambda^1}(u)) + \frac{u-1}{2} \log n + 2 \\
& \leq n \bar{F}_{\Lambda^1}(u) \log e + \frac{u-1}{2} \log n + 2
\end{aligned}$$

■

The next theorem complements the upper-bound on minimax regret for envelope classes (Theorem 4). It describes a general lower bound on minimax redundancy for envelope classes.

Theorem 5: Let f denote a non-increasing, summable envelope function. For any integer p , let $c(p) = \sum_{k=1}^p f(2k)$. Let $c(\infty) = \sum_{k \geq 1} f(2k)$. Assume furthermore that f satisfies property: $c(\infty) > 1$. Let $n \in \mathbb{N}_+$, $p \in \mathbb{N}_+$, $\epsilon > 0$ and $\lambda \in]0, 1[$ be such that $c(p) > 1$ and $(1 - \lambda)n \frac{f(2p)}{c(p)} > \frac{10}{\epsilon}$. Then

$$R^+(\Lambda_f^n) \geq C(p, n, \lambda, \epsilon) \sum_{i=1}^p \left(\frac{1}{2} \log \frac{n(1-\lambda)\pi f(2i)}{2c(p)e} - \epsilon \right), \quad (1)$$

where $C(p, n, \lambda, \epsilon) = \frac{1}{1 + \frac{c(p)}{\lambda^2 n f(2p)}} \left(1 - \frac{4}{\pi} \sqrt{\frac{10c(p)}{(1-\lambda)\epsilon n f(2p)}} \right)$.

Before proceeding to the proof, let us mention the following non-asymptotic bound from Xie and Barron (1997). Let m_n^* denote the Krichevsky-Trofimov distribution over $\{0, 1\}^n$. That is, for any $\mathbf{x} \in \{0, 1\}^n$, such that $n_1 = \sum_{i=1}^n x_i$ and $n_0 = n - n_1$

$$m_n^*(\mathbf{x}) = \pi \int_{[0,1]} \theta^{n_1-1/2} (1-\theta)^{n_0-1/2} d\theta.$$

Lemma 1: (Xie and Barron, 1997, Lemma 1) For any $\epsilon > 0$, there exists a $c(\epsilon)$ such that for $n > 2c(\epsilon)$ the following holds uniformly over $\theta \in [c(\epsilon)/n, 1 - c(\epsilon)/n]$:

$$\left| D(p_\theta^n, m_n^*) - \frac{1}{2} \log \frac{n}{2\pi e} - \log \pi \right| \leq \epsilon.$$

The bound $c(\varepsilon)$ can be chosen as small as $5/\varepsilon$.

Proof: Let f, n, p, ϵ and λ be as in the statement of the theorem. Let us first define a prior probability W on Λ_f^1 . For each integer i between 1 and p , let μ_i be defined as

$$\mu_i = \frac{f(2i)}{c(p)}.$$

Let $\boldsymbol{\theta} = (\theta_i)_{1 \leq i \leq p}$ be a collection of independent random variables each distributed according to a Beta distribution with parameters $(1/2, 1/2)$, hence the prior probability W has density w given by

$$w(\boldsymbol{\theta}) = \frac{1}{\pi^p} \prod_{i=1}^p \left(\theta_i^{-1/2} (1 - \theta_i)^{-1/2} \right).$$

The memoryless source parametrized by $\boldsymbol{\theta}$ is defined by the probability mass function $p_{\boldsymbol{\theta}}(2i - 1) = \theta_i \mu_i$ and $p_{\boldsymbol{\theta}}(2i) = (1 - \theta_i) \mu_i$ for $i : 1 \leq i \leq p$ and $p_{\boldsymbol{\theta}}(j) = 0$ for $j > 2i$. Thanks to the condition $c(p) > 1$, this probability mass function satisfies the envelope condition.

For $i \leq p$, let the random variable N_i (resp. N_i^0) be defined as the number of occurrences of $\{2i - 1, 2i\}$ (resp. $2i$) in the sequence \mathbf{x} . Let \mathbf{N} (resp. \mathbf{N}^0) denote the random vector N_1, \dots, N_p (resp. N_1^0, \dots, N_p^0). If a sequence \mathbf{x} from $\{1, \dots, 2p\}^n$ contains $n_i = N_i(\mathbf{x})$ symbols from $\{2i - 1, 2i\}$ for each $i \in \{1, \dots, p\}$, and if for each such i , the sequence contains $n_i^0 = N_i^0(\mathbf{x})$ (n_i^1) symbols equal to $2i - 1$ (resp. $2i$) then

$$P_{\boldsymbol{\theta}}^n(\mathbf{x}) = \prod_{i=1}^p \left(\mu_i^{n_i} \theta_i^{n_i^0} (1 - \theta_i)^{n_i^1} \right).$$

Note that when the source $\boldsymbol{\theta}$ is picked according to the prior W and the sequence $X_{1:n}$ picked according to $P_{\boldsymbol{\theta}}^n$, the random vector \mathbf{N} is multinomially distributed with parameters n and $(\mu_1, \mu_2, \dots, \mu_p)$, so the distribution of \mathbf{N} does not depend on the outcome of $\boldsymbol{\theta}$. Moreover, conditionally on \mathbf{N} , the conditional probability $P_{\boldsymbol{\theta}}^n \{ \cdot \mid \mathbf{N} \}$ is a product distribution:

$$P_{\boldsymbol{\theta}}^n(\mathbf{x} \mid \mathbf{N}) = \prod_{i=1}^p \left(\theta_i^{n_i^0} (1 - \theta_i)^{n_i^1} \right).$$

In statistical parlance, the random vectors \mathbf{N} and \mathbf{N}^0 form a sufficient statistic for $\boldsymbol{\theta}$.

Let Q^* denote the mixture distribution on \mathbb{N}_+^n induced by W :

$$Q^*(\mathbf{x}) = \mathbb{E}_W [P_{\boldsymbol{\theta}}^n(\mathbf{x})],$$

and, for each n , let m_n^* denote the Krichevsky-Trofimov mixture over $\{0, 1\}^n$, then

$$Q^*(\mathbf{x}) = \prod_{i=1}^p \left(\mu_i^{n_i} m_{n_i}^*(0^{n_i^0} 1^{n_i^1}) \right).$$

For a given value of \mathbf{N} , the conditional probability $Q^* \{ \cdot \mid \mathbf{N} \}$ is also a product distribution:

$$Q^* \{ \mathbf{x} \mid \mathbf{N} \} = \prod_{i=1}^p m_{N_i}^*(0^{n_i^0} 1^{n_i^1}),$$

so, we will be able to rely on:

$$\frac{P_{\theta}^n(\mathbf{x})}{Q^*(\mathbf{x})} = \frac{P_{\theta}^n(\mathbf{x} | \mathbf{N})}{Q^*(\mathbf{x} | \mathbf{N})}.$$

In the sequel, we let $P_{\theta}^n\{\cdot | \mathbf{N}\}$ ($Q^*\{\cdot | \mathbf{N}\}$) denote the conditional distribution of P_{θ}^n (Q_n^*) on strings with composition given by vector \mathbf{N} .

Now, the average redundancy of Q^* with respect to P_{θ}^n can be rewritten in a handy way.

$$\begin{aligned} \mathbb{E}_W [D(P_{\theta}^n, Q^*)] &= \mathbb{E}_W \left[\mathbb{E}_{P_{\theta}^n} \left[\log \frac{P_{\theta}^n(X_{1:n} | \mathbf{N})}{Q^*(X_{1:n} | \mathbf{N})} \right] \right] \\ &\quad \text{from the last equation,} \\ &= \mathbb{E}_W \left[\mathbb{E}_{P_{\theta}^n} \left[\mathbb{E}_{P_{\theta}^n} \left[\log \frac{P_{\theta}^n(X_{1:n} | \mathbf{N})}{Q^*(X_{1:n} | \mathbf{N})} | \mathbf{N} \right] \right] \right] \\ &= \mathbb{E}_W [\mathbb{E}_{P_{\theta}^n} [D(P_{\theta}^n(\cdot | \mathbf{N}), Q_n^*(\cdot | \mathbf{N}))]] \\ &= \mathbb{E}_W [\mathbb{E}_{\mathbf{N}} [D(P_{\theta}^n(\cdot | \mathbf{N}), Q^*(\cdot | \mathbf{N}))]] \\ &\quad \text{as the distribution of } \mathbf{N} \text{ does not depend on } \theta, \\ &= \mathbb{E}_{\mathbf{N}} [\mathbb{E}_W [D(P_{\theta}^n(\cdot | \mathbf{N}), Q^*(\cdot | \mathbf{N}))]] \\ &\quad \text{by Fubini's Theorem.} \end{aligned}$$

We may develop $D(P_{\theta}^n(\cdot | \mathbf{N}), Q^*(\cdot | \mathbf{N}))$ for a given value of $\mathbf{N} = (n_1, n_2, \dots, n_p)$. As both $P_{\theta}^n(\cdot | \mathbf{N})$ and $Q^*(\cdot | \mathbf{N})$ are product distributions on $\prod_{i=1}^p (\{2i-1, 2i\}^{n_i})$, we have

$$\begin{aligned} \mathbb{E}_W [D(P_{\theta}^n(\cdot | \mathbf{N}), Q^*(\cdot | \mathbf{N}))] &= \mathbb{E}_W \left[\sum_{i=1}^p D(P_{\theta_i}^{n_i}, m_{n_i}^*) \right] \\ &= \sum_{i=1}^p \mathbb{E}_{\theta_i} [D(P_{\theta_i}^{n_i}, m_{n_i}^*)]. \end{aligned}$$

The minimal average redundancy of Λ_f^n with respect to the mixing distribution W is thus finally given by:

$$\begin{aligned} \mathbb{E}_W [D(P_{\theta}^n, Q^*)] &= \mathbb{E}_{\mathbf{N}} \left[\sum_{i=1}^p \mathbb{E}_{\theta_i} [D(P_{\theta_i}^{n_i}, m_{n_i}^*)] \right] \\ &= \sum_{i=1}^p \mathbb{E}_{\mathbf{N}} [\mathbb{E}_{\theta_i} [D(P_{\theta_i}^{n_i}, m_{n_i}^*)]] \\ &= \sum_{i=1}^p \sum_{n_i=0}^n \binom{n}{n_i} \mu_i^{n_i} (1 - \mu_i)^{n-n_i} \mathbb{E}_{\theta_i} [D(P_{\theta_i}^{n_i}, m_{n_i}^*)]. \end{aligned} \quad (2)$$

Hence, the minimal redundancy of Λ_f^n with respect to prior probability W is a weighted average of redundancies of Krichevsky-Trofimov mixtures over binary strings with different lengths.

At some place, we will use the Chebychev-Cantelli inequality (see Devroye et al., 1996) which asserts that for a square-integrable random variable:

$$\Pr \{X \leq \mathbb{E}[X] - t\} \leq \frac{\text{Var}(X)}{\text{Var}(X) + t^2}.$$

Besides, note that $\forall \epsilon < \frac{1}{2}$,

$$\int_{\epsilon}^{1-\epsilon} \frac{dx}{\pi \sqrt{x(1-x)}} > 1 - \frac{4}{\pi} \sqrt{2\epsilon}. \quad (3)$$

Now, the proposition is derived by processing the right-hand-side of the last equation. Now, under condition $(1-\lambda) n \frac{f(2p)}{c(p)} > \frac{10}{\epsilon}$, we have $\forall i \leq p, n_i \geq (1-\lambda) n \mu_i \implies \frac{5}{n_i \epsilon} < \frac{1}{2}$. Hence,

$$\begin{aligned} & \mathbb{E}_W [D(P_{\theta}^n, Q^*)] \\ & \geq \sum_{i=1}^p \sum_{n_i \geq (1-\lambda) n \mu_i}^n \binom{n}{n_i} \mu_i^{n_i} (1-\mu_i)^{n-n_i} \int_0^1 \frac{D(P_{\theta_i}^{n_i}, m_{n_i}^*)}{\sqrt{\theta_i(1-\theta_i)}} d\theta_i \\ & \quad \text{by Proposition 1 from Xie and Barron (1997)} \\ & \geq \sum_{i=1}^p \sum_{n_i \geq (1-\lambda) n \mu_i}^n \binom{n}{n_i} \mu_i^{n_i} (1-\mu_i)^{n-n_i} \int_{\frac{5}{n_i \epsilon}}^{1-\frac{5}{n_i \epsilon}} \frac{\frac{1}{2} \log \frac{n_i}{2\pi e} + \log \pi - \epsilon}{\pi \sqrt{\theta_i(1-\theta_i)}} d\theta_i \\ & \quad \text{from (3)} \\ & \geq \sum_{i=1}^p \sum_{n_i \geq (1-\lambda) n \mu_i}^n \binom{n}{n_i} \mu_i^{n_i} (1-\mu_i)^{n-n_i} \left(1 - \frac{4}{\pi} \sqrt{\frac{10}{n_i \epsilon}}\right) \left(\frac{1}{2} \log \frac{n_i}{2\pi e} + \log \pi - \epsilon\right) \\ & \quad \text{invoking the Chebychef-Cantelli inequality,} \\ & \geq \sum_{i=1}^p \frac{1}{1 + \frac{1-\mu_i}{n \mu_i \lambda^2}} \left(1 - \frac{4}{\pi} \sqrt{\frac{10}{(1-\lambda) n \mu_i \epsilon}}\right) \left(\frac{1}{2} \log \frac{n(1-\lambda) \mu_i}{2\pi e} + \log \pi - \epsilon\right) \\ & \quad \text{using monotonicity assumption on } f \\ & \geq \frac{1}{1 + \frac{c(p)}{\lambda^2 n f(2p)}} \left(1 - \frac{4}{\pi} \sqrt{\frac{10c(p)}{(1-\lambda) \epsilon n f(2p)}}\right) \sum_{i=1}^p \left(\frac{1}{2} \log \frac{n(1-\lambda) f(2i)}{2c(p) \pi e} + \log \pi - \epsilon\right). \end{aligned}$$

■

IV. EXAMPLES OF ENVELOPE CLASSES

Theorems 3, 4 and 5 assert that the summability of the envelope defining a class of memoryless sources characterizes the (strong) universal compressibility of that class. However, it is not easy to figure out whether the bounds provided by the last two theorems are close to each other or not. In this Section, we investigate the case of envelopes which decline either like power laws or exponentially fast. In both cases, upper-bounds on minimax regret will follow directly from Theorem 4 and a straightforward optimization. Specific lower bounds on minimax redundancies are derived by mimicking the proof of Theorem 5, either faithfully as in the case of exponential envelopes or by developing an alternative prior as in the case of power-law envelopes.

A. Power-law envelope classes

Let us first agree on the classical notation: $\zeta(\alpha) = \sum_{k \geq 1} \frac{1}{k^\alpha}$, for $\alpha > 1$.

Theorem 6: Let α denote a real number larger than 1, and C be such that $C\zeta(\alpha) \geq 2^\alpha$. The source class $\Lambda_{C, -\alpha}$ is the envelope class associated with the decreasing function $f_{\alpha, C} : x \mapsto 1 \wedge \frac{C}{x^\alpha}$ for $C > 1$ and $\alpha > 1$.

Then:

1)

$$n^{1/\alpha} A(\alpha) \log [C\zeta(\alpha)] \leq R^+(\Lambda_{C,-\alpha}^n)$$

where

$$A(\alpha) = \frac{1}{\alpha} \int_1^\infty \frac{1}{u^{1-1/\alpha}} \left(1 - e^{-1/(\zeta(\alpha)u)}\right) du.$$

2)

$$R^*(\Lambda_{C,-\alpha}^n) \leq \left(\frac{2Cn}{\alpha-1}\right)^{1/\alpha} (\log n)^{1-1/\alpha} + O(1).$$

Remark 9: The gap between the lower-bound and the upper-bound is of order $(\log n)^{1-\frac{1}{\alpha}}$. We are not in a position to claim that one of the two bounds is tight, let alone which one is tight. Note however that as $\alpha \rightarrow \infty$ and $C = H^\alpha$, class $\Lambda_{C,-\alpha}$ converges to the class of memoryless sources on alphabet $\{1, \dots, H\}$ for which the minimax regret is $\frac{H-1}{2} \log n$. This is (up to a factor 2) what we obtain by taking the limits in our upper-bound of $R^*(\Lambda_{C,-\alpha}^n)$. On the other side, the limit of our lower-bound when α goes to 1 is infinite, which is also satisfying since it agrees with Theorem 3.

Remark 10: In contrast with various lower bounds derived using a similar methodology, the proof given here relies on a single prior probability distribution on the parameter space, it works for all values of n . It has been elaborated after helpful discussions with Laszló Györfi.

Note that the lower bound that can be derived from Theorem 5 is of the same order of magnitude $O(n^{1/\alpha})$ as the lower bound stated here (see Appendix II). The proof given here is completely elementary and does not rely on the subtle computations described in Xie and Barron (1997).

Proof: For the upper-bound on minimax regret, note that

$$\bar{F}_{\alpha,C}(u) = \sum_{k>u} 1 \wedge \frac{C}{k^\alpha} \leq \frac{C}{(\alpha-1)u^{\alpha-1}}.$$

Hence, choosing $u_n = \left(\frac{2Cn}{(\alpha-1)\log n}\right)^{\frac{1}{\alpha}}$, resorting to Theorem 4, we get:

$$R^*(\Lambda_{C,-\alpha}^n) \leq \left(\frac{2Cn}{\alpha-1}\right)^{1/\alpha} (\log n)^{1-1/\alpha} + O(1).$$

Let us now turn to the lower bound. We first define a set Θ of parameters such that $P_\theta^n \in \Lambda_{\alpha,C}^n$ for any $\theta \in \Theta$ and then we use the mutual information lower bound on redundancy. Let m be a positive integer such that $m^\alpha < C\zeta(\alpha)$. For all sufficiently large integer p , $m^\alpha \leq \sum_{k=1}^p \frac{C}{k^\alpha}$. Henceforth, let $c(p) = \sum_{k=1}^p \frac{C}{(km)^\alpha}$, so that the condition $m^\alpha < C\zeta(\alpha)$ translates into $c(p) \geq 1$.

The set $\{P_\theta, \theta \in \Theta\}$ consists of memoryless sources over the infinite alphabet \mathbb{N}_+ . Each parameter θ is a sequence of integers $\theta = (\theta_1, \theta_2, \dots)$. We take a prior distribution on Θ such that $(\theta_k)_k$ is a sequence of independent identically distributed random variables with uniform distribution on $\{1, \dots, m\}$. For any such θ , P_θ^1

is a probability distribution on \mathbb{N}_+ with support $\cup_{k \geq 1} \{(k-1)m + \theta_k\}$, namely:

$$P_{\theta}((k-1)m + \theta_k) = \frac{C}{c(p)(km)^\alpha} \quad \text{for } k = 1, \dots, p,$$

The condition $c_p \geq 1$ ensures that $P_{\theta}^1 \in \Lambda_{\alpha, C}^1$.

Now, the mutual information between parameter θ and source output $X_{1:n}$ is

$$I(\theta, X_{1:n}) = \sum_{k \geq 1} I(\theta_k, X_{1:n})$$

Let $N_k(\mathbf{x}) = 1$ if there exists some index $i \in \{1, \dots, n\}$ such that $\mathbf{x}_i \in [(k-1)m + 1, km]$, and 0 otherwise.

Note that the distribution of N_k does not depend on the value of θ . Thus we can write:

$$\begin{aligned} I(\theta_k, X_{1:n}) &= \sum_{j=1}^m \sum_{x \in \mathbb{N}_+^n} P(\theta_k = j, X_{1:n} = x) \log \frac{P(\theta_k = j | X_{1:n} = x)}{P(\theta_k = j)} \\ &= \sum_{j=1}^m \sum_{x: N_k(x)=0} P(\theta_k = j, X_{1:n} = x) \log \frac{1/m}{1/m} + \sum_{j=1}^m \sum_{x: N_k(x)=1} P(\theta_k = j, X_{1:n} = x) \log \frac{P(\theta_k = j | X_{1:n} = x)}{1/m} \\ &= 0 + \sum_{x: N_k(x)=1, P_{\theta}(x) > 0} \sum_{j=1}^m P(\theta_k = j, X_{1:n} = x) \log \frac{1}{1/m} \\ &= P(N_k(X_{1:n}) = 1) \log m. \end{aligned}$$

Hence,

$$\begin{aligned} I(\theta, X_{1:n}) &= \sum_{k \geq 1} P(N_k = 1) \log m \\ &= \mathbb{E}_{P_{\theta}}[Z_n] \log m, \end{aligned}$$

where $Z_n(\mathbf{x})$ denotes the number of distinct symbols in string \mathbf{x} (note that its distribution does not depend on the value of θ .)

As $Z_n = \sum_{k \geq 1} 1_{N_k(x)=1}$, the expectation translates into a sum

$$\mathbb{E}_{\theta}[Z_n] = \sum_{k=1}^{\infty} \left(1 - \left(1 - \frac{C}{c(p)(km)^\alpha} \right)^n \right)$$

which leads to:

$$R^+(\Lambda_{\alpha, C}^n) \geq \left(\sum_{k=1}^{\infty} \left(1 - \left(1 - \frac{C}{c(p)(km)^\alpha} \right)^n \right) \right) \times \log m.$$

In order to optimize the bound we choose $m = \lfloor C\zeta(\alpha) \rfloor$. The last sum can then be lower-bounded by:

$$\begin{aligned}
& \sum_{k=1}^{\infty} \left(1 - \left(1 - \frac{C m^\alpha}{C\zeta(\alpha)(k m)^\alpha} \right)^n \right) \\
& \geq \sum_{k=1}^{\infty} \left(1 - \exp \left(-\frac{n}{\zeta(\alpha) k^\alpha} \right) \right) \\
& \quad \text{as } 1 - x \leq \exp(-x) \\
& \geq \int_1^\infty \left(1 - \exp \left(-\frac{n}{\zeta(\alpha) x^\alpha} \right) \right) dx \\
& \geq \frac{n^{\frac{1}{\alpha}}}{\alpha} \int_1^\infty \frac{1}{u^{1-\frac{1}{\alpha}}} \left(1 - \exp \left(-\frac{1}{\zeta(\alpha) u} \right) \right) du.
\end{aligned}$$

■

For an alternative derivation of a similar lower-bound using Theorem 5, see Appendix II.

B. Exponential envelope classes

Theorems 4 and 5 provide almost matching bounds on the minimax redundancy for source classes defined by exponentially vanishing envelopes.

Theorem 7: Let C and α denote positive real numbers satisfying $C > e^{2\alpha}$. The class $\Lambda_{C e^{-\alpha \cdot}}$ is the envelope class associated with function $f_\alpha : x \mapsto 1 \wedge C e^{-\alpha x}$. Then

$$\frac{1}{8\alpha} \log^2 n (1 - o(1)) \leq R^+(\Lambda_{C e^{-\alpha \cdot}}^n) \leq R^*(\Lambda_{C e^{-\alpha \cdot}}^n) \leq \frac{1}{2\alpha} \log^2 n + O(1)$$

Proof: For the upper-bound, note that as $u \rightarrow \infty$,

$$\bar{F}_\alpha(u) = \sum_{k>u} 1 \wedge C e^{-\alpha k} \leq \frac{C}{1 - e^{-\alpha}} e^{-\alpha(u+1)}.$$

Hence, by choosing the optimal value $u_n = \frac{1}{\alpha} \log n$ in Theorem 4 we get:

$$R^*(\Lambda_{C e^{-\alpha \cdot}}^n) \leq \frac{1}{2\alpha} \log^2 n + O(1).$$

We will now prove the lower bound using Theorem 5. The constraint $C > e^{2\alpha}$ warrants that the sequence $c(p) = \sum_{k=1}^p f(2k) \geq C e^{-2\alpha} \frac{1 - e^{-2\alpha p}}{1 - e^{-2\alpha}}$ is larger than 1 for all p .

If we choose $p = \lfloor \frac{1}{2\alpha} (\log n - \log \log n) \rfloor$, then $n f(2p) > C n e^{-\log n + \log \log n - 2\alpha}$ goes to infinity with n . For $\epsilon = \lambda = \frac{1}{2}$, we get $C(p, n, \lambda, \epsilon) = 1 - o(1)$. Besides,

$$\begin{aligned}
\sum_{i=1}^p \left(\frac{1}{2} \log \frac{n(1-\lambda)C\pi e^{-2\alpha i}}{2c(p)e} - \epsilon \right) &= \frac{p}{2} \left(\log n + \log \frac{(1-\lambda)C\pi}{2c(p)e} - 2\epsilon \right) - \alpha \sum_{i=1}^p i \\
&= \left(\frac{1}{4\alpha} \log^2 n - \frac{\alpha}{2} \frac{1}{4\alpha^2} \log^2 n \right) (1 + o(1)) \\
&= \frac{1}{8\alpha} \log^2 n (1 + o(1)).
\end{aligned}$$

■

V. A CENSORING CODE FOR ENVELOPE CLASSES

The proofs of Theorems 4 and 5 suggest to handle separately small and large (allegedly infrequent) symbols. Such an algorithm should perform quite well as soon as the tail behavior of the envelope provides an adequate description of the sources in the class. The coding algorithm suggested by the proof of Theorem 4, which are based on the Shtarkov NML coder, are not computationally attractive. The design of the next algorithm (CensoringCode) is again guided by the proof of Proposition 4: it is parametrized by a sequence of cutoffs $(K_i)_{i \in \mathbb{N}}$ and handles the i^{th} symbol of the sequence to be compressed differently according to whether it is smaller or larger than cutoff K_i , in the latter situation, the symbol is said to be censored. The CensoringCode algorithm uses Elias penultimate code (Elias, 1975) to encode censored symbols and Krichevsky-Trofimov mixtures (Krichevsky and Trofimov, 1981) to encode the sequence of non-censored symbols padded with markers (zeros) to witness acts of censorship. The performance of this algorithm is evaluated on the power-law envelope class $\Lambda_{C, -\alpha}$, already investigated in Section IV. In this section, the parameters α and C are assumed to be known.

Let us first describe the algorithm more precisely. Given a non-decreasing sequence of cutoffs $(K_i)_{i \leq n}$, a string \mathbf{x} from \mathbb{N}_+^n defines two strings $\tilde{\mathbf{x}}$ and $\check{\mathbf{x}}$ in the following way. String $\tilde{\mathbf{x}}$ has length n and belongs to $\prod_{i=1}^n \mathcal{X}_i$, where $\mathcal{X}_i = \{0, \dots, K_i\}$:

$$\tilde{\mathbf{x}}_i = \begin{cases} \mathbf{x}_i & \text{if } \mathbf{x}_i \leq K_i \\ 0 & \text{otherwise (the symbol is censored),} \end{cases}$$

Meanwhile, string $\check{\mathbf{x}}$ is the subsequence of censored symbols, that is $(\mathbf{x}_i)_{\mathbf{x}_i > K_i, i \leq n}$.

The algorithm encodes \mathbf{x} as a pair of binary strings C1 and C2. The first one (C1) is obtained by applying Elias penultimate code to each symbol from $\tilde{\mathbf{x}}$, that is to each censored symbol. The second string (C2) is built by applying arithmetic coding to $\tilde{\mathbf{x}}$ using side-information from $\check{\mathbf{x}}$. Decoding C2 can be carried out using information gotten from decoding C1.

In order to describe the coding probability used to encode $\tilde{\mathbf{x}}$, we need a few more counters. For $j > 0$, let n_i^j be the number of occurrences of symbol j in $\mathbf{x}_{1:i}$ and let n_i^0 be the number of symbols larger than K_{i+1} in $\mathbf{x}_{1:i}$ (note that this not larger than the number of censored symbols in $\mathbf{x}_{1:i}$, and that the counters n_i^j can be recovered from $\tilde{\mathbf{x}}_{1:i}$ and $\check{\mathbf{x}}$). The conditional coding probability over alphabet $\mathcal{X}_{i+1} = \{0, \dots, K_{i+1}\}$ given $\tilde{\mathbf{x}}_{1:i}$ and $\check{\mathbf{x}}$ is derived from the Krichevsky-Trofimov mixture over \mathcal{X}_{i+1} . It is the posterior distribution corresponding to Jeffrey's prior on the $1 + K_{i+1}$ -dimensional probability simplex and counts n_i^j for j running from 0 to K_{i+1} :

$$Q(\tilde{X}_{i+1} = j \mid \tilde{X}_{1:i} = \tilde{\mathbf{x}}_{1:i}, \check{X} = \check{\mathbf{x}}) = \frac{n_i^j + \frac{1}{2}}{i + \frac{K_{i+1} + 1}{2}}.$$

The length of C2(\mathbf{x}) is (up to a quantity smaller than 1) given by

$$- \sum_{i=0}^{n-1} \log Q(\tilde{\mathbf{x}}_{i+1} \mid \mathbf{x}_{1:i}, \check{\mathbf{x}}) = - \log Q(\tilde{\mathbf{x}} \mid \check{\mathbf{x}}).$$

The following description of the coding probability will prove useful. For $1 \leq j \leq K_n$, let s^j be the number of censored occurrences of symbol j . Let n^j serve as a shorthand for n_n^j . Let $i(j)$ be n if $K_n < j$ or the largest

integer i such that K_i is smaller than j , then $s_j = n_{i(j)}^j$. The following holds

$$Q(\tilde{\mathbf{x}} \mid \mathbf{\tilde{x}}) = \left(\prod_{j=1}^{K_n} \frac{\Gamma(n^j + 1/2)}{\Gamma(s^j + 1/2)} \right) \left(\prod_{i:\mathbf{\tilde{x}}_i=0} (n_{i-1}^0 + 1/2) \right) \left(\prod_{i=0}^{n-1} \frac{1}{i + \frac{K_i+1}{2}} \right)$$

Note that the sequence $(n_i^0)_{i \leq n}$ is not necessarily non-decreasing.

A technical description of algorithm `CensoringCode` is given below. The procedure `EliasCode` takes as input an integer j and outputs a binary encoding of j using exactly $\ell(j)$ bits where ℓ is defined by: $\ell(j) = \lfloor \log j + 2 \log(1 + \log j) + 1 \rfloor$. The procedure `ArithCode` builds on the arithmetic coding methodology (Rissanen and Langdon, 1979). It is enough to remember that an arithmetic coder takes advantage of the fact that a coding probability Q is completely defined by the sequence of conditional distributions of the $i + 1$ th symbol given the past up to time i .

The proof of the upper-bound in Theorem 6, prompts us to choose $K_i = \lambda i^{\frac{1}{\alpha}}$ (it will appear at the end that the best choice is $\lambda = \left(\frac{2C}{\alpha-1} \right)^{\frac{1}{\alpha}}$).

Algorithm 1 `CensoringCode`

```

 $K \leftarrow 0$ 
 $counts \leftarrow [1/2, 1/2, \dots]$ 
for  $i$  from 1 to  $n$  do
   $cutoff \leftarrow \left\lfloor \left( 2 \frac{C_i}{\alpha-1} \right)^{1/\alpha} \right\rfloor$ 
  if  $cutoff > K$  then
    for  $j \leftarrow K + 1$  to  $cutoff$  do
       $counts[0] \leftarrow counts[0] - counts[j] + 1/2$ 
    end for
     $K \leftarrow cutoff$ 
  end if
  if  $x[i] \leq cutoff$  then
    ArithCode( $x[i], counts[0 : cutoff]$ )
  else
    ArithCode(0,  $counts[0 : cutoff]$ )
     $C1 \leftarrow C1 \cdot \text{EliasCode}(x[i])$ 
     $counts[0] \leftarrow counts[0] + 1$ 
  end if
   $counts[x[i]] \leftarrow counts[x[i]] + 1$ 
end for
 $C2 \leftarrow \text{ArithCode}()$ 
return  $C1 \cdot C2$ 

```

Theorem 8: Let C and α be positive reals. Let the sequence of cutoffs $(K_i)_{i \leq n}$ be given by

$$K_i = \left\lfloor \left(\frac{2C}{\alpha-1} i \right)^{1/\alpha} \right\rfloor.$$

The expected redundancy of procedure `CensoringCode` on the envelope class $\Lambda_{C, -\alpha}$ is not larger than

$$\left(\frac{2Cn}{\alpha-1} \right)^{\frac{1}{\alpha}} \log n (1 + o(1)).$$

Remark 11: The redundancy upper-bound in this Theorem is within a factor $\log n$ from the lower bound $O(n^{1/\alpha})$ from Theorem 6 .

The proof of the Theorem builds on the next two lemmas. The first lemma compares the length of $C2(\mathbf{x})$ with a tractable quantity. The second lemma upper-bounds the average length of $C1(\mathbf{x})$ by a quantity which is of the same order of magnitude as the upper-bound on redundancy we are looking for.

We need a few more definitions. Let \mathbf{y} be the string of length n over alphabet \mathcal{X}_n defined by:

$$\mathbf{y}_i = \begin{cases} \mathbf{x}_i & \text{if } \mathbf{x}_i \leq K_n; \\ 0 & \text{else.} \end{cases}$$

For $0 \leq j \leq K_n$, note that the previously defined shorthand n^j is the number of occurrences of symbol j in \mathbf{y} . The string \mathbf{y} is obtained from \mathbf{x} in the same way as $\tilde{\mathbf{x}}$ using the constant cutoff K_n .

Let m_n^* be the Krichevsky-Trofimov mixture over alphabet $\{0, \dots, K_n\}$:

$$m_n^*(\mathbf{y}) = \left(\prod_{j=0}^{K_n} \frac{\Gamma(n^j + \frac{1}{2})}{\Gamma(1/2)} \right) \frac{\Gamma(\frac{K_n+1}{2})}{\Gamma(n + \frac{K_n+1}{2})}.$$

String $\tilde{\mathbf{x}}$ seems easier to encode than \mathbf{y} since it is possible to recover $\tilde{\mathbf{x}}$ from \mathbf{y} . This observation does not however warrant automatically that the length of $C2(\mathbf{x})$ is not significantly larger than any reasonable codeword length for \mathbf{y} . Such a guarantee is provided by the following lemma.

Lemma 2: For every string $\mathbf{x} \in \mathbb{N}_+^n$, the length of the $C2(\mathbf{x})$ is not larger than $-\log m_n^*(\mathbf{y})$.

Proof: [Lemma 2] Let s^0 be the number of occurrences of 0 in \mathbf{y} , that is the number of symbols in \mathbf{x} that are larger than K_n . Let

$$T_0 = \prod_{i=1, \tilde{\mathbf{x}}_i=0}^n (n_{i-1}^0 + 1/2).$$

Then, the following holds:

$$\begin{aligned} T_0 &\stackrel{(a)}{=} \left(\prod_{i=1, \tilde{\mathbf{y}}_i=0}^n (n_{i-1}^0 + 1/2) \right) \prod_{j=1}^{K_n} \left(\prod_{i=1, \tilde{\mathbf{y}}_i=j}^{i(j)} (n_{i-1}^0 + 1/2) \right) \\ &\stackrel{(b)}{\geq} \left(\frac{\Gamma(s^0 + 1/2)}{\Gamma(1/2)} \right) \prod_{j=1}^{K_n} \left(\frac{\Gamma(s^j + 1/2)}{\Gamma(1/2)} \right), \end{aligned}$$

where (a) follows from the fact symbol \mathbf{x}_i is censored either because $\mathbf{x}_i > K_n$ (that is $\mathbf{y}_i = 0$) or because $\mathbf{x}_i = j \leq K_n$ and $i \leq i(j)$; (b) follows from the fact that for each $i \leq n$ such that $\mathbf{y}_i = 0$, $n_{i-1}^0 \geq \sum_{i' < i} \mathbf{1}_{\mathbf{x}_{i'} > K_n}$ while for each $j, 0 < j \leq K_n$, for each $i \leq i(j)$, $n_{i-1}^0 \geq n_{i-1}^j$.

From the last inequality, it follows that

$$\begin{aligned}
Q(\tilde{\mathbf{x}} \mid \tilde{\mathbf{x}}) &\geq \left(\prod_{j=1}^{K_n} \frac{\Gamma(n^j + 1/2)}{\Gamma(s^j + 1/2)} \right) \frac{\Gamma(s^0 + 1/2)}{\Gamma(1/2)} \prod_{j=1}^{K_n} \frac{\Gamma(s^j + 1/2)}{\Gamma(1/2)} \left(\prod_{i=0}^{n-1} \frac{1}{i + \frac{K_i+1}{2}} \right) \\
&\geq \left(\prod_{j=1}^{K_n} \frac{\Gamma(n^j + 1/2)}{\Gamma(s^j + 1/2)} \right) \frac{\Gamma(s^0 + 1/2)}{\Gamma(1/2)} \prod_{j=1}^{K_n} \frac{\Gamma(s^j + 1/2)}{\Gamma(1/2)} \left(\prod_{i=0}^{n-1} \frac{1}{i + \frac{K_n+1}{2}} \right) \\
&= m_n^*(\mathbf{y}),
\end{aligned}$$

where the last inequality holds since $(K_i)_i$ is a non-decreasing sequence. \blacksquare

The next lemma shows that the expected length of $\text{C1}(X_{1:n})$ is not larger than the upper-bound we are looking for.

Lemma 3: For every source $P \in \Lambda_{C,-\alpha}$, the expected length of the encoding of the censored symbols ($\text{C1}(X_{1:n})$) satisfies:

$$\mathbb{E}_P[|\text{C1}(X_{1:n})|] \leq \frac{C}{(\alpha-1)\lambda^{\alpha-1}} n^{\frac{1}{\alpha}} \log n (1 + o(1)).$$

Proof: [Lemma 3] Let $1 \leq a < b$ and $\beta > 0$. First note that:

$$\begin{aligned}
\int_a^b \frac{1}{x^\beta} dx &= \left[\frac{1}{(1-\beta)x^{\beta-1}} \right]_a^b, \\
\int_a^b \frac{\log x}{x^\beta} dx &= \left[\frac{\log x - \frac{1}{1-\beta}}{(1-\beta)x^{\beta-1}} \right]_a^b, \\
\int_a^b \frac{\log(1 + \log x)}{x^\beta} dx &= \left[\frac{\log(1 + \log x) - \frac{1}{1-\beta}}{(1-\beta)x^{\beta-1}} \right]_a^b.
\end{aligned}$$

The expected length of the first part of the code is thus:

$$\begin{aligned}
\mathbb{E}[|\text{C1}(X_1^n)|] &= \mathbb{E} \left[\sum_{j=1}^n \ell(X_j) \mathbb{1}_{X_j > K_j} \right] \\
&= \sum_{j=1}^n \sum_{x=K_j+1}^{\infty} \ell(x) P(x) \\
&\leq \sum_{j=1}^n \sum_{x=K_j+1}^{\infty} \ell(x) \frac{C}{x^\alpha} \\
&\leq C \sum_{j=1}^n \sum_{x=K_j+1}^{\infty} \frac{\log(x) + 2 \log(1 + \log x) + 1}{x^\alpha}.
\end{aligned}$$

Using the expressions for the integrals above, we get:

$$\begin{aligned}
\sum_{x=K_j+1}^{\infty} \frac{\log x + 2 \log(1 + \log x) + 1}{x^\alpha} &\leq \int_{K_j}^{\infty} \frac{\log x + 2 \log(1 + \log x) + 1}{x^\alpha} dx \\
&\leq \frac{\log K_j + 2 \log(1 + \log K_j) + \frac{4}{\alpha-1}}{(\alpha-1) K_j^{\alpha-1}}.
\end{aligned}$$

Thus, as $K_j = \lambda j^{1/\alpha}$, we substitute β by $1 - \frac{1}{\alpha}$ in Equations (4), (4) and (4) to obtain:

$$\begin{aligned}
\mathbb{E}[|\mathcal{C}1(X_1^n)|] &\leq \frac{C}{\alpha-1} \sum_{j=1}^n \frac{\log K_j + 2 \log(1 + \log K_j) + \frac{4}{\alpha-1}}{K_j^{\alpha-1}} \\
&= \frac{C}{\alpha-1} \sum_{j=1}^n \frac{\frac{1}{\alpha} \log j + \log \lambda + 2 \log(1 + \log(\lambda j^{1/\alpha})) + \frac{4}{\alpha-1}}{\lambda^{\alpha-1} j^{1-\frac{1}{\alpha}}} \\
&\leq \frac{C}{\alpha(\alpha-1)\lambda^{\alpha-1}} \left(C + \int_{x=2}^{n+1} \frac{(\log x + \alpha \log \lambda + 2\alpha \log(1 + \log(\lambda x^{1/\alpha})) + \frac{4\alpha}{\alpha-1})}{x^{1-\frac{1}{\alpha}}} dx \right) \\
&= \frac{C}{(\alpha-1)\lambda^{\alpha-1}} n^{\frac{1}{\alpha}} \log n (1 + o(1)).
\end{aligned}$$

■

We may now complete the proof of Theorem 8.

Proof: Remember that $\mathcal{X}_n = \{0, \dots, K_n\}$. If p is a probability mass function over alphabet \mathcal{X} , let $p^{\otimes n}$ be the probability mass function over \mathcal{X}^n defined by $p^{\otimes n}(\mathbf{x}) = \prod_{i=1}^n p(\mathbf{x}_i)$. Note that for every string $\mathbf{x} \in \mathbb{N}_+^n$,

$$\max_{p \in \mathfrak{M}_1(\mathcal{X}_n)} p^{\otimes n}(\mathbf{y}) \geq \max_{p \in \mathfrak{M}_1(\mathbb{N}_+)} p^{\otimes n}(\mathbf{x}) \geq \max_{P \in \Lambda_{C, -\alpha}} P^n(\mathbf{x}) = \hat{p}(\mathbf{x}).$$

Together with Lemma 2 and the bounds on the redundancy of the Krichevsky-Trofimov mixture (See Krichevsky and Trofimov, 1981), this implies:

$$|\mathcal{C}2(\mathbf{x})| \leq -\log \hat{p}(\mathbf{x}) + \frac{K_n}{2} \log n + O(1).$$

Let $L(\mathbf{x})$ be the length of the code produced by algorithm `CensoringCode` on the input string \mathbf{x} , then

$$\begin{aligned}
&\sup_{P \in \Lambda_{C, -\alpha}} \mathbb{E}_P [L(X_{1:n}) - \log 1/P^n(X_{1:n})] \\
&\leq \sup_{P \in \Lambda_{C, -\alpha}} \mathbb{E}_P [L(X_{1:n}) - \log 1/\hat{p}(X_{1:n})] \\
&\leq \sup_{P \in \Lambda_{C, -\alpha}^n} \mathbb{E}_P [|\mathcal{C}2(X_{1:n})| + \log \hat{p}(X_{1:n}) + |\mathcal{C}1(X_{1:n})|] \\
&\leq \sup_{\mathbf{x}} (|\mathcal{C}2(\mathbf{x})| + \log \hat{p}(\mathbf{x})) + \sup_{P \in \Lambda_{C, -\alpha}^n} \mathbb{E}_P [|\mathcal{C}1(X_{1:n})|] \\
&\leq \frac{\lambda n^{\frac{1}{\alpha}}}{2} \log n + \frac{C}{(\alpha-1)\lambda^{\alpha-1}} n^{\frac{1}{\alpha}} \log n (1 + o(1)).
\end{aligned}$$

The optimal value is $\lambda = \left(\frac{2C}{\alpha-1}\right)^{\frac{1}{\alpha}}$, for which we get:

$$R^+(Q^n, \Lambda_{C, -\alpha}^n) \leq \left(\frac{2Cn}{\alpha-1}\right)^{\frac{1}{\alpha}} \log n (1 + o(1)).$$

■

Note that the proof of the theorem provides an upper-bound on the expected regret of the censoring code.

VI. ADAPTIVE ALGORITHMS

The performance of `CensoringCode` depends on the fit of the cutoffs sequence to the tail behavior of the envelope. From the proof of Theorem 8, it should be clear that if `CensoringCode` is fed with a source which marginal is light-tailed, it will be unable to take advantage of this, and will suffer from excessive redundancy.

In this section, a sequence $(Q^n)_n$ of coding probabilities is said to be *approximately asymptotically adaptive* with respect to a collection $(\Lambda_m)_{m \in \mathcal{M}}$ of source classes if for each $P \in \cup_{m \in \mathcal{M}} \Lambda_m$, for each Λ_m such that $P \in \Lambda_m$:

$$D(P^n, Q^n) / R^+(\Lambda_m^n) \in O(\log n).$$

Such a definition makes sense, since we are considering massive source classes which minimax redundancies are large with respect to the logarithm function.

A. Pattern coding

First, the use of *pattern coding* Orlitsky et al. (2004), Shamir (2006) leads to an almost minimax adaptive procedure for small values of α , that is heavy-tailed distributions. Let us introduce the notion of pattern using the example of string $\mathbf{x} = \text{"abracadabra"}$, which is made of $n = 11$ characters. The information it conveys can be separated in two blocks:

- 1) a *dictionary* $\Delta = \Delta(\mathbf{x})$: the sequence of distinct symbols occurring in \mathbf{x} in order of appearance (in the example, $\Delta = (a, b, r, c, d)$).
- 2) a *pattern* $\psi = \psi(\mathbf{x})$ where ψ_i is the rank of \mathbf{x}_i in the dictionary Δ (here, $\psi = 1231415123$).

Now, consider the algorithm coding message \mathbf{x} by transmitting successively

- 1) the dictionary $\Delta_n = \Delta(\mathbf{x})$ (by concatenating the Elias codes for successive symbols);
- 2) and the pattern $\Psi_n = \psi(\mathbf{x})$, using a procedure for coding patterns as suggested by Orlitsky et al. (2004) or Shamir (2006). Henceforth, the latter procedure is called pattern coding.

Theorem 9: Let Q^n denote the coding probability associated with the coding algorithm which consists in applying Elias penultimate coding to the dictionary $\Delta(\mathbf{x})$ of a string \mathbf{x} from \mathbb{N}_+^n and then pattern coding to the pattern $\psi(\mathbf{x})$.

Then for any α such that $1 \leq \alpha \leq 5/2$, there exists a constant K depending on α and C such that

$$R^+(Q^n, \Lambda_{C, -\alpha}^n) \leq K n^{1/\alpha} \log n$$

Proof: For a given value of C and α , the Elias encoding of the dictionary uses on average

$$\mathbb{E}[|\Delta_n|] = K' n^{\frac{1}{\alpha}} \log n$$

bits (as proved in Appendix IV), for some constant K' depending on α and C .

If our pattern coder reaches (approximately) the minimax pattern redundancy

$$R_{\Psi}^+(\Psi_{1:n}) = \inf_{q \in \mathfrak{M}_1(\mathbb{N}_+^n)} \sup_{P \in \mathfrak{M}_1(\mathbb{N}_+)} \mathbb{E}_P \left[\log \frac{P^{\otimes n}(\Psi_{1:n})}{q(\Psi_{1:n})} \right],$$

the encoding of the pattern uses on average

$$H(\Psi_{1:n}) + R_{\Psi}^+(\Psi_{1:n}) \leq H(X_{1:n}) + R_{\Psi}^+(\Psi_{1:n}) \text{ bits.}$$

But in Orlitsky et al. (2004), the authors show that $R_{\Psi}^+(\Psi_{1:n})$ is upper-bounded by $O(\sqrt{n})$ and even $O(n^{\frac{2}{5}})$ according to Shamir (2004) (actually, these bounds are even satisfied by the minimax individual pattern redundancy). Hence, this code is adaptive, up to a factor $\log n$, in the range $1 < \alpha \leq \frac{5}{2}$. ■

This remarkably simple method is however expected to have a poor performance when α is large. Indeed, Garivier (2006) proves that $R_{\Psi}^+(\Psi_{1:n})$ is lower-bounded by $1.84 \left(\frac{n}{\log n} \right)^{\frac{1}{3}}$, which indicates that pattern coding is probably suboptimal as soon as α is larger than 3.

B. An (approximately) adaptive censoring code

Given the limited scope of the pattern coding method, we will attempt to turn the censoring code into an adaptive method, that is to tune the cutoff sequence so as to model the source statistics. As the cutoffs are chosen in such a way that they model the tail-heaviness of the source, we are facing a tail-heaviness estimation problem. In order to focus on the most important issues we do not attempt to develop a sequential algorithm. The $n + 1$ th cutoff \hat{K}_{n+1} is chosen according to the number of *distinct* symbols $Z_n(\mathbf{x})$ in \mathbf{x} .

This is a reasonable method if the probability mass function defining the source statistics P^1 actually decays like $\frac{1}{k^\alpha}$. Unfortunately, sparse distributions consistent with $\Lambda_{-\alpha}$ may lead this project astray. If, for example, $(Y_n)_n$ is a sequence of geometrically distributed random variables, and if $X_n = \left\lfloor 2^{\frac{Y_n}{\alpha}} \right\rfloor$, then the distribution of the X_n just fits in $\Lambda_{C, -\alpha}$ but obviously $Z_n(X_{1:n}) = Z_n(Y_{1:n}) = O(\log n)$.

Thus, rather than attempting to handle $\cup_{\alpha>0} \Lambda_{-\alpha}$, we focus on subclasses $\cup_{\alpha>0} \mathcal{W}_\alpha$, where

$$\mathcal{W}_\alpha = \left\{ P : P \in \Lambda_{-\alpha}, 0 < \liminf_k k^\alpha P^1(k) \leq \limsup_k k^\alpha P^1(k) < \infty \right\}.$$

The rationale for tuning cutoff \hat{K}_n using Z_n comes from the following two propositions.

Proposition 7: For every memoryless source $P \in \mathcal{W}_\alpha$, there exist constants c_1 and c_2 such that for all positive integer n ,

$$c_1 n^{1/\alpha} \leq \mathbb{E}_{P^n}[Z_n] \leq c_2 n^{1/\alpha}.$$

Proposition 8: The number of distinct symbols Z_n output by a memoryless source satisfies a Bernstein inequality:

$$\Pr \{Z_n \leq \mathbb{E}[Z_n]\} \leq e^{-\frac{\mathbb{E}[Z_n]}{8}}. \quad (4)$$

Proof: Note that Z_n is a function of n independent random variables. Moreover, Z_n is a configuration function as defined by Talagrand (1995) since $Z_n(\mathbf{x})$ is the size of a maximum subsequence of \mathbf{x} satisfying an hereditary property (all its symbols are pairwise distinct). Using the main theorem in Boucheron et al. (2000), this is enough to conclude. ■

Noting that $Z_n \geq 1$, we can derive the following inequality that will prove useful later on:

$$\begin{aligned} \mathbb{E} \left[\frac{1}{Z_n^{\alpha-1}} \right] &= \mathbb{E} \left[\frac{1}{Z_n^{\alpha-1}} \mathbb{1}_{Z_n > \frac{1}{2} \mathbb{E}[Z_n]} \right] + \mathbb{E}_P \left[\frac{1}{Z_n^{\alpha-1}} \mathbb{1}_{Z_n \leq \frac{1}{2} \mathbb{E}[Z_n]} \right] \\ &\leq \frac{1}{\left(\frac{1}{2} \mathbb{E}[Z_n]\right)^{\alpha-1}} + \Pr \left(Z_n \leq \frac{1}{2} \mathbb{E}[Z_n] \right). \end{aligned} \quad (5)$$

We consider here a modified version of `CensoringCode` that operates similarly, except that

- 1) the string \mathbf{x} is first scanned completely to determine $Z_n(\mathbf{x})$;
- 2) the constant cutoff $\hat{K}_n = \mu Z_n$ is used for all symbols \mathbf{x}_i , $1 \leq i \leq n$, where μ is some positive constant.
- 3) the value of \hat{K}_n is encoded using Elias penultimate code and transmitted before `C1` and `C2`.

Note that this version of the algorithm is not sequential because of the initial scanning.

Algorithm 2 AdaptiveCensoringCode

```

cutoff  $\leftarrow \mu Z_n(\mathbf{x})$  {Determination of the constant cutoff}
counts  $\leftarrow [1/2, 1/2, \dots]$ 
for  $i$  from 1 to  $n$  do
  if  $x[i] \leq \textit{cutoff}$  then
    ArithCode( $x[i]$ , counts[0 : cutoff])
  else
    ArithCode(0, counts[0 : cutoff])
     $\text{C1} \leftarrow \text{C1} \cdot \text{EliasCode}(x[i])$ 
    counts[0]  $\leftarrow \textit{counts}[0] + 1
  end if
  counts[ $x[i]$ ]  $\leftarrow \textit{counts}[ $x[i]$ ] + 1
end for
 $\text{C2} \leftarrow \text{ArithCode}()$ 
return  $\text{C1} \cdot \text{C2}$$$ 
```

We may now assert.

Theorem 10: The algorithm AdaptiveCensoringCode is approximately asymptotically adaptive with respect to $\bigcup_{\alpha > 0} \mathcal{W}_\alpha$.

Proof: Let us again denote by $\text{C1}(\mathbf{x})$ and $\text{C2}(\mathbf{x})$ the two parts of the code-string associated with \mathbf{x} .

Let \hat{L} be the codelength of the output of algorithm AdaptiveCensoringCode.

For any source P :

$$\begin{aligned} \mathbb{E}_P \left[\hat{L}(X_{1:n}) \right] - H(X_{1:n}) &= \mathbb{E}_P \left[\ell(\hat{K}_n) + |\text{C1}(X_{1:n})| + |\text{C2}(X_{1:n})| \right] - n \sum_{k=1}^{\infty} P^1(k) \log \frac{1}{P_1(k)} \\ &\leq \mathbb{E}_P \left[\ell(\hat{K}_n) \right] + \mathbb{E}_P [|\text{C1}(X_{1:n})|] + \mathbb{E}_P \left[|\text{C2}(X_{1:n})| - n \sum_{k=1}^{\hat{K}_n} P^1(k) \log \frac{1}{P_1(k)} \right]. \end{aligned}$$

As function ℓ is increasing and equivalent to log at infinity, the first summand is obviously $o\left(\mathbb{E}_P \left[\hat{K}_n \right]\right)$. Moreover,

if $P \in \mathcal{W}_\alpha$ there exists C such that $P^1(k) \leq \frac{C}{k^\alpha}$ and the second summand satisfies:

$$\begin{aligned}
\mathbb{E}_P [|\mathcal{C}1(X_{1:n})|] &= \mathbb{E}_P \left[\sum_{k \geq \hat{K}_n+1} P^1(k) \ell(k) \right] \\
&\leq nC \mathbb{E}_P \left[\int_{\hat{K}_n}^{\infty} \frac{\ell(x)}{x^\alpha} dx \right] \\
&= nC \mathbb{E}_P \left[\frac{1}{\hat{K}_n^{\alpha-1}} \int_1^{\infty} \frac{\ell(\hat{K}_n u)}{u^\alpha} du \right] \\
&\leq nC \mathbb{E}_P \left[\frac{1}{\hat{K}_n^{\alpha-1}} \right] \int_1^{\infty} \frac{\log(nu)}{u^\alpha} du (1 + o(1)) \\
&= O\left(n^{\frac{1}{\alpha}} \log n\right)
\end{aligned}$$

by Proposition (7) and Inequality (5).

By Theorem 2, every string $x \in \mathbb{N}_+^n$ satisfies

$$|\mathcal{C}2(x)| - n \sum_{k=1}^{\hat{K}_n} P^1(k) \log \frac{1}{P_1(k)} \leq \frac{\hat{K}_n}{2} \log n + 2.$$

Hence, the third summand is upper-bounded as:

$$\begin{aligned}
\mathbb{E}_P \left[|\mathcal{C}2(X_{1:n})| - n \sum_{k=1}^{\hat{K}_n} P^1(k) \log \frac{1}{P_1(k)} \right] &\leq \frac{\mathbb{E}_P [\hat{K}_n]}{2} \log n + 2 \\
&= O\left(n^{\frac{1}{\alpha}} \log n\right)
\end{aligned}$$

which finishes to prove the theorem. ■

REFERENCES

- A. Barron, J. Rissanen, and B. Yu. The minimum description length principle in coding and modeling. *IEEE Trans. Inform. Theory*, 44(6):2743–2760, 1998. ISSN 0018-9448.
- S. Boucheron, G. Lugosi, and P. Massart. A sharp concentration inequality with applications. *Random Struct. & Algorithms*, 16:277–292, 2000.
- O. Catoni. *Statistical learning theory and stochastic optimization*, volume 1851 of *Lecture Notes in Mathematics*. Springer-Verlag, 2004. Ecole d’Ete de Probabilites de Saint-Flour XXXI.
- N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, 2006. ISBN 0521841089.
- B. Clarke and A. Barron. Information-theoretic asymptotics of bayes methods. *IEEE Trans. Inform. Theory*, 36: 453–471, 1990.
- B. Clarke and A. Barron. Jeffrey’s prior is asymptotically least favorable under entropy risk. *J. Stat. Planning and Inference*, 41:37–60, 1994.

- T. Cover and J. Thomas. *Elements of information theory*. John Wiley & sons, 1991.
- D. Cox. An analysis of bayesian inference for nonparametric regression. *Annals of Statistics*, 21:903–923, 1993.
- I. Csiszár. Class notes on information theory and statistics. University of Maryland, 1990.
- I. Csiszár and P. Shields. Redundancy rates for renewal and other processes. *IEEE Trans. Inform. Theory*, 42(6): 2065–2072, 1996.
- L. D. Davisson. Universal noiseless coding. *IEEE Trans. Information Theory*, IT-19:783–795, 1973. ISSN 0018-9448.
- A. Dembo and O. Zeitouni. *Large deviation techniques and applications*. Springer, 1998.
- L. Devroye, L. Györfi, and G. Lugosi. *A probabilistic theory of pattern recognition*. Springer., 1996.
- P. Elias. Universal codeword sets and representations of the integers. *IEEE Trans. Information Theory*, IT-21: 194–203, 1975. ISSN 0018-9448.
- M. Feder, N. Merhav, and M. Gutman. Universal prediction of individual sequences. *IEEE Trans. Inform. Theory*, 38(4):1258–1270, 1992. ISSN 0018-9448.
- D. Freedman. On the Bernstein von Mises theorem with infinite dimensional parameters. *Annals of Statistics*, 27: 1119–1140, 1999.
- R. G. Gallager. *Information theory and reliable communication*. John Wiley & sons, 1968.
- A. Garivier. A lower bound for the maximin redundancy in pattern coding. Technical report, Université Paris-Sud, 2006.
- J. K. Ghosh and R. V. Ramamoorthi. *Bayesian nonparametrics*. Springer Series in Statistics. Springer-Verlag, New York, 2003. ISBN 0-387-95537-2.
- L. Györfi, I. Páli, and E. C. van der Meulen. There is no universal source code for an infinite source alphabet. *IEEE Trans. Inform. Theory*, 40(1):267–271, 1994. ISSN 0018-9448.
- D. Haussler. A general minimax result for relative entropy. *IEEE Trans. Inform. Theory*, 43(4):1276–1280, 1997. ISSN 0018-9448.
- D. Haussler and M. Opper. Mutual information, metric entropy and cumulative relative entropy risk. *Ann. Statist.*, 25(6):2451–2492, 1997. ISSN 0090-5364.
- J. C. Kieffer. A unified approach to weak universal source coding. *IEEE Trans. Inform. Theory*, 24(6):674–682, 1978. ISSN 0018-9448.
- R. E. Krichevsky and V. K. Trofimov. The performance of universal encoding. *IEEE Trans. Inform. Theory*, 27(2): 199–207, 1981. ISSN 0018-9448.
- A. Orlitsky and N. P. Santhanam. Speaking of infinity. *IEEE Trans. Inform. Theory*, 50(10):2215–2230, 2004. ISSN 0018-9448.
- A. Orlitsky, N. P. Santhanam, and J. Zhang. Universal compression of memoryless sources over unknown alphabets. *IEEE Trans. Inform. Theory*, 50(7):1469–1481, 2004. ISSN 0018-9448.
- J. Rissanen and G. G. Langdon, Jr. Arithmetic coding. *IBM J. Res. Develop.*, 23(2):149–162, 1979. ISSN 0018-8646.
- G. Shamir. On the MDL principle for i.i.d. sources with large alphabets. *IEEE Trans. Inform. Theory*, 52(5):

- 1939–1955, 2006. ISSN 0018-9448.
- G. I. Shamir. A new redundancy bound for universal lossless compression of unknown alphabets. In *Proceedings of The 38th Annual Conference on Information Sciences and Systems-CISS*, pages 1175–1179, Princeton, New-Jersey, U.S.A., 2004.
- P. Shields. Universal redundancy rates do not exist. *IEEE Trans. Inform. Theory*, 39:520–524, 1993.
- Y. Shtarkov. Universal sequential coding of messages. *Probl. Inform. Transmission*, 23:3–17, 1987.
- M. Sion. On general minimax theorems. *Pacific J. Math.*, 8:171–176, 1958. ISSN 0030-8730.
- W. Szpankowski. On asymptotics of certain recurrences arising in universal coding. *Probl. Inf. Transm.*, 34(2): 142–146, 1998.
- M. Talagrand. Concentration of measure and isoperimetric inequalities in product spaces. *Publ. Math., Inst. Hautes Etud. Sci.*, 81:73–205, 1995.
- A. van der Vaart. *Asymptotic statistics*. Cambridge University Press, 1998.
- E. Whittaker and G. Watson. *A course of modern analysis*. Cambridge Mathematical Library. Cambridge University Press, Cambridge, 1996. ISBN 0-521-58807-3. An introduction to the general theory of infinite processes and of analytic functions; with an account of the principal transcendental functions, Reprint of the fourth (1927) edition.
- Q. Xie and A. R. Barron. Minimax redundancy for the class of memoryless sources. *IEEE Trans. Inform. Theory*, 43:646–656, 1997.
- Q. Xie and A. R. Barron. Asymptotic minimax regret for data compression, gambling and prediction. *IEEE Trans. Inform. Theory*, 46:431–445, 2000.

APPENDIX I

UPPER-BOUND ON MINIMAX REGRET

This sections contains the proof of the last inequality in Theorem 2.

The minimax regret is not larger than the maximum regret of the Krichevsky-Trofimov mixture over k -ary alphabet over strings of length n . The latter is classically upper-bounded by

$$\log \left(\frac{\Gamma(n + \frac{k}{2})\Gamma(\frac{1}{2})}{\Gamma(n + \frac{1}{2})\Gamma(\frac{k}{2})} \right),$$

as proved for example in (Csiszár, 1990).

Now thanks to the Stirling approximation to the Gamma function (See Whittaker and Watson, 1996, Chapter XII) asserts that for any $x > 0$, there exists $\beta \in [0, 1]$ such that

$$\Gamma(x) = x^{x-\frac{1}{2}} e^{-x} \sqrt{2\pi} e^{\frac{\beta}{12x}}.$$

The announced upper-bound follows in a straightforward way by plugging this approximation into the preceding upper-bound.

APPENDIX II

LOWER BOUND ON REDUNDANCY FOR POWER-LAW ENVELOPES

In this appendix we derive a lower-bound for power-law envelopes using Theorem 5. Let α denote a real larger than 1. Let C be such that $C^{1/\alpha} > 4$. As the envelope function is defined by $f(i) = 1 \wedge C/i^\alpha$, the constant $c(\infty) = \sum_{i \geq 1} f(2i)$ satisfies

$$\frac{\alpha}{\alpha-1} \frac{C^{1/\alpha}}{2} - 1 \leq c(\infty) \leq \frac{C^{1/\alpha}}{2} + \frac{C}{(\alpha-1)2^\alpha} \left(\frac{C^{1/\alpha}}{2} \right)^{1-\alpha}.$$

The condition on C and α warrants that, for sufficiently large p , we have $c(p) > 1$ (this is indeed true for $p > C^{1/\alpha}$).

We choose $p = an^{\frac{1}{\alpha}}$ for a small enough to have

$$\frac{(1-\lambda)C\epsilon}{(2a)^{\frac{1}{\alpha}} c(\infty)} > 10,$$

so that condition $(1-\lambda)n \frac{f(2p)}{c(p)} > \frac{10}{\epsilon}$ is satisfied for n large enough. Then

$$R^+(\Lambda_f^n) \geq C(p, n, \lambda, \epsilon) \sum_{i=1}^p \left(\frac{1}{2} \log \frac{n(1-\lambda)\pi f(2i)}{2c(p)e} - \epsilon \right),$$

where $C(p, n, \lambda, \epsilon) = \frac{1}{1 + \frac{(2a)^{\frac{1}{\alpha}} c(\infty)}{C\lambda^2}} \left(1 - \frac{4}{\pi} \sqrt{\frac{10c(\infty)(2a)^\alpha}{(1-\lambda)C\epsilon}} \right)$, and

$$\begin{aligned} \sum_{i=1}^p \left(\frac{1}{2} \log \frac{n(1-\lambda)\pi f(2i)}{2c(p)e} - \epsilon \right) &= \frac{p}{2} \log n - \frac{\alpha}{2} \sum_{i=1}^p \log i + \left(\frac{1}{2} \log \frac{(1-\lambda)\pi C}{2^{1+\alpha} c(\infty)e} - \epsilon \right) p \\ &= \frac{p}{2} \log n - \frac{\alpha}{2} (p \log p - p + o(p)) + \left(\frac{1}{2} \log \frac{(1-\lambda)\pi C}{2^{1+\alpha} c(\infty)e} - \epsilon \right) p \\ &= \frac{an^{\frac{1}{\alpha}}}{2} \log n - \frac{\alpha}{2} \left(an^{\frac{1}{\alpha}} \log a + \frac{a}{\alpha} n^{\frac{1}{\alpha}} \log n - an^{\frac{1}{\alpha}} + o\left(n^{\frac{1}{\alpha}}\right) \right) \\ &\quad + \left(\frac{1}{2} \log \frac{(1-\lambda)\pi C}{2^{1+\alpha} c(\infty)e} - \epsilon \right) an^{\frac{1}{\alpha}} \\ &= \left(\frac{\alpha}{2} (1 - \log a) + \frac{1}{2} \log \frac{(1-\lambda)\pi C}{2^{1+\alpha} c(\infty)e} - \epsilon + o(1) \right) an^{\frac{1}{\alpha}}. \end{aligned}$$

For a small enough, this gives the existence of a positive constant η such that $R^+(\Lambda_f^n) \geq \eta n^{\frac{1}{\alpha}}$.

APPENDIX III

PROOF OF LEMMA 7

Suppose that there exist k_0 , c and C such that for all $k \geq k_0$, $\frac{c}{k^\alpha} \leq p_k \leq \frac{C}{k^\alpha}$.

For $0 \leq x \leq \frac{1}{2}$, it holds that $-(2 \log 2)x \leq \log(1-x) \leq -x$ and thus

$$e^{-(2 \log 2)nx} \leq (1-x)^n \leq e^{-nx}.$$

Hence (as $p_k \leq \frac{1}{2}$ for all $k \geq 2$) :

$$\begin{aligned} \sum_{k=k_0}^{\infty} \left(1 - \left(1 - \frac{C}{k^\alpha}\right)^n\right) &\leq \mathbb{E}[Z_n] \leq \sum_{k=1}^{\infty} \left(1 - \left(1 - \frac{C}{k^\alpha}\right)^n\right) \\ \sum_{k=k_0}^{\infty} \left(1 - e^{-\frac{Cn}{k^\alpha}}\right) &\leq \mathbb{E}[Z_n] \leq 1 + \sum_{k=2}^{\infty} \left(1 - e^{-\frac{(2 \log 2) C n}{k^\alpha}}\right) \\ \int_{k_0}^{\infty} \left(1 - e^{-\frac{Cn}{x^\alpha}}\right) dx &\leq \mathbb{E}[Z_n] \leq 1 + \int_1^{\infty} \left(1 - e^{-\frac{(2 \log 2) C n}{x^\alpha}}\right) dx. \end{aligned}$$

But, for any $t, K > 0$, it holds that

$$\int_t^{\infty} \left(1 - e^{-\frac{K n}{x^\alpha}}\right) dx = \frac{(K n)^{1/\alpha}}{\alpha} \int_0^{\frac{K n}{t^\alpha}} \frac{1 - e^{-u}}{u^{1+1/\alpha}} du.$$

Thus, by noting that integral

$$A(\alpha) = \int_0^{\infty} \frac{1 - e^{-u}}{u^{1+1/\alpha}} du,$$

is convergent, we get

$$\frac{c^{1/\alpha} A(\alpha)}{\alpha} n^{1/\alpha} (1 - o(1)) \leq \mathbb{E}[Z_n] \leq \frac{((2 \log 2) C)^{1/\alpha} A(\alpha)}{\alpha} n^{1/\alpha}.$$

APPENDIX IV

EXPECTED SIZE OF DICTIONARY ENCODING

Assume that the probability mass function (p_k) satisfies $p_k \leq \frac{C}{k^\alpha}$ for $C > 0$ and all $k \geq 0$. Then, using Elias penultimate code for the first occurrence of each symbol in $X_{1:n}$, the expected length of the binary encoding of the dictionary can be upper-bounded in the following way. Let U_k be equal to 1 if symbol k occurs in $X_{1:n}$, and equal to 0 otherwise.

$$\begin{aligned} \mathbb{E}[|\Delta_n|] &= \mathbb{E}\left[\sum_{k=1}^{\infty} U_k \ell(k)\right] \\ &= \sum_{k=1}^{\infty} \mathbb{E}[U_k \ell(k)] \\ &\leq \sum_{k=1}^{\infty} \left(1 - \left(1 - \frac{C}{k^\alpha}\right)^n\right) \ell(k) \\ &\leq 2 \left(1 + \sum_{k=2}^{\infty} \left(1 - e^{-\frac{(2 \log 2) C n}{k^\alpha}}\right) \log k\right) \\ &\leq 2 \left(1 + \int_1^{\infty} \left(1 - e^{-\frac{(2 \log 2) C n}{x^\alpha}}\right) \log x dx\right) \\ &\leq 2 \left(\frac{((2 \log 2) C n)^{1/\alpha}}{\alpha^2} \int_0^{(2 \log 2) C n} \frac{1 - e^{-u}}{u^{1+1/\alpha}} \log\left(\frac{(2 \log 2) C n}{u}\right) du\right) \\ &\leq T \frac{((2 \log 2) C n)^{1/\alpha}}{\alpha^2} \log n \int_0^{\infty} \frac{1 - e^{-u}}{u^{1+1/\alpha}} du \text{ bits} \end{aligned}$$

for some positive constant T .

Acknowledgment

The authors wish to thank Professor László Györfi for useful discussion and helpful comments.