
A SHORT INTRODUCTION TO INFORMATION THEORY AND BAYESIAN MODEL AVERAGING

OLIVIER CATONI

March 22, 2012

1. INFORMATION THEORY AND LOSSLESS CODES

1.1. BINARY CODES. Let us consider a finite alphabet A and a random sequence $(X_n)_{n=1}^\infty$ taking its values in A . (The alphabet can be any finite set here.)

Let $\{0, 1\}^* = \bigcup_{n=1}^\infty \{0, 1\}^n$ be the set of finite binary sequences.

DEFINITION 1.1 *Given some block length n , a binary code c is an injective map from A^n to $\{0, 1\}^n$. The mean length of c is*

$$\mathbb{E}\{\ell[c(X^n)]\},$$

where the expectation is taken with respect to the distribution of the sequence $X^n \stackrel{\text{def}}{=} (X_1, \dots, X_n)$ and where the length function ℓ is defined as $\ell(w) = k$ for any $w \in \{0, 1\}^k$.

Minimizing the mean code length under various assumptions on the source and the set of authorized codes is the main subject of lossless coding theory.

1.2. OPTIMAL CODE LENGTH FOR A KNOWN SOURCE. In the case when \mathbb{P}_{X^n} (the distribution of $X^n = (X_1, \dots, X_n)$) is known and c is arbitrary, minimizing the code length is achieved by sorting A^n in order of decreasing probabilities. More precisely, let us write $A^n = \{b_i, i = 1, \dots, d^n\}$, where $d = |A|$ is the size of the alphabet and where the blocks of n letters are indexed by order of decreasing probabilities:

$$\mathbb{P}_{X^n}(b_i) \leq \mathbb{P}_{X^n}(b_{i-1}), i = 2, \dots, d^n.$$

Let us sort $\{0, 1\}^*$ by order of increasing lengths, writing

$$\{0, 1\}^* = \{w(j), j \in \mathbb{N} \setminus \{0\}\},$$

CNRS – UMR 8553, Département de Mathématiques et Applications, Ecole Normale Supérieure, 45, rue d’Ulm, F75230 Paris cedex 05, and INRIA Paris-Rocquencourt – CLASSIC team.

where $w(1) = 0$, $w(2) = 1$, $w(3) = 00$, $w(4) = 01$, \dots and more generally

$$w\left(2^\ell + \sum_{k=1}^{\ell} d_{\ell-k} 2^k - 1\right) = (d_k)_{k=1}^\ell,$$

where $d_k \in \{0, 1\}$ (this means that $w(i)$ is the binary representation of $i + 1$ read from left to right, after removing the leftmost bit, which is always equal to one). Let us remark that $\ell[w(i)] = \lfloor \log_2(i + 1) \rfloor$.

PROPOSITION 1.1 *The code $c(b_i) = w(i)$ minimizes the mean code length among all binary codes.*

PROOF. Let c' be some other code. Its mean length can be written as

$$\mathbb{E}\left\{\ell[c'(X^n)]\right\} = \sum_{k=1}^{\infty} \mathbb{P}\left\{\ell[c'(X^n)] \geq k\right\}. \quad (1)$$

On the other hand, since $i \mapsto \ell[w(i)]$ is non-decreasing, for any code length k ,

$$\{i \in \mathbb{N} \setminus \{0\} : \ell[w(i)] < k\} = \llbracket 1, i_k \rrbracket,$$

where $i_k = 2^k - 2$ is the number of binary codes of length less than k , and $\llbracket 1, i_k \rrbracket \stackrel{\text{def}}{=} \{1, \dots, i_k\}$. Thus

$$\mathbb{P}\left\{\ell[c'(X^n)] < k\right\} = \mathbb{P}_{X^n}\left\{c'^{-1}[w(\llbracket 1, i_k \rrbracket)]\right\}.$$

As c' is an injective map, $|c'^{-1}[w(\llbracket 1, i_k \rrbracket)]| \leq i_k$, and can thus be written as $\{b_{j_1}, \dots, b_{j_m}\}$, where $j_1 < j_2 < \dots < j_m$ and $m \leq i_k$. Thus

$$\begin{aligned} \mathbb{P}\left\{\ell[c'(X^n)] < k\right\} &= \mathbb{P}_{X^n}(b_{j_1}) + \dots + \mathbb{P}_{X^n}(b_{j_m}) \\ &\leq \mathbb{P}_{X^n}(b_1) + \dots + \mathbb{P}_{X^n}(b_m) \\ &\leq \mathbb{P}_{X^n}(b_1) + \dots + \mathbb{P}_{X^n}(b_{i_k}) = \mathbb{P}\left\{\ell[c(X^n)] < k\right\}, \end{aligned}$$

since $s \leq j_s$, $i \mapsto \mathbb{P}_{X^n}(b_i)$ is non-increasing and $m \leq i_k$. This proves that

$$\mathbb{P}_{X^n}\left\{\ell[c'(X^n)] \geq k\right\} \geq \mathbb{P}_{X^n}\left\{\ell[c(X^n)] \geq k\right\},$$

and therefore according to equation (1) that the mean code length of c is optimal. \square

The mean length of optimal codes is related to the Shannon entropy $H(\mathbb{P}_{X^n})$ of the distribution of X^n .

DEFINITION 1.2 *The Shannon entropy $H(p)$ of a probability measure $p \in \mathcal{M}_+^1(\mathcal{X})$ on a finite set \mathcal{X} is defined as*

$$H(p) = - \sum_{x \in \mathcal{X}} p(x) \log_2[p(x)],$$

where $\log_2(z) = \log(z)/\log(2)$. *The entropy is measured in bits, for reasons that will become clear later.*

Another quantity of interest is the Kullback Leibler divergence, or relative entropy, between two probability measures P and Q .

DEFINITION 1.3 *Let P and $Q \in \mathcal{M}_+^1(\mathcal{X})$ be two probability measures defined on some measurable space \mathcal{X} (that needs not be finite). The Kullback Leibler divergence of P with respect to Q is defined as*

$$\mathcal{K}(P, Q) = \begin{cases} \int \log\left(\frac{dP}{dQ}\right) dP & \text{when } P \ll Q, \\ +\infty & \text{otherwise.} \end{cases}$$

Let us remark that when $P \ll Q$, $\log\left(\frac{dP}{dQ}\right)$ has an integrable negative part, because

$$\int \log\left(\frac{dP}{dQ}\right)_- dP = \int \log\left(\frac{dP}{dQ}\right)_- \frac{dP}{dQ} dQ < \infty,$$

where $z_- = \max\{-z, 0\}$, due to the fact that $z \mapsto z \log(z) : \mathbb{R}_+ \rightarrow \mathbb{R}$ is bounded from below (by $-1/e$). Thus the integral $\int \log\left(\frac{dP}{dQ}\right) dP$ is defined as a generalized integral, taking values in $\mathbb{R} \cup \{+\infty\}$.

Exercise 1 *Give an example, where $P \ll Q$ and $\mathcal{K}(P, Q) = +\infty$.*

PROPOSITION 1.2 *The Kullback Leibler divergence is a non-negative function.*

PROOF. When $P \ll Q$, we can write the divergence as

$$\mathcal{K}(P, Q) = \int 1 - \frac{dP}{dQ} + \frac{dP}{dQ} \log\left(\frac{dP}{dQ}\right) dQ.$$

Since the function $z \mapsto 1 - z + z \log(z) : \mathbb{R}_+ \rightarrow \mathbb{R}$ has a positive range, it shows that $\mathcal{K}(P, Q) \in \mathbb{R}_+ \cup \{+\infty\}$. \square

Exercise 2 Show that for any $P_1, P_2, Q \in \mathcal{M}_+^1(\mathcal{X})$, any $\lambda \in [0, 1]$,

$$\begin{aligned}\mathcal{K}(\lambda P_1 + (1 - \lambda)P_2, Q) &\leq \lambda \mathcal{K}(P_1, Q) + (1 - \lambda)\mathcal{K}(P_2, Q), \\ \mathcal{K}(Q, \lambda P_1 + (1 - \lambda)P_2) &\leq \lambda \mathcal{K}(Q, P_1) + (1 - \lambda)\mathcal{K}(Q, P_2).\end{aligned}$$

PROPOSITION 1.3 The map $p \mapsto H(p) : \mathcal{M}_+^1(\mathcal{X}) \rightarrow \mathbb{R}_+$ is concave.

PROOF. The function $z \mapsto z \log_2(z)$ is convex. \square

PROPOSITION 1.4 The Shannon entropy is sub-additive: For any random sequence X^{n+m} ,

$$H(\mathbb{P}_{X^{n+m}}) \leq H(\mathbb{P}_{X^n}) + H(\mathbb{P}_{X_{n+1}^{n+m}}),$$

where $X_{n+1}^{n+m} \stackrel{\text{def}}{=} (X_{n+1}, \dots, X_{n+m})$.

PROOF. It is enough to prove that for any couple (X, Y) of random variables,

$$H(\mathbb{P}_{X,Y}) \leq H(\mathbb{P}_X) + H(\mathbb{P}_Y).$$

We can then remark that

$$\begin{aligned}H(\mathbb{P}_{X,Y}) &= - \int \log_2[\mathbb{P}_{X,Y}(x, y)] d\mathbb{P}_{X,Y}(x, y) \\ &= - \int \log_2[\mathbb{P}_X(x)] d\mathbb{P}_X(x) - \int \log_2[\mathbb{P}_{Y|X=x}(y)] d\mathbb{P}_{Y|X=x}(y) d\mathbb{P}_X(x) \\ &= H(\mathbb{P}_X) + \int H(\mathbb{P}_{Y|X=x}) d\mathbb{P}_X(x) \\ &\leq H(\mathbb{P}_X) + H\left(\int \mathbb{P}_{Y|X=x} d\mathbb{P}_X(x)\right) = H(\mathbb{P}_X) + H(\mathbb{P}_Y),\end{aligned}$$

where we have used the fact that H is concave. \square

DEFINITION 1.4 A random sequence $(X_n)_{n=1}^{+\infty}$ is said to be stationary when, for any $n, m \in \mathbb{N} \setminus \{0\}$, $\mathbb{P}_{X_{n+1}^{n+m}} = \mathbb{P}_{X^m}$.

PROPOSITION 1.5 If $(X_n)_{n=1}^{+\infty} \stackrel{\text{def}}{=} X^{+\infty}$ is a stationary source, the following limit exists

$$\lim_{n \rightarrow +\infty} \frac{1}{n} H(\mathbb{P}_{X^n}) = \inf_{n=1, \dots, +\infty} \frac{1}{n} H(\mathbb{P}_{X^n}) \stackrel{\text{def}}{=} \overline{H}(\mathbb{P}_{X^{+\infty}}).$$

It is called the Shannon entropy of $X^{+\infty}$.

PROOF. Let $h(n) = n^{-1}H(\mathbb{P}_{X^n})$. From the previous proposition $h(n+m) \leq \frac{nh(n) + mh(m)}{n+m}$. Let $n = pq + r$ where $0 \leq r < p$ be the result of the Euclidean division of n by p . We obtain

$$\begin{aligned} h(n) &\leq \frac{pqh(pq) + rh(r)}{pq+r} \leq \frac{pqh(p) + rh(r)}{pq+r} \\ &\leq h(p) + \frac{p-1}{n} \max\{h(1), \dots, h(p-1)\}. \end{aligned}$$

Thus $\limsup_n h(n) \leq h(p)$, for any $p \in \mathbb{N} \setminus \{0\}$. This implies that $\limsup_n h(n) \leq \inf_n h(n) \leq \liminf_n h(n)$, proving that $\lim_n h(n)$ exists and is equal to $\inf_n h(n)$. \square

PROPOSITION 1.6 *For any source, any optimal code c ,*

$$\begin{aligned} (1 - n^{-1})[H(\mathbb{P}_{X^n}) - \log_2(n-1)] - 1 \\ \leq \sup_{\alpha > 1} \alpha^{-1}[H(\mathbb{P}_{X^n}) + \log_2(\alpha-1)] - 1 \\ \leq \mathbb{E}_{X^n}[\ell(c(X^n))] \leq H(\mathbb{P}_{X^n}) + 1. \end{aligned}$$

Consequently, if the source is stationary,

$$\lim_{n \rightarrow +\infty} \frac{1}{n} \mathbb{E}[\ell(c(X^n))] = \overline{H}(\mathbb{P}_{X^{+\infty}}).$$

PROOF. Let $\mathbb{P}_{X^n}(b_i) \stackrel{\text{def}}{=} p(b_i)$. As

$$ip(b_i) \leq \sum_{j=1}^i p(b_j) \leq 1,$$

$p(b_i) \leq i^{-1}$. Since all optimal codes have the same mean length, we can choose the optimal code of Proposition 1.1 (page 2). We get

$$\begin{aligned} \mathbb{E}[\ell(c(X^n))] &= \sum_{i=1}^{d^n} p(b_i) \lceil \log_2(i+1) \rceil \\ &\leq \sum_i p(b_i) \log_2(p(b_i)^{-1} + 1) \leq 1 - \sum_i p(b_i) \log_2(p(b_i)) = H(p) + 1. \end{aligned}$$

On the other hand,

$$\begin{aligned} \mathbb{E}[\ell(c(X^n))] &\geq \sum_i p(b_i) [\log_2(i+1) - 1] \\ &\geq -\alpha^{-1} \sum_i p(b_i) \log_2\left(\frac{\alpha-1}{(i+1)^\alpha}\right) - 1 + \frac{\log_2(\alpha-1)}{\alpha}. \end{aligned}$$

Since $\sum_i \frac{\alpha-1}{(i+1)^\alpha} \leq (\alpha-1) \int_1^{+\infty} \frac{dz}{z^\alpha} = 1$, there is a probability measure $q \in \mathcal{M}_+^1(D^n)$ such that $q(b_i) \geq \frac{\alpha-1}{(i+1)^\alpha}$. We can then remark that

$$-\sum_i p(b_i) \log_2\left(\frac{\alpha-1}{(i+1)^\alpha}\right) - H(p) \geq \log(2)^{-1} \mathcal{K}(p, q) \geq 0.$$

This proves the second inequality of the proposition. The first one is obtained by choosing $\alpha = (1 - n^{-1})^{-1}$. \square

1.3. INSTANTANEOUS CODES. When transmitting $c(X^n)$ through a channel, it is useful for the receiver to know without delay when the message ends. More precisely, if $c(A^n) = W$, the receiver would like to dispose of a function $f : \{0, 1\}^* \rightarrow \{0, 1\}$ such that $\mathbb{1}(i = \ell(w)) = f(w(1), \dots, w(i))$. If such a function exists, it is necessarily equal to $\mathbb{1}(w(1, \dots, i) \in W)$, so that in this case

$$\mathbb{1}(i = \ell(w)) = \mathbb{1}(w(1, \dots, i) \in W), \quad w \in W, 1 \leq i \leq \ell(w),$$

which can also be written as

$$w(1, \dots, i) \notin W, \quad w \in W, 1 \leq i < \ell(w).$$

A code having this property is called an instantaneous code, and the previous reasoning proves

PROPOSITION 1.7 *The binary code $c : A^n \rightarrow \{0, 1\}^*$ is an instantaneous code if and only if it is a prefix code, which means that no codeword is the prefix of another codeword. In this case, the range $c(A^n) = W \subset \{0, 1\}^*$ of c is called a prefix set.*

PROPOSITION 1.8 (KRAFT INEQUALITY) *For any prefix set $W \subset \{0, 1\}^*$,*

$$\sum_{w \in W} 2^{-\ell(w)} \leq 1.$$

PROOF. Consider $a(w) = \sum_{i=1}^{\ell(w)} w(i) 2^{-i}$ and $b(w) = a(w) + 2^{-\ell(w)}$. Let us write $w \prec w'$ when $w(1, \dots, \ell(w)) = w'(1, \dots, \ell(w))$, that is when w is a prefix of w' . To each finite binary word $w \in \{0, 1\}^*$ corresponds an interval $I(w) = [a(w), b(w)[\subset [0, 1[$. It is easy to see that $I(w') \subset I(w)$ when $w \prec w'$ and that $I(w') \cap I(w) = \emptyset$ otherwise. Thus W is a prefix set if and only if the intervals $I(w)$ are disjoint. In this case, considering the Lebesgue measure λ , we see that $1 = \lambda([0, 1]) \geq \sum_{w \in W} \lambda(I(w)) = \sum_{w \in W} 2^{-\ell(w)}$. \square

PROPOSITION 1.9 (INVERSE KRAFT INEQUALITY) *If*

$$\sum_{i=1}^T 2^{-r_i} \leq 1, \quad (r_i, i = 1, \dots, T) \in \mathbb{N}^T,$$

then there is a prefix set $\{w_1, \dots, w_T\}$ such that $\ell(w_i) = r_i$.

PROOF. Let us sort r_i so that $r_{i-1} \leq r_i$, $i = 2, \dots, T$. Let

$$\alpha_i = \sum_{j=1}^{i-1} 2^{-r_j} = \sum_{k=1}^{r_i} w_i(k) 2^{-k}, \quad \text{where } w_i \in \{0, 1\}^{r_i}$$

is the binary representation of α_i up to precision 2^{-r_i} . As $a(w_i) = \alpha_i$ and $b(w_i) = \alpha_{i+1}$, it is clear that $I(w_i) \cap I(w_j) = \emptyset$, for any $i \neq j$, proving that $\{w_i, i = 1, \dots, T\}$ is indeed a prefix set. \square

DEFINITION 1.5 (COMPLETE PREFIX CODE AND PREFIX SET) *A prefix set W is said to be complete if no other prefix set W' is such that $W \subsetneq W'$. A prefix code is complete if it uses a complete prefix set of codewords.*

PROPOSITION 1.10 (KRAFT EQUALITY) *A prefix set $W \subset \{0, 1\}^*$ is complete if and only if*

$$\sum_{w \in W} 2^{-\ell(w)} = 1.$$

PROOF. We are going to show that W is *not* complete if and only if $\sum_{w \in W} 2^{-\ell(w)} < 1$. If W is not complete, $W \subsetneq W'$, where W' is another prefix set, so that $\sum_{w \in W} 2^{-\ell(w)} < \sum_{w \in W'} 2^{-\ell(w)} \leq 1$, and $\sum_{w \in W} 2^{-\ell(w)} < 1$. On the other hand, if

$$\sum_{w \in W} 2^{-\ell(w)} < 1, \quad \text{then } \bigcup_{w \in W} I(w) \subsetneq [0, 1[.$$

Let $L = \max\{\ell(w), w \in W\}$ and consider the partition of $[0, 1[$ made by the intervals $\{J(k) = [k2^{-L}, (k+1)2^{-L}[, k = 0, \dots, 2^L - 1\}$. Since for any $w \in W$, $I(w) = \bigcup_{k \in K(w)} J(k)$, we see that $\bigcup_{w \in W} K(w) \subsetneq \{0, \dots, 2^L - 1\}$. Therefore, there is $k \in \{0, \dots, 2^L - 1\}$ such that $J(k) \not\subseteq \bigcup_{w \in W} I(w)$. Define $w' \in \{0, 1\}^L$ by the formula $k2^{-L} = \sum_{i=1}^L w'(i)2^{-i}$. Since $I(w') = J(k)$, $W \cup \{w'\}$ forms a prefix set, hence W is not complete. \square

PROPOSITION 1.11 *For any prefix code c , and any random source $X^n \in A^n$,*

$$\mathbb{E}[\ell(c(X^n))] \geq H(\mathbb{P}_{X^n}).$$

Moreover, there exists a prefix code c such that

$$\mathbb{E}[\ell(c(X^n))] < H(\mathbb{P}_{X^n}) + 1.$$

PROOF. Due to the Kraft inequality, we can find a probability measure $Q \in \mathcal{M}_+^1(A^n)$ such that $Q(x^n) \geq 2^{-\ell(c(x^n))}$, for any $x^n \in A^n$. We can then remark that

$$\mathbb{E}[\ell(c(X^n))] - H(\mathbb{P}_{X^n}) \geq \mathcal{K}(P, Q) / \log(2) \geq 0,$$

as proved in Proposition 1.2 (page 3). As for the second part of the proposition, let us put

$$r(x^n) = \lceil -\log_2[\mathbb{P}_{X^n}(x^n)] \rceil.$$

Since $\sum_{x^n \in A^n} 2^{-r(x^n)} \leq 1$, there is, according to Proposition 1.9 (page 7), a prefix code c such that $\ell(c(x^n)) = r(x^n) < -\log_2(\mathbb{P}_{X^n}(x^n)) + 1$, for any $x^n \in A^n$, and consequently such that $\mathbb{E}[\ell(c(X^n))] < H(\mathbb{P}_{X^n}) + 1$. \square

PROPOSITION 1.12 (HUFFMAN CODE) *Let $p \in \mathcal{M}_+^1(\{1, \dots, d\})$ be a probability measure on the d first integers. Let i and $j \leq d$ be such that $p(i)$ and $p(j)$ are the smallest, meaning that*

$$\max\{p(i), p(j)\} \leq \min\{p(k), k \neq i, k \neq j\}.$$

Let c be some optimal prefix code for the probability vector

$$(p(k), k \notin \{i, j\}, p(i) + p(j))$$

of length $d - 1$, that is some prefix code with minimum mean code length. Let $c' : \{1, \dots, d\} \rightarrow \{0, 1\}^*$ be the prefix binary code defined as

$$c'(k) = \begin{cases} c(k), & k \notin \{i, j\}, \\ (c(i, j), 0), & k = i, \\ (c(i, j), 1), & k = j. \end{cases}$$

Then, c' is an optimal code for the source with probability p .

PROOF. Let $c'' : \{1, \dots, d\} \rightarrow \{0, 1\}^*$ be some optimal prefix code and let $W = c''(\{1, \dots, d\})$. Let us choose $w \in W$ such that $\ell(w) = \max\{\ell(v), v \in W\}$. Let $w' = (w(1), \dots, \ell(w) - 1, 1 - w(\ell(w)))$ be its brother. If $w' \notin W$, then $w(1, \dots, \ell(w) - 1)$ would be the prefix of w only, so that it would be possible to replace w with it and to decrease the mean code length of c'' , leading to a contradiction. Thus necessarily $w' \in W$ also. We can assume that $c''(i) = w$ and $c''(j) = w'$, indeed, if this is not the case, we can exchange w and $c''(i)$ and w' and $c''(j)$ without increasing the mean length of c'' . We can then consider the prefix code c''' on $\{k \notin \{i, j\}, 1 \leq k \leq d\} \cup \{(i, j)\}$ defined as

$$c'''(k) = \begin{cases} c''(k), & k \notin \{i, j\}, \\ w(1, \dots, \ell(w) - 1), & k = (i, j). \end{cases}$$

It is easy to see that

$$\begin{aligned} \mathbb{E}[\ell(c'')] &= \mathbb{E}[\ell(c''')] + p(i) + p(j), \\ \mathbb{E}[\ell(c')] &= \mathbb{E}[\ell(c)] + p(i) + p(j). \end{aligned}$$

Since c is optimal, $\mathbb{E}[\ell(c)] \leq \mathbb{E}[\ell(c''')]$, showing that $\mathbb{E}[\ell(c')] \leq \mathbb{E}[\ell(c'')]$, so that c' is also optimal. \square

Exercise 3 Find an optimal code for the probability vector $(1/3, 1/3, 1/6, 1/6)$.

1.4. ARITHMETIC CODES. The previous section shows that for prefix codes as well, finding an optimal code requires to sort A^n according to decreasing probabilities. If n is large, this may induce intensive computations. Arithmetic codes are almost optimal and do not share this weakness.

DEFINITION 1.6 A non ordered vector of probabilities $p(1, \dots, d)$ being given, let us define

$$\xi_i = \sum_{j=1}^{i-1} p(j), \quad i = 1, \dots, d + 1.$$

Let us then consider the intervals $J(i) = [\xi_i, \xi_{i+1}[$ and define

$$w_i \in \arg \max_w \{\lambda(I(w)); w \in \{0, 1\}^*, I(w) \subset J(i)\}.$$

The code $c(i) = w_i$ is called an arithmetic, or Shannon-Fano-Elias code, and satisfies $\mathbb{E}[\ell(c)] < H(p) + 2$.

PROOF. The code c is prefix because the intervals $J(i)$ being non overlapping, this is also the case for $I(w)$, $w \in c(\{1, \dots, d\})$. Let us define $L = \lceil 1 - \log_2[p(i)] \rceil$ and let

$$w = \arg \min_w \{a(w), w \in \{0, 1\}^*, \ell(w) = L, a(w) \geq \xi_i\}.$$

Since

$$b(w) = a(w) + 2^{-L} < \xi_i + 2^{-(L-1)} \leq \xi_i + p(i) = \xi_{i+1},$$

$I(w) \subset J(i)$ and therefore $\ell[c(i)] \leq \ell(w) < 2 - \log_2(p(i))$, implying that $\mathbb{E}[\ell(c)] < H(p) + 2$. \square Arithmetic coding of $x^n \in A^n$ is fast when A^n is sorted in lexicographic order. Indeed,

$$\begin{aligned} \xi(x^n) &= \sum_{k=1}^n \sum_{y < x_k} p(x^{k-1}, y) \\ &= \sum_{k=1}^n \sum_{y < x_k} \prod_{j=1}^{k-1} p(x_j | x^{j-1}) p(y | x^{k-1}), \end{aligned}$$

which can be computed from at most dn conditional probabilities performing at most dn additions and dn multiplications.

Exercise 4 Write a program that draws a probability vector $p \in \mathcal{M}_+^1(\{1, \dots, d\})$ at random according to the uniform probability measure on the simplex. (Hint: draw (X_1, \dots, X_{d-1}) i.i.d. according to the uniform measure on the interval $[0, 1]$, set $X_0 = 0$, $X_d = 1$, consider the order statistics $X_{(i)}$ such that $\{X_i\} = \{X_{(i)}\}$ and $X_{(0)} \leq \dots \leq X_{(d)}$ and put $p_i = X_{(i)} - X_{(i-1)}$.) Compute an optimal prefix code for p (following Huffman's algorithm). Compute also an arithmetic prefix code. Print the mean code lengths of the two codes, as well as the Shannon entropy of p , for various draws of p and various values of d .

REFERENCES

- [1] O. Catoni. *Statistical Learning Theory and Stochastic Optimization, Lectures on Probability Theory and Statistics, École d'Été de Probabilités de Saint-Flour XXXI – 2001*, volume 1851 of *Lecture Notes in Mathematics*. Springer, 2004. Pages 1–269.
- [2] T. Cover and J. Thomas. *Elements of Information Theory*. Wiley and Sons, New York, second edition, 2006.