
A SHORT INTRODUCTION TO INFORMATION THEORY AND BAYESIAN MODEL AVERAGING

OLIVIER CATONI

May 3, 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 .

2. UNIVERSAL CODES AND BAYESIAN MODEL AVERAGING

We have seen that prefix codes were closely related to probability distributions, due to the Kraft inequality. From now on, we will identify the two, and consider any probability measure on A^n as an *ideal code*. In the previous section, we described some ways of coding efficiently a source X^n whose distribution was known. Here, we are going to study the case of an unknown distribution.

DEFINITION 2.1 *The performance of an ideal code $Q \in \mathcal{M}_+^1(A^n)$ applied to a block X^n of length n of a source with probability law \mathbb{P}_{X^n} , is measured, in bits, by its redundancy*

$$\mathcal{R}(\mathbb{P}_{X^n}, Q) = \log(2)^{-1} \mathcal{K}(\mathbb{P}_{X^n}, Q) = \mathbb{E}_{X^n} \left\{ -\log_2 [Q(X^n)] \right\} - H(\mathbb{P}_{X^n})$$

More precisely, the redundancy measures the difference between the mean length of a prefix code built according to Q and the optimal code built using \mathbb{P}_{X^n} , up to the discretization errors of 1 (Huffman) or 2 (Shannon-Fano-Elias) bits introduced by the actual coding algorithms.

2.1. BAYES REDUNDANCY. One important approach to choose Q when \mathbb{P}_{X^n} is unknown is provided by the Bayesian approach.

DEFINITION 2.2 *Let π be some probability measure on some measurable set of parameters Θ . Let $\{P_\theta \in \mathcal{M}_+^1(A^n); \theta \in \Theta\}$ be a family of measures on A^n , where A is some finite set (extensions to more general measurable sets being possible). Let us assume that $\theta \mapsto P_\theta(x^n)$ is measurable for any $x^n \in A^n$. The mean Bayesian redundancy with respect to π is defined as*

$$\mathcal{R}_\pi(Q) = \int \mathcal{R}(P_\theta, Q) \, d\pi(\theta).$$

Thus, the Bayesian redundancy measures the average loss of efficiency of the coding distribution Q , when the source distribution is drawn from $\{P_\theta; \theta \in \Theta\}$ according to the “prior” probability measure π .

PROPOSITION 2.1 *The infimum $\inf_{Q \in \mathcal{M}_+^1(A^n)} \mathcal{R}_\pi(Q)$ is reached for only one coding distribution Q , called the bayesian mixture ideal code, and defined as*

$$Q = \int P_\theta \, d\pi(\theta) \stackrel{\text{def}}{=} P_\pi.$$

PROOF. The Bayesian redundancy is composed of

$$\begin{aligned}\mathcal{R}_\pi(Q) &= \int \mathcal{R}(P_\theta, P_\pi) d\pi(\theta) + \mathcal{R}(P_\pi, Q) \\ &= \mathcal{R}_\pi(P_\pi) + \mathcal{R}(P_\pi, Q).\end{aligned}$$

The conclusion comes from the fact that $\mathcal{R}(P, Q) > 0$ for $P \neq Q$ and $\mathcal{R}(P, P) = 0$, because it is proportional to the Kullback Leibler divergence. \square

DEFINITION 2.3 *The mutual information of the joint distribution of two random variables X and Y defined on the same probability space is*

$$\begin{aligned}\mathcal{J}(\mathbb{P}_{X,Y}) &= \log(2)^{-1} \mathcal{K}(\mathbb{P}_{X,Y}, \mathbb{P}_X \otimes \mathbb{P}_Y) \\ &= \log(2)^{-1} \mathbb{E}_X [\mathcal{K}(\mathbb{P}_{Y|X}, \mathbb{P}_Y)].\end{aligned}$$

PROPOSITION 2.2 *The optimal Bayesian redundancy, achieved by the Bayesian ideal code, is equal to the mutual information between the parameter θ and the source (X_1, \dots, X_n) , when their joint distribution is defined as*

$$\mathbb{P}_{(\theta, X^n)}(B) = \int \mathbf{1}[(\theta, x^n) \in B] d\pi(\theta) dP_\theta(x^n),$$

for any measurable set B . In other words,

$$\mathcal{R}_\pi(P_\pi) = \mathcal{J}(\mathbb{P}_{\theta, X^n}).$$

PROOF. This is a straightforward consequence of the definitions. The mutual information between X and Y tells how many bits can be saved on the transmission of Y when X is known both to the sender and the receiver of the message, as is clear from the second expression of the mutual information. In the case of the Bayesian redundancy, we see that the average increase of code length due to the fact that the parameter of the source has to be learnt, depends on how much information about the unknown parameter θ is contained in the source X^n . \square

2.2. MINIMAX REDUNDANCY. In this approach we want to make sure to do our best in the worst case. In the setting of Definition 2.2 (page 11), let us define the worst case redundancy as

$$\mathcal{R}(Q) = \sup_{\theta \in \Theta} \mathcal{R}(P_\theta, Q).$$

DEFINITION 2.4 *The ideal code $\widehat{Q} \in \mathcal{M}_+^1(A^n)$ is said to be a minimax coding distribution when*

$$\mathcal{R}(\widehat{Q}) = \inf_{Q \in \mathcal{M}_+^1(A^n)} \mathcal{R}(Q)$$

PROPOSITION 2.3 (LEAST FAVORABLE PRIOR) *Under the same hypotheses as in Definition 2.2 (page 11),*

$$\begin{aligned} \operatorname{ess\,inf}_{d\pi(\theta)} \mathcal{R}(P_\theta, P_\pi) &\leq \int \mathcal{R}(P_\theta, P_\pi) d\pi(\theta) \\ &\leq \inf_Q \sup_\theta \mathcal{R}(P_\theta, Q) \leq \sup_{\theta \in \Theta} \mathcal{R}(P_\theta, P_\pi). \end{aligned}$$

If moreover, for some $\widehat{\pi} \in \mathcal{M}_+^1(\Theta)$, $\mathcal{R}(P_\theta, P_{\widehat{\pi}}) = \sup_\theta \mathcal{R}(P_\theta, P_{\widehat{\pi}})$, $\widehat{\pi}$ almost surely, then $\widehat{\pi}$ is called a least favorable prior. In this case $P_{\widehat{\pi}}$ is the unique minimax coding distribution. Moreover $\widehat{\pi}$ is solution of

$$\mathcal{R}_{\widehat{\pi}}(P_{\widehat{\pi}}) = \sup_{\pi \in \mathcal{M}_+^1(\Theta)} \mathcal{R}_\pi(P_\pi). \quad (2)$$

On the other hand, if $\widehat{\pi}$ satisfies equation (2), then $P_{\widehat{\pi}}$ is the unique minimax coding distribution.

REMARK 2.1 *Equation (2) does have a solution under mild assumptions, we refer to [1][theorem 1.2.1 page 17] for further details on this question.*

PROOF. The chain of inequalities at the beginning of the proposition is a consequence of Proposition 3.1, saying that

$$\int \mathcal{R}(P_\theta, P_\pi) d\pi(\theta) = \inf_{Q \in \mathcal{M}_+^1(X)} \int \mathcal{R}(P_\theta, Q) d\pi(\theta).$$

In the case when $\mathcal{R}(P_\theta, P_{\widehat{\pi}}) = \sup_\theta \mathcal{R}(P_\theta, P_{\widehat{\pi}})$ almost surely, all the inequalities become equalities, showing that $P_{\widehat{\pi}}$ is a minimax coding distribution. To show that it is unique, let us consider another minimax distribution \widehat{Q} . It would satisfy necessarily

$$\sup_{\theta \in \Theta} \mathcal{R}(P_\theta, \widehat{Q}) = \int \mathcal{R}(P_\theta, P_{\widehat{\pi}}) d\widehat{\pi}(\theta),$$

and therefore

$$\int \mathcal{R}(P_\theta, \widehat{Q}) d\widehat{\pi}(\theta) \leq \int \mathcal{R}(P_\theta, P_{\widehat{\pi}}) d\widehat{\pi}(\theta).$$

On the other hand

$$\int \mathcal{R}(P_{\widehat{\theta}}, \widehat{Q}) d\widehat{\pi}(\theta) = \int \mathcal{R}(P_\theta, P_{\widehat{\pi}}) d\widehat{\pi}(\theta) + \mathcal{R}(P_{\widehat{\pi}}, \widehat{Q}),$$

implying that $\mathcal{R}(P_{\hat{\pi}}, \hat{Q}) = 0$, and consequently that $\hat{Q} = P_{\hat{\pi}}$. Finally, equation (2) is a consequence of the equality

$$\int \mathcal{R}(P_{\theta}, P_{\hat{\pi}}) d\hat{\pi}(\theta) = \inf_Q \sup_{\theta} \mathcal{R}(P_{\theta}, Q) \geq \int \mathcal{R}(P_{\theta}, P_{\pi}) d\pi(\theta).$$

On the other hand, let us now assume that $\hat{\pi}$ is solution of equation (2). In this case, we can write

$$\mathcal{R}_{\hat{\pi}}(P_{\hat{\pi}}) - \mathcal{R}_{\pi}(P_{\pi}) = \int \mathcal{R}(P_{\theta}, P_{\hat{\pi}}) (d\hat{\pi} - d\pi)(\theta) + \mathcal{R}(P_{\pi}, P_{\hat{\pi}}).$$

Applying this inequality to $\pi = \lambda\nu + (1 - \lambda)\hat{\pi}$, for all $\lambda \in [0, 1]$, and $\nu \in \mathcal{M}_+^1(\Theta)$, we see that the right-hand side should have a positive derivative at $\lambda = 0$, showing that

$$\int \mathcal{R}(P_{\theta}, P_{\hat{\pi}}) (d\hat{\pi} - d\nu)(\theta) \geq 0.$$

Considering for any $\theta \in \Theta$, $\nu = \delta_{\theta}$, the Dirac mass at θ , we deduce that

$$\int \mathcal{R}(P_{\theta}, P_{\hat{\pi}}) d\hat{\pi}(\theta) = \sup_{\theta \in \Theta} \mathcal{R}(P_{\theta}, P_{\hat{\pi}}),$$

showing from the first part of the proposition that $P_{\hat{\pi}}$ is indeed the minimax coding distribution in this case. \square

2.3. ORACLE INEQUALITIES. Another point of view is to consider a family of coding distributions, $\{Q_{\theta} : \theta \in \Theta\}$, instead of considering a family of source distributions. Given some prior probability measure $\pi \in \mathcal{M}_+^1(\Theta)$, it is easy to compare the mean length of the code mixture

$$Q_{\pi} \stackrel{\text{def}}{=} \int Q_{\theta} d\pi(\theta),$$

with the mean code length of any Q_{θ} , $\theta \in \Theta$. Indeed,

PROPOSITION 2.4 *For any $(x_1, \dots, x_n) \in A^n$,*

$$\begin{aligned} -\log_2 [Q_{\pi}(x^n)] &\leq \inf_{\rho \in \mathcal{M}_+^1(\Theta)} \int -\log_2 [Q_{\theta}(x^n)] d\rho(\theta) + \mathcal{R}(\rho, \pi) \\ &\leq \inf_{\theta \in \Theta} -\log_2 [Q_{\theta}(x^n)] - \log_2 [\pi(\{\theta\})]. \end{aligned}$$

Consequently, for any source distribution $P \in \mathcal{M}_+^1(A^n)$,

$$\begin{aligned} \mathcal{R}(P, Q_\pi) &\leq \inf_{\rho \in \mathcal{M}_+^1(\Theta)} \int \mathcal{R}(P, Q_\theta) d\rho(\theta) + \mathcal{R}(\rho, \pi) \\ &\leq \inf_{\theta \in \Theta} \mathcal{R}(P, Q_\theta) - \log_2[\pi(\{\theta\})]. \end{aligned}$$

PROOF. If $\mathcal{R}(\rho, \pi) = +\infty$, there is nothing to prove. Let us assume consequently that $\mathcal{R}(\rho, \pi) < \infty$. Let $\pi = \pi_s + \pi_a$ where $\pi_a \ll \rho$ and π_s and ρ are singular. It is easy to see that $\frac{d\pi_a}{d\rho}(\theta) = \frac{d\rho}{d\pi}(\theta)^{-1}$, ρ almost surely, therefore

$$\begin{aligned} \int Q_\theta(x^n) d\pi(\theta) d\pi(\theta) &\geq \int Q_\theta(x^n) d\pi_a(\theta) \\ &= \int \exp\left\{ \log[Q_\theta(x^n)] + \log\left[\frac{d\pi_a}{d\rho}(\theta)\right] \right\} d\rho(\theta) \\ &= \int \exp\left\{ \log[Q_\theta(x^n)] - \log\left[\frac{d\rho}{d\pi}(\theta)\right] \right\} d\rho(\theta) \\ &\geq \exp\left\{ \int \log[Q_\theta(x^n)] d\rho(\theta) - \mathcal{K}(\rho, \pi) \right\}, \end{aligned}$$

by Jensen's inequality. \square

COROLLARY 2.5 (DOUBLE MIXTURE CODES AND BAYESIAN MODEL AVERAGING) *Consider a family of models of ideal codes $\{Q_\theta \in \mathcal{M}_+^1(A^n); \theta \in \Theta_k\}$, $k \in \mathbb{N}$, where Θ_k are measurable parameter spaces. Let $\mu \in \mathcal{M}_+^1(\mathbb{N})$ and for each $k \in \mathbb{N}$ let $\nu_k \in \mathcal{M}_+^1(\Theta_k)$ be some prior measure on Θ_k . Let us consider $\Theta = \bigcup_{k \in \mathbb{N}} (\{k\} \times \Theta_k)$. Let the prior distribution π on Θ be defined for any measurable set $B \subset \Theta$ by the formula*

$$\pi(B) = \int \mathbf{1}[(k, \theta) \in B] d\mu(k) d\nu_k(\theta).$$

The coding distribution Q_π is called a double mixture code. The sequential prediction method that estimates $\mathbb{P}_{X_i|X^{i-1}}$ by

$$Q_\pi(x_i|x^{i-1}) \stackrel{\text{def}}{=} Q_\pi(x^i)/Q_\pi(x^{i-1})$$

is called Bayesian Model Averaging. It is such that

$$\sum_{i=1}^n -\log_2[Q_\pi(x_i|x^{i-1})] \leq \inf_{k \in \mathbb{N}} -\log_2 Q_{\nu_k}(x_i|x^{i-1}) - \log_2[\mu(k)].$$

Consequently for any source $X^n \in A^n$,

$$\begin{aligned} & \sum_{i=1}^n \mathbb{E} \left\{ \mathcal{K} \left[\mathbb{P}_{X_i|X^{i-1}}, Q_\pi(\cdot|X^{i-1}) \right] \right\} \\ & \leq \inf_{k \in \mathbb{N}} \left\{ \sum_{i=1}^n \mathbb{E} \left\{ \mathcal{K} \left[\mathbb{P}_{X_i|X^{i-1}}, Q_{\nu_k}(\cdot|X^{i-1}) \right] \right\} - \log[\mu(k)] \right\}. \end{aligned}$$

PROOF. It is a consequence of the previous proposition and of the identity

$$\sum_{i=1}^n -\log_2[Q_\pi(x_i|x^{i-1})] = -\log_2[Q_\pi(x^n)].$$

Let us remark that we could have replaced in this proof the family $\{Q_{\nu_k}, k \in \mathbb{N}\}$ by any family of ideal codes $\{Q_k, k \in \mathbb{N}\}$. \square

This is a result about sequential prediction. We consider estimators $Q_k(\cdot|X^{i-1})$ of the conditional distributions $\mathbb{P}_{X_i|X^{i-1}}$ and are provided with an aggregated estimator $Q_\pi(\cdot|X^{i-1})$ that performs almost as well as the best choice of k , when performance is measured by the sum of the expected Kullback divergences over all sample sizes in the range $1, \dots, n$. The sequential nature of the result comes from the fact that the criterion is a sum over all sample sizes : this is an example of cumulated risk function. The following two subsections are devoted to an important example of double mixture codes : the context tree weighting algorithm, where we are going to compute everything explicitly.

2.4. MIXTURES OF I.I.D. DISTRIBUTIONS. Let $\Theta = \mathcal{M}_+^1(A)$ and consider the Lebesgue measure λ on Θ . Let the prior $\nu \ll \lambda$ be defined by its density

$$\frac{d\nu}{d\lambda}(\theta) = \frac{\Gamma(\frac{d}{2})}{\sqrt{d}\Gamma(\frac{1}{2})^d} \prod_{x \in A} \theta(x)^{-1/2}, \quad (3)$$

where $\Gamma(z) = \int_{\mathbb{R}_+} t^{z-1} \exp(-t) dt = (z-1)\Gamma(z-1)$ is the usual Γ function.

LEMMA 2.6 *The mixture code based on the prior ν and the family of product measures $Q_\theta(x^n) = \prod_{i=1}^n \theta(x_i)$ is such that*

$$Q_\nu(x^n) = \frac{\Gamma(\frac{d}{2}) \prod_{y \in A} \Gamma(n\bar{\mathbb{P}}_{x^n}(y) + \frac{1}{2})}{\Gamma(\frac{1}{2})^d \Gamma(n + \frac{d}{2})} = \frac{\prod_{y \in A} \prod_{j=1}^{n\bar{\mathbb{P}}_{x^n}(y)} (j - \frac{1}{2})}{\prod_{j=1}^n (j + \frac{d}{2} - 1)},$$

where

$$\bar{\mathbb{P}}_{x^n} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$$

is the empirical probability measure of the sequence (x_1, \dots, x_n) , so that $n\bar{\mathbb{P}}_{x^n}(y) = \sum_{i=1}^n \mathbf{1}(x_i = y)$.

PROOF. Let us remark that

$$Q_\theta(x^n) = \prod_{y \in A} \theta(y)^{n\bar{\mathbb{P}}_{x^n}(y)},$$

and thus

$$Q_\nu(x^n) = \frac{\Gamma(\frac{d}{2})}{\sqrt{d}\Gamma(\frac{1}{2})^d} \int_{\Theta} \prod_{y \in A} \theta(y)^{n\bar{\mathbb{P}}_{x^n}(y) - \frac{1}{2}} d\lambda(\theta). \quad (4)$$

On the other hand, for any $\alpha \in]-1, +\infty[A$,

$$\int_{\Theta} \prod_{y \in A} \theta(y)^{\alpha_y} d\lambda(\theta) = \frac{\sqrt{d} \prod_{y \in A} \Gamma(\alpha_y + 1)}{\Gamma(\sum_{y \in A} \alpha_y + d)}. \quad (5)$$

Indeed, the change of variables $t_y = s\theta(y)$, where $t_y \in \mathbb{R}_+$, $\theta \in \mathcal{M}_+^1(A)$ and $s = \sum_{y \in A} t_y \in \mathbb{R}_+$ shows that

$$\int_{\mathbb{R}_+^d} \prod_{y \in A} t_y^{\alpha_y} \exp(-t_y) dt_y = \int_{\mathbb{R}_+} \int_{\Theta} \prod_{y \in A} \theta(y)^{\alpha_y} s^{\alpha_y} \exp(-s) s^{d-1} d\lambda(\theta) \frac{ds}{\sqrt{d}}.$$

Putting together equations (4) and (5) gives the desired result. \square

PROPOSITION 2.7 (KRICHEVSKI, TROFIMOV) *For any $x^n \in A^n$,*

$$Q_\nu(x^n) \geq \frac{2^{-nH(\bar{\mathbb{P}}_{x^n})}}{dn^{(d-1)/2}} = d^{-1} n^{-(d-1)/2} \sup_{\theta \in \Theta} Q_\theta(x^n).$$

PROOF. Let us put for any $a = (a_i)_{i=1}^d \in \mathbb{N}^d$

$$\Delta(a) = \frac{\Gamma(\frac{d}{2}) \prod_{i=1}^d \Gamma(a_i + \frac{1}{2}) (\sum_{i=1}^d a_i)^{\sum_{i=1}^d a_i + (d-1)/2}}{\Gamma(\frac{1}{2})^d \Gamma(\sum_{i=1}^d a_i + \frac{d}{2}) \prod_{i=1}^d a_i^{a_i}}$$

We have to show that $\Delta(a) \geq d^{-1}$. Let us notice first that $\Delta((1, 0, \dots, 0)) = d^{-1}$, and that Δ is invariant under any permutation of the vector $(a_i)_{i=1}^d$. It

is therefore enough to check that $\Delta(a) \geq \Delta((a_1 - 1, a_2, \dots, a_d))$. Let us put $s = \sum_{i=1}^d a_i$ and $t = a_1$. With these notations,

$$\Delta(a) = \Delta(a_1 - 1, a_2^d) \frac{(t - \frac{1}{2})(t - 1)^{t-1} s^{s + \frac{d-1}{2}}}{(s + \frac{d}{2} - 1)t^t(s - 1)^{s-1 + \frac{d-1}{2}}}.$$

To end the proof, we have to check that

$$\frac{(t - \frac{1}{2})(t - 1)^{t-1} s^{s + \frac{d-1}{2}}}{(s + \frac{d}{2} - 1)t^t(s - 1)^{s-1 + \frac{d-1}{2}}} \geq 1, \quad t \geq 1, s \geq 2.$$

This can also be written as $g(s) \geq f(t)$, where

$$\begin{aligned} g(s) &= -(s + \frac{d-3}{2}) \log\left(\frac{s-1}{s}\right) - \log\left(\frac{s + \frac{d-3}{2}}{s}\right), \\ f(t) &= t \log(t) - (t-1) \log(t-1) - \log\left(t - \frac{1}{2}\right). \end{aligned}$$

Using the fact that $\log(1+z) \leq z$, we see that $g(s) \geq (s + \frac{d-3}{2})s^{-1} - \frac{d-3}{2s} = 1$. On the other hand, making a Taylor expansion of $z \mapsto z \log(z)$, we see that

$$(z+u) \log(z+u) = z \log(z) + u[\log(z) + 1] + \int_0^u (u-v) \frac{dv}{z+v}.$$

Applying this to $z = t - \frac{1}{2}$ and $u \in \{-\frac{1}{2}, \frac{1}{2}\}$, gives

$$\begin{aligned} &t \log(t) - (t-1) \log(t-1) - \log\left(t - \frac{1}{2}\right) - 1 \\ &= \int_0^{\frac{1}{2}} \left(\frac{1}{2} - v\right) \left(\frac{1}{t - \frac{1}{2} + v} - \frac{1}{t - \frac{1}{2} - v}\right) dv \leq 0, \end{aligned}$$

proving that $f(t) \leq 1$. \square

COROLLARY 2.8 *The redundancy of the KT ideal code is such that*

$$\mathcal{R}(P, Q_\nu) \leq \inf_{\theta \in \Theta} \mathcal{R}(P, Q_\theta) + \frac{d-1}{2} \log_2(n) + \log_2(d).$$

2.5. UNIVERSAL CODES AND CONTEXT TREE WEIGHTING. Let \mathcal{D} be the set of complete suffix sets of A^* and \mathcal{D}_L the subset of complete suffix sets of length not greater than L . More explicitly, let us introduce the notations $A^*w = \{w'w, w' \in A^*\}$ and

$$\overline{D} = D \cup \{w \in A^*; A^*w \cap D \neq \emptyset\},$$

the set of suffixes of D . We can define \mathcal{D}_L formally as

$$\mathcal{D}_L = \{D \subset A^*; D \cap A^*D = \emptyset, A^*\overline{D} = A^*D, \text{ and } \max_{w \in D} \ell(w) \leq L\}.$$

Let $D \in \mathcal{D}_L$ be a complete suffix set and let $f_D : A^{\mathbb{Z}^-} \rightarrow D$ be defined by

$$f_D(x_{-\infty}^0) = x_{-k}^0,$$

where k is the only index such that $x_{-k}^0 \in D$. Given a past context $x_{-\infty}^0$, we want to define a probability measure on $x_1^{+\infty}$. A stationary context tree distribution is a probability measure defined by

$$\mathbb{P}(X_n = x_n | X_{-\infty}^{n-1} = x_{-\infty}^{n-1}) = \mathbb{P}[X_1 = x_n | f_D(X_{-\infty}^0) = f_D(x_{-\infty}^{n-1})].$$

Its parameter set is (D, Θ_D) where $\Theta_D = \{\theta_w \in \mathcal{M}_+^1(A); w \in D\} = \mathcal{M}_+^1(A)^D$. Consider the parameter space

$$\Theta = \bigcup_{D \in \mathcal{D}_L} \{D\} \times \Theta_D.$$

Let us define on Θ the prior probability measure π by

$$\pi(B) = \int \mathbb{1}[(D, \theta) \in B] d\mu(D) \prod_{w \in D} d\nu(\theta_w),$$

where ν is the Krichevski Trofimov prior on $\mathcal{M}_+^1(A)$ and where for some real parameter $\alpha \in]0, 1[$

$$\mu(D) = \alpha^{(|D|-1)/(d-1)}(1-\alpha)^{|D \setminus A^L|}, \quad D \in \mathcal{D}_L, \quad (6)$$

is the measure of a Galton Watson process on complete suffix sets with offspring probability α . More precisely, the distribution of \overline{D} under μ is given by the formula

$$\begin{aligned} \mu(\overline{D} = \overline{B}) &= \prod_{j=1}^L \mu(\overline{D} \cap A^j = \overline{B} \cap A^j | \overline{D} \cap \overline{A^{j-1}} = \overline{B} \cap \overline{A^{j-1}}) \\ &\stackrel{\text{def}}{=} \prod_{j=1}^L \prod_{w \in \overline{B} \cap A^{j-1}} [\alpha \mathbb{1}(A^*w \subset \overline{B}) + (1-\alpha) \mathbb{1}(A^*w \cap \overline{B} = \emptyset)] \\ &= \alpha^{|\overline{B} \setminus B|+1} (1-\alpha)^{|B \setminus A^L|}, \quad B \in \mathcal{D}_L. \end{aligned}$$

Equation (6) is a consequence of the fact that

$$|\overline{B} \setminus B| = (|B| - 1)/(d - 1) - 1,$$

for any $B \in \mathcal{D}_L$, a fact that can easily be checked by induction.

PROPOSITION 2.9 *Let us be given some sequence $x_{1-L}^n \in A^{n+L}$. Let us consider the counters*

$$\begin{aligned} a_w^y(n) &= \sum_{k=1}^n \mathbb{1}(x_{k-\ell(w)}^{k-1} = w, x_k = y), \quad w \in A^* \cup \{\emptyset\}, y \in A, \\ b_w(n) &= \sum_{y \in A} a_w^y(n), \\ K_w(n) &= \frac{\Gamma(\frac{d}{2}) \prod_{y \in A} \Gamma(a_w^y(n) + \frac{1}{2})}{\Gamma(\frac{1}{2})^d \Gamma(b_w(n) + \frac{d}{2})}, \end{aligned}$$

The mixture code Q_π can be computed as

$$Q_\pi(x_1^n | x_{1-L}^0) = \sum_{D \in \mathcal{D}_L} \mu(D) \prod_{w \in D} K_w(n).$$

This computation can be made by induction according to the following scheme

$$\begin{aligned} p_w(n) &= K_w(n), \quad w \in A^L, \\ p_w(n) &= (1 - \alpha)K_w(n) + \alpha \prod_{y \in A} p_{yw}(n), \quad w \in A^* \cup \{\emptyset\}, \ell(w) < L, \\ p_\emptyset(n) &= Q_\pi(x_1^n | x_{1-L}^0). \end{aligned}$$

PROOF. For any $D \in \mathcal{D}_L$,

$$\begin{aligned} \int_{\Theta_D} Q_{D,\theta}(x_1^n | x_{1-L}^0) \prod_{w \in D} d\nu(\theta_w) \\ = \int_{\Theta_D} \prod_{w \in D} \prod_{y \in A} \theta_w(y)^{a_w^y(n)} \prod_{w \in D} d\nu(\theta_w) = \prod_{w \in D} K_w(n), \end{aligned}$$

so that

$$Q_\pi(x_1^n | x_{1-L}^0) = \int_{\mathcal{D}_L} \prod_{w \in D} K_w(n) d\mu(D).$$

Let

$$p_w(n) = \int_{\mathcal{D}_L} \prod_{w' \in D \cap (A^*w \cup \{w\})} K_{w'}(n) d\mu(D | w \in \bar{D} \cup \{\emptyset\}),$$

so that $p_\emptyset(n) = Q_\pi(x_1^n | x_{1-L}^0)$. We see from the definition of μ that

$$p_w(n) = \mu(w \in D | w \in \bar{D}) K_w(n)$$

$$\begin{aligned}
& + \mu(A^*w \subset \overline{D} \mid w \in \overline{D}) \prod_{y \in A} \int_{\mathcal{D}_L} \prod_{w' \in D \cap (A^*yw \cup \{yw\})} K_{w'}(n) d\mu(D \mid yw \in \overline{D}) \\
& = (1 - \alpha)K_w(n) + \alpha \prod_{y \in A} p_{yw}(n), \quad w \in A^* \cup \{\emptyset\}, \ell(w) < L,
\end{aligned}$$

and that $p_w(n) = K_w(n)$ for any $w \in A^L$. \square

Exercise 5 Show that the computation of the context tree weighting ideal code can be updated online according to the following rules :

$$\begin{aligned}
& a_w^y(0) = b_w(y) = 0, \\
& K_w(0) = 1, \\
& p_w(0) = 1, \\
& \quad \vdots \\
& a_w^y(n) = \begin{cases} a_w^y(n-1) + 1, & \text{when } w = x_{n-\ell(w)}^{n-1} \text{ and } y = x_n, \\ a_w^y(n-1), & \text{otherwise,} \end{cases} \\
& b_w(n) = \begin{cases} b_w(n-1) + 1, & \text{when } w = x_{n-\ell(w)}^{n-1}, \\ b_w(n-1), & \text{otherwise,} \end{cases} \\
& K_w(n) = \begin{cases} K_w(n-1) \frac{a_w^{x_n}(n) - \frac{1}{2}}{b_w(n) + \frac{d}{2} - 1}, & \text{when } w = x_{n-\ell(w)}^{n-1}, \\ K_w(n-1), & \text{otherwise,} \end{cases} \\
& p_w(n) = \begin{cases} p_w(n-1), & \text{when } w \neq x_{n-\ell(w)}^{n-1}, \\ K_w(n), & \text{when } w = x_{n-L}^{n-1}, \\ (1 - \alpha)K_w(n) + \alpha \prod_{y \in A} p_{yw}(n), & \text{when } w = x_{n-\ell}^{n-1}, \text{ with } \ell < L. \end{cases}
\end{aligned}$$

Conclude that the number of operations needed to compute $Q_\pi(x_1^n \mid x_{1-L}^0)$ for a given value of $x_{1-L}^n \in A^{n+L}$ grows only linearly with n .

PROPOSITION 2.10 For any $x_{1-L}^n \in A^{n+L}$,

$$\begin{aligned}
-\log_2[Q_\pi(x_1^n \mid x_{1-L}^0)] & \leq \inf_{D \in \mathcal{D}_L} \inf_{\theta \in \Theta_D} \left\{ -\log_2[Q_{D,\theta}(x_1^n \mid x_{1-L}^0)] \right. \\
& \quad + \frac{|D|(d-1)}{2} \log\left(\frac{n}{|D|}\right) + |D| \log_2(d) \\
& \quad \left. - \frac{|D|-1}{d-1} \log_2(\alpha) - |D| \log_2(1 - \alpha) \right\}.
\end{aligned}$$

As a consequence for any process X_{1-L}^n , and any $x_{1-L}^0 \in A^L$,

$$\begin{aligned}
& \mathcal{R}[\mathbb{P}_{X_1^n | X_{1-L}^0 = x_{1-L}^0}, Q_\pi(\cdot | x_{1-L}^0)] \\
& \leq \inf_{D \in \mathcal{D}_L} \inf_{\theta \in \Theta_D} \left\{ \mathcal{R}[\mathbb{P}_{X_1^n | X_{1-L}^0 = x_{1-L}^0}, Q_{D,\theta}(\cdot | x_{1-L}^0)] \right. \\
& \quad \left. + \frac{|D|(d-1)}{2} \log\left(\frac{n}{|D|}\right) + |D| \log_2(d) \right. \\
& \quad \left. - \frac{|D|-1}{d-1} \log_2(\alpha) - |D| \log_2(1-\alpha) \right\}.
\end{aligned}$$

PROOF. Let $\gamma(z) = \frac{d-1}{2} \log_2(z) + \log_2(d)$. According to Proposition 2.7 (page 17),

$$\begin{aligned}
& -\log_2 \left[\int_{\Theta_D} Q_{D,\theta}(x_1^n | x_{1-L}^0) \prod_{w \in D} d\nu(\theta_w) \right] \\
& \leq \inf_{\theta \in \Theta_D} -\log_2 [Q_{D,\theta}(x_1^n | x_{1-L}^0)] + \sum_{w \in D} \gamma[b_w(n)] \\
& \leq \inf_{\theta \in \Theta_D} -\log_2 [Q_{D,\theta}(x_1^n | x_{1-L}^0)] + |D| \gamma \left(\frac{1}{|D|} \sum_{w \in D} b_w(n) \right),
\end{aligned}$$

where we have used the fact that γ is concave. As $\sum_{w \in D} b_w(n) = n$, we get

$$\begin{aligned}
-\log_2 [Q_\pi(x_1^n | x_{1-L}^0)] & \leq \inf_{D \in \mathcal{D}_L} \inf_{\theta \in \Theta_D} -\log_2 [Q_{D,\theta}(x_1^n | x_{1-L}^0)] \\
& \quad + |D| \gamma \left(\frac{n}{|D|} \right) - \log_2 [\mu(D)],
\end{aligned}$$

as stated in the proposition. \square

Exercise 6 (A simulation study of a switch code) *Let us consider the alphabet $A = \{1, \dots, 4\}^2$ of size 16. Let us define $q \in \mathcal{M}_+^1(A)$ as*

$$q(i, j) = \frac{4(i-1) + j}{8 \times 17}$$

and let $q_1(i) = \sum_{j=1}^4 q(i, j)$ and $q_2(j) = \sum_{i=1}^4 q(i, j)$ be its two marginal distributions. Let

$$p(i, j) = \frac{4}{5} q(i, j) + \frac{1}{5} q_1(i) q_2(j), \quad i, j \in \{1, \dots, 4\}.$$

For any $\theta_1, \theta_2 \in \Theta_1 \stackrel{\text{def}}{=} \mathcal{M}_+^1(\{1, \dots, 4\})$, let

$$Q_{\theta_1, \theta_2}(x^n) = \prod_{i=1}^n \theta_1(x_{i,1}) \theta_2(x_{i,2}), \quad x^n \in A^n,$$

where $x_i = (x_{i,1}, x_{i,2}) \in \{1, \dots, 4\}^2$. In the same way, for any $\theta \in \Theta_2 \stackrel{\text{def}}{=} \mathcal{M}_+^1(A)$, let

$$Q_\theta(x^n) = \prod_{i=1}^n \theta(x_i).$$

Let ν_1 be the Krichevski Trofimov prior distribution on Θ_1 , given by equation (3, page 16), and ν_2 the Krichevski Trofimov prior distribution on Θ_2 . Let us consider the following ideal codes

$$\begin{aligned} Q_1(x^n) &= \int_{\Theta_1 \times \Theta_1} Q_{\theta_1, \theta_2}(x^n) d\nu_1(\theta_1) d\nu_1(\theta_2), \\ Q_2(x^n) &= \int_{\Theta_2} Q_\theta(x^n) d\nu_2(\theta), \\ Q_m(x^n) &= \frac{1}{2}Q_1(x^n) + \frac{1}{2}Q_2(x^n), \\ Q_s(x^n) &= \begin{cases} Q_1(x^n), & \text{when } n < m, \\ Q_1(x^m)Q_2(x^n)/Q_2(x^m), & \text{when } n \geq m. \end{cases} \end{aligned}$$

Take $m = 2000$, and write a program to compute for $n \in 200\mathbb{N}$, $n \leq 10000$ the empirical mean ideal code length

$$L(Q, n) = -\frac{1}{100} \sum_{j=1}^{100} \log_2 [Q(X_1^n(j))],$$

where $X_1^n(j)$ are independent trials of the i.i.d. random process distributed according to $p^{\otimes n}$, and where $Q \in \{Q_1, Q_2, Q_m, Q_s\}$. Plot the curves

$$\begin{aligned} n &\mapsto L(Q_2, n) - L(Q_1, n), \\ n &\mapsto L(Q_m, n) - L(Q_1, n), \\ n &\mapsto L(Q_s, n) - L(Q_1, n). \end{aligned}$$

Make your comments (about the weak points and possible improvements of mixture codes).

Programming hint : all you need is to simulate the empirical measure of X_1^n , something you can do in Octave, using the function

```
tabulate(rand(1,n), df),
```

where **df** is the distribution function of your probability vector (to be computed with **cumsum**).

3. PAC-BAYES BOUNDS FOR SUPERVISED CLASSIFICATION

3.1. INTRODUCTION. PAC-Bayes theory was first developed in the framework of supervised classification (see [6, 7, 8, 9, 5]) and subsequently extended to other settings. We will not deal with these extensions in the present notes, but instead focus on supervised classification, a setting that plays a central role in statistical learning theory and requires specific techniques of proofs.

In this section, we are given some i.i.d. sample $(W_i)_{i=1}^n \in \mathcal{W}^n$, where \mathcal{W} is a measurable space, and some binary measurable loss function $L : \mathcal{W} \times \Theta \rightarrow \{0, 1\}$, where Θ is a measurable parameter space. Our aim is to minimize with respect to $\theta \in \Theta$ the expected loss

$$\int L(w, \theta) d\mathbb{P}(w),$$

where \mathbb{P} is the marginal distribution of the observed sample $(W_i)_{i=1}^n$. More precisely, assuming that \mathbb{P} is unknown, we would like to find an estimator $\hat{\theta}(W_1^n)$ depending on the observed sample W_1^n such that the excess risk

$$\int L(w, \hat{\theta}) d\mathbb{P}(w) - \inf_{\theta \in \Theta} \int L(w, \theta) d\mathbb{P}(w)$$

is small. The previous quantity is random, since $\hat{\theta}$ depends on the random sample W_1^n . Therefore, how small it is can be understood in different ways. Here we will focus on the *deviations* of the excess risk. Accordingly, we will look for estimators providing a small risk with a probability close to one.

A typical example of such a problem is provided by supervised classification. In this setting $\mathcal{W} = \mathcal{X} \times \mathcal{Y}$, where \mathcal{Y} is a finite set, $W_i = (X_i, Y_i)$, where (X_i, Y_i) are input-output pairs, a family of measurable classification rules $\{f_\theta : \mathcal{X} \rightarrow \mathcal{Y}; \theta \in \Theta\}$ is considered and the loss function $L(w, \theta)$ is defined as the classification error

$$L[(x, y), \theta] = \mathbf{1}[f_\theta(x) \neq y].$$

Accordingly the aim is to minimize the expected classification error

$$\mathbb{P}_{X,Y}[f_\theta(X) \neq Y]$$

in view of a sample $(X_i, Y_i)_{i=1}^n$ of observations.

Let us remark that the point of view adopted in this section is different from the point of view of the previous section in some important ways. With the tools of the last section, we could have considered the loss function

$$L_c[(x, y), q_\theta] = -\log_2[q_\theta(y | x)],$$

where $q_\theta(y|x)$ is some family of conditional distributions indexed by θ .

We could then have obtained results concerned with

$$\sum_{i=1}^n L_c \left[(X_i, Y_i), \int q_\theta d\rho_i(\theta) \right],$$

where $\rho_i \in \mathcal{M}_+^1(\Theta)$ is a random probability measure defined with the help of some prior probability measure $\pi \in \mathcal{M}_+^1(\Theta)$ by its density

$$\frac{d\rho_i}{d\pi}(\theta) = \frac{\prod_{i=1}^n q_\theta(Y_i|X_i)}{\int \prod_{i=1}^n q_{\theta'}(Y_i|X_i) d\pi(\theta')}.$$

Namely, the previous section provides the following almost sure upper bound:

$$\sum_{i=1}^n L_c \left[(X_i, Y_i), \int q_\theta d\rho_i(\theta) \right] \leq \inf_{\theta \in \Theta} \sum_{i=1}^n L_c [(X_i, Y_i), q_\theta] - \log_2 [\pi(\{\theta\})]. \quad (7)$$

Importantly, the previous section is concerned with a *cumulated* notion of risk, whereas this one will deal with an *instantaneous* notion of risk. For the cumulated risk, we could obtain almost sure results, this will not be possible any more for the instantaneous risk we are considering here. We will have instead to make an hypothesis about the probabilistic nature of the observations, assuming that they are independent, and obtain results holding with a probability close to one but not equal to one.

Exercise 7 *To make the link with the previous section more specific, we can focus on*

$$q_\theta(y|x) = \frac{\exp\{-\beta \mathbf{1}[f_\theta(x) \neq y]\}}{[1 + (|Y| - 1) \exp(-\beta)]}.$$

Show that

$$\begin{aligned} [1 - \exp(-\beta)] \sum_{i=1}^n \int \mathbf{1}[f_\theta(X_i) \neq Y_i] d\rho_i(\theta) \\ \leq \inf_{\theta \in \Theta} \beta \sum_{i=1}^n \mathbf{1}[f_\theta(X_i) \neq Y_i] - \log_2 [\pi(\{\theta\})]. \end{aligned}$$

(Hint : consider a lower bound of the left-hand side of equation (7, page 25) and an upper bound of its right-hand side, using the fact that $\log(1+z) \leq z$.) Conclude that, almost surely

$$\begin{aligned} \sum_{i=1}^n \int L[(X_i, Y_i), \theta] d\rho_i(\theta) \\ \leq \inf_{\theta \in \Theta} \frac{1}{1 - \beta/2} \sum_{i=1}^n L[(X_i, Y_i), \theta] - \frac{\log_2[\pi(\{\theta\})]}{\beta(1 - \beta/2)}. \end{aligned}$$

(Hint: use the inequality $\exp(-\beta) \leq 1 - \beta + \beta^2/2$, $\beta \geq 0$.) How would you choose the value of the parameter β ?

3.2. DEVIATION BOUNDS FOR SUMS OF BERNOULLI RANDOM VARIABLES. Given some parameter $\lambda \in \mathbb{R}$, let us consider the (normalized) log-Laplace transform of the Bernoulli distribution :

$$\Phi_\lambda(p) \stackrel{\text{def}}{=} -\frac{1}{\lambda} \log[1 - p + p \exp(-\lambda)].$$

Let us also consider the Kullback-Leibler divergence of two Bernoulli distributions

$$K(q, p) \stackrel{\text{def}}{=} q \log\left(\frac{q}{p}\right) + (1 - q) \log\left(\frac{1 - q}{1 - p}\right).$$

In the sequel $\bar{\mathbb{P}}$ will be the empirical measure

$$\bar{\mathbb{P}} = \frac{1}{n} \sum_{i=1}^n \delta_{W_i}$$

of an i.i.d. sample $(W_i)_{i=1}^n$ drawn from $\mathbb{P}^{\otimes n} \in \mathcal{M}_+^1(\mathcal{W}^n)$. We will use a short notation for integrals, putting for any $\rho, \pi \in \mathcal{M}_+^1(\Theta)$ and any integrable function $f \in \mathbb{L}_1(\mathcal{W} \times \Theta^2, \mathbb{P} \otimes \pi \otimes \rho)$

$$f(\mathbb{P}, \rho, \pi) = \int f(w, \theta, \theta') d\mathbb{P}(w) d\rho(\theta) d\pi(\theta'),$$

so that for instance $L(\mathbb{P}, \rho) = \int L(w, \theta) d\mathbb{P}(w) d\rho(\theta)$.

Let us recall first Chernoff's bound.

PROPOSITION 3.1 For any fixed value of the parameter $\theta \in \Theta$, the identity

$$\int \exp[-\lambda L(\bar{\mathbb{P}}, \theta)] d\mathbb{P}^{\otimes n} = \exp\left\{-\lambda \Phi_\lambda[L(\mathbb{P}, \theta)]\right\}$$

shows that with probability at least $1 - \epsilon$,

$$L(\mathbb{P}, \theta) \leq B_+[L(\bar{\mathbb{P}}, \theta), \log(\epsilon^{-1})/n],$$

$$\begin{aligned} \text{where } B_+(q, \delta) &= \inf_{\lambda \in \mathbb{R}_+} \Phi_\lambda^{-1} \left(q + \frac{\delta}{\lambda} \right) \\ &= \sup \left\{ p \in [0, 1] : K(q, p) \leq \delta \right\}, \quad q \in [0, 1], \delta \in \mathbb{R}_+. \end{aligned}$$

Moreover

$$-\delta q \leq B_+(q, \delta) - q - \sqrt{2\delta q(1-q)} \leq 2\delta(1-q).$$

In the same way, the identity

$$\int \exp[\lambda L(\bar{\mathbb{P}}, \theta)] d\mathbb{P}^{\otimes n} = \exp \left\{ \lambda \Phi_{-\lambda} [L(\mathbb{P}, \theta)] \right\}$$

shows that with probability at least $1 - \epsilon$

$$L(\bar{\mathbb{P}}, \theta) \leq B_- [L(\mathbb{P}, \theta), \log(\epsilon^{-1})/n],$$

$$\begin{aligned} \text{where } B_-(q, \delta) &= \inf_{\lambda \in \mathbb{R}_+} \Phi_{-\lambda}(q) + \frac{\delta}{\lambda} \\ &= \sup \left\{ p \in [0, 1] : K(p, q) \leq \delta \right\}, \quad q \in [0, 1], \delta \in \mathbb{R}_+, \end{aligned}$$

and

$$-\delta q \leq B_-(q, \delta) - q - \sqrt{2\delta q(1-q)} \leq 2\delta(1-q).$$

Let us mention here some important identity.

PROPOSITION 3.2 *For any probability measures π and ρ on some measurable space, such that $\mathcal{K}(\rho, \pi) < \infty$, and any bounded measurable function h , let us define the transformed probability measure $\pi_{\exp(h)} \ll \pi$ by its density*

$$\frac{d\pi_{\exp(h)}}{d\pi} = \frac{\exp(h)}{Z},$$

where $Z = \int \exp(h) d\pi$. Let us moreover define

$$\mathbf{Var}(h d\pi) = \int (h - \int h d\pi)^2 d\pi.$$

The expectations with respect to ρ and π of h and the log-Laplace transform of h are linked by the identities

$$\int h d\rho - \mathcal{K}(\rho, \pi) + \mathcal{K}(\rho, \pi_{\exp(h)}) = \log \left[\int \exp(h) d\pi \right] \quad (8)$$

$$= \int h d\pi + \int_0^1 (1 - \alpha) \mathbf{Var} [h d\pi_{\exp(\alpha h)}] d\alpha. \quad (9)$$

PROOF. The first identity is a straightforward consequence of the definitions of $\pi_{\exp(h)}$ and of the Kullback-Leibler divergence function. The second one is the Taylor expansion of order one with integral remainder of the function

$$f(\alpha) = \log \left[\int \exp(\alpha h) d\pi \right],$$

which says that $f(1) = f(0) + f'(0) + \int_0^1 (1-\alpha) f''(\alpha) d\alpha$. \square

Exercise 8 Prove that $f \in \mathcal{C}^\infty$. Hint : write

$$h^k \exp(\alpha h) = h^k + \int_0^{+\infty} \mathbf{1}(\gamma \leq \alpha) h^{k+1} \exp(\gamma h) d\gamma$$

and use Fubini's theorem to show that $\alpha \mapsto \int h^k \exp(\alpha h) d\pi$ belongs to \mathcal{C}^1 and compute its derivative.

Let us come now to the proof of Proposition 3.1 (page 26). Chernoff's inequality reads

$$\Phi_\lambda [L(\mathbb{P}, \theta)] - \frac{\log(\epsilon^{-1})}{n\lambda} \leq L(\bar{\mathbb{P}}, \theta),$$

where the inequality holds with probability at least $1 - \epsilon$. Since the left-hand side is non-random, it can be optimized in λ , giving

$$L(\mathbb{P}, \theta) \leq B_+ [L(\bar{\mathbb{P}}, \theta), \log(\epsilon^{-1})/n].$$

Exercise 9 Prove this statement in more details. For any integer $k > 1$, consider the event

$$A_k = \left\{ \sup_{\lambda \in \mathbb{R}_+} F(\lambda) - k^{-1} > L(\bar{\mathbb{P}}, \theta) \right\},$$

where $F(\lambda) = \Phi_\lambda [L(\mathbb{P}, \theta)] - \frac{\log(\epsilon^{-1})}{n\lambda}$. Show that $\mathbb{P}^{\otimes n}(A_k) \leq \epsilon$ by choosing some suitable value of λ . Remark that $A_k \subset A_{k+1}$ and conclude that $\mathbb{P}^{\otimes n}(\cup_k A_k) \leq \epsilon$.

Since

$$\lim_{\lambda \rightarrow +\infty} \Phi_\lambda^{-1} \left(q + \frac{\delta}{\lambda} \right) = \lim_{\lambda \rightarrow +\infty} \frac{1 - \exp(-\lambda q - \delta)}{1 - \exp(-\lambda)} \leq 1,$$

$$B_+(q, \delta) \leq 1.$$

Applying equation (8, page 27) to Bernoulli distributions gives

$$\lambda \Phi_\lambda(p) = \lambda q + K(q, p) - K(q, p_\lambda)$$

where

$$p_\lambda = \frac{p}{p + (1-p)\exp(\lambda)}.$$

This shows that

$$\begin{aligned} B_+(q, \delta) &= \sup \left\{ p \in [0, 1] : \Phi_\lambda(p) \leq q + \frac{\delta}{\lambda}, \lambda \in \mathbb{R}_+ \right\} \\ &= \sup \left\{ p \in [q, 1[: K(q, p) \leq \delta + K(q, p_\lambda), \lambda \in \mathbb{R}_+ \right\} \\ &= \sup \left\{ p \in [q, 1[: K(q, p) \leq \delta \right\} \\ &= \sup \left\{ p \in [0, 1] : K(q, p) \leq \delta \right\}, \end{aligned}$$

because when $q \leq p < 1$ then $\lambda = \log\left(\frac{q^{-1}-1}{p^{-1}-1}\right) \in \mathbb{R}_+$, $q = p_\lambda$ and therefore $K(q, p_\lambda) = 0$.

Let us remark now that $\frac{\partial^2}{\partial x^2} K(x, p) = x^{-1}(1-x)^{-1}$. Thus if $p \geq q \geq 1/2$, then

$$K(q, p) \geq \frac{(p-q)^2}{2q(1-q)},$$

so that if $K(q, p) \leq \delta$, then

$$p \leq q + \sqrt{2\delta q(1-q)}.$$

Now if $q \leq 1/2$ and $p \geq q$ then

$$K(q, p) \geq \left\{ \begin{array}{ll} \frac{(p-q)^2}{2p(1-p)}, & p \leq 1/2 \\ \frac{(p-q)^2}{2(p-q)^2}, & p \geq 1/2 \end{array} \right\} \geq \frac{(p-q)^2}{2p(1-p)},$$

so that if $K(q, p) \leq \delta$, then

$$(p-q)^2 \leq 2\delta p(1-q),$$

implying that

$$p - q \leq \delta(1-q) + \sqrt{2\delta q(1-q) + \delta^2(1-q)^2} \leq \sqrt{2\delta q(1-q) + 2\delta(1-q)}.$$

On the other hand,

$$K(q, p) \leq \frac{(p-q)^2}{2 \min\{q(1-q), p(1-p)\}} \leq \frac{(p-q)^2}{2q(1-p)},$$

thus when $K(q, p) = \delta$ with $p > q$, then

$$(p - q)^2 \geq 2\delta q(1 - p),$$

implying that

$$p - q \geq -\delta q + \sqrt{2\delta q(1 - q) + \delta^2 q^2} \geq \sqrt{2\delta q(1 - q)} - \delta q.$$

Exercise 10 *The second part of Proposition 3.1 (page 26) is proved in the same way and left as an exercise.*

3.3. PAC-BAYES BOUNDS. We are now going to make Proposition 3.1 uniform with respect to θ . The PAC-Bayes approach to this is to randomize θ , so we will consider now joint distributions on $(W_1, \dots, W_n, \theta)$, where the distribution of (W_1, \dots, W_n) is still $\mathbb{P}^{\otimes n}$ and the conditional distribution of θ given the sample is given by some transition probability kernel $\rho : \mathcal{W}^n \rightarrow \mathcal{M}_+^1(\Theta)$, called in this context a posterior distribution*. This posterior distribution ρ will be compared with a prior (meaning non-random) probability measure $\pi \in \mathcal{M}_+^1(\Theta)$.

PROPOSITION 3.3 *Let us introduce the notation*

$$B_\Lambda(q, \delta) = \inf_{\lambda \in \Lambda} \Phi_\lambda^{-1} \left(q + \frac{\delta}{\lambda} \right).$$

For any prior probability measure $\pi \in \mathcal{M}_+^1(\Theta)$ and any $\lambda \in \mathbb{R}_+$,

$$\int \exp \left[\sup_{\rho \in \mathcal{M}_+^1(\Theta)} n\lambda \left\{ \Phi_\lambda [L(\mathbb{P}, \rho)] - L(\bar{\mathbb{P}}, \rho) \right\} - \mathcal{K}(\rho, \pi) \right] d\mathbb{P}^{\otimes n} \leq 1, \quad (10)$$

and therefore for any finite set $\Lambda \subset \mathbb{R}_+$, with probability at least $1 - \epsilon$, for any $\rho \in \mathcal{M}_+^1(\Theta)$,

$$L(\mathbb{P}, \rho) \leq B_\Lambda \left(L(\bar{\mathbb{P}}, \rho), \frac{\mathcal{K}(\rho, \pi) + \log(|\Lambda|/\epsilon)}{n} \right),$$

PROOF. The exponential moment inequality (10) is a consequence of equation (8, page 27), showing that

*We will assume that ρ is a regular conditional probability kernel, meaning that for any measurable set A the map $(w_1, \dots, w_n) \mapsto \rho(w_1, \dots, w_n)(A)$ is assumed to be measurable. We will also assume that the σ -algebra we consider on Θ is generated by a countable family of subsets. See [1][page 50] for more details

$$\begin{aligned} \exp \left\{ \sup_{\rho \in \mathcal{M}_+^1(\Theta)} n\lambda \int \left\{ \Phi_\lambda[L(\mathbb{P}, \theta)] - L(\bar{\mathbb{P}}, \theta) \right\} d\rho(\theta) - \mathcal{K}(\rho, \pi) \right\} \\ \leq \int \exp \left[n\lambda \left\{ \Phi_\lambda[L(\mathbb{P}, \theta)] - L(\bar{\mathbb{P}}, \theta) \right\} \right] d\pi(\theta), \end{aligned}$$

and of the fact that Φ_λ is convex, showing that $\Phi_\lambda[L(\mathbb{P}, \rho)] \leq \int \Phi_\lambda[L(\mathbb{P}, \theta)] d\rho(\theta)$. The deviation inequality follows as usual. \square

We cannot take the infimum on $\lambda \in \mathbb{R}_+$ as in Proposition 3.1 (page 26), because we can no more cast our deviation inequality in such a way that λ appears on some non-random side of the inequality. Nevertheless, we can get a more explicit bound from some specific choice of the set Λ .

PROPOSITION 3.4 *Let us define the least increasing upper bound of the variance of a Bernoulli distribution of parameter $p \in [0, 1]$ as*

$$\bar{v}(p) = \begin{cases} p(1-p), & p \leq 1/2, \\ 1/4, & \text{otherwise.} \end{cases}$$

Let us choose some positive integer parameter m and let us put

$$t = \frac{1}{4} \log \left(\frac{n}{8 \log[(m+1)/\epsilon]} \right).$$

With probability at least $1 - \epsilon$, for any $\rho \in \mathcal{M}_+^1(\Theta)$,

$$L(\mathbb{P}, \rho) \leq L(\bar{\mathbb{P}}, \rho) + B_m[L(\bar{\mathbb{P}}, \rho), \mathcal{K}(\rho, \pi), \epsilon],$$

where

$$\begin{aligned} B_m(q, e, \epsilon) &= \max \left\{ \sqrt{\frac{2\bar{v}(q)\{e + \log[(m+1)/\epsilon]\}}{n}} \cosh(t/m) \right. \\ &\quad \left. + \frac{2(1-q)\{e + \log[(m+1)/\epsilon]\}}{n} \cosh(t/m)^2, \right. \\ &\quad \left. \frac{2\{e + \log[(m+1)/\epsilon]\}}{n} \right\} \\ &\leq \sqrt{\frac{2\bar{v}(q)\{e + \log[(m+1)/\epsilon]\}}{n}} \cosh(t/m) \\ &\quad + \frac{2\{e + \log[(m+1)/\epsilon]\}}{n} \cosh(t/m)^2. \end{aligned}$$

Moreover, as soon as $n \geq 5$,

$$B_{\lfloor \log(n)^2 \rfloor - 1}(q, e, \epsilon) \leq B(q, e, \epsilon) \stackrel{\text{def}}{=} \sqrt{\frac{2\bar{v}(q)\{e + \log[\log(n)^2/\epsilon]\}}{n}} \cosh[\log(n)^{-1}] + \frac{2\{e + \log[\log(n)^2/\epsilon]\}}{n} \cosh[\log(n)^{-1}]^2, \quad (11)$$

so that with probability at least $1 - \epsilon$, for any $\rho \in \mathcal{M}_+^1(\Theta)$,

$$L(\mathbb{P}, \rho) \leq L(\bar{\mathbb{P}}, \rho) + \sqrt{\frac{2\bar{v}[L(\bar{\mathbb{P}}, \rho)]\{\mathcal{K}(\rho, \pi) + \log[\log(n)^2/\epsilon]\}}{n}} \cosh[\log(n)^{-1}] + \frac{2\{\mathcal{K}(\rho, \pi) + \log[\log(n)^2/\epsilon]\}}{n} \cosh[\log(n)^{-1}]^2.$$

PROOF. Let us put

$$\begin{aligned} q &= L(\bar{\mathbb{P}}, \rho), \\ \delta &= \frac{\mathcal{K}(\rho, \pi) + \log[(m+1)/\epsilon]}{n}, \\ \lambda_{\min} &= \sqrt{\frac{8 \log[(m+1)/\epsilon]}{n}}, \\ \Lambda &= \left\{ \lambda_{\min}^{1-k/m}, k = 0, \dots, m \right\}, \\ p &= B_\Lambda(q, \delta) = \inf_{\lambda \in \Lambda} \Phi_\lambda^{-1}\left(q + \frac{\delta}{\lambda}\right), \\ \hat{\lambda} &= \sqrt{\frac{2\delta}{\bar{v}(p)}}. \end{aligned}$$

According to equation (9, page 27) applied to Bernoulli distributions, for any $\lambda \in \Lambda$,

$$\Phi_\lambda(p) = p - \frac{1}{\lambda} \int_0^\lambda (\lambda - \alpha) p_\alpha (1 - p_\alpha) d\alpha \leq q + \frac{\delta}{\lambda}.$$

As moreover $p_\alpha \leq p$,

$$p - q \leq \inf_{\lambda \in \Lambda} \frac{\lambda \bar{v}(p)}{2} + \frac{\delta}{\lambda} = \inf_{\lambda \in \Lambda} \sqrt{2\delta \bar{v}(p)} \cosh\left[\log\left(\frac{\hat{\lambda}}{\lambda}\right)\right].$$

As $\bar{v}(p) \leq 1/4$ and $\delta \geq \frac{\log[(m+1)/\epsilon]}{n}$,

$$\sqrt{\frac{2\delta}{\bar{v}(p)}} = \hat{\lambda} \geq \lambda_{\min} = \sqrt{\frac{8 \log[(m+1)/\epsilon]}{n}}.$$

Therefore either $\lambda_{\min} \leq \hat{\lambda} \leq 1$, or $\hat{\lambda} > 1$. Let us consider these two cases separately.

If $\lambda_{\min} = \min \Lambda \leq \hat{\lambda} \leq \max \Lambda = 1$, then $\log(\hat{\lambda})$ is at distance at most t/m from some $\log(\lambda)$ where $\lambda \in \Lambda$, because $\log(\Lambda)$ is a grid with constant steps of size $2t/m$. Thus

$$p - q \leq \sqrt{2\delta\bar{v}(p)} \cosh(t/m).$$

If moreover $q \leq 1/2$, then $\bar{v}(p) \leq p(1-q)$, so that we obtain a quadratic inequality in p , whose solution is less than

$$p \leq q + \sqrt{2\delta q(1-q)} \cosh(t/m) + 2\delta(1-q) \cosh(t/m)^2.$$

If on the contrary $q \geq 1/2$, then $\bar{v}(p) = \bar{v}(q) = 1/4$ and

$$p \leq q + \sqrt{2\delta\bar{v}(q)} \cosh(t/m),$$

so that in both cases

$$p - q \leq \sqrt{2\delta\bar{v}(q)} \cosh(t/m) + 2\delta(1-q) \cosh(t/m)^2. \quad (12)$$

Let us consider now the case when $\hat{\lambda} > 1$. In this case

$$p - q \leq \sqrt{2\delta\bar{v}(p)} \hat{\lambda} = 2\delta.$$

In conclusion, applying Proposition 3.3 (page 30) we see that with probability at least $1 - \epsilon$, for any posterior distribution ρ ,

$$L(\mathbb{P}, \rho) \leq p \leq q + \max\left\{2\delta, \sqrt{2\delta\bar{v}(q)} \cosh(t/m) + 2\delta(1-q) \cosh(t/m)^2\right\},$$

which is precisely the statement to be proved.

In the special case when $m = \lfloor \log(n)^2 \rfloor - 1 \geq \log(n)^2 - 2$,

$$\frac{t}{m} \leq \frac{1}{4 \lfloor \log(n)^2 - 2 \rfloor} \log\left(\frac{n}{8 \log \lfloor \log(n)^2 - 1 \rfloor}\right) \leq \log(n)^{-1}$$

as soon as the last inequality holds, that is as soon as $n \geq \exp(\sqrt{2}) \simeq 4.11$ to make $\log(n)^2 - 2$ positive and

$$3 \log(n)^2 - 8 + \log(n) \log\left\{8 \log \lfloor \log(n)^2 - 1 \rfloor\right\} \geq 0,$$

which holds true for any $n \geq 5$, as can be checked numerically. \square

 4. LINEAR CLASSIFICATION AND SUPPORT VECTOR MACHINES

We are going in this section to consider more specifically the case of linear binary classification. In this setting $\mathcal{W} = \mathcal{X} \times \mathcal{Y} = \mathbb{R}^d \times \{-1, +1\}$, $w = (x, y)$, where $x \in \mathbb{R}^d$ and $y \in \{-1, +1\}$, $\Theta = \mathbb{R}^d$, and

$$L(w, \theta) = \mathbf{1}[\langle \theta, x \rangle y \leq 0].$$

Although we will stick in this presentation to the case when \mathcal{X} is a vector space of finite dimension, the results also apply to support vector machines, where the pattern space is some arbitrary space mapped to a Hilbert space \mathcal{H} by some implicit mapping $\Psi : \mathcal{X} \rightarrow \mathcal{H}$, $\Theta = \mathcal{H}$ and $L(w, \theta) = \mathbf{1}(\langle \theta, \Psi(x) \rangle y \leq 0)$. It turns out that classification algorithms do not need to manipulate \mathcal{H} itself, but only to compute scalar products of the form $k(x_1, x_2) = \langle \Psi(x_1), \Psi(x_2) \rangle$, defining a symmetric positive kernel k on the original pattern space \mathcal{X} . The converse is also true, any positive symmetric kernel k can be represented as a scalar product in some mapped Hilbert space (this is the Moore-Aronszajn theorem). Often used kernels on \mathbb{R}^d are

$$\begin{aligned} k(x_1, x_2) &= (1 + \langle x_1, x_2 \rangle)^s, \text{ for which } \dim \mathcal{H} < \infty, \\ k(x_1, x_2) &= \exp(-\|x_1 - x_2\|^2), \text{ for which } \dim \mathcal{H} = +\infty. \end{aligned}$$

In the following, we will work in \mathbb{R}^d , which covers only the case when $\dim \mathcal{H} < \infty$, but extensions would be possible.

Let us consider, after [5, 8] as prior probability measure π the centered Gaussian measure with covariance $\beta^{-1} \mathbf{Id}$, so that

$$\frac{d\pi}{d\theta}(\theta) = \left(\frac{\beta}{2\pi}\right)^{d/2} \exp\left(-\frac{\beta\|\theta\|^2}{2}\right).$$

Let us also consider the function

$$\begin{aligned} \varphi(x) &= \frac{1}{\sqrt{2\pi}} \int_x^{+\infty} \exp(-t^2/2) dt, \quad x \in \mathbb{R} \\ &\leq \min\left\{\frac{1}{x\sqrt{2\pi}}, \frac{1}{2}\right\} \exp\left(-\frac{x^2}{2}\right), \quad x \in \mathbb{R}_+. \end{aligned}$$

Let π_θ be the measure π shifted by θ , defined by the identity

$$\int h(\theta') d\pi_\theta(\theta') = \int h(\theta + \theta') d\pi(\theta').$$

In this case

$$\mathcal{K}(\pi_\theta, \pi) = \frac{\beta}{2} \|\theta\|^2,$$

and

$$L(w, \pi_\theta) = \varphi[\sqrt{\beta} \|x\|^{-1} \langle \theta, x \rangle y].$$

Thus the randomized loss function has an explicit expression : randomization replaces the indicator function of the negative real line by a smooth approximation. As we are eventually interested in $L(w, \theta)$, we will shift things a little bit, considering along with the classification error function L some *error with margin*

$$M(w, \theta) = \mathbf{1}[y \|x\|^{-1} \langle \theta, x \rangle \leq 1].$$

Unlike $L(w, \theta)$ which is independent of the norm of θ , the margin error $M(w, \theta)$ depends on $\|\theta\|$, counting a classification error each time x is at distance less than $\|x\|/\|\theta\|$ from the boundary $\{x' : \langle \theta, x' \rangle = 0\}$, so that the error with margin region is the complement of the open cone $\{x \in \mathbb{R}^d; y \langle \theta, x \rangle > \|x\|\}$.

Let us compute the randomized margin error

$$M(w, \pi_\theta) = \varphi\left\{\sqrt{\beta} [y \|x\|^{-1} \langle \theta, x \rangle - 1]\right\}.$$

It satisfies the inequality

$$M(w, \pi_\theta) \geq \varphi(-\sqrt{\beta}) L(w, \theta) = [1 - \varphi(\sqrt{\beta})] L(w, \theta).$$

Applying previous results we obtain

PROPOSITION 4.1 *With probability at least $1 - \epsilon$, for any $\theta \in \mathbb{R}^d$,*

$$L(\mathbb{P}, \theta) \leq [1 - \varphi(\sqrt{\beta})]^{-1} M(\mathbb{P}, \pi_\theta) \leq C_1(\theta),$$

where

$$C_1(\theta) = [1 - \varphi(\sqrt{\beta})]^{-1} B\left(M(\bar{\mathbb{P}}, \pi_\theta), \frac{\beta \|\theta\|^2}{2}, \epsilon\right),$$

the bound B being defined by equation (11, page 32).

We can now minimize this empirical upper-bound to define an estimator. Let us consider some estimator $\hat{\theta}$ such that

$$C_1(\hat{\theta}) \leq \inf_{\theta \in \mathbb{R}^d} C_1(\theta) + \zeta.$$

Then for any fixed parameter θ_* , $C_1(\theta) \leq C_1(\theta_*) + \zeta$. On the other hand, with probability at least $1 - \epsilon$

$$M(\bar{\mathbb{P}}, \pi_{\theta_*}) \leq B_- \left(M(\mathbb{P}, \pi_{\theta_*}), \frac{\log(\epsilon^{-1})}{n} \right).$$

Indeed

$$\begin{aligned} & \int \exp \left\{ n\lambda [M(\bar{\mathbb{P}}, \pi_{\theta_*}) - \Phi_{-\lambda}[M(\mathbb{P}, \pi_{\theta_*})]] \right\} d\mathbb{P}^{\otimes n} \\ & \leq \int \exp \left\{ n\lambda \int \left\{ M(\bar{\mathbb{P}}, \theta) - \Phi_{-\lambda}[M(\mathbb{P}, \theta)] \right\} d\pi_{\theta_*}(\theta) \right\} d\mathbb{P}^{\otimes n} \leq 1, \end{aligned}$$

because $p \mapsto -\Phi_{-\lambda}(p)$ is convex. As a consequence

PROPOSITION 4.2 *With probability at least $1 - 2\epsilon$,*

$$\begin{aligned} L(\mathbb{P}, \hat{\theta}) & \leq \\ & \inf_{\theta_* \in \Theta} [1 - \varphi(\sqrt{\beta})]^{-1} B \left(B_- \left(M(\mathbb{P}, \pi_{\theta_*}), \frac{\log(\epsilon^{-1})}{n} \right), \frac{\beta \|\theta_*\|^2}{2}, \epsilon \right) + \zeta. \end{aligned}$$

It is also possible to state a result in terms of empirical margins. Indeed

$$M(w, \pi_\theta) \leq M(w, \theta/2) + \varphi(\sqrt{\beta}).$$

Thus with probability at least $1 - \epsilon$, for any $\theta \in \mathbb{R}^d$,

$$L(\mathbb{P}, \theta) \leq C_2(\theta),$$

where

$$C_2(\theta) = [1 - \varphi(\sqrt{\beta})]^{-1} B \left(M(\bar{\mathbb{P}}, \theta/2) + \varphi(\sqrt{\beta}), \frac{\beta \|\theta\|^2}{2}, \epsilon \right).$$

However, C_1 and C_2 are non-convex criteria, faster minimization algorithms are available for the usual SVN loss function, for which it is also possible to derive some generalization bound. Indeed

$$M(w, \pi_\theta) = \varphi[\sqrt{\beta}(\langle \theta, x \rangle y - 1)] \leq (2 - \langle \theta, x \rangle y)_+ + \varphi(\sqrt{\beta}).$$

Thus we also have, using this time Proposition 3.3 (page 30)

PROPOSITION 4.3 *With probability at least $1 - \epsilon$, for any $\theta \in \mathbb{R}^d$,*

$$\begin{aligned} L(\mathbb{P}, \theta) &\leq [1 - \varphi(\sqrt{\beta})]^{-1} B_\Lambda \left(\int (2 - \langle \theta, x \rangle y)_+ d\bar{\mathbb{P}}(x, y) + \varphi(\sqrt{\beta}), \right. \\ &\quad \left. \frac{\beta \|\theta\|^2 + 2 \log(|\Lambda|/\epsilon)}{2n} \right) \\ &= [1 - \varphi(\sqrt{\beta})]^{-1} \inf_{\lambda \in \Lambda} \Phi_\lambda^{-1} \left[C_3(\lambda, \theta) + \varphi(\sqrt{\beta}) + \frac{\log(|\Lambda|/\epsilon)}{n\lambda} \right], \end{aligned}$$

where

$$C_3(\lambda, \theta) = \int (2 - \langle \theta, x \rangle y)_+ d\bar{\mathbb{P}}(x, y) + \frac{\beta \|\theta\|^2}{2n\lambda}.$$

The loss function $C_3(\lambda, \theta)$ is the most employed learning criterion for support vector machines, and is called the box constraint. It is convex in θ . There are fast algorithms to compute $\inf_\theta C_3(\lambda, \theta)$ for any fixed value of λ . Here we get an empirical criterion which could be used to optimize also the value of λ .

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] O. Catoni. *PAC-Bayesian Supervised Classification: The Thermodynamics of Statistical Learning*, volume 56 of *IMS Lecture Notes Monograph Series*. Institute of Mathematical Statistics, 2007. Pages i-xii, 1-163.
- [3] T. Cover and J. Thomas. *Elements of Information Theory*. Wiley and Sons, New York, second edition, 2006.
- [4] Pascal Germain, Alexandre Lacasse, François Laviolette, and Mario Marchand. Pac-bayesian learning of linear classifiers. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML '09*, pages 353–360, New York, NY, USA, 2009. ACM.
- [5] J. Langford and J. Shawe-Taylor. PAC-bayes & margins. In *Advances in Neural Information Processing Systems*, pages 423–430, 2002.

-
- [6] D. A. McAllester. PAC-Bayesian model averaging. In *Proceedings of the 12th annual conference on Computational Learning Theory*. Morgan Kaufmann, 1999.
- [7] D. A. McAllester. PAC-Bayesian stochastic model selection. *Mach. Learn.*, 51(1):5–21, April 2003.
- [8] David Mcallester. Simplified pac-bayesian margin bounds. In *In COLT*, pages 203–215, 2003.
- [9] M. Seeger. PAC-Bayesian generalization error bounds for gaussian process classification. Informatics report series EDI-INF-RR-0094, Division of Informatics, University of Edinburgh, 2002.
- [10] T. van Erven, P. D. Grünwald, and S. de Rooij. Catching up faster by switching sooner : A predictive approach to adaptive estimation with an application to the AIC-BIC dilemma. *JRSS B*, 2011.
- [11] F.M.J. Willems, Y.M. Shtarkov, and T.J. Tjalkens. The context-tree weighting method: basic properties. *IEEE Trans. Inform. Theory*, 41(3):653–664, 1995.