

Compression-based methods for nonparametric on-line prediction, regression, classification and density estimation of time series *

Boris Ryabko

Siberian State University of Telecommunications and Informatics,
Institute of Computational Technologies of Siberian Branch of Russian Academy of Sciences,
Novosibirsk, Russia. e-mail: boris@ryabko.net

Abstract

Jorma Rissanen has discovered some deep connections between universal coding (or universal data compression) and mathematical statistics. In particular, the MDL principle has been one of the most powerful methods of modern mathematical statistics. In this paper we apply Rissanen's approach and ideas to some statistical problems concerned with time series. We address the problem of nonparametric estimation of characteristics of stationary and ergodic time series. We consider finite-alphabet as well as real-valued time series and the following four problems: i) estimation of the (limiting) probability $P(u_0 \dots u_s)$ for every s and each sequence $u_0 \dots u_s$ of letters over the process alphabet (or estimation of the density $p(x_0, \dots, x_s)$ for real-valued time series), ii) so-called on-line prediction, where the conditional probability $P(x_{t+1}|x_1 x_2 \dots x_t)$ (or the conditional density $p(x_{t+1}|x_1 x_2 \dots x_t)$) should be estimated (when $x_1 x_2 \dots x_t$ is known), iii) regression and iv) classification (or so-called problems with side information). We show that so-called archivers (or data compressors) can be used as a tool for solving these problems. In particular, it is proven that any universal code (or universal data compressor) can be used as a basis for constructing asymptotically optimal methods for the above problems. (By definition, a universal code can "compress" any sequence generated by a stationary and ergodic source asymptotically to the Shannon entropy of the source.)

AMS subject classification: 60G10, 60J10, 62G07, 62G08, 62M20, 94A29.

keywords: time series, nonparametric estimation, prediction, universal coding, data compression, on-line prediction, Shannon entropy, stationary and ergodic process, regression.

1 Introduction

We consider a stationary and ergodic source which generates sequences $x_1 x_2 \dots$ of elements (letters) from some set (alphabet) A , which is either finite or real-valued. It is supposed that the probability distribution (or distribution of limiting probabilities) $P(x_1 = a_{i_1}, x_2 = a_{i_2}, \dots, x_t = a_{i_t})$ (or the density $p(x_1, x_2, \dots, x_t)$) is unknown, but we are given either one sample $x_1 \dots x_t$ or several (r) independent samples $x^1 = x_1^1 \dots x_{t_1}^1, \dots, x^r = x_1^r \dots x_{t_r}^r$ generated by the source. (Generally speaking, they cannot be combined into one sample for a stationary and ergodic source as it can be done for an i.i.d. one.) Of course, if someone knows the probability distribution (or the density) he has all information about the source and can solve all problems in the best

*Research was supported by Russian Foundation for Basic Research (grant no. 06-07-89025.)

way. Hence, precise estimations of the probability distribution and the density can be used for prediction, regression estimation, etc. In this paper we use the following scheme: we consider the problems of estimation of the probability distribution or density. Then we show how the solution can be applied to other problems, paying the main attention to the problem of prediction, because of its practical applications and importance for probability theory, information theory, statistics and other theoretical sciences, see [1, 16, 17, 20, 28, 29, 31, 34, 35, 36, 38, 41, 46]. We show that universal codes (or data compressors) can be applied directly to the problems of estimation, prediction, regression and classification. It is not surprising because for any stationary and ergodic source P generating letters from a finite alphabet and any universal code U the following equality is valid with probability 1:

$$\lim_{t \rightarrow \infty} t^{-1}(-\log P(x_1 \dots x_t) - |U(x_1 \dots x_t)|) = 0,$$

where $x_1 \dots x_t$ is generated by P . (Here and below $\log = \log_2$, $|v|$ is the length of v , if v is a word, and the number of elements of v if v is a set.) So, in fact, the length of the universal code ($|U(x_1 \dots x_t)|$) can be used as an estimate of the logarithm of the unknown probability and, obviously, $2^{-|U(x_1 \dots x_t)|}$ can be considered as the estimation of $P(x_1 \dots x_t)$. In fact, a universal code can be viewed as a non-parametrical estimation of (limiting) probabilities for stationary and ergodic sources. This was recognized shortly after the discovery of universal codes (for the set of stationary and ergodic processes with finite alphabets [40]) and universal codes were applied for solving prediction problem [35, 41].

We would like to emphasize that, on the one hand, all results are obtained in the framework of classical probability theory and mathematical statistics and, on the other hand, everyday methods of data compression (or archivers) can be used as a tool for density estimation, prediction and other problems, because they are practical realizations of universal codes. It is worth noting that modern data compressors are based on deep theoretical results of source coding theory (see, e.g., [11, 17, 21, 24, 33, 34, 35, 37, 48]) and have demonstrated high efficiency in practice as compressors of texts (zip, arj, rar, etc.), DNA sequences [24] and many other types of real data. In fact, archivers can find many kinds of latent regularities, that is why they look like a promising tool for prediction and other problems. Moreover, recently universal codes and archivers were effectively applied to some problems which are very far from data compression: first, their applications created a new and rapidly growing line of investigations in clustering and classification (see [4, 5] and references therein) and, second, universal codes were used as a basis for non-parametric tests for the main statistical hypotheses concerned with stationary and ergodic time series [44, 45]. The outline of the paper is as follows. Section 2 contains description of the Laplace predictor and its generalizations, a review of known results and description of one universal code. Sections 3 and 4 are devoted to processes with finite and real-valued alphabets, correspondingly.

2 Predictors and universal data compressors

2.1 The Laplace measure and on-line prediction for i.i.d. processes

We consider a source with unknown statistics which generates sequences $x_1 x_2 \dots$ of letters from some set (or alphabet) A . It will be convenient at first to describe briefly the prediction problem. Let the source generate a message $x_1 \dots x_{t-1} x_t$, $x_i \in A$ for all i , and the next letter x_{t+1} needs to be predicted. This can be traced back to Laplace who considered the problem how to estimate the probability that the sun will rise tomorrow, given that it has risen every day since Creation

(see [12]). In our notation the alphabet A contains two letters: 0 ("the sun rises") and 1 ("the sun does not rise"), t is the number of days since Creation, $x_1 \dots x_{t-1} x_t = 00 \dots 0$.

Laplace suggested the following predictor:

$$L_0(a|x_1 \dots x_t) = (\nu_{x_1 \dots x_t}(a) + 1)/(t + |A|), \quad (1)$$

where $\nu_{x_1 \dots x_t}(a)$ denotes the count of letter a occurring in the word $x_1 \dots x_{t-1} x_t$. For example, if $A = \{0, 1\}$, $x_1 \dots x_5 = 01010$, then the Laplace prediction is as follows: $L_0(x_6 = 0|01010) = (3 + 1)/(5 + 2) = 4/7$, $L_0(x_6 = 1|01010) = (2 + 1)/(5 + 2) = 3/7$. In other words, $3/7$ and $4/7$ are estimates of the unknown probabilities $P(x_{t+1} = 0|x_1 \dots x_t = 01010)$ and $P(x_{t+1} = 1|x_1 \dots x_t = 01010)$. (It is worth noting that the estimated probability to encounter zero after observing a binary string that contains only zeros is not one.)

We can see that Laplace considered prediction as a set of estimations of unknown (conditional) probabilities. This approach to the problem of prediction was developed in [41] and now is often called on-line prediction or universal prediction [1, 20, 28, 31]. As we mentioned above, it seems natural to consider conditional probabilities to be the best prediction, because they contain all information about the future behavior of the stochastic process. Moreover, this approach is deeply connected with game-theoretical interpretation of prediction (see [18, 43]) and, in fact, all obtained results can be easily transferred from one model to the other.

Any predictor γ defines a measure by the following equation

$$\gamma(x_1 \dots x_t) = \prod_{i=1}^t \gamma(x_i|x_1 \dots x_{i-1}). \quad (2)$$

For example, $L_0(0101) = \frac{1}{2} \frac{1}{3} \frac{1}{2} \frac{2}{5} = \frac{1}{30}$. And, vice versa, any measure γ (or estimate of the measure) defines a predictor: $\gamma(x_i|x_1 \dots x_{i-1}) = \gamma(x_1 \dots x_{i-1} x_i) / \gamma(x_1 \dots x_{i-1})$. The same is true for a density (and its estimate): a predictor is defined by conditional density and, vice versa, the density is equal to the product of conditional densities:

$$p(x_i|x_1 \dots x_{i-1}) = \frac{p(x_1 \dots x_{i-1} x_i)}{p(x_1 \dots x_{i-1})}, \quad p(x_1 \dots x_t) = \prod_{i=1}^t p(x_i|x_1 \dots x_{i-1}).$$

The next natural question is how to estimate the precision of a prediction and a probability estimation. Mainly we will estimate the error of prediction by the Kullback-Leibler (KL) divergence between a distribution P and its estimate as follows:

$$\rho_{\gamma, P}(x_1 \dots x_t) = \sum_{a \in A} P(a|x_1 \dots x_t) \log \frac{P(a|x_1 \dots x_t)}{\gamma(a|x_1 \dots x_t)}, \quad (3)$$

where γ is the estimate of an unknown conditional probability. It is well-known that for any distributions P and γ the KL divergence is nonnegative and equals 0 if and only if $P(a) = \gamma(a)$ for all a , see, e.g., [15]. The following inequality (Pinsker's inequality)

$$\sum_{a \in A} P(a) \log \frac{P(a)}{Q(a)} \geq \frac{\log e}{2} \|P - Q\|^2 \quad (4)$$

connects the KL divergence with the so-called variation distance

$$\|P - Q\| = \sum_{a \in A} |P(a) - Q(a)|,$$

where P and Q are distributions over A , see [7]. For fixed t , $\rho_{\gamma,P}(\cdot)$ is a random variable, because x_1, x_2, \dots, x_t are random variables. We define the average error at time t by

$$\rho^t(P\|\gamma) = E(\rho_{\gamma,P}(\cdot)) = \sum_{x_1 \dots x_t \in A^t} P(x_1 \dots x_t) \rho_{\gamma,P}(x_1 \dots x_t). \quad (5)$$

It is shown in [42] that the error of Laplace predictor L_0 goes to 0 for any i.i.d. source P . More precisely, it is proven that

$$\rho^t(P\|L_0) \leq (|A| - 1) \log e / (t + 1) \quad (6)$$

for any source P , [42], see also [46]. So, we can see from this inequality that the average error of the Laplace predictor L_0 (estimated either by the KL divergence or the variation distance) goes to zero for any unknown i.i.d. source, when the sample size t grows. Moreover, it can be easily shown that the error (3) (and the corresponding variation distance) goes to zero with probability 1, when t goes to infinity. Obviously, such a property is very desirable for any predictor and for larger classes of sources, like Markov, stationary and ergodic, etc. However, it is proven in [41] (see also [1]) that such predictors do not exist for the class of all stationary and ergodic sources (generating letters from a given finite alphabet). More precisely, for any predictor γ there exists a source P and $\delta > 0$ such that with probability 1 $\rho_{\gamma,P}(x_1 \dots x_t) \geq \delta$ infinitely often when $t \rightarrow \infty$. So, the error of any predictor may not go to 0, if the predictor is applied to an arbitrary stationary and ergodic source, that is why it is difficult to use (3) and (5) to compare different predictors.

On the other hand, it is shown in [41], that there exists a predictor R , such that the following Cesaro average $t^{-1} \sum_{i=1}^t \rho_{R,P}(x_1 \dots x_t)$ goes to 0 (with probability 1) for any stationary and ergodic source P , where t goes to infinity. That is why we will focus our attention on such averages and by analogy with (5) we define

$$\bar{\rho}_{\gamma,P}(x_1 \dots x_t) = t^{-1} (\log(P(x_1 \dots x_t)/\gamma(x_1 \dots x_t))) \quad (7)$$

and

$$\bar{\rho}_t(\gamma, P) = t^{-1} \sum_{x_1 \dots x_t \in A^t} P(x_1 \dots x_t) \log(P(x_1 \dots x_t)/\gamma(x_1 \dots x_t)), \quad (8)$$

where, as before, $\gamma(x_1 \dots x_t) = \prod_{i=1}^t \gamma(x_i|x_1 \dots x_{i-1})$.

From these definitions and (6) we obtain the following estimation of the error of the Laplace predictor L_0 for any i.i.d. source:

$$\bar{\rho}_t(L_0, P) < ((|A| - 1) \log t + c)/t, \quad (9)$$

where c is a certain constant. So, we can see that the average error of the Laplace predictor goes to zero for any i.i.d. source (which generates letters from a known finite alphabet). As a matter of fact, the Laplace probability $L_0(x_1 \dots x_t)$ is a consistent estimate of the unknown probability $P(x_1 \dots x_t)$.

A natural problem is to find a predictor whose error is minimal (for i.i.d. sources). This problem was considered and solved by Krichevsky in [25], see also [26]. He suggested the following predictor:

$$K_0(a|x_1 \dots x_t) = (\nu_{x_1 \dots x_t}(a) + 1/2)/(t + |A|/2), \quad (10)$$

where, as before, $\nu_{x_1 \dots x_t}(a)$ is the count of letter a occurring in the word $x_1 \dots x_t$. We can see that the Krichevsky predictor is quite close to Laplace's one (1). For example, if $A = \{0, 1\}$, $x_1 \dots x_5 = 01010$, then $K_0(x_6 = 0|01010) = (3+1/2)/(5+1) = 7/12$, $K_0(x_6 = 1|01010) = (2+1/2)/(5+1) = 5/12$ and $K_0(01010) = \frac{1}{2} \frac{1}{4} \frac{1}{2} \frac{3}{8} \frac{1}{2} = \frac{3}{256}$.

The Krichevsky measure K_0 can be presented as follows:

$$K_0(x_1 \dots x_t) = \prod_{i=1}^t \frac{\nu_{x_1 \dots x_{i-1}}(x_i) + 1/2}{i - 1 + |A|/2} = \frac{\prod_{a \in A} (\prod_{j=1}^{\nu_{x_1 \dots x_t}(a)} (j - 1/2))}{\prod_{i=0}^{t-1} (i + |A|/2)}. \quad (11)$$

It is known that

$$(r + 1/2)((r + 1) + 1/2) \dots (s - 1/2) = \frac{\Gamma(s + 1/2)}{\Gamma(r + 1/2)}, \quad (12)$$

where $\Gamma(\cdot)$ is the gamma function (see, e.g., [22] for definition). So, (11) can be presented as follows:

$$K_0(x_1 \dots x_t) = \frac{\prod_{a \in A} (\Gamma(\nu_{x_1 \dots x_t}(a) + 1/2) / \Gamma(1/2))}{\Gamma(t + |A|/2) / \Gamma(|A|/2)}. \quad (13)$$

For this predictor

$$\bar{\rho}_t(K_0, P) < ((|A| - 1) \log t + c) / (2t), \quad (14)$$

where c is a constant, and, moreover, in a certain sense this average error is minimal: for any predictor γ there exists such a source p^* that

$$\bar{\rho}_t(\gamma, p^*) \geq ((|A| - 1) \log t + c) / (2t),$$

see [25, 26].

2.2 Consistent estimations and on-line predictors for Markov and ergodic processes

Now we briefly describe consistent estimations of unknown probabilities and efficient on-line predictors for general stochastic processes (or sources of information). Denote by A^t and A^* the set of all words of length t over A and the set of all finite words over A correspondingly ($A^* = \bigcup_{i=1}^{\infty} A^i$). By $M_{\infty}(A)$ we denote the set of all stationary and ergodic sources which generate letters from A and let $M_0(A) \subset M_{\infty}(A)$ be the set of all i.i.d. processes. Let $M_m(A) \subset M_{\infty}(A)$ be the set of Markov sources of order (or with memory, or connectivity) not larger than m , $m \geq 0$. Let $M^*(A) = \bigcup_{i=0}^{\infty} M_i(A)$ be the set of all finite-memory sources.

The Laplace and Krichevsky predictors can be extended to general Markov processes. The trick is to view a Markov source $P \in M_m(A)$ as resulting from $|A|^m$ i.i.d. sources. We illustrate this idea by an example from [46]. So assume that $A = \{O, I\}$, $m = 2$ and assume that the source $P \in M_2(A)$ has generated the sequence

OOIOIHOOIHIOIO.

We represent this sequence by the following four subsequences:

*** I ** ** I ** ** **,*
*** * O * I ** * I ** * O,*
*** ** I * O ** ** I *,*
*** ** ** O ** * IO **.*

These four subsequences contain letters which follow *OO*, *OI*, *IO* and *II*, respectively. By definition, $P \in M_m(A)$ if $P(a|x_1 \dots x_t) = P(a|x_{t-m+1} \dots x_t)$, for all $0 < m \leq t$, all $a \in A$ and all $x_1 \dots x_t \in A^t$. Therefore, each of the four generated subsequences may be considered as

generated by a Bernoulli source. Further, it is possible to reconstruct the original sequence if we know the four ($= |A|^m$) subsequences and the two ($= m$) first letters of the original sequence.

Any predictor γ for i.i.d. sources can be applied to Markov sources. Indeed, in order to predict, it is enough to store in the memory $|A|^m$ sequences, one corresponding to each word in A^m . Thus, in the example, the letter x_3 which follows OO is predicted based on the Bernoulli method γ corresponding to the x_1x_2 -subsequence ($= OO$), then x_4 is predicted based on the Bernoulli method corresponding to x_2x_3 , i.e. to the OI -subsequence, and so forth. When this scheme is applied along with either L_0 or K_0 we denote the obtained predictors as L_m and K_m , correspondingly, and define the probabilities for the first m letters as follows: $L_m(x_1) = L_m(x_2) = \dots = L_m(x_m) = 1/|A|$, $K_m(x_1) = K_m(x_2) = \dots = K_m(x_m) = 1/|A|$. For example, having taken into account (13), we can present the Krichevsky predictors for $M_m(A)$ as follows:

$$K_m(x_1 \dots x_t) = \begin{cases} \frac{1}{|A|^t}, & \text{if } t \leq m, \\ \frac{1}{|A|^m} \prod_{v \in A^m} \frac{\prod_{a \in A} ((\Gamma(\nu_x(va)+1/2) / \Gamma(1/2))}{(\Gamma(\bar{\nu}_x(v)+|A|/2) / \Gamma(|A|/2))}, & \text{if } t > m, \end{cases} \quad (15)$$

where $\bar{\nu}_x(v) = \sum_{a \in A} \nu_x(va)$, $x = x_1 \dots x_t$; see [25] and references therein. It is worth noting that representation (12) can be more convenient for carrying out calculations. Let us consider an example. For the word $OOIOIIOOIIIIOIO$ considered in the previous example, we obtain $K_2(OOIOIIOOIIIIOIO) = 2^{-2} \frac{1}{2} \frac{3}{4} \frac{1}{2} \frac{1}{4} \frac{3}{8} \frac{1}{2} \frac{1}{4} \frac{1}{2} \frac{1}{4} \frac{1}{2}$.

Let us define the measure R , which is a consistent estimator of probabilities for the class of all stationary and ergodic processes with a finite alphabet. First we define a probability distribution $\{\omega = \omega_1, \omega_2, \dots\}$ on integers $\{1, 2, \dots\}$ by

$$\omega_1 = 1 - 1/\log 3, \dots, \omega_i = 1/\log(i+1) - 1/\log(i+2), \dots \quad (16)$$

(In what follows we will use this distribution, but results described below are obviously true for any distribution with nonzero probabilities.) The measure R is defined as follows:

$$R(x_1 \dots x_t) = \sum_{i=0}^{\infty} \omega_{i+1} K_i(x_1 \dots x_t). \quad (17)$$

It is worth noting that this construction can be applied to the Laplace measure (if we use L_i instead of K_i) and any other family of measures.

The main properties of the measure R are connected with the Shannon entropy, which is defined as follows

$$H(P) = \lim_{m \rightarrow \infty} -\frac{1}{m} \sum_{v \in A^m} P(v) \log P(v). \quad (18)$$

Theorem 1 ([41]). *For any stationary and ergodic source P the following equalities are valid:*

$$i) \lim_{t \rightarrow \infty} \frac{1}{t} \log(1/R(x_1 \dots x_t)) = H(P)$$

with probability 1,

$$ii) \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{u \in A^t} P(u) \log(1/R(u)) = H(P).$$

2.3 Nonparametric estimations and data compression

One of the goals of the paper is to show how practically used data compressors can be employed as a tool for nonparametric estimation, prediction and other problems. That is why a short description of universal data compressors (or universal codes) will be given here.

A data compression method (or code) φ is defined as a set of mappings φ_n such that $\varphi_n : A^n \rightarrow \{0, 1\}^*$, $n = 1, 2, \dots$ and for each pair of different words $x, y \in A^n$ $\varphi_n(x) \neq \varphi_n(y)$. It is also required that each sequence $\varphi_n(u_1)\varphi_n(u_2)\dots\varphi_n(u_r)$, $r \geq 1$, of encoded words from the set A^n , $n \geq 1$, could be uniquely decoded into $u_1u_2\dots u_r$. Such codes are called uniquely decodable. For example, let $A = \{a, b\}$, the code $\psi_1(a) = 0, \psi_1(b) = 00$, obviously, is not uniquely decodable. It is well known that if a code φ is uniquely decodable then the lengths of the codewords satisfy the following inequality (Kraft's inequality): $\sum_{u \in A^n} 2^{-|\varphi_n(u)|} \leq 1$, see, e.g., [15]. It will be convenient to reformulate this property as follows:

Claim 1. *Let φ be a uniquely decodable code over an alphabet A . Then for any integer n there exists a measure μ_φ on A^n such that*

$$-\log \mu_\varphi(u) \leq |\varphi(u)| \quad (19)$$

for any u from A^n .

(Obviously, Claim 1 is true for the measure $\mu_\varphi(u) = 2^{-|\varphi(u)|} / \sum_{u \in A^n} 2^{-|\varphi(u)|}$). In what follows we call uniquely decodable codes just "codes".

It is worth noting that, in fact, any measure μ defines a code for which the length of the codeword associated with a word u is (close to) $-\log \mu(u)$.

Now we consider universal codes. By definition, a code U is universal if for any stationary and ergodic source P the following equalities are valid:

$$\lim_{t \rightarrow \infty} |U(x_1 \dots x_t)|/t = H(P) \quad (20)$$

with probability 1, and

$$\lim_{t \rightarrow \infty} E(|U(x_1 \dots x_t)|)/t = H(P), \quad (21)$$

where $H(P)$ is the Shannon entropy of P , $E(f)$ is a mean value of f . In fact, (21) and (20) are valid for known universal codes, but there exist codes for which only one equality is valid.

3 Finite-alphabet processes

3.1 The estimation of (limiting) probabilities

The following theorem shows how universal codes can be applied for probability estimations.

Theorem 2. *Let U be a universal code and*

$$\mu_U(u) = 2^{-|U(u)|} / \sum_{v \in A^{|u|}} 2^{-|U(v)|}. \quad (22)$$

Then, for any stationary and ergodic source P the following equalities are valid:

$$i) \lim_{t \rightarrow \infty} \frac{1}{t} (-\log P(x_1 \dots x_t) - (-\log \mu_U(x_1 \dots x_t))) = 0$$

with probability 1,

$$ii) \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{u \in A^t} P(u) \log(P(u)/\mu_U(u)) = 0,$$

Proof. The proof is based on Shannon-MacMillan-Breiman Theorem which states that for any stationary and ergodic source P

$$\lim_{t \rightarrow \infty} -\log P(x_1 \dots x_t)/t = H(P)$$

with probability 1, see [3, 15]. From this equality and (20) we obtain the statement i). The second statement follows from the definition of Shannon entropy (18) and (21). \square

So, we can see that, in a certain sense, the measure μ_U is a consistent (nonparametric) estimate of the (unknown) measure P .

Nowadays there are many efficient universal codes (and universal predictors connected with them), see [11, 21, 26, 34, 35, 40, 48], which can be applied to estimation. For example, the above described measure R is based on the code from [40, 41] and can be applied for probability estimation. More precisely, Theorem 2 (and the following theorems) are true for R , if we replace μ_U by R .

It is important to note that the measure R has some additional properties, which can be useful for applications. The following theorem describes these properties (whereas all other theorems are valid for all universal codes and corresponding measures, including the measure R).

Theorem 3. *For any Markov process P with memory k*

i) the error of the probability estimator, which is based on the measure R , is upper-bounded as follows:

$$\frac{1}{t} \sum_{u \in A^t} P(u) \log(P(u)/R(u)) \leq \frac{(|A| - 1)|A|^k \log t}{2t} + O\left(\frac{1}{t}\right),$$

ii) in a certain sense the error of R is asymptotically minimal: for any measure μ there exists a k -memory Markov process p_μ such that

$$\frac{1}{t} \sum_{u \in A^t} p_\mu(u) \log(p_\mu(u)/\mu(u)) \geq \frac{(|A| - 1)|A|^k \log t}{2t} + O\left(\frac{1}{t}\right),$$

iii) Let Θ be a set of stationary and ergodic processes such that there exists a measure μ_Θ for which the estimation error of the probability goes to 0 uniformly:

$$\lim_{t \rightarrow \infty} \sup_{P \in \Theta} \left(\frac{1}{t} \sum_{u \in A^t} P(u) \log(P(u)/\mu_\Theta(u)) \right) = 0.$$

Then the error of estimator, which is based on the measure R , goes to 0 uniformly too:

$$\lim_{t \rightarrow \infty} \sup_{P \in \Theta} \left(\frac{1}{t} \sum_{u \in A^t} P(u) \log(P(u)/R(u)) \right) = 0.$$

The proof can be found in [40, 41].

3.2 Prediction

As we mentioned above, any universal code U can be applied for prediction. Namely, the measure μ_U (22) can be used for prediction as the following conditional probability:

$$\mu_U(x_{t+1}|x_1\dots x_t) = \mu_U(x_1\dots x_t x_{t+1})/\mu_U(x_1\dots x_t). \quad (23)$$

Theorem 4. *Let U be a universal code and P be any stationary and ergodic process. Then*

$$i) \lim_{t \rightarrow \infty} \frac{1}{t} \{E(\log \frac{P(x_1)}{\mu_U(x_1)}) + E(\log \frac{P(x_2|x_1)}{\mu_U(x_2|x_1)}) + \dots + E(\log \frac{P(x_t|x_1\dots x_{t-1})}{\mu_U(x_t|x_1\dots x_{t-1})})\} = 0,$$

$$ii) \lim_{t \rightarrow \infty} E(\frac{1}{t} \sum_{i=0}^{t-1} (P(x_{i+1}|x_1\dots x_i) - \mu_U(x_{i+1}|x_1\dots x_i))^2) = 0,$$

and

$$iii) \lim_{t \rightarrow \infty} E(\frac{1}{t} \sum_{i=0}^{t-1} |P(x_{i+1}|x_1\dots x_i) - \mu_U(x_{i+1}|x_1\dots x_i)|) = 0.$$

Proof. i) immediately follows from the second statement of the previous theorem and properties of log. The statement ii) can be proven as follows:

$$\begin{aligned} & \lim_{t \rightarrow \infty} E(\frac{1}{t} \sum_{i=0}^{t-1} (P(x_{i+1}|x_1\dots x_i) - \mu_U(x_{i+1}|x_1\dots x_i))^2) = \\ & \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{i=0}^{t-1} \sum_{x_1\dots x_i \in A^i} P(x_1\dots x_i) (\sum_{a \in A} |P(a|x_1\dots x_i) - \mu_U(a|x_1\dots x_i)|)^2 \leq \\ & \lim_{t \rightarrow \infty} \frac{const}{t} \sum_{i=0}^{t-1} \sum_{x_1\dots x_i \in A^i} P(x_1\dots x_i) \sum_{a \in A} P(a|x_1\dots x_i) \log \frac{P(a|x_1\dots x_i)}{\mu_U(a|x_1\dots x_i)} = \\ & \lim_{t \rightarrow \infty} (\frac{const}{t} \sum_{x_1\dots x_t \in A^t} P(x_1\dots x_t) \log(P(x_1\dots x_t)/\mu(x_1\dots x_t))). \end{aligned}$$

Here the first inequality is obvious, the second follows from the Pinsker's inequality (4), the others from properties of expectation and log. iii) can be derived from ii) and the Jensen inequality for the function x^2 . \square

Comment 1. The measure R described above has one additional property if it is used for prediction. Namely, for any Markov process P ($P \in M^*(A)$) the following is true:

$$\lim_{t \rightarrow \infty} \log \frac{P(x_{t+1}|x_1\dots x_t)}{R(x_{t+1}|x_1\dots x_t)} = 0$$

with probability 1, where $R(x_{t+1}|x_1\dots x_t) = R(x_1\dots x_t x_{t+1})/R(x_1\dots x_t)$; see [42].

Comment 2. In fact, the statements ii) and iii) are equivalent, because one of them follows from the other. For details see Lemma 2 in [47].

3.3 Problems with side information

Now we consider so-called problems with side information, which are described as follows: there is a stationary and ergodic source, whose alphabet A is presented as a product $A = X \times Y$. We are given a sequence $(x_1, y_1), \dots, (x_{t-1}, y_{t-1})$ and so-called side information y_t . The goal is to predict, or estimate, x_t . This problem arises in statistical decision theory, pattern recognition, and machine learning, see [29]. Obviously, if someone knows the conditional probabilities $P(x_t | (x_1, y_1), \dots, (x_{t-1}, y_{t-1}), y_t)$ for all $x_t \in X$, he has all information about x_t , available before x_t is known. That is why we will look for the best (or, at least, good) estimators for this conditional probabilities. Our solution will be based on results obtained in the parts 3.1 and 3.2. More precisely, for any universal code U and the corresponding measure μ_U (22) we define the following estimate for the problem with side information:

$$\mu_U(x_t | (x_1, y_1), \dots, (x_{t-1}, y_{t-1}), y_t) = \frac{\mu_U((x_1, y_1), \dots, (x_{t-1}, y_{t-1}), (x_t, y_t))}{\sum_{x_t \in X} \mu_U((x_1, y_1), \dots, (x_{t-1}, y_{t-1}), (x_t, y_t))}.$$

Theorem 5. *Let U be a universal code and P any stationary and ergodic process. Then*

$$\begin{aligned} \text{i) } \lim_{t \rightarrow \infty} \frac{1}{t} \{ & E(\log \frac{P(x_1 | y_1)}{\mu_U(x_1 | y_1)}) + E(\log \frac{P(x_2 | (x_1, y_1), y_2)}{\mu_U(x_2 | (x_1, y_1), y_2)}) + \dots \\ & + E(\log \frac{P(x_t | (x_1, y_1), \dots, (x_{t-1}, y_{t-1}), y_t)}{\mu_U(x_t | (x_1, y_1), \dots, (x_{t-1}, y_{t-1}), y_t)}) \} = 0, \\ \text{ii) } \lim_{t \rightarrow \infty} E(\frac{1}{t} \sum_{i=0}^{t-1} & (P(x_{i+1} | (x_1, y_1), \dots, (x_i, y_i), y_{i+1})) - \\ & \mu_U(x_{i+1} | (x_1, y_1), \dots, (x_i, y_i), y_{i+1}))^2) = 0, \end{aligned}$$

and

$$\begin{aligned} \text{iii) } \lim_{t \rightarrow \infty} E(\frac{1}{t} \sum_{i=0}^{t-1} & |P(x_{i+1} | (x_1, y_1), \dots, (x_i, y_i), y_{i+1})) - \\ & \mu_U(x_{i+1} | (x_1, y_1), \dots, (x_i, y_i), y_{i+1})|) = 0. \end{aligned}$$

Proof. The following inequality follows from the nonnegativity of the KL divergency (see (4)), whereas equality is obvious.

$$\begin{aligned} E(\log \frac{P(x_1 | y_1)}{\mu_U(x_1 | y_1)}) + E(\log \frac{P(x_2 | (x_1, y_1), y_2)}{\mu_U(x_2 | (x_1, y_1), y_2)}) + \dots & \leq E(\log \frac{P(y_1)}{\mu_U(y_1)}) \\ + E(\log \frac{P(x_1 | y_1)}{\mu_U(x_1 | y_1)}) + E(\log \frac{P(y_2 | (x_1, y_1)}{\mu_U(y_2 | (x_1, y_1))}) + E(\log \frac{P(x_2 | (x_1, y_1), y_2)}{\mu_U(x_2 | (x_1, y_1), y_2)}) + \dots & \\ = E(\log \frac{P(x_1, y_1)}{\mu_U(x_1, y_1)}) + E(\log \frac{P((x_2, y_2) | (x_1, y_1))}{\mu_U((x_2, y_2) | (x_1, y_1))}) + \dots & \end{aligned}$$

Now we can apply the first statement of Theorem 4 to the last sum as follows:

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{1}{t} E(\log \frac{P(x_1, y_1)}{\mu_U(x_1, y_1)}) + E(\log \frac{P((x_2, y_2) | (x_1, y_1))}{\mu_U((x_2, y_2) | (x_1, y_1))}) + \dots & \\ E(\log \frac{P((x_t, y_t) | (x_1, y_1) \dots (x_{t-1}, y_{t-1}))}{\mu_U((x_t, y_t) | (x_1, y_1) \dots (x_{t-1}, y_{t-1}))}) = 0. & \end{aligned}$$

From this equality and last inequality we obtain the proof of i). The proof of the second statement can be obtained from the similar representation for ii) and the second statement of Theorem 4. iii) can be derived from ii) and the Jensen inequality for the function x^2 . \square

3.4 The case of several independent samples

Now we extend our consideration to the case where the sample is presented as several independent samples $x^1 = x_1^1 \dots x_{t_1}^1$, $x^2 = x_1^2 \dots x_{t_2}^2, \dots, x^r = x_1^r \dots x_{t_r}^r$ generated by a source. More precisely, we will suppose that all sequences were independently created by one stationary and ergodic source. (As it was mentioned above, it is impossible just to combine all samples into one, if the source is not i.i.d.) We denote this sample by $x^1 \diamond x^2 \diamond \dots \diamond x^r$ and define $\nu_{x^1 \diamond x^2 \diamond \dots \diamond x^r}(v) = \sum_{i=1}^r \nu_{x^i}(v)$. For example, if $x^1 = 0010$, $x^2 = 011$, then $\nu_{x^1 \diamond x^2}(00) = 1$. The definition of K_m and R can be extended to this case:

$$K_m(x^1 \diamond x^2 \diamond \dots \diamond x^r) = \left(\prod_{i=1}^r |A|^{-\min\{m, t_i\}} \right) \prod_{v \in A^m} \frac{\prod_{a \in A} ((\Gamma(\nu_{x^1 \diamond x^2 \diamond \dots \diamond x^r}(va) + 1/2) / \Gamma(1/2))}{(\Gamma(\bar{\nu}_{x^1 \diamond x^2 \diamond \dots \diamond x^r}(v) + |A|/2) / \Gamma(|A|/2))}, \quad (24)$$

whereas the definition of R is the same (see (17)). (Here, as before, $\bar{\nu}_{x^1 \diamond x^2 \diamond \dots \diamond x^r}(v) = \sum_{a \in A} \nu_{x^1 \diamond x^2 \diamond \dots \diamond x^r}(va)$. Note, that $\bar{\nu}_{x^1 \diamond x^2 \diamond \dots \diamond x^r}() = \sum_{i=1}^r t_i$ if $m = 0$.)

The following example is intended to show the difference between the case of many samples and one. Let there be two independent samples $y = y_1 \dots y_4 = 0101$ and $x = x_1 \dots x_3 = 101$, generated by a stationary and ergodic source with the alphabet $\{0, 1\}$. One wants to estimate the (limiting) probabilities $P(z_1 z_2)$, $z_1, z_2 \in \{0, 1\}$ (here $z_1 z_2 \dots$ can be considered as an independent sequence, generated by the source) and predict $x_4 x_5$ (i.e. estimate conditional probability $P(x_4 x_5 | x_1 \dots x_3 = 101, y_1 \dots y_4 = 0101)$). For solving both problems we will use the measure R (see (17)). First we consider the case where $P(z_1 z_2)$ is to be estimated without knowledge of sequences x and y . From (11) and (15) we obtain:

$$K_0(00) = K_0(11) = \frac{1/2}{1} \frac{3/2}{1+1} = 3/8, \quad K_0(01) = K_0(10) = \frac{1/2}{1+0} \frac{1/2}{1+1} = 1/8,$$

$$K_i(00) = K_i(01) = K_i(10) = K_i(11) = 1/4; \quad , \quad i \geq 1.$$

Having taken into account the definitions of ω_i (16) and the measure R (17), we can calculate $R(z_1 z_2)$ as follows:

$$R(00) = \omega_1 K_0(00) + \omega_2 K_1(00) + \dots = (1 - 1/\log 3) 3/8 + (1/\log 3 - 1/\log 4) 1/4 + (1/\log 4 - 1/\log 5) 1/4 + \dots = (1 - 1/\log 3) 3/8 + (1/\log 3) 1/4 \approx 0.296.$$

Analogously, $R(01) = R(10) \approx 0.204$, $R(11) \approx 0.296$.

Let us now estimate the probability $P(z_1 z_2)$ taking into account that there are two independent samples $y = y_1 \dots y_4 = 0101$ and $x = x_1 \dots x_3 = 101$. First of all we note that such estimates are based on the formula for conditional probabilities:

$$R(z|x \diamond y) = R(x \diamond y \diamond z) / R(x \diamond y).$$

First we estimate the frequencies :

$$\nu_{0101 \diamond 101}(0) = 3, \nu_{0101 \diamond 101}(1) = 4, \nu_{0101 \diamond 101}(00) = \nu_{0101 \diamond 101}(11) = 0, \nu_{0101 \diamond 101}(01) = 3,$$

$$\nu_{0101 \diamond 101}(10) = 2, \nu_{0101 \diamond 101}(010) = 1, \nu_{0101 \diamond 101}(101) = 2, \nu_{0101 \diamond 101}(0101) = 1,$$

whereas frequencies of all other three-letters and four-letters words are 0. Then we calculate :

$$K_0(0101 \diamond 101) = \frac{1}{2} \frac{3}{4} \frac{5}{6} \frac{1}{8} \frac{3}{10} \frac{5}{12} \frac{1}{14} \approx 0.00244, \quad K_1(0101 \diamond 101) = (2^{-1})^2 \frac{1}{2} \frac{3}{4} \frac{5}{6} \frac{1}{8} \frac{3}{10} \frac{5}{12} \frac{1}{14} \approx 0.00061$$

$$\begin{aligned} &\approx 0.0293, \quad K_2(0101 \diamond 101) \approx 0.01172, \quad K_i(0101 \diamond 101) = 2^{-7}, \quad i \geq 3, \\ &R(0101 \diamond 101) = \omega_1 K_0(0101 \diamond 101) + \omega_2 K_1(0101 \diamond 101) + \dots \approx \\ &0.369 \ 0.00244 + 0.131 \ 0.0293 + 0.06932 \ 0.01172 + 2^{-7} / \log 5 \approx 0.0089. \end{aligned}$$

In order to avoid repetitions, we estimate only one probability $P(z_1 z_2 = 01)$. Carrying out similar calculations, we obtain

$$\begin{aligned} R(0101 \diamond 101 \diamond 01) &\approx 0.00292, \\ R(z_1 z_2 = 01 | y_1 \dots y_4 = 0101, x_1 \dots x_3 = 101) &= \\ R(0101 \diamond 101 \diamond 01) / R(0101 \diamond 101) &\approx 0.32812. \end{aligned}$$

If we compare this value and the estimation $R(01) \approx 0.204$, which is not based on the knowledge of samples x and y , we can see that the measure R uses additional information quite naturally (indeed, 01 is quite frequent in $y = y_1 \dots y_4 = 0101$ and $x = x_1 \dots x_3 = 101$).

Such generalization can be applied for many universal codes, but, generally speaking, there exist codes U for which $U(x^1 \diamond x^2)$ is not defined and, hence, the measure $\mu_U(x_1 \diamond x_2)$ is not defined. That is why we will describe properties of R , but do not describe properties of universal codes in general. For the measure R all asymptotic properties are the same for the cases of one sample and several samples. More precisely, the following statement is true:

Claim 2. *Let x^1, x^2, \dots, x^r be independent sequences generated by a stationary and ergodic source and t be a total length of those sequences ($t = \sum_{i=1}^r |x^i|$). Then, if $t \rightarrow \infty$, (and r is fixed) the statements of Theorems 1–5 are valid, when applied to $x^1 \diamond x^2 \diamond \dots \diamond x^r$ instead of $x_1 \dots x_t$. (In Theorems 2, 4, 5 μ_U should be changed to R .)*

The proofs are analogous to the proofs of Theorems 1–5.

4 Real-valued time series

Here we will consider problems of the density estimation and prediction for a stationary process with densities. We have seen that Shannon-MacMillan-Breiman theorem played a key role in the case of finite-alphabet processes. In this part we will use its generalization to the processes with densities. This result was proved by Barron [2] and was an extension of the $L1$ convergence obtained in [19, 30, 32]. First we describe considered processes with some properties needed for a fulfilment of the generalized Shannon-MacMillan-Breiman theorem.

Let (Ω, F, P) be a probability space and let X_1, X_2, \dots be a stochastic process with each X_t taking values in a standard Borel space. As in [2], suppose that the joint distribution P_n for (X_1, X_2, \dots, X_n) has a probability density function $p(x_1 x_2 \dots x_n)$ with respect to a sigma-finite measure M_n . Assume that the sequence of dominating measures M_n is Markovian of order $m \geq 0$ with a stationary transition measure. Familiar cases for M_n are Lebesgue measure and counting measure. Let $p(x_{n+1} | x_1 \dots x_n)$ denote the conditional density given by the ratio $p(x_1 \dots x_{n+1}) / p(x_1 \dots x_n)$ for $n > 1$. It is known that for stationary and ergodic processes there exists a so-called relative entropy rate h defined by

$$h = \lim_{n \rightarrow \infty} -E(\log p(x_{n+1} | x_1 \dots x_n)), \quad (25)$$

where E denotes expectation with respect to P . The following generalization of the Shannon-MacMillan-Breiman theorem is obtained by Barron in [2]:

Claim 3. If $\{X_n\}$ is a stationary ergodic process with density $p(x_1 \dots x_n) = dP_n/dM_n$ and $h_n < 1$ for some $n \geq m$, the sequence of relative entropy densities

$$-(1/n) \log p(x_1 \dots x_n)$$

converges almost surely to the relative entropy rate, i.e.,

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log p(x_1 \dots x_t) = h. \quad (26)$$

with probability 1 (according to P).

Now we return to the estimation problems. Let $\{\Pi_n\}$, $n \geq 1$, be an increasing sequence of finite partitions of Ω that asymptotically generates the Borel sigma-field on F , and let $x^{[k]}$ denote the element of Π_k that contains the point x . (Informally, if Ω is an interval, $x^{[k]}$ is obtained by quantizing x to k bits of precision). For integers s and n we define the following approximation of the density

$$p^s(x_1, \dots, x_n) = P(x_1^{[s]}, \dots, x_n^{[s]})/M_n(x_1^{[s]} \dots x_n^{[s]}). \quad (27)$$

We also consider

$$h_s = \lim_{n \rightarrow \infty} E(\log p^s(x_{n+1}|x_1, \dots, x_n)). \quad (28)$$

Applying Claim 3 to the density $p^s(x_1, \dots, x_t)$, we obtain that a.s.

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log p^s(x_1, \dots, x_t) = h_s. \quad (29)$$

Let U be a universal code, which is defined for any finite alphabet. In order to describe the density estimate we will use the distribution ω ; see (16). Now we define the corresponding density r_U as follows:

$$r_U(x_1 \dots x_t) = \sum_{i=0}^{\infty} \omega_i \mu_U(x_1^{[i]} \dots x_t^{[i]}) / M_t(x_1^{[i]} \dots x_t^{[i]}), \quad (30)$$

where the measure μ_U is defined by (22). (It is supposed here that the code $U(x_1^{[i]} \dots x_t^{[i]})$ is defined for the alphabet, which contains $|\Pi_i|$ letters.)

It turns out that, in a certain sense, the density $r_U(x_1 \dots x_t)$ estimates the unknown density $p(x_1, \dots, x_t)$.

Theorem 6. Let X_t be a stationary ergodic process with densities $p(x_1 \dots x_t) = dP_t/dM_t$ such that

$$\lim_{s \rightarrow \infty} h_s = h < \infty, \quad (31)$$

where h and h_s are relative entropy rates, see (25), (28). Then

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{p(x_1 \dots x_t)}{r_U(x_1 \dots x_t)} = 0 \quad (32)$$

with probability 1 and

$$\lim_{t \rightarrow \infty} \frac{1}{t} E(\log \frac{p(x_1 \dots x_t)}{r_U(x_1 \dots x_t)}) = 0. \quad (33)$$

Proof. First we note that for any integer s the following obvious equality is true: $r_U(x_1 \dots x_t) = \omega_s \mu_U(x_1^{[s]} \dots x_t^{[s]}) / M_t(x_1^{[s]} \dots x_t^{[s]}) (1 + \delta)$ for some $\delta > 0$. From this equality, (22) and (32) we immediately obtain that a.s.

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{p(x_1 \dots x_t)}{r_U(x_1 \dots x_t)} \leq \lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{p(x_1 \dots x_t)}{2^{-|U(x_1^{[s]} \dots x_t^{[s]})|} / M_t(x_1^{[s]} \dots x_t^{[s]})}. \quad (34)$$

The right part can be presented as follows:

$$\begin{aligned} & \lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{p(x_1 \dots x_t)}{2^{-|U(x_1^{[s]} \dots x_t^{[s]})|} / M_t(x_1^{[s]} \dots x_t^{[s]})} \\ &= \lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{p^s(x_1 \dots x_t) M_t(x_1^{[s]} \dots x_t^{[s]})}{2^{-|U(x_1^{[s]} \dots x_t^{[s]})|}} + \lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{p(x_1 \dots x_t)}{p^s(x_1 \dots x_t)}. \end{aligned} \quad (35)$$

Having taken into account that U is the universal code, (27) and Theorem 1, we can see that the first term equals to zero. From (26) and (29) we can see that a.s. the second term is equal to $h_s - h$. This equality is valid for any integer s and, according to (31), the second term equals to zero too, and we obtain (33). The first statement is proven.

From (34) and (35) we can see that

$$\begin{aligned} E \log \frac{p(x_1 \dots x_t)}{r_U(x_1 \dots x_t)} &\leq E \log \frac{p_t^s(x_1, \dots, x_t) M_t(x_1^{[s]} \dots x_t^{[s]})}{2^{-|U(x_1^{[s]} \dots x_t^{[s]})|}} \\ &+ E \log \frac{p(x_1 \dots x_t)}{p^s(x_1, \dots, x_t)}. \end{aligned} \quad (36)$$

The first term is the average redundancy of the universal code for a finite- alphabet source, hence, according to Theorem 1, it tends to 0. The second term tends to $h_s - h$ for any s and from (31) we can see that it is equal to zero. The second statement is proven. \square

We have seen that the requirement (31) plays an important role in the proof. A natural question is whether there exist processes for which (31) is valid. The answer is positive. For example, let Ω be an interval $[-1, 1]$, M_n be Lebesgue measure and a considered process is Markovian with conditional density

$$p(x|y) = \begin{cases} 1/2 + \alpha \operatorname{sign}(y), & \text{if } x < 0 \\ 1/2 - \alpha \operatorname{sign}(y), & \text{if } x \geq 0, \end{cases}$$

where $\alpha \in (0, 1)$ is a parameter and

$$\operatorname{sign}(y) = \begin{cases} -1, & \text{if } y < 0, \\ 1, & \text{if } y \geq 0. \end{cases}$$

It is easy to see that (31) is true for any $\alpha \in (0, 1/2)$.

The following theorem describes properties of conditional probabilities $r_U(x|x_1 \dots x_m) = r_U(x_1 \dots x_m x) / r_U(x_1 \dots x_m)$ which, in turn, is connected with the prediction problem. We will see that the conditional density $r_U(x|x_1 \dots x_m)$ is a reasonable estimation of $p(x|x_1 \dots x_m)$.

Theorem 7. *Let f be an integrable function whose absolute value is bounded by a certain constant \bar{M} . Then the following equalities are valid:*

$$i) \lim_{t \rightarrow \infty} \frac{1}{t} E \left(\sum_{m=0}^{t-1} \left(\int f(x) p(x|x_1 \dots x_m) dM_m - \int f(x) r_U(x|x_1 \dots x_m) dM_m \right)^2 \right) = 0, \quad (37)$$

$$ii) \lim_{t \rightarrow \infty} \frac{1}{t} E \left(\sum_{m=0}^{t-1} \left| \int f(x) p(x|x_1 \dots x_m) dM_m - \int f(x) r_U(x|x_1 \dots x_m) dM_m \right| \right) = 0.$$

Proof. The last inequality of the following chain follows from the Pinsker's one, whereas all others are obvious.

$$\begin{aligned} & \left(\int f(x) p(x|x_1 \dots x_m) dM_m - \int f(x) r_U(x|x_1 \dots x_m) dM_m \right)^2 = \\ & \quad \left(\int f(x) (p(x|x_1 \dots x_m) - r_U(x|x_1 \dots x_m)) dM_m \right)^2 \\ & \leq \bar{M}^2 \left(\int (p(x|x_1 \dots x_m) - r_U(x|x_1 \dots x_m)) dM_m \right)^2 \\ & \leq \bar{M}^2 \left(\int |p(x|x_1 \dots x_m) - r_U(x|x_1 \dots x_m)| dM_m \right)^2 \leq \\ & \quad \text{const} \int p(x|x_1 \dots x_m) \log(p(x|x_1 \dots x_m)/r_U(x|x_1 \dots x_m)) dM_m. \end{aligned}$$

From these inequalities we obtain:

$$\begin{aligned} & \sum_{m=0}^{t-1} E \left(\int f(x) p(x|x_1 \dots x_m) dM_m - \int f(x) r_U(x|x_1 \dots x_m) dM_m \right)^2 \leq \quad (38) \\ & \sum_{m=0}^{t-1} \text{const} E \left(\int p(x|x_1 \dots x_m) \log(p(x|x_1 \dots x_m)/r_U(x|x_1 \dots x_m)) dM_m \right). \end{aligned}$$

The last term can be presented as follows:

$$\begin{aligned} & \sum_{m=0}^{t-1} E \left(\int p(x|x_1 \dots x_m) \log(p(x|x_1 \dots x_m)/r_U(x|x_1 \dots x_m)) dM_m \right) = \\ & \sum_{m=0}^{t-1} \int p(x_1 \dots x_m) \int p(x|x_1 \dots x_m) \log(p(x|x_1 \dots x_m)/r_U(x|x_1 \dots x_m)) dM_1 dM_m = \\ & \quad \int p(x_1 \dots x_t) \log(p(x_1 \dots x_t)/r_U(x_1 \dots x_t)) dM_t. \end{aligned}$$

From this equality, (38) and (33) we obtain (37). ii) can be derived from (38) and the Jensen inequality for x^2 . \square

References

- [1] P.Algoet, "Universal Schemes for Learning the Best Nonlinear Predictor Given the Infinite Past and Side Information", *IEEE Trans. Inform. Theory*, vol. 45, pp. 1165–1185, 1999.
- [2] A.R. Barron, "The strong ergodic theorem for dencities: generalized Shannon-McMillan-Breiman theorem", *The annals of Probability*, vol. 13, pp.1292–1303, 1985.
- [3] P. Billingsley, *Ergodic theory and information*, John Wiley & Sons, 1965.
- [4] R. Cilibrasi and P.M.B. Vitanyi, "Clustering by Compression," *IEEE Transactions on Information Theory*, vol. 51, 2005.

- [5] R. Cilibrasi, R. de Wolf and P.M.B. Vitanyi, "Algorithmic Clustering of Music," *Computer Music Journal*, vol. 28, pp.49–67, 2004.
- [6] T.M. Cover and J.A. Thomas, *Elements of Information Theory*, John Wiley and sons, 1991.
- [7] I. Csiszár, J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*, Budapesht, Akadémiai Kiadó, 1981.
- [8] I. Csiszár, P. Shields, "The consistency of the BIC Markov order estimation", *Annals of Statistics*, vol. 6, pp. 1601–1619, 2000.
- [9] G.A. Darbellay, I. Vajda, "Entropy expressions for multivariate continuous distributions," in. Research Report no 1920, UTIA, Academy of Science, Prague (library@utia.cas.cz), 1998.
- [10] G.A. Darbellay, I. Vajda, "Estimation of the mutual information with data-dependent partitions," *IEEE Trans. Inform. Theory*, vol. 48, pp. 1061–1081, 1999.
- [11] M. Effros, K. Visweswariah, S.R. Kulkarni and S. Verdu, "Universal lossless source coding with the Burrows Wheeler transform," *IEEE Trans. Inform. Theory*, vol. 45, pp. 1315–1321, 1999.
- [12] W. Feller, *An Introduction to Probability Theory and Its Applications, vol.1*, John Wiley & Sons, New York, 1970.
- [13] L. Finesso, Chuang-Chun Liu, P. Narayan, "The optimal error exponent for Markov order estimation," *IEEE Trans. on Information Theory*, vol. 42, pp. 1488–1497, 1996.
- [14] B.M. Fitingof, "Optimal encoding for unknown and changing statistica of messages," *Problems of Information Transmission*, vol.2, n. 2, pp. 3–11, 1966.
- [15] R.G. Gallager, *Information Theory and Reliable Communication*, John Wiley & Sons, New York, 1968.
- [16] L. Gyorfí, G. Morvai and S.J. Yakowitz, "Limits to consistent on-line forecasting for ergodic time series," *IEEE Transactions on nformation Theory*, vol. 44, pp.886– 892, 1998.
- [17] P. Jacquet, W. Szpankowski and L. Apostol, "Universal predictor based on pattern matching," *IEEE Trans. Inform. Theory*, vol.48, pp. 1462–1472, 2002.
- [18] J.L. Kelly, "A new interpretation of information rate," *Bell System Tech. J.*, vol. 35, pp. 917–926, 1956.
- [19] J. Kieffer, "A simple proof of the Moy-Perez generalization of the Shannon- MacMillan theorem," *Pacific J. Math.*, vol.51, pp. 203–206, 1974.
- [20] J. Kieffer, *Prediction and Information Theory*, Preprint, 1998. (available at <ftp://oz.ee.umn.edu/users/kieffer/papers/prediction.pdf/>)
- [21] J.C. Kieffer and En-Hui Yang, "Grammar-based codes: a new class of universal lossless source codes," *IEEE Transactions on Information Theory*, vol.46, pp.737 – 754, 2000.
- [22] D.E. Knuth, *The art of computer programming, Vol.2*. Addison Wesley, 1981.
- [23] A.N. Kolmogorov, "Three approaches to the quantitative definition of information," *Problems of Inform. Transmission*, vol. 1, pp.3–11, 1965.

- [24] G. Korodi, I. Tabus, J. Rissanen and J. Astola, "DNA sequence compression - based on the normalized maximum likelihood model," *IEEE Signal Processing Magazine*, vol. 24, pp.47 – 53, 2007.
- [25] R. Krichevsky, "A relation between the plausibility of information about a source and encoding redundancy," *Problems Inform. Transmission*, vol. 4, n.3, pp. 48–57, 1968.
- [26] R. Krichevsky, *Universal Compression and Retrieval*. Kluwer Academic Publishers, 1993.
- [27] S. Kullback, *Information Theory and Statistics*, Wiley, New York, 1959.
- [28] D.S. Modha and E. Masry, "Memory-universal prediction of stationary random processes," *IEEE Trans. Inform. Theory*, vol. 44, pp. 117–133, 1998.
- [29] G. Morvai, S.J. Yakowitz and P.H. Algoet, "Weakly convergent nonparametric forecasting of stationary time series," *IEEE Trans. Inform. Theory*, vol. 43, pp.483 – 498, 1997.
- [30] S.C. Moy, "Generalisations of Shannon-MacMillan theorem," *Pacific J. Math.*, vol. 11, pp.705–714, 1961.
- [31] A.B. Nobel, "On optimal sequential prediction," *IEEE Trans. Inform. Theory*, vol. 49, pp. 83–98, 2003.
- [32] A. Perez, "Extensions of Shannon-MacMillan's limit theorem to more general stochastic processes," in *Trans. Third Prague Conference on Information Theory, Statistical Decision Functions, and Random Processes*, . Czechoslovak Academy of Sciences, Prague, 1964, pp. 545–574.
- [33] J. Rissanen, "Generalized Kraft inequality and arithmetic coding," *IBM J. Res. Dev.*, vol. 20, n. 5, pp. 198–203, 1976.
- [34] J. Rissanen, "Modeling by shortest data description," *Automatica*, vol.14, pp. 465–471, 1978.
- [35] J. Rissanen, "Universal coding, information, prediction, and estimation", *IEEE Trans. Inform. Theory*, vol. 30, pp. 629–636, 1984.
- [36] J. Rissanen, *Stochastic Complexity in Statistical Inquiry*, World Scientific Publishing Co., Singapore, 1989.
- [37] J. Rissanen, "Fisher information and stochastic complexity," *IEEE Trans. Inform. Theory*, vol.42, n.1, pp. 40–47 1996.
- [38] J. Rissanen, *Information and complexity in statistical modeling*, Springer Verlag, 2007.
- [39] A. Rukhin and others. *A statistical test suite for random and pseudorandom number generators for cryptographic applications*, NIST Special Publication 800- 22 (with revision dated May,15,2001). <http://csrc.nist.gov/rng/SP800-22b.pdf>
- [40] B.Ya. Ryabko, "Twice-universal coding," *Problems of Information Transmission*, vol.20, n.3, pp. 173–177, 1984.
- [41] B.Ya. Ryabko, "Prediction of random sequences and universal coding," *Problems of Inform. Transmission*, vol. 24, n.2, pp. 87–96, 1988.

- [42] B.Ya. Ryabko, "A fast adaptive coding algorithm," *Problems of Inform. Transmission*, vol. 26, pp. 305–317, 1990.
- [43] B. Ya. Ryabko, "The complexity and effectiveness of prediction algorithms," *J. Complexity*, vol. 10, n.3, pp. 281–295, 1994.
- [44] B. Ryabko, J. Astola and A. Gammernan, "Application of Kolmogorov complexity and universal codes to identity testing and nonparametric testing of serial independence for time series," *Theoretical Computer Science*, vol. 359, pp.440–448, 2006.
- [45] B. Ya. Ryabko and V.A. Monarev, "Using information theory approach to randomness testing," *Journal of Statistical Planning and Inference*, vol. 133, n.1, pp. 95–110, 2005.
- [46] B. Ryabko and F. Topsoe, "On Asymptotically Optimal Methods of Prediction and Adaptive Coding for Markov Sources," *Journal of Complexity*, vol. 18, n.1, pp. 224–241, 2002.
- [47] D. Ryabko and M. Hutter, "Sequence prediction for non-stationary processes." In Proceedings IEEE International Symposium on Information Theory, 2006. pp. 2346-2350. (see also <http://arxiv.org/pdf/cs.LG/0606077>)
- [48] S. A. Savari, "A probabilistic approach to some asymptotics in noiseless communication," *IEEE Transactions on Information Theory*, vol. 46, pp. 1246-1262, 2000.
- [49] C. E. Shannon, "A mathematical theory of communication," *Bell Sys. Tech. J.*, vol. 27, pp. 379–423 and pp.623–656, 1948.
- [50] C.E. Shannon, "Communication theory of secrecy systems," *Bell Sys. Tech. J.*, vol. 28, pp. 656–715, 1948.
- [51] P.C. Shields, "The interactions between ergodic theory and information theory," *IEEE Transactions on Information Theory*, vol. 44, pp.2079–2093, 1998.
- [52] W. Szpankowski, *Average case analysis of algorithms on sequences*, John Wiley and Sons, New York, 2001.
- [53] P.M.B. Vitanyi and M. Li, "Minimum description length induction, Bayesianism, and Kolmogorov complexity," *IEEE Trans. Inform. Theory*, vol.46, pp. 446–464, 2000.