

# The nonprobabilistic approach to learning the best prediction

Boris Ryabko

## Abstract

The problem of predicting a sequence  $x_1, x_2, \dots$  where each  $x_i$  belongs to a finite alphabet  $A$  is considered. Each letter  $x_{t+1}$  is predicted using information on the word  $x_1, x_2, \dots, x_t$  only. We use the game theoretical interpretation which can be traced to Laplace where there exists a gambler who tries to estimate probabilities for the letter  $x_{t+1}$  in order to maximize his capital. The optimal method of prediction is described for the case when the sequence  $x_1, x_2, \dots$  is generated by a stationary and ergodic source. It turns out that the optimal method is based only on estimations of conditional probabilities. In particular, it means that if we work in the framework of the ergodic and stationary source model, we cannot consider pattern recognition and other complex and interesting tools, because even the optimal method does not need them. That is why we suggest a so-called nonprobabilistic approach which is not based on the stationary and ergodic source model and show that complex algorithms of prediction can be considered in the framework of this approach.

The new approach is to consider a set of all infinite sequences (over a given finite alphabet) and estimate the size of sets of predictable sequences with the help of the Hausdorff dimension. This approach enables us first, to show that there exist large sets of well predictable sequences which have zero measure for each stationary and ergodic measure. (In fact, it means that such sets are invisible in the framework of the ergodic and stationary source model and shows the necessity of the new approach.) Second, it is shown that there exist quite large sets of such sequences that can be predicted well by complex algorithms which use not only estimations of conditional probabilities.

*Index terms:* prediction, learning, stationary and ergodic source, Hausdorff dimension, Kolmogorov complexity, Turing machine.

## 1 Introduction

Presently, the problem of prediction is investigated by many researches because its practical applications and importance for probability theory, machine intelligence, learning and other theoretical sciences: see, for example, [1]. We shall investigate a model of prediction of time series which can be traced to Laplace (cf. [4] where the problem is referred to as the problem of succession). Namely, we consider the finite alphabet  $A = \{a_1, \dots, a_m\}$ ,  $m \geq 2$ , and an infinite sequence  $x_1, x_2, \dots$ , where  $x_i \in A$ . Let us assume that gambler  $I$  has capital  $V_t$  at moment  $t = 0, 1, \dots$ ;  $V_0 = 1$ . Each moment  $t$  that the gambler divides the capital  $V_t$  into  $|A|$  parts equals  $V_t P_I(a/x_1, \dots, x_t)$ ,  $a \in A$ . (The Stakes on the  $x_{t+1} = a \in A$ ).

$P_I(a/x_1, \dots, x_t)$  reflects the  $I$ - gambler's confidence in the appearance of  $a \in A$  at the moment  $t + 1$ . (Here and in the sequel  $|X|$  denotes the number of letters, if  $X$  is an alphabet or a set, and the length of  $X$ , if the  $X$  is a word).

The gambler  $I$  learns  $x_{t+1}$  at the moment  $t + 1$  and his capital becomes equal to  $|A|V_t P_I(x_{t+1}/x_1, \dots, x_t)$ . (i.e., the stake on the right letter is increasing  $|A|$  times).

This model was suggested in [6] and used by many authors, see, for example, [3], [9]. It seems that this model is the simplest and can be considered as the useful tool for making first steps towards an understanding of the connection between complexity of prediction algorithms and their efficiency.

Denote the win by the method (or gambler)  $I$  with  $W_I(x_1, \dots, x_t)$  :

$$W_I(x_1, \dots, x_t) = \prod_{i=1}^t (|A| P_I(x_{i+1}/x_1, \dots, x_i)) \quad (1)$$

In this case the gambler divides his capital into  $|A|$  parts each moment (stakes on the letter of  $A$ ), and at the moment  $(t + 1)$  his capital should be equal to

$$W_I(x_1, \dots, x_t) (|A| P_I(x_{t+1}/x_1, \dots, x_t)),$$

i.e., the stake on the letter that actually appears, increases  $|A|$  times.

Our goal is to find the algorithms maximizing (1). From the mathematical point of view it is convenient to study the logarithm of  $W_I(x_1, \dots, x_t)$  divided by  $t$ . Let

$$L_I(x_1, \dots, x_t) = \log(W_I(x_1, \dots, x_t))/t.$$

(Here and below  $\log x = \log_2 x$ ). Then

$$L_I(x_1, \dots, x_t) = \log |A| + \sum_{i=1}^t \log P_I(x_i/x_1, \dots, x_{i-1})/t. \quad (2)$$

We note that, if at the word  $x_1, \dots, x_t$  strategy  $I$  recommends staking all the capital  $V_0$  ( $V_0 = 1$ ) on  $x_1$  (that is unknown), then the capital  $V_1 = |A|$  on  $x_2$  and so on. This results in a maximal win equal to  $|A|^t$  after  $t$  games.

If the gambler breaks up all his capital into  $|A|$  equal parts each time and stakes those parts on the letters of  $A$ , then the win is equal to 0 and the capital remains the same:  $V_0 = V_1 = \dots = V_t = 1$ . In general,  $L_I(x_1, \dots, x_t)$  is the mean value of the exponent of the capital increase over the first  $t$  games:  $V_t = |A|^{2^{tL_I(x_1, \dots, x_t)}}$ . This value characterizes quantitatively the efficiency of prediction  $I$ .

It is natural to look for the gambler's strategy which gives the maximal asymptotic value of  $L_I(x_1, \dots, x_t)$  in (2). In the next section we consider this problem when it is known that the sequence  $x_1 \dots x_t$  is generated by an ergodic and stationary source but parameters of the source are not known (ESS model sometimes called as a problem of universal prediction). The main result of that part is the asymptotically optimal (for the ESS model) method of prediction named  $\rho$ . This method can find each periodicity. For example, let  $A = \{0, 1\}$  and  $x_1 \dots x_4 = 0101$ . The described below optimal method  $\rho$  predicts

$$P_\rho(0/0101) = 33/52, P_\rho(1/0101) = 19/52.$$

Thus, the stake on 0 after 0101 is almost twice as high as the stake on 1.

The method  $\rho$  is based on estimations of conditional probabilities. In particular, it means that the optimal (for the ESS) method  $\rho$  does not use even simple regularities. As an example, consider two infinite sequences

$$x^1 = 1010010001000010\dots, x^2 = 0100011011000001010011\dots \quad (3)$$

The structure of the first sequence is obvious. In order to understand regularities in  $x^2$  it is enough to figure out that it is the sequence of the all lexicographically ordered binary words  $(0, 1, 00, 01, 10, 11, 000, 001, \dots)$ . It is easy to show that the method  $\rho$  does not use regularities of the  $x^1$  and  $x^2$ , and that is why its prediction is not efficient. In fact, the sequences like  $x^1$  and  $x^2$  are invisible in the framework of the ESS model because they belong to the set which measure is equal to zero according to any ergodic and stationary measure.

So, on the one hand the method  $\rho$  is optimal in framework of the ESS model, on the other hand, it cannot find even quite simple regularities and use them for predicting. Moreover, it is shown in the third part of the paper, that the invisible set is quite large. In order to measure the 'size' of sets we use the Hausdorff dimension. (Hausdorff dimension, best known as a powerful tool of fractal geometry, has been known for over fifty years, to be closely related to information theory and prediction ,see, for example, [2, 9, 10].)

Thus, we know that the set of sequences which cannot be investigated in the framework of the ESS model is large. But how many sequences from this set can be predicted quite well ? This question is considered in the last part of the paper. It is shown that there exists a large subset of sequences from the 'invisible' set which can be predicted well by algorithms realizable by Turing machine. It turns out that the efficiency of the best algorithms is closely related with Kolmogorov complexity.

## 2 A method which is asymptotically optimal for ergodic and stationary sources

We shall describe the optimal method  $\rho$  which is based on results from universal coding theory [9]. First, we give some definitions. Let  $A^n$ ,  $A^*$  and  $A^\infty$  be the sets of all  $n$ -length words, all finite words and all infinite words over the alphabet  $A$ , respectively. Let  $u_1u_2\dots u_n$  and  $v_1v_2\dots v_k$  be two words of  $A^*$  and  $k \leq n$ . By  $\tau_v(u)$  we denote the rate of the word  $v$  occurring in the sequence  $u_1u_2\dots u_k, u_2u_3\dots u_{k+1}, u_3u_4\dots u_{k+2}, \dots, u_{n-k+1}\dots u_n$ . For example,  $\tau_{00}(000100) = 3$ . For  $k = 0, 1, 2, \dots$  denote the mapping  $\rho_k$  that assigns to each word  $u \in A^*$  the value

$$p_k(u) = \begin{cases} |A|^{-|u|}, & \text{for } |u| \leq k \\ \left(\frac{\Gamma(|A|/2)}{\Gamma(1/2)|A|}\right)^{|A|^k} \cdot \frac{1}{|A|^k} \prod_{\alpha \in A^k} \frac{\prod_{a \in A} \Gamma(\tau_{\alpha a}(u) + 1/2)}{\Gamma(\tau_\alpha^*(u) + |A|/2)}, & \text{for } |u| > k \end{cases}$$

where  $\tau_\alpha^* = \sum_{a \in A} \tau_{\alpha a}(u)$  and  $\Gamma()$  is a gamma function. Now, define the probability distribution  $\lambda$  on the set of nonnegative integers  $0, 1, 2, \dots$  using the code from [7]. Let  $\log^{(0)}(x) = x$ ,  $\log^{(i)}(x) = \log_2(\log^{(i-1)}(x))$  for  $i \geq 1$  and  $m(x) = i$ , so that  $0 \leq \log^{(i)}(x) \leq 1$ . For  $n = 0, 1, 2, \dots$ , define

$$w(n) = \sum_{i=1}^{m(n)} \lceil \log^{(i)}(n) \rceil + m(n) + 1, \lambda(n) = 2^{-w(n)}.$$

It is known that

$$-\log \lambda(n) = \log n + O(\log \log n)$$

when  $n \rightarrow \infty$  [9].

We define the measure  $\rho$  by

$$\rho = \sum_{k=0}^{\infty} \lambda(k) \rho_k(u).$$

The method of prediction is defined by

$$P_\rho(a/x_1x_2 \dots x_t) = \rho(x_1x_2 \dots x_t a) / \rho(x_1x_2 \dots x_t).$$

As an example, consider the prediction computation on  $x = 0101$ ,  $A = \{0, 1\}$  by the method  $\rho$ . We have

$$\rho(01010) = 33/2^{10}, \rho(01011) = 19/2^{10}.$$

Hence, the prediction follows:

$$P_\rho(0/0101) = 33/52, P_\rho(1/0101) = 19/52.$$

Thus, the stake on 0 after 0101 is almost twice as high as the stake on 1.

From the definition of  $\rho_k$  and  $\rho$  we can see that the prediction method  $\rho$  uses the rates of the subwords occurring in the sequence  $x_1x_2 \dots x_t$ . In fact, those rates are the estimations of the conditional probabilities.

The theorem below shows that the method  $\rho$  is close to being optimal when  $t \rightarrow \infty$ .

**Theorem 1** *For any stationary and ergodic process over  $A^\infty$  with probability equal to 1,*

$$\lim_{t \rightarrow \infty} L_\rho(x_1x_2 \dots x_t) = \log |A| - h$$

where  $h$  is the Shannon entropy of the process. (The limit does exist with probability equal to 1.)

On the other hand, any prediction method  $\alpha$  is asymptotically no better than  $\rho$ , so the following inequality is valid with probability equal to 1:

$$\lim_{t \rightarrow \infty} L_\alpha(x_1x_2 \dots x_t) \leq \log |A| - h$$

The proof is given in [9]. The definition of the Shannon entropy can be found, for example, in [5].

### 3 The new approach

Let us consider the performance of the prediction method  $\rho$  when it is applied to the sequences  $x^1$  and  $x^2$ , which are defined by (3). In both cases  $\rho$  does not recognize regularities. When  $\rho$  is applied to  $x^1$  and  $t \rightarrow \infty$ , the gambler will stake almost all money on '0' and will not be able to predict appearances of 1's because, informally speaking, the sequence  $x^1$  does not contain subwords which are repeated periodically. It is not the case when  $\rho$  is applied to  $x^2$  because it is known that any subword  $u \in \{0, 1\}^n$ ,  $n \leq 1$  has the

frequency of occurrence  $2^{-n}$  in  $x^2$ . That is why the gambler will stake around a half of his capital on '0' and the other half on '1', when  $t \rightarrow \infty$ , and his capital remains the same. On the other hand, if the gambler finds the regularity his capital will grow exponentially.

So, on the one hand, Theorem 1 claims that the prediction method  $\rho$  is asymptotically optimal when it is applied to a stationary and ergodic source. On the other hand, we have seen that  $\rho$  is completely inefficient when it is applied to sequences with regularities. In order to get over the paradox we suggest a new approach. First we give some new definitions. For simplicity sake, we consider only the case of  $A = \{0, 1\}$  in this section, but all results can be easily extended to the general case. Let  $\mu$  be a stationary and ergodic measure on  $\{0, 1\}^\infty$  and  $x \in \{0, 1\}^\infty$ . The sequence  $x$  is defined to be  $\mu$ -typical if for any word  $u \in \{0, 1\}^*$

$$\lim_{n \rightarrow \infty} (\tau_u(x|_1^n) / (n - |u| + 1)) = \mu(u)$$

where  $x|_1^n = x_1, \dots, x_n$ . Informally, it means that the frequency of occurrences of the word  $u$  in the  $x$  is equal to the probability of  $u$  according to the measure  $\mu$ . Let  $T$  be the set all typical sequences:

$$T = \{x \in \{0, 1\}^\infty : x \text{ is typical for some stationary ergodic measure } \}.$$

First, we show that the set of untypical sequences which is defined as follows:

$$T^* = \{0, 1\}^\infty - T$$

is 'invisible' and quite large.

**Theorem 2** *i)  $\mu(T^*) = 0$  for any ergodic and stationary measure  $\mu$ .  
ii)  $\dim_H(T^*) = 1$*

Proof of the Theorem 2 as well as the definition of the Hausdorff dimension is given in Appendix. (See also [2] for the definition.)

From i) we can see that the set  $T^*$  has the measure 0 for each ergodic stationary measure. That is why  $T^*$  is invisible in the framework of the ESS model. On the other hand, it is known that  $\dim_H(S) \leq 1$  for each subset  $S$  from  $\{0, 1\}^\infty$ . So, the size of the invisible set  $T^*$  is maximal.

But how many well predictable sequences does the set  $T^*$  contain? If there does not exist many such sequences we may not look for new models and approaches. The informal answer is following: the set of well predictable but invisible sequences is as large as the set of predictable and 'visible' sequences. Theorem 3 gives a more formal answer on that question. But first, we have to specify the notation of a prediction method. Since this moment we will consider only algorithmically realizable prediction methods because only such methods are interesting from practical point of view. (Algorithmically realizable prediction methods can be defined by using each 'standard' notation of the algorithm. For example, we may think that such algorithms are described as the Turing machines).

**Theorem 3** *For every  $\alpha \in [0, 1]$  there exists the set  $W_\alpha$  such that*

- i) there exists the prediction method  $\gamma$  for which  $\lim_{n \rightarrow \infty} L_\gamma(x|_1^n) \geq 1 - \alpha$  for all  $x \in W_\alpha$*
- ii)  $\dim_H(W_\alpha) = \alpha$*
- iii)  $\mu(W_\alpha) = 0$  for every stationary and ergodic measure  $\mu$ .*

The proof is given in Appendix.

We define  $U_\alpha$  as the set of all sequences  $x$  such that there exists an algorithmically realizable prediction method  $\delta(x)$  for which  $\lim_{n \rightarrow \infty} L_{\delta(x)}(x|_1^n) \geq 1 - \alpha$ .

The set  $U_\alpha$  is worth considering because it contains all well predictable sequences. It turns out that this set is closely connected with Kolmogorov complexity (the definition of Kolmogorov complexity can be found, for instance, in [?].)

**Theorem 4** *i) For each  $x \in U_\alpha$*

$$\lim_{n \rightarrow \infty} KC(x|_1^n)/n \leq \alpha$$

where  $KC(u)$  is the Kolmogorov complexity of the word  $u$ .

*ii)*

$$\dim_H(U_\alpha) = \alpha.$$

The proof may be easily obtained from results of the papers [8, 9].

By definition,  $U_\alpha$  contains all well predictable sequences. Hence, it contains the set  $W_\alpha$  which is also well predictable and invisible in the framework of the ESS model. So, Theorems 3 and 4 show that the invisible set  $W_\alpha$  has the maximal size.

## 4 Appendix

*The definition of the Hausdorff dimension.*

A set  $S$  consisting of subsets in  $A^*$  is called a  $\rho$ -cover of  $Y \subset A^\infty$  for  $\rho > 0$  if

1.  $x \in Y$  has  $\sigma \in S$  as its prefix; and
2.  $2^{-|\sigma|} \leq \rho$  for  $\forall \sigma \in S$ .

Let  $C(Y, \rho)$  denote the set of such  $S$ 's, and

$$l^{(\alpha)}(Y, \rho) = \inf_{S \in C(Y, \rho)} \sum_{\sigma \in S} 2^{-\alpha|\sigma|}.$$

Then, the Hausdorff dimension of the set  $Y$  is defined by

$$\dim_H(Y) = \inf\{\alpha : \lim_{\rho \rightarrow 0} l^{(\alpha)}(Y, \rho) = 0\} = \sup\{\alpha : \lim_{\rho \rightarrow 0} l^{(\alpha)}(Y, \rho) = \infty\}.$$

*Proof of Theorem 3.*

It is easy to see that for every  $\alpha \in [0, 1]$  there exists  $\pi$  for which  $-(\pi \log \pi + (1 - \pi) \log(1 - \pi)) = \alpha$  where, by definition,  $-(\pi \log \pi + (1 - \pi) \log(1 - \pi))$  is the Shannon entropy  $h(\pi)$ . We consider two Bernoulli sources  $\mu_1, \mu_2$  over  $\{0, 1\}^\infty$  such that

$$\mu_1(0) = \pi, \mu_1(1) = 1 - \pi, \mu_2(0) = 1 - \pi, \mu_2(1) = \pi$$

and let  $B_1$  and  $B_2$  be the sets of all  $\mu_1$ -typical sequences and  $\mu_2$ -typical sequences, respectively. In order to define the set  $W_\alpha$  we take any sequences  $b^1 = b_1^1 b_2^1 \dots \in B_1$  and  $b^2 = b_1^2 b_2^2 \dots \in B_2$  and define the sequence  $w \in W_\alpha$  as follows:

$$w = b_1^1 b_2^2 b_3^2 b_4^1 b_5^1 b_6^1 b_7^1 b_8^2 b_9^2 b_{10}^2 \dots b_{15}^2 b_{16}^1 b_{17}^1 \dots b_{31}^1 b_{32}^2 \dots$$

It is known that the Hausdorff dimension of the set of all typical sequences of an ergodic stationary source is equal to the Shannon entropy of the source [2]. So,  $\dim_H(B_1) = \dim_H(B_2) = h(\pi) = -(\pi \log \pi + (1-\pi) \log(1-\pi)) = \alpha$ . On the other hand, in [9] it is shown that there exists a strategy  $B$  such that for each  $x \in B_1$  the equality  $\lim_{n \rightarrow \infty} L_B(x|_1^n) = -(\pi \log \pi + (1-\pi) \log(1-\pi)) = \alpha$  is valid. (This strategy is quite simple; we should stake the share of capital  $\pi$  on '0' and the share of capital  $1-\pi$  on '1'.) It is obvious how to change this strategy in order to obtain the strategy  $\gamma$  from the statement i) of the theorem. On the other hand, it is intuitively clear that  $W_\alpha$  can be transformed into  $B_1$  (or  $B_2$ ) without any compression and expansion. That is why all three sets have the same Hausdorff dimension  $\alpha$  and we obtain the statement ii). In order to prove iii) we should take into account that all sequences in  $W_\alpha$  are untypical for each stationary and ergodic measure simply because  $\lim_{n \rightarrow \infty} (\tau_0(x|_1^n)/(n - |u| + 1))$  does not exist. (According to the definition of typical sequences all such limits should exist.) Theorem is proved.

*Proof of Theorem 2.* It is known that for any ergodic stationary measure  $\mu$  the following equality is valid :

$$\mu(\text{the set of all } \mu\text{-typical sequences}) = 1,$$

see [2]. From this equality and the definition of  $T$  we can see that  $\mu(T) = 1$ . Hence,  $\mu(T^*) = \mu(\{0, 1\}^\infty - T) = 0$ . In order to prove the second statement of the theorem it is enough to note that  $T$  contains the sets  $W_\alpha$  which are described in Theorem 3. So, we obtain the statement ii) from Theorem 3 if the parameter  $\alpha \rightarrow 1$ . (That is why the proof of Theorem 3 is given before the proof of Theorem 2.)

## References

- [1] P. Algoet. *Universal Schemes for Learning the Best Nonlinear Predictor Given the Infinit Past and Side Information* , IEEE Trans. Inform. Theory, v. 45, n.4, pp. 1165-1185, 1999.
- [2] P. Billingsley, *Ergodic theory and information*, John Wiley & Sons (1965).
- [3] T.Cover,J.Thomas.*Elements of Information Theory*. New York: Wiley & Sons, 1991.
- [4] W. Feller, *An Introduction to Probability Theory and Its Applications*. New York: Wiley & Sons, 1970.
- [5] Gallager R. G., *Information Theory and Reliable Communication*. John Wiley & Sons, New York, 1968.
- [6] Kelly J.L., *A new interpretation of information rate*, Bell System Tech. J., v. 35, pp. 917-926, 1956.
- [7] V.Levinshstein, *The redundancy and deceleration of a separative encoding of natural numbers*. In: 'Problems of Cybernetics', v.20, Moscow, pp.173-179.
- [8] B. Y. Ryabko, "Noiseless coding of combinatorial sources, Hausdorff dimension, and Kolmogorov Complexity" , Problems of Information Transmission, vol. 22, No. 1, pp. 16-26 (1986).
- [9] B.Ya.Ryabko,"The complexity and effectiveness of prediction algorithms." *J. of Complexity*, v.10 (1994), no. 3, pp.281–295.
- [10] Ryabko B., Suzuki J., Topsoe, F.,*Hausdorff dimension as a new dimension in source coding and predicting*, In:Proceedings of the 1999 IEEE Information Theory and Communications Workshop, 1999. pp. 66 -68