

A Preadapted Universal Switch Distribution for Testing Hilberg's Conjecture

Łukasz Dębowski*

Abstract

Hilberg's conjecture about natural language states that the mutual information between two adjacent long blocks of text grows like a power of the block length. The exponent in this statement can be upper bounded using the pointwise mutual information estimate computed for a carefully chosen code. The bound is the better, the lower the compression rate is but there is a requirement that the code be universal. So as to improve a received upper bound for Hilberg's exponent, in this paper, we introduce two novel universal codes, called the plain switch distribution and the preadapted switch distribution. Generally speaking, switch distributions are certain mixtures of adaptive Markov chains of varying orders with some additional communication to avoid so called catch-up phenomenon. The advantage of these distributions is that they both achieve a low compression rate and are guaranteed to be universal. Using the switch distributions we obtain that a sample of a text in English is non-Markovian with Hilberg's exponent being ≤ 0.83 , which improves over the previous bound ≤ 0.94 obtained using the Lempel-Ziv code.

Keywords: universal coding, natural language, Hilberg's conjecture

*Ł. Dębowski is with the Institute of Computer Science, Polish Academy of Sciences, ul. Jana Kazimierza 5, 01-248 Warszawa, Poland (e-mail: ldebowsk@ipipan.waw.pl).

I Introduction

Hilberg's conjecture is a hypothesis concerning natural language which states that the mutual information between two adjacent long blocks of text grows very fast, namely as a power of the block length [1, 2, 3, 4, 5, 6, 7]. There are two important information-theoretic results concerning this conjecture. On the one hand, Hilberg's hypothesis can be linked with the idea that texts in natural language refer to large amounts of randomly accessed knowledge in a repetitive way [8, 9]. On the other hand, Hilberg's hypothesis can be linked with the fact that the number of distinct words in a text grows as a power of the text length [10, 8], the fact known as Herdan's or Heaps' law [11, 12]. These two results make Hilberg's conjecture interesting and worth direct empirical testing.

To present Hilberg's conjecture formally, let us introduce some notations. Consider a probability space (Ω, \mathcal{F}, Q) with $\Omega = \{1, 2, \dots, D\}^{\mathbb{Z}}$, random variables $X_k : \Omega \ni (x_i)_{i \in \mathbb{Z}} \mapsto x_k \in \{1, 2, \dots, D\}$, and distribution Q which is stationary on $(X_i)_{i \in \mathbb{Z}}$ but not necessarily ergodic. Blocks of symbols or variables are denoted as $X_n^m = (X_i)_{n \leq i \leq m}$ with X_n^m being the empty block for $m < n$. Moreover, for a random variable X we introduce a random variable $Q(X)$ which takes value $Q(X = x)$ for $X = x$. The pointwise entropy of variable X is the random variable

$$H^Q(X) = -\log Q(X) \quad (1)$$

whereas the pointwise mutual information between X and Y is

$$I^Q(X; Y) = -\log Q(X) - \log Q(Y) + \log Q(X, Y). \quad (2)$$

Having this in mind, Hilberg's conjecture states that

$$I^Q(X_1^n; X_{n+1}^{2n}) \propto n^\beta, \quad \beta \in (0, 1). \quad (3)$$

Hilberg [1] supposed that $\beta \approx 0.5$ holds for texts in English but his estimate was very rough, based on Shannon's psycholinguistic experiment [13].

It is an interesting open question how much the exponent β varies across different texts and whether it is possibly a text-independent language universal. This question can be connected to some fundamental limitations of human memory and attention. Consequently, in this paper, we want to improve the method of upper bounding Hilberg's exponent β proposed in [14], which is based on universal coding [15] or universal distributions [16]. The focus of the present paper is to develop a better tool for estimating mutual information than used so far and to demonstrate how it works with an instance of empirical language data. Whereas we propose some method of upper bounding exponent β and showing that texts in natural language are non-Markovian, it remains to find a computational method of lower bounding Hilberg's exponent β .

Before we show how to upper bound exponent β using a universal distribution, it is important to note that the discussion of Hilberg's hypothesis is intimately connected with the question whether the natural language production is nonergodic and what its plausible ergodic decomposition is. According to the ergodic decomposition theorems [17, 18, 8], any stationary measure Q equals the expectation $\mathbf{E}_Q F$, where $F = Q(\cdot | \mathcal{I})$ is the random ergodic measure for measure Q and \mathcal{I} is the shift-invariant algebra. There exist some stationary

nonergodic measures Q , called Santa Fe processes, for which mutual information $I^Q(X_1^n; X_{n+1}^{2n})$ grows according to a power law but the main contribution of the mutual information comes from identifying the random ergodic measure F given the block X_1^n [8, 9]. In that case mutual information $I^F(X_1^n; X_{n+1}^{2n})$ for the random ergodic measure F itself is negligibly small. The Santa Fe processes are not irrelevant for our discussion. In their original construction they were intended as some idealized models for the transmission of knowledge in natural language. Thus when estimating the exponent in Hilberg's conjecture we have first to decide whether we do it for measure Q or for measure F . The methods for estimating $I^Q(X_1^n; X_{n+1}^{2n})$ and $I^F(X_1^n; X_{n+1}^{2n})$ are very different. We suppose that estimating the mutual information for the nonergodic measure Q is closer to the original intention of Hilberg although some philosophical problem remains 'how many' ergodic components we actually admit (e.g. do we assume that Q models texts in a given register of a particular language or in any register of any natural language). In contrast to the random ergodic measure F , the possibly nonergodic measure Q is not identifiable given a single realization $(X_i)_{i \in \mathbb{Z}}$. Despite that, it is somewhat baffling that some nontrivial upper bound for Hilberg's exponent β , a property of measure Q , can be learned from a single realization $(X_i)_{i \in \mathbb{Z}}$ if the growth of mutual information is uniform in Q .

As we have indicated, our method of upper bounding Hilberg's exponent β is based on universal coding. Here we say that a distribution P is weakly universal if for every stationary distribution Q we have

$$\lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{E}_Q H^P(X_1^n) = h_Q, \quad (4)$$

where the entropy rate h_Q is

$$h_Q := \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{E}_Q H^Q(X_1^n) = \inf_{k \in \mathbb{N}} \mathbf{E}_Q [-\log Q(X_{k+1}|X_1^k)]. \quad (5)$$

On the other hand, the distribution P is called strongly universal if for every stationary ergodic distribution Q we have Q -almost surely

$$\limsup_{n \rightarrow \infty} \frac{1}{n} H^P(X_1^n) \leq h_Q. \quad (6)$$

Strongly universal distributions are weakly universal under mild conditions [19]. Dębowski [14] proposed to investigate the empirical law of form

$$I^P(X_1^n; X_{n+1}^{2n}) \propto n^\gamma, \quad \gamma \in (0, 1), \quad (7)$$

where P is a strongly universal distribution. Relationship (7), which can be called the codewise Hilberg conjecture, has been checked experimentally for the Lempel-Ziv code on a sample of 10 texts in English and it holds surprisingly uniformly with $\gamma = 0.94$ [14]. The same estimate $\gamma = 0.94$ has been obtained for 21 other texts in German and French (work under review).

Are laws (3) and (7) related? In fact, if they hold uniformly for large n then exponents β and γ can be linked. To see it, the following lemma is helpful.

Lemma 1 ([8]) *Consider a function $G : \mathbb{N} \rightarrow \mathbb{R}$ such that $\lim_k G(k)/k = 0$ and $G(n) \geq 0$ for all but finitely many n . For infinitely many n , we have $2G(n) - G(2n) \geq 0$.*

If P is weakly universal, the above statement is satisfied for Kullback-Leibler divergence

$$G(n) = \mathbf{E}_Q [H^P(X_1^n) - H^Q(X_1^n)] . \quad (8)$$

Hence we obtain that

$$\mathbf{E}_Q I^P(X_1^n; X_{n+1}^{2n}) \geq \mathbf{E}_Q I^Q(X_1^n; X_{n+1}^{2n}) \quad (9)$$

holds for infinitely many n . Thus if relationships (3) and (7) hold uniformly for large n then

$$\gamma \geq \beta. \quad (10)$$

In other words, the smaller γ we observe for a text (or methodologically better, for a large sample of different texts), the better bound it gives for β . It can also be easily shown that the bound is the tighter, the smaller compression rate $H^P(X_1^n)/n$ is, with the sole provision that distribution P be weakly universal. Results of our experiment suggest that this requirement is essential.

Thus the question of upper bounding Hilberg's exponent β boils down, if we assume uniform information growth in (3), to finding appropriate universal distributions. Many methods have been proposed for compression of texts in natural language, e.g.: Lempel-Ziv (LZ) code [15], n -gram models [20, 21, 22], prediction by partial match (PPM) [23], context tree weighting (CTW) [24], probabilistic suffix trees (PST) [25], grammar-based codes [26], PAQ codes [27], and switch distributions [28]. These compression schemes can be divided into two classes: (a) preadapted distributions, which are trained on large corpora and achieve low compression rate—as low as 0.88 bpc (bits per character) for WinRK 3.1.2,¹ and (b) adaptive distributions, which are not pre-trained and achieve larger compression rate but are proven to be universal. Whereas the distributions proposed so far belong either to class (a) or (b), for upper bounding Hilberg's exponent, we need a distribution that would combine the advantages of classes (a) and (b), namely low compression rate and universality.

In this paper we propose and investigate two novel universal distributions, one of which is not preadapted and the other is preadapted. The point of our departure is a modification of the switch distributions proposed in [28, 29]. The idea of a switch distribution is to use a mixture of adaptive Markov chains of varying orders but, at each data point, the probabilities are partly transferred among different orders. In this way, lower order Markov chains are used to compress the data exclusively until enough information is gathered to predict new outcomes with higher order chains. This avoids so called catch-up phenomenon and leads to much better compression than while there is no transmission of probabilities among different Markov chain orders [28]. If we combine the idea of the switch distribution with smoothing proposed in [30, p. 111] and the idea of a universal distribution called R measure, proposed in [16], we obtain another new universal compression scheme, which is efficiently computable. This scheme will be called the plain switch distribution. It is not preadapted yet. The preadapted switch distribution is obtained by initializing the Markov chains with frequencies coming from a large corpus and letting them gradually adapt to the compressed source. It will be shown that the

¹<http://www.maximumcompression.com/data/text.php>

preadapted switch distributions is also universal. For the considered input text the nonpreadapted and the preadapted switch distributions achieve almost the same ultimate compression rate 2.21 bpc, approximately twice smaller than for the LZ code. This figure is not so favorable as for the WinRK 3.1.2 but we have a guarantee that the switch distributions are universal.

Once we have constructed the universal switch distributions, we can use them for upper bounding Hilberg's exponent. In the previous paper [14], the LZ code was used for a sample of texts in English which yielded $\gamma = 0.94$. Here using the plain switch distribution we obtain a slightly tighter bound $\gamma = 0.83$. Surprisingly, the preadapted switch distribution yields almost the same compression rate for long blocks as the plain switch distribution and does not give a tighter bound for γ . Differences in the estimates of γ may also stem from differences in data representation. In [14] the alphabet of $D = 27$ symbols was used. Here we use $D = 256$ and obtain $\gamma = 0.89$ for the LZ code. It is important to underline that meaningful estimates of γ can be only obtained using universal distributions. As we show, if a nonuniversal distribution is used, the pointwise mutual information can be very low despite a good-looking compression rate. To a certain extent this also applies to the preadapted switch distribution, where the pointwise mutual information is low for short blocks.

There remains a question what the estimates of the codewise Hilberg exponent γ tell about the true Hilberg exponent β for texts in natural language. In particular, how large can the difference between γ and β be in case of the considered universal codes? Let us recall that for memoryless sources, i.e., IID processes with $\beta = 0$, the pointwise mutual information of the LZ code is of order $I^P(X_1^n; X_{n+1}^{2n}) \propto n/\log n$ [31] so empirically we should observe $\gamma \approx 1$. Thus the difference between γ and β can be close to 1. Moreover, observing $\gamma \approx 1$ for the LZ code we cannot reject the hypothesis that the source is IID. In contrast, for stationary Markov chains, the pointwise mutual information of a Bayesian mixture of Markov chains, which is the building block for the switch distribution, cf. [32], is only of order $I^P(X_1^n; X_{n+1}^{2n}) \propto \log n$ [33] so empirically we should observe $\gamma \approx 0$ in that case. The same property carries over to the switch distributions introduced in this paper. Thus, observing $\gamma > 0$ for a switch distribution, we are compelled to reject the hypothesis that the source is a Markov chain. Possibly, $\beta > 0$ may hold in that case. In other words, given the presented empirical data, Hilberg's conjecture may be true but we need still some stronger evidence in favor of this hypothesis, such as a nontrivial lower bound for exponent β (work under review).

The further organization of the paper is as follows. In Section II, we present the plain switch distribution. In Section III, we discuss the preadapted switch distribution. In Section IV, we investigate the codewise Hilberg's conjecture (7) experimentally using the introduced distributions.

II The plain switch distribution

The frequency of substring $w_1^k \in \{1, 2, \dots, D\}^k$ in string $z_1^n \in \{1, 2, \dots, D\}^n$ will be denoted as

$$c(w_1^k | z_1^n) = \sum_{i=0}^{n-k} \mathbf{1}\{w_1^k = z_{i+1}^{i+k}\}. \quad (11)$$

The plain switch distribution is defined as follows:

Definition 1 (plain switch distribution) Define conditional probabilities $B(x_{n+1}|x_1^n, -1) = D^{-1}$ and

$$B(x_{n+1}|x_1^n, k) = \frac{c(x_{n+1-k}^{n+1}|x_1^n) + B(x_{n+1}|x_1^n, k-1)}{c(x_{n+1-k}^n|x_1^{n-1}) + 1}. \quad (12)$$

Let coefficients $p_n \in (0, 1)$, where $n = 0, 1, 2, \dots$, satisfy $\prod_{n=0}^{\infty} p_n > 0$. Put also $q_n = 1 - p_n$. We define the partial switch distribution $P(x_1^n, k)$ by conditions

$$P(x_1, -1) = p_0 B(x_1 | -1), \quad (13)$$

$$P(x_1, 0) = q_0 B(x_1 | 0), \quad (14)$$

$$P(x_1^n, k) = 0 \text{ for } k < -1 \text{ or } k \geq n, \quad (15)$$

$$P(x_1^{n+1}, k) = [p_n P(x_1^n, k) + q_n P(x_1^n, k-1)] B(x_{n+1}|x_1^n, k) \\ \text{for } n \geq 1 \text{ and } -1 \leq k \leq n. \quad (16)$$

The total probability for block x_1^n according to the switch distribution is

$$P(x_1^n) = \sum_{k=-1}^{n-1} P(x_1^n, k). \quad (17)$$

The scheme of computing $P(x_1^n)$ is depicted in Figure 1.

Remark 1: Condition $\prod_{n=0}^{\infty} p_n > 0$ holds for instance if we fix

$$p_n = \exp [-(n+1)^{-\alpha}], \quad \alpha > 1. \quad (18)$$

Value α is a parameter.

Remark 2: Probability $B(x_{n+1}|x_1^n, k)$ defines an adaptive k -th order Markov model. Probability $P(x_1^n, k)$ represents the mass of the adaptive k -th order Markov model modified by communication with models of lower orders. The motivation for this communication, carried out in formula (16), is that lower order Markov models should be solely used for compression until enough data are collected to predict new outcomes with higher order Markov models, cf., the catch-up phenomenon described in [28]. Distribution $P(x_1^n)$ is a special case of the general scheme of switch distributions considered by [28, 29] to overcome the catch-up phenomenon. In contrast to the models discussed in [28, 29], the switch distribution considered here is universal in the sense made precise in the Introduction and still can be efficiently computed.

Remark 3: Probability $B(x_{n+1}|x_1^n, k)$ of an adaptive k -th order Markov model so as to be defined for all x_1^n is smoothed using probability $B(x_{n+1}|x_1^n, k-1)$ of an adaptive $(k-1)$ -th order Markov model in a way inspired by [30, p. 111]. Thus we add the probability of a lower order $B(x_{n+1}|x_1^n, k-1)$ in the numerator and 1 in the denominator rather than using the Laplace rule or the Krichevski-Trofimov rule, i.e., adding $\alpha \in (0, 1]$ in the numerator and αD in the denominator as in, e.g., [16]. We have checked that this trick works better for natural language data than the Laplace or Krichevski-Trofimov rules.

Remark 4: The infinite mixture of probabilities $B(x_{n+1}|x_1^n, k)$ for orders $k = 0, 1, 2, \dots$, smoothed using the Laplace or Krichevski-Trofimov rules, was

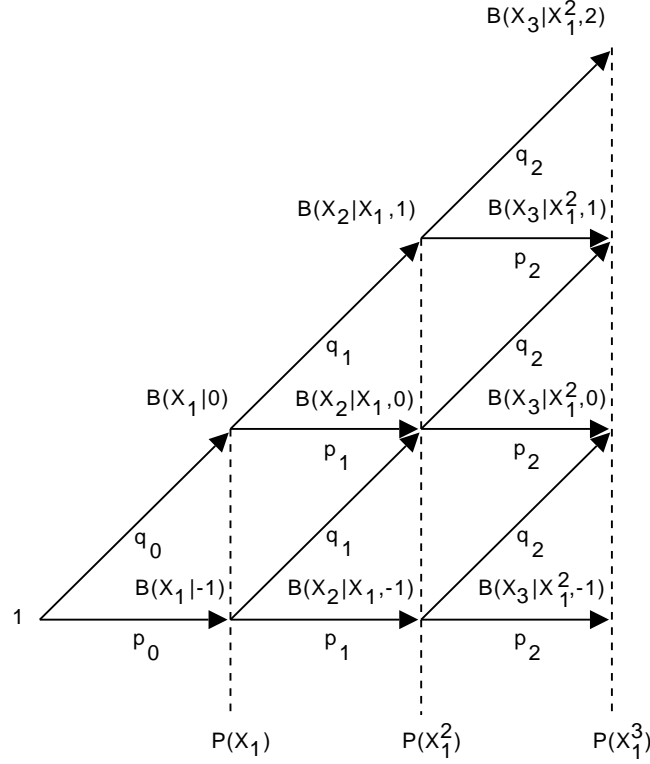


Figure 1: The scheme of computing $P(x_1^3)$.

investigated in [16] under the name of R measure and shown to be a universal distribution. Our construction is quite similar in spirit but avoids the catch-up phenomenon.

Now we will show that the switch distribution (17) is both strongly and weakly universal. First we need this simple fact:

Lemma 2 *Introduce notation*

$$B(x_l^n | x_{l-k}^{l-1}, k) = \prod_{i=l}^n B(x_i | x_1^{i-1}, k). \quad (19)$$

The switch distribution satisfies the following:

i) *there exists a constant $\delta_{-1} > 0$ such that for all $n \geq 1$ we have*

$$P(x_1^n) \geq \delta_{-1} D^{-n}, \quad (20)$$

ii) *for each $k \geq 0$ there exists a constant $\delta_k > 0$ such that for all $n \geq k + 1$ we have*

$$P(x_1^n) \geq \delta_k B(x_{k+1}^n | x_1^k, k). \quad (21)$$

Proof: For $n \geq 1$ we have

$$P(x_1^n) \geq \left(\prod_{i=0}^{n-1} p_i B(x_{i+1}|x_1^i, -1) \right) \geq \delta_{-1} P(x_1^n | -1). \quad (22)$$

where $\delta_{-1} = \prod_{i=0}^{\infty} p_i > 0$. Thus we have claim (i). On the other hand, for $k \geq 0$ and $n \geq k+1$ we obtain

$$\begin{aligned} P(x_1^n) &\geq \left(\prod_{i=0}^k q_i B(x_{i+1}|x_1^i, i) \right) \left(\prod_{i=k+1}^{n-1} p_i B(x_{i+1}|x_1^i, k) \right) \\ &= \left(\prod_{i=0}^k q_i \right) D^{-k} \left(\prod_{i=k+1}^{n-1} p_i \right) B(x_{k+1}^n | x_1^k, k) \\ &\geq \delta_k B(x_{k+1}^n | x_1^k, k), \end{aligned} \quad (23)$$

where

$$\delta_k = \left(\prod_{i=0}^k q_i \right) D^{-k} \left(\prod_{i=k+1}^{\infty} p_i \right) > 0. \quad (24)$$

Hence the claim (ii) follows. \square

Combining Lemma 2(ii) with the ergodic theorem we obtain the proof of universality.

Theorem 1 *The switch distribution is strongly and weakly universal.*

Proof: Let Q be a stationary ergodic distribution. Since the alphabet of X_i is finite, by the ergodic theorem differences $B(X_n | X_1^{n-1}, k) - Q(X_n | X_{n-k-1}^{n-1})$ converge to 0 Q -almost surely. Hence

$$\lim_{n \rightarrow \infty} \frac{1}{n} [-\log B(X_{k+1}^n | X_1^k, k)] = \lim_{n \rightarrow \infty} \frac{1}{n} \left[- \sum_{i=k+1}^n \log Q(X_i | X_{i-k-1}^{i-1}) \right]. \quad (25)$$

Applying the ergodic theorem again, we obtain

$$\lim_{n \rightarrow \infty} \frac{1}{n} \left[- \sum_{i=k+1}^n \log Q(X_i | X_{i-k-1}^{i-1}) \right] = \mathbf{E}_Q [-\log Q(X_{k+1} | X_1^k)]. \quad (26)$$

Now we combine these facts with Lemma 2(ii), which yields

$$\limsup_{n \rightarrow \infty} \frac{1}{n} [-\log P(X_1^n)] \leq \inf_{k \in \mathbb{N}} \lim_{n \rightarrow \infty} \frac{1}{n} [-\log B(X_{k+1}^n | X_1^k, k)] \quad (27)$$

$$= \inf_{k \in \mathbb{N}} \mathbf{E}_Q [-\log Q(X_{k+1} | X_1^k)] = h_Q. \quad (28)$$

Hence the distribution P is strongly universal. Moreover, as shown in [19], the claim of Lemma 2(i) and the strong universality are sufficient conditions that distribution P be weakly universal. \square

A naive implementation of the switch distribution $P(x_1^n)$ has the time complexity $O(n^4)$ for the following reason: There are $O(n^2)$ calls of $B(x_{l+1}|x_1^l, k)$ where $k, l \leq n$ and in a naive implementation each $B(x_{l+1}|x_1^l, k)$ has time complexity $O(kl)$. This is, however, a very careless approach and usually we can do much better. Let us denote the maximal length of a substring that appears at least twice in a string z_1^n as

$$\mathbf{L}(z_1^n) := \max \{k : \exists w_1^k : c(w_1^k|z_1^n) > 1\}. \quad (29)$$

For brevity, $\mathbf{L}(z_1^n)$ will be called the depth of z_1^n .

Theorem 2 *The value of the switch distribution $P(x_1^n)$ can be computed in time $O(ns)$ where $s = \mathbf{L}(x_1^n)$ is the depth of x_1^n .*

Remark: The depth $\mathbf{L}(X_1^n)$ is bounded by $O(\log n)$ for a large class of processes called finite-energy processes. They can be obtained by dithering ergodic processes with an IID noise [34]. For texts in natural language, an experiment indicates that the depth $\mathbf{L}(x_1^n)$ is of order $O((\log n)^\alpha)$, where $\alpha < 4$ [35].

Proof: We can knock down the complexity of an individual call of $B(x_{l+1}|x_1^l, k)$ to a constant if we store the frequencies of substrings tested in formula (12) and we increment them on line. Some further important savings can be done if we know the depth $s = \mathbf{L}(x_1^n)$. The value of s can be computed in time $O(n)$ by building the suffix tree of x_1^n [36]. Once we have that s , let us observe that

$$B(x_{l+1}|x_1^l, k) = B(x_{l+1}|x_1^l, s) \quad (30)$$

holds for all $k > s$.

Thus we can flush all probabilities $P(x_1^n, k)$ for $k > s$ into a dummy variable $P(x_1^n, \bullet)$. In the following, without affecting the value of $P(x_1^n)$, the recursion (15)–(16) can be altered to

$$P(x_1^n, k) = 0 \text{ for } k < -1 \text{ or } k \geq n, \quad (31)$$

$$P(x_1^n, \bullet) = 0 \text{ for } n < s + 1, \quad (32)$$

$$P(x_1^{n+1}, k) = [p_n P(x_1^n, k) + q_n P(x_1^n, k - 1)] B(x_{n+1}|x_1^n, k) \quad \text{for } n \geq 1 \text{ and } -1 \leq k \leq \min(n, s), \quad (33)$$

$$P(x_1^{n+1}, \bullet) = [P(x_1^n, \bullet) + q_n P(x_1^n, s)] B(x_{n+1}|x_1^n, s) \text{ for } n \geq s + 1. \quad (34)$$

The formula for the total probability becomes

$$P(x_1^n) = \sum_{k=-1}^s P(x_1^n, k) + P(x_1^n, \bullet). \quad (35)$$

Hence the time complexity of $P(x_1^n)$ is of order $O(ns)$. \square

The space complexity of the switch distribution can also be reduced by observing that in order to compute $B(x_{l+1}|x_1^l, k)$ we only need to store the frequencies of substrings w that appear in x_1^n at least twice and the frequencies of their extensions wa , where $a \in \{1, 2, \dots, D\}$. These strings can be also found while building the suffix tree of x_1^n .

Parameter s in the algorithm (31)–(35) will be called the depth of the switch distribution. Without a significant change of $P(x_1^n)$, the depth of the switch

distribution can be chosen as much smaller than the depth of string x_1^n . This fact can be also used for the further speed-up of computation. Fixing the depth, however, leads asymptotically to the Q -almost sure bound

$$\lim_{n \rightarrow \infty} \frac{1}{n} [-\log P(X_1^n)] = \lim_{n \rightarrow \infty} \frac{1}{n} [-\log B(X_{s+1}^n | X_1^s, s)] \quad (36)$$

$$= \mathbf{E}_Q [-\log Q(X_{s+1} | X_1^s)] \quad (37)$$

if Q is ergodic. Conditional entropy $\mathbf{E}_Q [-\log Q(X_{s+1} | X_1^s)]$ is greater than h_Q .

III The preadapted switch distribution

Often we want to predict or compress data x_1^n that are generated by a class of complex unknown distributions Q that partly resemble the empirical distribution of another, much larger data y_1^j . Such a case arises in particular in the compression of texts in natural language. Then using a universal distribution such as the plain switch distribution need not be the best approach, since this distribution has to learn all frequencies of substrings from the data x_1^n . A competing approach is to use frequencies of substrings from the larger data y_1^j . This can yield a better compression rate for finite data x_1^n . The problem of using a fixed empirical distribution of y_1^j is, however, that it is not universal. The source of the problem lies in using non-adaptive substring frequencies. A simple solution for this problem is to initialize the substring frequencies with the frequencies coming from y_1^j and let them gradually adapt to x_1^n . In this way we obtain a preadapted universal compression scheme. One can suppose that this scheme may compress better than both the plain switch distribution and the empirical distribution of y_1^j .

Let us clarify this idea.

Definition 2 (fixed switch distribution) Let y_1^j be a fixed sequence, called the training data. Define conditional probabilities $B(x_{n+1} | x_1^n, -1) = D^{-1}$ and

$$B(x_{n+1} | x_1^n, k) = \frac{c(x_{n+1-k}^{n+1} | y_1^j) + B(x_{n+1} | x_1^n, k-1)}{c(x_{n+1-k}^n | y_1^{j-1}) + 1}. \quad (38)$$

Using these $B(x_{n+1} | x_1^n, k)$, we define the fixed switch distribution P via formulae (13)–(17).

For short blocks x_1^n , the fixed switch distribution can achieve much lower compression rate than the plain switch distribution but it is not universal. To obtain a universal distribution which combines the advantages of the fixed switch distribution and the plain switch distribution, we may consider a compromise between expressions (12) and (38). This can be done easily as follows.

Definition 3 (preadapted switch distribution) Let y_1^j be a fixed sequence, called the training data. Define conditional probabilities $B(x_{n+1} | x_1^n, -1) = D^{-1}$ and

$$B(x_{n+1} | x_1^n, k) = \frac{c(x_{n+1-k}^{n+1} | y_1^j x_1^n) + B(x_{n+1} | x_1^n, k-1)}{c(x_{n+1-k}^n | y_1^j x_1^{n-1}) + 1}. \quad (39)$$

Using these $B(x_{n+1} | x_1^n, k)$, we define the preadapted switch distribution P via formulae (13)–(17).

As in the plain case, we can show that the preadapted switch distribution is universal and efficiently computable. The proof of universality relies on the observation that the influence of training data y_1^j on the probability of long blocks x_1^n is asymptotically negligible.

Theorem 3 *The preadapted switch distribution is strongly and weakly universal.*

Proof: Analogously to the plain switch distribution, the preadapted switch distribution satisfies the analogue of 2. Having this fact in mind, we can prove the universality. Let Q be a stationary ergodic distribution. Since the alphabet of X_i is finite, by the ergodic theorem differences $B(X_n|X_1^{n-1}, k) - Q(X_n|X_{n-k-1}^{n-1})$ converge to 0 Q -almost surely. The further reasoning proceeds like the proof of Theorem 1. \square

Theorem 4 *The value of the preadapted switch distribution $P(x_1^n)$ can be computed in time $O((j+n)s)$ where $s = \mathbf{L}(y_1^j x_1^n)$.*

Proof: The complexity of an individual call of $B(x_{l+1}|x_1^l, k)$ can be reduced to a constant if we record the frequencies of substrings tested in formula (39) and we increment them on line. Initializing these frequencies takes time $O(js)$. Let us also observe that

$$B(x_{l+1}|x_1^l, k) = B(x_{l+1}|x_1^l, s) \quad (40)$$

holds for all $k > s$. Thus without affecting the value of $P(x_1^n)$, the algorithm (15)–(16) can be changed to (31)–(34) and the formula for the total probability becomes (35). Thus the time complexity of $P(x_1^n)$ is of order $O((j+n)s)$. \square

The space complexity of the preadapted switch distribution can also be reduced by noticing that in order to compute $B(x_{l+1}|x_1^l, k)$ we only have to record the frequencies of substrings w that appear in $y_1^j x_1^n$ at least twice and the frequencies of their extensions wa , where $a \in \{1, 2, \dots, D\}$.

IV Measuring codewise Hilberg exponent γ

Here we describe a simple experiment that we have performed using the three switch distributions and the Lempel-Ziv code. As the training data we have taken *The Complete Memoirs* by J. Casanova (6,719,801 characters), and as the compressed text—*Gulliver's Travels* by J. Swift (579,438 characters). Both texts were downloaded from the Project Gutenberg.² The alphabet size was set as $D = 256$. The switch distributions were computed using transition probabilities p_n of form (18) with $\alpha = 1.001$ since we observed that the lower the α is the better compression is achieved. Moreover, we have used algorithm (31)–(35) with fixed depth $s = 7$ since more than 99.99% of the probability mass in the observed cases concentrated in $P(x_1^n, k)$ with $k \leq 4$. Hence it was a safe approximation. The Lempel-Ziv code was computed by our own implementation for the ASCII encoding of the text. The results are presented in Tables 1 and 2 and Figures 2 and 3.

²<http://www.gutenberg.org/>

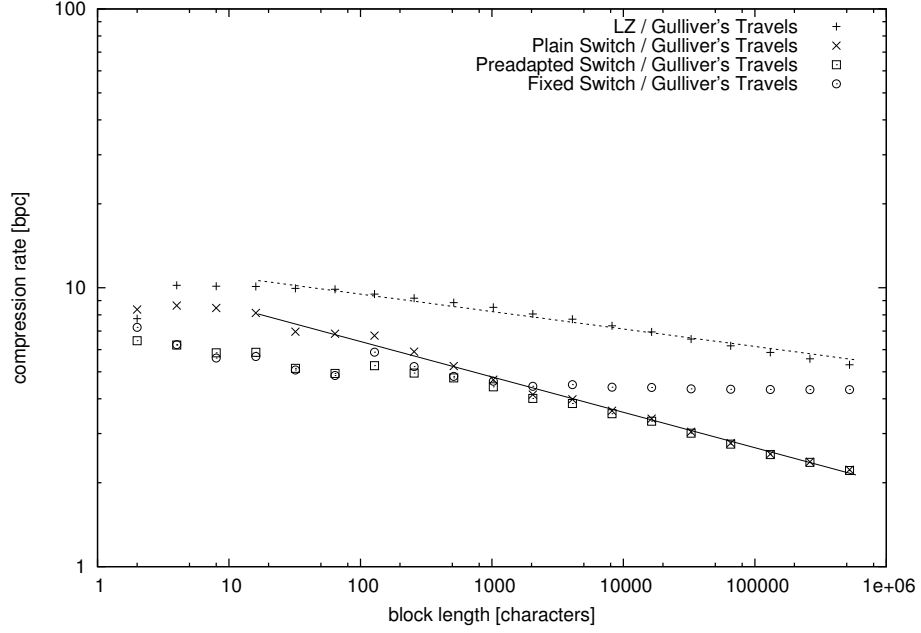


Figure 2: Compression rates for the switch distributions and the LZ code. The solid line is the least square regression $y = 11.51n^{-0.127}$, computed for the plain switch distribution. The dotted line is the least square regression $y = 12.66n^{-0.0625}$, computed for the LZ code.

n	$H^P(x_1^n)/n$ [bpc]			
	LZ	plain switch	preadapted switch	fixed switch
2	7.7459	8.3547	6.4605	7.2124
4	10.2089	8.6367	6.2363	6.2506
8	10.1347	8.4657	5.847	5.5963
16	10.1122	8.1227	5.8676	5.6687
32	9.9482	6.9604	5.1395	5.0712
64	9.8816	6.8478	4.9319	4.8489
128	9.4894	6.7355	5.254	5.8753
256	9.1817	5.9023	4.9481	5.2167
512	8.8427	5.2381	4.7506	4.8141
1024	8.5069	4.6831	4.414	4.626
2048	8.0525	4.1411	4.0127	4.4325
4096	7.7158	3.9809	3.8476	4.4953
8192	7.3084	3.6209	3.5361	4.4023
16384	6.9471	3.3941	3.3238	4.3935
32768	6.5467	3.0459	3.0114	4.3422
65536	6.1909	2.7745	2.7504	4.327
131072	5.865	2.5342	2.5223	4.3188
262144	5.5665	2.3759	2.3664	4.3142
524288	5.2928	2.2252	2.213	4.3126

Table 1: Compression rates for the switch distributions and the LZ code.

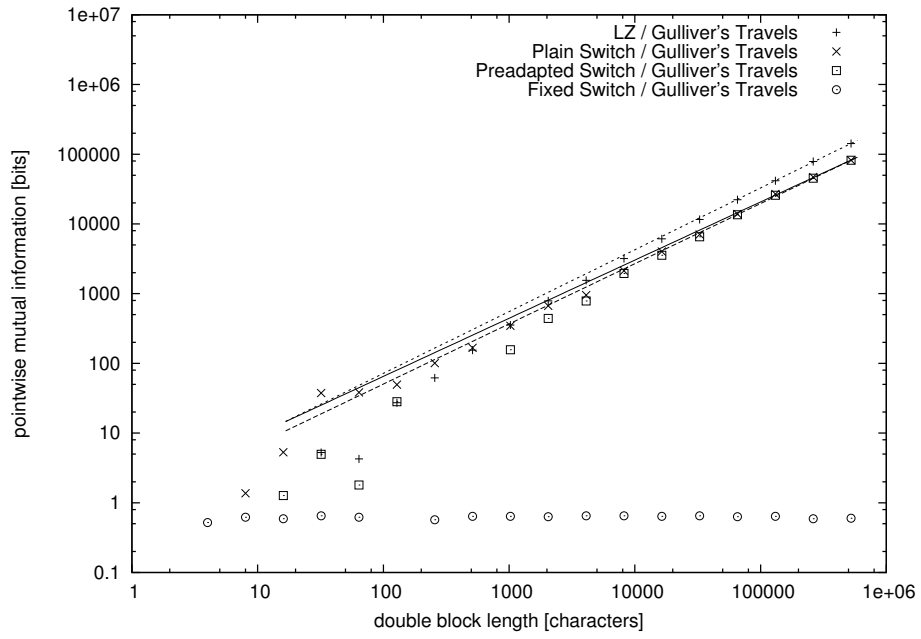


Figure 3: Pointwise mutual information for the switch distributions and the LZ code. The solid line is the least square regression $y = 1.395n^{0.834}$, computed for the plain switch distribution. The dashed line is the least square regression $y = 0.946n^{0.863}$, computed for the preadapted switch distribution. The dotted line is the least square regression $y = 1.209n^{0.887}$, computed for the LZ code.

n	$I^P(x_1^{n/2}; x_{n/2+1}^n)$ [bits]			
	LZ	plain switch	preadapted switch	fixed switch
2	-1.32	-0.71	-0.64	-2.14
4	-6.04	-1.13	-0.67	0.52
8	-3.99	1.37	-0.5	0.62
16	-4.96	5.3	1.27	0.59
32	5.25	37.5	4.96	0.65
64	4.26	38.26	1.8	0.62
128	27.27	49.6	28.22	-2.47
256	61.92	100.78	-3.89	0.57
512	155.1	166.76	-29.5	0.64
1024	353.89	345.54	156.62	0.64
2048	789.56	668.01	441.52	0.63
4096	1554.31	954.54	786.68	0.65
8192	3187.28	2128.53	1945.46	0.65
16384	6119.95	4017.21	3558.77	0.64
32768	11608.28	7062.68	6551.43	0.65
65536	22241.83	13877.79	13549.91	0.63
131072	41621.75	26852.4	25727.25	0.64
262144	78530.25	47113.91	45313.69	0.59
524288	142330.87	81859.85	81873.95	0.6

Table 2: Pointwise mutual information for the switch distributions and the LZ code.

In Figure 2 and Table 1, the quality of compression can be compared for the particular distributions. Among the universal schemes, the best compression is given by the preadapted switch distribution followed by the plain switch distribution followed by the LZ code. However, our hope that the preadapted switch distribution will significantly beat the plain switch distribution has not been fully confirmed. Indeed for short blocks the preadapted switch distribution mimics the behavior of the fixed switch distribution and performs much better than the plain switch distribution. Alas, for long blocks the difference between the two universal switch distributions becomes negligible. Ultimately, both universal switch distributions compress the text twice better than the LZ code. For no universal code we can observe the ultimate stabilization of the compression rate. On the other hand, the fixed switch distribution, which is not universal, stabilizes ultimately at the constant rate of 4.31 bpc.

The stabilization of the fixed switch distribution is clearly visible in Figure 3 and Table 2, which concern the pointwise mutual information. Namely, we can see that pointwise mutual information for the fixed switch distribution does not grow, whereas for the other distributions, which are universal, the pointwise mutual information grows rather fast. The tightest bound for the pointwise mutual information is obtained in the case of the plain switch distribution, which gives the exponent $\gamma = 0.83$ for the codewise Hilberg conjecture (7). It is surprising that the pointwise mutual information for the two other universal distributions grows almost at the same rate, despite the large difference of compression rates between the universal switch distributions and the LZ code.

As we have indicated in the Introduction, observing $\gamma > 0$ for the switch

distribution, we have to reject the hypothesis that the source is a Markov chain. Possibly, Hilberg’s conjecture (3) may be true but we need still some stronger evidence for this hypothesis such as a lower bound for the exponent β .

Acknowledgment

The author wishes to thank Jacek Koronacki and Jan Miłniczuk for a discussion and an anonymous referee for suggesting some relevant references.

References

- [1] W. Hilberg, “Der bekannte Grenzwert der redundanzfreien Information in Texten — eine Fehlinterpretation der Shannonschen Experimente?” *Frequenz*, vol. 44, pp. 243–248, 1990.
- [2] W. Ebeling and G. Nicolis, “Entropy of symbolic sequences: the role of correlations,” *Europhys. Lett.*, vol. 14, pp. 191–196, 1991.
- [3] —, “Word frequency and entropy of symbolic sequences: a dynamical perspective,” *Chaos Sol. Fract.*, vol. 2, pp. 635–650, 1992.
- [4] W. Ebeling and T. Pöschel, “Entropy and long-range correlations in literary English,” *Europhys. Lett.*, vol. 26, pp. 241–246, 1994.
- [5] W. Bialek, I. Nemenman, and N. Tishby, “Predictability, complexity and learning,” *Neural Comput.*, vol. 13, p. 2409, 2001.
- [6] —, “Complexity through nonextensivity,” *Physica A*, vol. 302, pp. 89–99, 2001.
- [7] J. P. Crutchfield and D. P. Feldman, “Regularities unseen, randomness observed: The entropy convergence hierarchy,” *Chaos*, vol. 15, pp. 25–54, 2003.
- [8] Ł. Dębowski, “On the vocabulary of grammar-based codes and the logical consistency of texts,” *IEEE Trans. Inform. Theory*, vol. 57, pp. 4589–4599, 2011.
- [9] —, “Mixing, ergodic, and nonergodic processes with rapidly growing information between blocks,” *IEEE Trans. Inform. Theory*, vol. 58, pp. 3392–3401, 2012.
- [10] —, “On Hilberg’s law and its links with Guiraud’s law,” *J. Quantit. Linguist.*, vol. 13, pp. 81–109, 2006.
- [11] G. Herdan, *Quantitative Linguistics*. Butterworths, 1964.
- [12] H. S. Heaps, *Information Retrieval—Computational and Theoretical Aspects*. Academic Press, 1978.
- [13] C. Shannon, “Prediction and entropy of printed English,” *Bell Syst. Tech. J.*, vol. 30, pp. 50–64, 1951.

- [14] Ł. Dębowski, “Empirical evidence for Hilberg’s conjecture in single-author texts,” in *Methods and Applications of Quantitative Linguistics—Selected papers of the 8th International Conference on Quantitative Linguistics (QUALICO)*, I. Obradović, E. Kelih, and R. Köhler, Eds. Belgrade: Academic Mind, 2013, pp. 143–151.
- [15] J. Ziv and A. Lempel, “A universal algorithm for sequential data compression,” *IEEE Trans. Inform. Theory*, vol. 23, pp. 337–343, 1977.
- [16] B. Ryabko, “Applications of universal source coding to statistical analysis of time series,” in *Selected Topics in Information and Coding Theory*, ser. Series on Coding and Cryptology, I. Woungang, S. Misra, and S. C. Misra, Eds. World Scientific Publishing, 2010.
- [17] R. M. Gray and L. D. Davisson, “The ergodic decomposition of stationary discrete random processes,” *IEEE Trans. Inform. Theory*, vol. 20, pp. 625–636, 1974.
- [18] O. Kallenberg, *Foundations of Modern Probability*. Springer, 1997.
- [19] T. Weissman, “Not all universal source codes are pointwise universal,” 2004, <http://web.stanford.edu/~tsachy/pdf/files/Not%20All%20Universal%20Source%20Codes%20are%20Pointwise%20Universal.pdf>.
- [20] P. F. Brown, S. A. D. Pietra, V. J. D. Pietra, J. C. Lai, and R. L. Mercer, “An estimate of an upper bound for the entropy of English,” *Comput. Linguist.*, vol. 18, pp. 31–40, 1992.
- [21] F. Jelinek, *Statistical Methods for Speech Recognition*. The MIT Press, 1997.
- [22] C. D. Manning and H. Schütze, *Foundations of Statistical Natural Language Processing*. The MIT Press, 1999.
- [23] J. G. Cleary and I. H. Witten, “Data compression using adaptive coding and partial string matching,” *IEEE Trans. Comm.*, vol. 32, pp. 396–402, 1984.
- [24] F. M. J. Willems, Y. M. Shtarkov, and T. J. Tjalkens, “The context tree weighting method: Basic properties,” *IEEE Trans. Inform. Theory*, vol. 41, pp. 653–664, 1995.
- [25] D. Ron, Y. Singer, and N. Tishby, “The power of amnesia: Learning probabilistic automata with variable memory length,” *Machine Learn.*, vol. 25, pp. 117–149, 1996.
- [26] J. C. Kieffer and E. Yang, “Grammar-based codes: A new class of universal lossless source codes,” *IEEE Trans. Inform. Theory*, vol. 46, pp. 737–754, 2000.
- [27] D. S. an Giovanni Motta, *Handbook of Data Compression*. Springer, 2009.
- [28] T. van Erven, P. Grünwald, and S. de Rooij, “Catching up faster in Bayesian model selection and model averaging,” in *Advances in Neural Information Processing Systems 20 (NIPS 2007)*, 2007.

- [29] W. M. Koolen and S. de Rooij, “Universal codes from switching strategies,” *IEEE Trans. Inform. Theory*, vol. 59, pp. 7168–7185, 2013.
- [30] D. Hindle and M. Rooth, “Structural ambiguity and lexical relations,” *Comput. Linguist.*, vol. 19, pp. 103–120, 1993.
- [31] G. Louchard and W. Szpankowski, “On the average redundancy rate of the Lempel-Ziv code,” *IEEE Trans. Inform. Theory*, vol. 43, pp. 2–8, 1997.
- [32] P. D. Grünwald, *The Minimum Description Length Principle*. The MIT Press, 2007.
- [33] K. Atteson, “The asymptotic redundancy of Bayes rules for Markov chains,” *IEEE Trans. Inform. Theory*, vol. 45, pp. 2104–2109, 1999.
- [34] P. C. Shields, “String matching bounds via coding,” *Ann. Probab.*, vol. 25, pp. 329–336, 1997.
- [35] Ł. Dębowski, “Maximal lengths of repeat in English prose,” in *Synergetic Linguistics. Text and Language as Dynamic System*, S. Naumann, P. Grzybek, R. Vulcanović, and G. Altmann, Eds. Wien: Praesens Verlag, 2012, pp. 23–30.
- [36] E. Ukkonen, “On-line construction of suffix trees,” *Algorithmica*, vol. 14, pp. 249–260, 1995.