

# Information Theory and Statistics

## Lecture 4: Lempel-Ziv code

Łukasz Dębowski  
ldebowsk@ipipan.waw.pl

Ph. D. Programme 2013/2014

# Entropy rate is the limiting compression rate

## Theorem

For a stationary process  $(\mathbf{X}_i)_{i=-\infty}^{\infty}$ , let  $L_n$  denote the minimal expected compression rate of a uniquely decodable code  $\mathbf{B}_n : \mathbb{X}^n \rightarrow \{0, 1\}^*$  for the block of  $n$  variables. That is,

$$L_n := \min_{\mathbf{B}_n} \frac{1}{n} \mathbf{E} |\mathbf{B}_n(\mathbf{X}_1, \dots, \mathbf{X}_n)|.$$

We claim that  $\lim_{n \rightarrow \infty} L_n = h$ .

## Proof

Assume that  $\mathbf{B}_n$  is the Shannon-Fano code for the block  $(\mathbf{X}_1, \dots, \mathbf{X}_n)$ . Then  $\mathbf{H}(\mathbf{X}_1^n) \leq nL_n \leq \mathbf{E} |\mathbf{B}_n(\mathbf{X}_1, \dots, \mathbf{X}_n)| \leq \mathbf{H}(\mathbf{X}_1^n) + 1$ . Hence the claim follows.

# The problem of universal compression

- To compute the Shannon-Fano code we need to know the probability distribution of the block.
- Such a situation is unlikely in practical applications of data compression, where we have no prior information about the probability distribution of blocks.
- Fortunately, as an important corollary of the ergodic theorem, there exist universal codes whose compression rates tend to the entropy rate for any stationary process.

# Universal codes

## Definition (weakly universal code)

A uniquely decodable code  $\mathbf{B} : \mathbb{X}^* \rightarrow \{0, 1\}^*$  is called *weakly universal* if for any stationary process (not necessarily ergodic) we have

$$\lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{E} |\mathbf{B}(\mathbf{X}_1^n)| = h.$$

## Definition (strongly universal code)

A uniquely decodable code  $\mathbf{B} : \mathbb{X}^* \rightarrow \{0, 1\}^*$  is called *strongly universal* if for any stationary ergodic process inequality

$$\limsup_{n \rightarrow \infty} \frac{1}{n} |\mathbf{B}(\mathbf{X}_1^n)| \leq h.$$

holds with probability 1.

# Strongly universal codes are better

## Theorem

Let code  $\mathbf{B}$  be strongly universal. If there exists a constant  $\mathbf{K}$  such that

$$|\mathbf{B}(x_1^n)| \leq \mathbf{K}n$$

for each string  $x_1^n$  then code  $\mathbf{B}$  is weakly universal.

# Compression and statistics

The problem of universal compression falls under the scope of statistics. Indeed, the interest of statisticians lies in identifying parameters of a stochastic process basing on the data typical for that process. Entropy rate of an ergodic process is an example of such a parameter. When we have a universal code, we may estimate the entropy rate as the compression rate.

# Lempel-Ziv code

The code was derived by Abraham Lempel (1936–) and Jacob Ziv (1931–) in 1977 and is partly implemented in `gzip` and `compress`.

## Definition (LZ code)

For simplicity of the algorithm description we assume that the compressed data are binary sequences, that is  $\mathbb{X} = \{0, 1\}$ . The Lempel-Ziv compression algorithm is as follows.

- 1 The compressed sequence is parsed into a sequence of shortest phrases that have not appeared before (except for the last phrase). For example, the sequence **001010010011100...** is split into phrases **0, 01, 010, 0100, 1, 11, 00, ...**
- 2 In the following, each phrase is described using a binary index of the longest prefix that appeared earlier and a single bit that follows that prefix. For the considered sequence, this representation is as follows: **(0, 0)(1, 1)(10, 0)(11, 0)(0, 1)(101, 1)(1, 0)**.

# The length of the LZ code

- Let  $C_n$  be the number of phrases in the compressed block  $\mathbf{X}_1^n$ . If we know  $C_n$ , we need  $\log C_n$  bits to identify the prefix index for each phrase and 1 bit to describe the following bit. Thus the LZ code uses  $|\mathbf{B}(\mathbf{X}_1^n)| = C_n [\log C_n + O(1)]$  bits in total.
- A splitting of a sequence into distinct phrases will be called a *distinct parsing* of the sequence.

## Theorem

Let  $(\mathbf{X}_i)_{i=-\infty}^{\infty}$  be a stationary ergodic process and let  $C_n$  be the number of phrases in a distinct parsing of block  $(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)$ . With probability 1 we have

$$\limsup_{n \rightarrow \infty} \frac{C_n [\log C_n + O(1)]}{n} \leq h.$$

*Remark:* Hence the LZ code is strongly universal. It can be also shown that the LZ code is weakly universal.



# The first lemma

## Lemma

*The number of phrases  $C_n$  in any distinct parsing of block  $(X_1, X_2, \dots, X_n)$  satisfies inequality*

$$\lim_{n \rightarrow \infty} \frac{C_n \log n}{n} \leq 1.$$

# Proof of the first lemma

Let  $n_k = \sum_{j=1}^k j2^j = (k-1)2^{k+1} + 2$  be the sum of lengths of distinct phrases that are not longer than  $k$ . The number of phrases  $C_n$  in a distinct parsing will be maximal if the phrases are as short as possible. For  $n_k \leq n < n_{k+1}$  this happens if we take all phrases of length  $\leq k$  and  $\delta/(k+1)$  phrases of length  $k+1$ , where  $\delta = n - n_k$ . Then

$$C_n \leq \sum_{j=1}^k 2^j + \frac{\delta}{k+1} \leq \frac{n_k}{k-1} + \frac{\delta}{k+1} \leq \frac{n}{k-1}.$$

In the following we will provide a bound for  $k$  given  $n$ . We have  $n \geq n_k = (k-1)2^{k+1} + 2 \geq 2^k$ , so  $k \leq \log n$ . Moreover  $n < n_{k+1} = k2^{k+2} + 2 \leq (\log n + 2)2^{k+2}$ . Hence

$$k + 2 > \log \frac{n}{\log n + 2}.$$

Further transformations yield  $k - 1 > \log n - \log(\log n + 2) - 3$ . Hence we obtain the claim.

# Ziv inequality

Let  $\mathbf{P}^k$  denote the measure of the  $k$ -th order Markov approximation of the process  $(\mathbf{X}_i)_{i=-\infty}^{\infty}$ . That is

$$\mathbf{P}^k(\mathbf{X}_{-k+1}^n | \mathbf{X}_{-k+1}^0) := \prod_{i=1}^n \mathbf{P}(\mathbf{X}_i | \mathbf{X}_{i-k}^{i-1}).$$

Moreover, assume that sequence  $(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)$  is parsed into  $\mathbf{C}_n$  distinct phrases  $(\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_{\mathbf{C}_n})$ . Let  $\mathbf{W}_i$  denote the  $k$  bits preceding  $\mathbf{Y}_i$ . Next, let  $\mathbf{C}_n^{l,w}$  denote the number of phrases  $\mathbf{Y}_i$  that have length  $l$  and context  $\mathbf{W}_i = \mathbf{w}$ .

## Lemma (Ziv inequality)

*We have inequality*

$$-\log \mathbf{P}^k(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n | \mathbf{W}_1) \geq \sum_{l,w} \mathbf{C}_n^{l,w} \log \mathbf{C}_n^{l,w}.$$

# Proof of Ziv inequality

## Proof

Observe that

$$\begin{aligned} -\log P^k(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n | \mathbf{W}_1) &= -\sum_{j=1}^{C_n} \log P(\mathbf{Y}_j | \mathbf{W}_j) \\ &= -\sum_{l,w} C_n^{lw} \cdot \frac{1}{C_n^{lw}} \sum_{j: |\mathbf{Y}_j|=l, \mathbf{W}_j=w} \log P^k(\mathbf{Y}_j | \mathbf{W}_j) \\ &\geq -\sum_{l,w} C_n^{lw} \log \left( \frac{1}{C_n^{lw}} \sum_{j: |\mathbf{Y}_j|=l, \mathbf{W}_j=w} P^k(\mathbf{Y}_j | \mathbf{W}_j) \right), \end{aligned}$$

where the inequality follows from the Jensen inequality because the logarithm function is concave. Because the phrases  $\mathbf{Y}_j$  under the sum are distinct, we have  $\sum_{j: |\mathbf{Y}_j|=l, \mathbf{W}_j=w} P^k(\mathbf{Y}_j | \mathbf{W}_j) \leq 1$ . Hence the claim follows.

# Third lemma

## Lemma

Let  $L$  be a nonnegative random variable taking values in integers and having expectation  $EL$ . Then entropy  $H(L)$  is bounded by inequality

$$H(L) \leq (EL + 1) \log (EL + 1) - EL \log EL.$$

The proof of this lemma will be discussed after the lecture on maximum entropy modeling as an easy exercise.

# Proof that LZ code is universal

Let  $\mathbf{L}$  and  $\mathbf{W}$  be random variables such that

$$P(\mathbf{L} = l, \mathbf{W} = w) = \frac{C_n^{lw}}{C_n}.$$

The expectation of  $\mathbf{L}$  is

$$E \mathbf{L} = \sum_{l,w} \frac{l C_n^{lw}}{C_n} = \frac{n}{C_n}.$$

Hence by the third lemma, we obtain

$$\begin{aligned} H(\mathbf{L}) &\leq (E \mathbf{L} + 1) \log (E \mathbf{L} + 1) - E \mathbf{L} \log E \mathbf{L} \\ &= \log \frac{n}{C_n} + \left( \frac{n}{C_n} + 1 \right) \log \left( \frac{C_n}{n} + 1 \right). \end{aligned}$$

# Proof that LZ code is universal (continued)

On the other hand,  $H(W) \leq k$ , so

$$\begin{aligned} H(L, W) &\leq H(L) + H(W) \\ &\leq \log \frac{n}{C_n} + \left( \frac{n}{C_n} + 1 \right) \log \left( \frac{C_n}{n} + 1 \right) + k. \end{aligned}$$

Then by the first lemma, we have

$$\lim_{n \rightarrow \infty} \frac{C_n}{n} H(L, W) = 0.$$

# Proof that LZ code is universal (finished)

Now using the first lemma again, the Ziv inequality, and the ergodic theorem, we obtain

$$\begin{aligned}
 \limsup_{n \rightarrow \infty} \frac{C_n [\log C_n + O(1)]}{n} &= \limsup_{n \rightarrow \infty} \left( \frac{C_n \log C_n}{n} - \frac{C_n}{n} H(L, W) \right) \\
 &= \limsup_{n \rightarrow \infty} \frac{1}{n} \left( C_n \log C_n + C_n \sum_{l,w} \frac{C_n^{lw}}{C_n} \log \frac{C_n^{lw}}{C_n} \right) \\
 &= \limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{l,w} C_n^{lw} \log C_n^{lw} \leq - \lim_{n \rightarrow \infty} \frac{1}{n} \log P^k(X_1^n | X_{-k+1}^0) \\
 &= - \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \log P(X_i | X_{i-k}^{i-1}) = H(X_i | X_{i-k}^{i-1}).
 \end{aligned}$$

with probability 1. This inequality holds for any  $k$ . Considering  $k \rightarrow \infty$ , we obtain the claim.



# Motivation for the R measure

- We can estimate the entropy rate of a stationary process as the length of the Lempel-Ziv code for a sequence of symbols drawn for the process divided by the sequence length.
- This method of estimation is far from satisfactory since the estimate of the entropy rate converges very slow as a function of the sequence length.
- We will present a method which seems to be better in case of empirical sources such as the natural language.

# R measure

Let the frequency of substring  $\mathbf{w}_1^k$  in string  $\mathbf{z}_1^n \in \{0, 1, \dots, D-1\}^n$  be

$$c(\mathbf{w}_1^k | \mathbf{z}_1^n) = \sum_{i=0}^{n-k} \mathbf{1}\{\mathbf{w}_1^k = \mathbf{z}_{i+1}^{i+k}\}.$$

## Definition (R measure)

Define conditional probabilities  $\mathbf{B}(\mathbf{x}_{n+1} | \mathbf{x}_1^n, -1) = D^{-1}$  and

$$\mathbf{B}(\mathbf{x}_{n+1} | \mathbf{x}_1^n, k) = \frac{c(\mathbf{x}_{n+1-k}^{n+1} | \mathbf{x}_1^n) + \mathbf{B}(\mathbf{x}_{n+1} | \mathbf{x}_1^n, k-1)}{c(\mathbf{x}_{n+1-k}^n | \mathbf{x}_1^{n-1}) + 1}.$$

We write  $\mathbf{B}(\mathbf{x}_1^n, k) = \prod_{i=1}^n \mathbf{B}(x_i | x_1^{i-1}, k)$ . Let  $\mathbf{p}_k \in (0, 1)$  satisfy  $\sum_{k=-1}^{\infty} \mathbf{p}_k = 1$ . The R measure is

$$\mathbf{Q}(\mathbf{x}_1^n) = \sum_{k=-1}^{\infty} \mathbf{p}_k \mathbf{B}(\mathbf{x}_1^n, k).$$

# Universality of R measure

A probability distribution  $Q$  is called *weakly universal* if for any stationary process (not necessarily ergodic) we have

$$\lim_{n \rightarrow \infty} \frac{1}{n} E [-\log Q(X_1^n)] = h.$$

A probability distribution  $Q$  is called *strongly universal* if for any stationary ergodic process inequality

$$\limsup_{n \rightarrow \infty} \frac{1}{n} [-\log Q(X_1^n)] \leq h.$$

holds with probability 1.

## Theorem

*The R measure is both strongly and weakly universal.*

# Proof of universality

Let  $\mathbf{P}$  be a stationary ergodic distribution. Since the alphabet of  $\mathbf{X}_i$  is finite, by the ergodic theorem differences  $\mathbf{B}(\mathbf{X}_n|\mathbf{X}_1^{n-1}, k) - \mathbf{P}(\mathbf{X}_n|\mathbf{X}_{n-k-1}^{n-1})$  converge to 0 with  $\mathbf{P}$ -probability 1. Hence

$$\lim_{n \rightarrow \infty} \frac{1}{n} [-\log \mathbf{B}(\mathbf{X}_1^n, k)] = \lim_{n \rightarrow \infty} \frac{1}{n} \left[ - \sum_{i=k+1}^n \log \mathbf{P}(\mathbf{X}_i | \mathbf{X}_{i-k-1}^{i-1}) \right].$$

Applying the ergodic theorem again, we obtain

$$\lim_{n \rightarrow \infty} \frac{1}{n} \left[ - \sum_{i=k+1}^n \log \mathbf{P}(\mathbf{X}_i | \mathbf{X}_{i-k-1}^{i-1}) \right] = \mathbf{E} \left[ -\log \mathbf{P}(\mathbf{X}_{k+1} | \mathbf{X}_1^k) \right].$$

Hence

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{1}{n} [-\log \mathbf{Q}(\mathbf{X}_1^n)] &\leq \inf_{k \in \mathbb{N}} \lim_{n \rightarrow \infty} \frac{1}{n} [-\log \mathbf{B}(\mathbf{X}_1^n, k)] \\ &= \inf_{k \in \mathbb{N}} \mathbf{E} \left[ -\log \mathbf{P}(\mathbf{X}_{k+1} | \mathbf{X}_1^k) \right] = h. \end{aligned}$$

Hence the distribution  $\mathbf{Q}$  is strongly universal. Since

$-\log \mathbf{Q}(\mathbf{X}_1^n) \leq -\log \mathbf{p}_{-1} + n \log \mathbf{D}$ , distribution  $\mathbf{Q}$  is also weakly universal.

# The R measure is effectively computable

Denote the maximal length of a substring that appears at least twice in  $z_1^n$  as

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

For  $k > L(x_1^n)$ ,

$$B(x_1^n, k) = B(x_1^n, k - 1).$$

Hence the R measure is

$$Q(x_1^n) = \sum_{k=-1}^{L(x_1^n)} p_k B(x_1^n, k) + \left( 1 - \sum_{k=-1}^{L(x_1^n)} p_k \right) B(x_1^n, L(x_1^n)).$$