Lempel-Ziv Coding

Prof. Ja-Ling Wu

Department of Computer Science and Information Engineering National Taiwan University

J.Ziv and A.Lempel, "Compression of individual sequences by variable rate coding, "IEEE Trans. Information Theory, 1978.

Algorithm:

- The source sequence is sequentially parsed into strings that have not appeared so far. For example, 1011010100010..., will be parsed as 1,0,11,0 1,010,00,10,...
- After every comma, we look along the input sequence until we come to the shortest string that has not been marked off before. Since this is the shortest such string, all its prefixes must have occurred easier. In particular, the string consisting of all but the last bit of this string must have occurred earlier.

We code this phrase by giving the location of the prefix and the value of the last bit.

Let C(n) be the number of phrases in the parsing of the input n-sequence. We need log C(n) bits to describe the location of the prefix to the phrase and 1 bit to describe the last bit.

In the above example, the code sequence is: (000,1), (000,0), (001,1), (010,1), (100,0), (010,0), (001,0) where the first number of each pair gives the index of the prefix and the 2nd number gives the last bit of the phrase.

Decoding the coded sequence is straight forward and we can recover the source sequence without error!! Two passes-algorithm:

pass-1: parse the string and calculate C(n) Dictionary construction

pass-2: calculate the points and produce the code string

Pattern Matching

- The two-passes algorithm allots an equal number of bits to all the pointers. This is not necessary, since the range of the pointers is smaller at the initial portion of the string.
 - T-A. Welch, A technique for high-performance data Compression, Computer, 1984.
 - Bell, Cleary, Witlen, Text Compression, Prentice-Hall, 1990.

Definition:

A parsing S of a binary string $x_1, x_2,...,x_n$ is a division of the string into phrases, separated by commas. A distinct parsing is a parsing such that no two phrases are identical.

The L-Z algorithm described above gives a <u>distinct</u> parsing of the source sequence. Let <u>C(n)</u> denote the number of phrases in the L-Z parsing of a sequence of length n.

The compressed sequence (after applying the L-Z algorithm) consists of a list of <u>c(n) pairs</u> of numbers, each pair consisting of a <u>pointer</u> to the previous

log c(n)

occurrence of the prefix of the phrase and the <u>last bit</u> of the phrase.

the total length of the compressed sequence is : $C(n)(\log c(n)+1)$

we will show that

$$\frac{C(n)(\log C(n)+1)}{n} \to H(\chi)$$

for a stationary ergodic sequence $X_1, X_2, ..., X_n$.

Lemma 1 (Lempel and Ziv):

The number of phrase c(n) in a distinct parsing of a binary sequence X1,X2,...,Xn satisfies

$$C(n) \le \frac{n}{(1 - \varepsilon_n) \log n}$$

where $\varepsilon_n \rightarrow 0$ as $n \rightarrow \infty$.

Let
$$n_k = \sum_{j=1}^k j \cdot 2^j = (k-1)2^{k+1} + 2$$

- : the sum of the lengths of all distinct strings of length less than or equal to k.
- C(n): the no. of phrases in a distinct parsing of a sequence of length n this number is maximized when all the phrases are as short as possible.

 $n=n_k \rightarrow all$ the phrases are of length $\leq k$, thus

$$C(n) \le \sum_{j=1}^{k} 2^{j} = 2^{K+1} - 2 \le 2^{k+1} = \frac{n_k - 2}{k - 1} \le \frac{n_k}{k - 1}$$

If $n_k \le n < n_{k+1}$, we write $n = n_k + \Delta$, then $\Delta = n - n_k < n_{k+1} - n_k = (k+1)2^{k+1}$

Then the parsing into shortest phrases has each of the phases of length \leq k and $\stackrel{\Delta}{\longrightarrow}$ phrases of the length k+1.

$$C(n) = C(n_k) + \frac{\Delta}{k+1} \le \frac{n_k}{K-1} + \frac{\Delta}{k+1} \le \frac{n_k + \Delta}{k-1} = \frac{n}{K-1}$$

$$\Rightarrow C(n) \le \frac{n}{K-1}$$

We now bound the size of k for a given n.

Let
$$n_k \le n < n_{k+1}$$
. Then

 $n \ge n_k = (k-1)2^{k+1} + 2 \ge 2^k$

→ $k \le \log n$
 $n \le n_{k+1} = k \cdot 2^{k+2} + 2 \le (k+2)2^{k+2}$
 $\le [(\log n+2) \cdot 2^{k+2}]$

⇒ $k+2 \ge \log \frac{n}{\log n+2}$

for all $n \ge 4$

rall
$$n \ge 4$$

 $k-1 = (k+2)-3 \ge \log n - \log(\log n + 2) - 3$
 $= \left(1 - \frac{\log(\log n + 2) + 3}{\log n}\right) \log n$
 $\ge \left(1 - \frac{\log(2\log n) + 3}{\log n}\right) \log n$
 $= \left(1 - \frac{\log(\log n) + 4}{\log n}\right) \log n$
 $= (1 - \varepsilon_n) \log n$
 $\Rightarrow C(n) \le \frac{n}{K-1} = \frac{n}{(1 - \varepsilon_n) \log n}$
where $\varepsilon_n = \frac{\log(\log n) + 4}{\log n} \to 0$ as $n \to \infty$

Lemma 2:

Let Z be a positive integer valued r.v. with mean μ . Then the entropy H(z) is bound by H(z) \leq (μ +1) log(μ +1) - μ log μ pf : H W.

Let $\{X_i\}_{i=-\infty}^{\infty}$ be a stationary ergodic process with pmf $p(x_1,x_2,...x_n)$. For fixed integer k, defined the k^{th} order Markov approximation to P as

$$Q_{k}\left(x_{-(K-1)},...,x_{0},x_{1},...,x_{n}\right)\underline{\underline{\Delta}}P\left(x_{-(k-1)}^{0}\right)\prod_{j=1}^{n}P\left(x_{j}|x_{j-k}^{j-1}\right)$$

where $x_i^j \underline{\underline{\Delta}}(x_i, x_{i+1}, ..., x_j)$, $i \leq j$, and the initial state $x_{-(k-1)}^0$ will be part of the specification of Q_k .

Since $P(X_n|X_{n-k}^{n-1})$ is itself an ergodic process, we have

$$-\frac{1}{n}\log Q_{K}(X_{1}, X_{2}, ..., X_{n}|X_{-(k-1)}^{0}) = -\frac{1}{n}\sum_{j=1}^{n}\log P(X_{j}|X_{j-k}^{j-1})$$

$$\to -E\log P(X_{j}|X_{j-k}^{j-1})\underline{\Delta}H(X_{j}|X_{j-k}^{j-1})$$

We will bound the rate of the L-Z code by the entropy rate of the k-th order Markov approximation for all k. The entropy rate of the Markov approximation $H(X_j|X_{j-k}^{j-1})$ converges to the entropy rate of the process as $k \to \infty$ and this will prove the result.

S'pose $X_{-(k-1)}^n = x_{-(k-1)}^n$, and s'pose that x_1^n is parsed into C distinct phrases, y_1, y_2, \ldots, y_c . Let ν_i be the index of the start of the i-th phrase, i.e., $y_i = x_{\nu_i}^{\nu_{i+1}-1}$. For each i=1,2,...,c, defines $S_i = x_{\nu_i-k}^{\nu_i-1}$. Thus, S_i is the k bits of x preceding y_i , of course, $S_1 = x_{-(k-1)}^0$

Let C_{ls} be the number of phrases y_i with length l and preceding state S_i =S for l=1,2,... and $s \in X^k$, we then have $\sum_{l,s} C_{ls} = C$

and
$$\sum_{l,s}^{s} lC_{ls} = n$$

Lemma 3: (Ziv's inequality)

For any distinct parsing (in particular, the L-Z parsing) of the string $x_1, x_2, ..., x_n$, we have

$$\log Q_k(x_1, x_2, ..., x_n | s_1) \le -\sum_{l,s} C_{ls} \log C_{ls}$$

Note that the right hand side does not depend on Q_k.

proof: we write

$$Q_{k}(x_{1}, x_{2},..., x_{n}|s_{1}) = Q(y_{1}, y_{2},..., y_{c}|s_{1})$$

$$= \prod_{i=1}^{C} P(y_{i}|s_{i})$$

or
$$\log Q_k(x_1, x_2, ..., x_n | s_1) = \sum_{i=1}^{C} \log P(y_i | s_i)$$

$$= \sum_{l,s} \sum_{i:|y_i|=l, s_i=s} \log P(y_i | s_i)$$

$$= \sum_{l,s} C_{ls} \sum_{\substack{i:|y_i|=l \ s_i=S}} \frac{1}{C_{ls}} \log P(y_i | s_i)$$

$$\leq \sum_{l,s} C_{ls} \log \left(\sum_{\substack{i:|y_i|=l\\|s_i|=S}} \frac{1}{C_{ls}} P(y_i|s_i) \right)$$

Now since the y_i are distinct, we have $\sum_{i:|y_i|=l \atop |s_i|=S} P(y_i|s_i) \le 1$ $\Rightarrow \log Q_K(x_1, x_2, ..., x_n|s_1) \le \sum_{l:s} C_{ls} \log \frac{1}{C_{ls}}$

$$\Rightarrow \log Q_K(x_1, x_2, ..., x_n | s_1) \le \sum_{l,s} C_{ls} \log \frac{1}{C_{ls}}$$

Theorem:

Let $\{X_n\}$ be a stationary ergodic process with entropy rate H(X), and let C(n) be the number of phrases in a distinct parsing of a sample of length n from this process. Then

 $\lim_{n\to\infty} \sup \frac{C(n)\log C(n)}{n} \le H(X)$

with probability 1.

$$\begin{aligned} &\text{pf}: \log Q_k \left(x_1, x_2, ..., x_n \middle| s_1 \right) \leq -\sum_{ls} C_{ls} \log \frac{C_{ls} \cdot C}{C} \\ &= -C \log C - C \sum_{l,s} \frac{C_{ls}}{C} \log \frac{C_{ls}}{C} \\ &\text{writting} \quad \Pi_{ls} = \frac{C_{ls}}{C} \; , \quad \text{we have} \\ &\sum_{l,s} \Pi_{l,s} = 1, \quad \sum_{l,s} l \Pi_{l,s} = \frac{n}{C} \end{aligned}$$

We now define r.v.'s U,V, such that

$$Pr(U=I,V=s) = \prod_{I,s}$$

Thus E(U)=n/c, and

$$\log Q_k(x_1, x_2, ..., x_n \mid s_1) \le CH(U, V) - c \log c$$

$$\Rightarrow -\frac{1}{n}\log Q_K\left(x_1, x_2, ..., x_n \middle| s_1\right) \ge \frac{C}{n}\log C - \frac{C}{n}H\left(U, V\right)$$

Since $H(U,V) \leq H(U) + H(V)^n$

and
$$H(V) \leq \log |X|^k = k$$

From Lemma 2, we have

$$H(U) \le (EU+1)\log(EU+1) - EU\log EU$$

$$= \left(\frac{n}{C}+1\right)\log\left(\frac{n}{C}+1\right) - \frac{n}{C}\log\frac{n}{C}$$

$$= \log\frac{n}{C} + \left(\frac{n}{C}+1\right)\log\left(\frac{C}{n}+1\right)$$

Thus,
$$\frac{C}{n}H(U,V) \le \frac{C}{n}k + \frac{C}{n}\log\frac{n}{C} + O(1)$$
.

For a given n, the maximum of $\frac{C}{n}\log\frac{n}{C}$ is attained for the maximum value of C $\left(\text{for } \frac{C}{n} \leq \frac{1}{e} \right)$. But from Lemma 1,

$$C \le \frac{n}{\log n} (1 + O(1))$$
. Thus

$$\frac{C}{n}\log\frac{n}{C} \le O\left(\frac{\log\log n}{\log n}\right)$$

and therefore $\frac{C}{n}H(U,V) \to 0$ as $n \to \infty$

Therefore,

$$\frac{C(n)\log C(n)}{n} \le \frac{-1}{n}\log Q_k(x_1, x_2, ..., x_n | s_1) + \varepsilon_k(n)$$

where ε k(n) \rightarrow 0 as n $\rightarrow\infty$. Hence, with probability 1,

$$\lim_{n \to \infty} \sup \frac{C(n) \log C(n)}{n} \le \lim_{n \to \infty} -\frac{1}{n} \log Q_k \left(X_1, X_2, ..., X_n \middle| X_{-(k-1)}^0 \right)$$

$$= H\left(X_0 \middle| X_{-1}, ..., X_{-k} \right) \to H(X) \quad as \quad k \to \infty.$$

<u>Theorem</u>: Let $\{X_i\}_{i=-\infty}^{\infty}$ be a stationary ergodic stochastic process. Let I(X1,X2,...Xn) be the L-Z codeword length associated with X1,X2,...,Xn. Then

$$\lim_{n\to\infty} \frac{1}{n} l(X_1, X_2, ..., X_n) \le H(X) \quad \text{with probability 1.}$$

$$\begin{aligned} & \text{pf}: \ l(X_1, X_2, ..., X_n) = C(n)(\log C(n) + 1) \\ & \text{by Lemma 1,} \quad \limsup \frac{C(n)}{n} = 0 \\ & \Rightarrow \limsup_{n \to \infty} \frac{l(X_1, X_2, ..., X_n)}{n} = \limsup_{n \to \infty} \left(\frac{C(n)\log C(n)}{n} + \frac{C(n)}{n} \right) \\ & \leq H(X) \quad \text{, with probability 1.} \end{aligned}$$

The L-Z code is a universal code, in which, the code does not depend on the distribution of the source!!

Optimal Variable Length-to-Fixed Length Code

Algorithm:

step 1. 永遠分最大的 node

step 2. 分到有 2^l leaf nodes 則停

Example: Source data input: A B C C B A A A C

$$P_{A} = 0.7$$

Α	1011	
В	1010	
С	1001	
AA	1000	A B, C, C, B, A A A C
AB	0111	
AC	0110	0111 1001 1001 1010 0000
AAA	0101	
AAB	0100	Tunstall '77
AAC	0011	GIT ph.D. Thesis
AAAA	0010	
AAAB	0001	
AAAC	0000	