

Universal Lossless coding

Lempel-Ziv Coding

Outline

- Shortcomings of Block-to-variable-length codes
- Variable-length-to-block coding
- Lempel-Ziv coding

April 19, 2007

1

Shortcomings of Block-to-Variable-length Codes

1. Complexity

A blocklength K code for a source with M letters has M^K codewords.

Example: $M = 26$, $K = 5 =$ avg. word length: $M^K \approx 12 \times 10^6$

2. Dependence upon statistical model

- a. What if source probabilities are not known?
- b. What if source probabilities "change" over time?
- c. What if there are no probabilities?

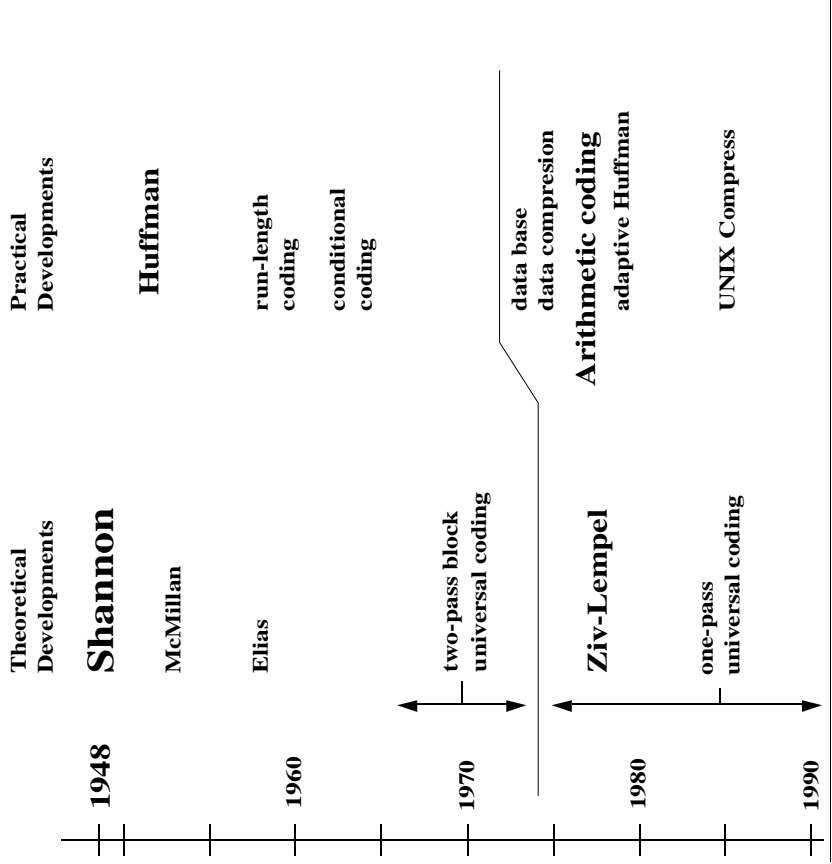
3. Goals

Find low complexity codes that are *universal*, meaning that when applied to any stationary ergodic source they should attain rate close to the entropy-rate of that source.

April 19, 2007

2

Historical Review



April 19, 2007

3

Lempel-Ziv Coding

- In these notes we introduce Lempel-Ziv codes, which is the most widely used family of universal codes.
 - There are many variations.
 - All are basically a kind of adaptive variable-length-to-block coding. Variable-length to block coding is introduced on the next pages.
 - LZ-78 and LZ-77 are the two broad categories of LZ codes.
- LZ-78 refers to codes following the approach described in

J. Ziv and A. Lempel, "Compression of individual sequences via variable-rate coding," *IEEE Trans. Inform. Theory*, vol. 24, pp. 530-536, Sept. 1978.

LZ-77 refers to codes following the approach described in

J. Ziv and A. Lempel, "A universal algorithm for sequential data compression," *IEEE Trans. Inform. Theory*, vol. 23, pp. 337-343, May. 1977.

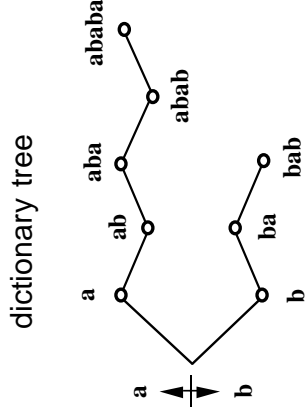
April 19, 2007

4

Variable-Length-to-Block Coding

aka Dictionary Coding

dictionary of source words	binary codewords w_i
a	000
b	001
ab	010
aba	011
ba	100
bab	101
abab	110
ababa	111



"Greedy" Encoding Rule:

Find longest sourceword w that prefixes remaining symbols; output codeword.

Example: ababbababababab is parsed into ababbabbababababab

$$\text{Rate: } R = \frac{\text{codeword len.}}{\text{avg. sourceword len.}} = \frac{L}{\sum_i p_i n_i} = \frac{5 \times 3}{18} = \frac{5}{6} \text{ bit/symb. for this example}$$

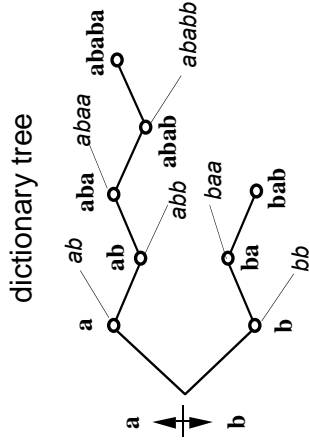
where p_i = probability that i th source word is chosen, and n_i is its length

April 19, 2007

5

Variable-Length-to-Block Coding

dictiry of source words	prob- ability p_i	source length n_i	code word w_i
a	.050	1	000
b	.050	1	001
ab	.045	2	010
aba	.041	3	011
ba	.045	2	100
bab	.405	3	101
abab	.036	4	110
ababa	.328	5	111



For binary alphabet, the probability of a given source word equals the probability of the "stopping word".

The above probabilities assume greedy encoding and a Markov source with $p(a) = p(b) = 0.5$, $p(b|a) = p(a|b) = 0.9$.

$$\text{Rate: } R = \frac{\text{codeword len.}}{\text{avg. sourceword len.}} = \frac{L}{\sum_i p_i n_i} = \frac{3}{3.4} = 0.88 \text{ bits/symbol}$$

The above probabilities assume greedy encoding and (approximately) a Markov model with $p(a) = p(b) = 0.5$, $p(b|a) = p(a|b) = 0.9$. $H_\infty = 0.47$.

April 19, 2007

6

Basic idea: A good code has $p_i \approx 2^{-L}$ for most i .

LZ-78: Lempel-Ziv Coding An Adaptive/Universal Block to VL Code

example sequence to encode: a b a b a b a b a b a b a b a b a b

encoding table

source word	index	code word
a	1	00..000
b	2	00..001
	3	00..010
	4	00..011
	5	00..100
	6	00..101
	7	00..110
	8	00..111

	2^L	11..111

Rate: $R = \frac{L+1}{\text{length}(w)+1}$ bits/symbol (drop "+1" when table full)

For the above sequence: let n denote the number of source words encoded. Then

$$R \approx \frac{4(L+1)}{n}, \text{ when } n \leq 2^L, \text{ and Rate} \approx \frac{(L+1)}{2^{L-1}}, \text{ when } n > 2^L.$$

Encoding Algorithm

1. Initialize encoding table with each letter of source alphabet as a source word, in alphabetical order, and with 2^L codewords of length L in usual order.
2. Find longest sourceword \underline{w} that prefixes remaining symbols.
3. Send binary codeword for \underline{w} and next symbol α (using some binary code)
4. Mark $\underline{w}\alpha$ as "already encoded".
5. If table not full, add $\underline{w}\alpha$ to encoder table.
6. Go to 2.

Decoding LZW

- Assume decoder knows the source alphabet and L.
- Decoding table is initialized with each letter of the alphabet as a source word, in alphabetical order, and with 2^L codewords of length L in usual order.
- As each w is received, decoder augments the table in the same way as the encoder has done, except that it leaves a "blank" for α , which it will not know until it receives the next w. Specifically, α is the first letter of the next w.
- For the example of the previous page, the encoded symbols are:

1 2 3 4 5 6 7 8 ...

decoding table		
index	source word	code word
1	a	00..000
2	b	00..001
3		00..010
4		00..011
5		00..100
6		00..101
7		00..110
8		00..111
2^L		11..111

April 19, 2007

11

Other variants: do not limit the size of table, nor fix the length of the codewords.

April 19, 2007

12

