

problems with Huffman codes

changing ensembles

- Huffman is optimal for iid sources

the extra bit: we know Huffman gives $H(X) \leq \mathbb{E}[L_C(X)] \leq H(X) + 1$

a	0.001	00000
b	0.001	00001
c	0.990	1
d	0.001	00010
e	0.001	00011
f	0.001	0100
g	0.001	0101
h	0.001	0110
i	0.001	0111
j	0.001	0010
k	0.001	0011

$$\mathbb{E}[\text{length}] = 1.034$$

$$H(X) = 0.114$$

$$\mathbb{E}[L]/H(X) = 9$$

the guessing game

M A J O R I T Y O F A M E R I C A N S -
17 1 12 2 0 0 0 0 1 1 3 1 2 1 0 0 0 0 0

F O U R I N T E N - H A T E M A T H
20 5 2 1 1 1 1 1 3 1 3 1 1 1 1 1

Code - Huffman Code of (17. 1. 12. 2. 0. . . .)

how to model data sources

- 1) iid sources, known distr
- 2) iid sources, unknown distr
(universal codes)
- 3) Oracle-based models (probabilistic model)
(arithmetic coding) $X_1, X_2, \dots, X_n, P(x_1), P(x_2|x_1), P(x_3|x_1, x_2) \dots$
- 4) Unknown probabilistic model
(Lempel-Ziv, Dictionary codes)

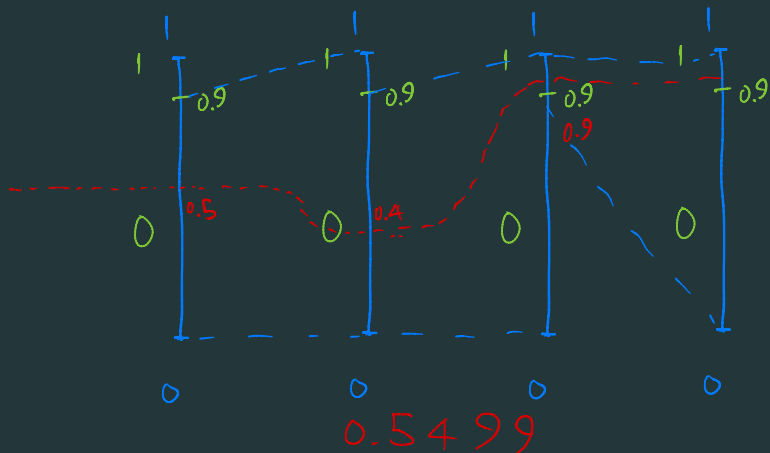
two approaches to stream coding

- Arithmetic coding
 - 'Encodes entire document as a single number'
 - Requires probabilistic model
(Separation of modeling & coding)
- Lempel-Ziv coding (Dictionary coding)
 - Generate and store a 'dictionary'
 - Codeword consists of dictionary + message

arithmetic coding

Want to compress $X_1 X_2 X_3 \dots X_n$, $X_i \sim \text{Ber}(0.1)$

• $D = 00110000000$



arithmetic coding

Idea - After encoding x_1, x_2 as $0.c_1c_2c_3$,
partition $[0,1]$ according to $P[x_3 | x_1, x_2]$



Average codeword length $\approx \log_2 \frac{1}{P(x_1, x_2, \dots, x_n)}$

$P[x_3 | x_1, x_2]$

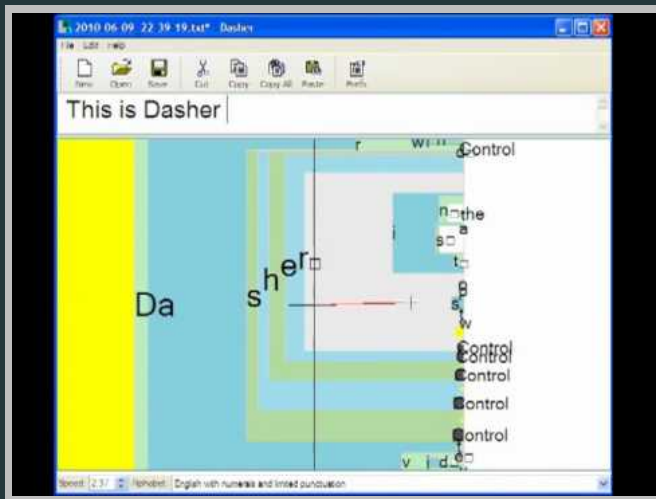
Fact - Symbol

x_i heads

$$\approx h(x_i | x_1, x_2, \dots, x_{i-1})$$

$$= \log_2 \frac{1}{P(x_i | x_1, x_2, \dots, x_{i-1})} \text{ bits}$$

application of arithmetic coding beyond compression



<https://www.youtube.com/watch?v=nr3s4613DX8>

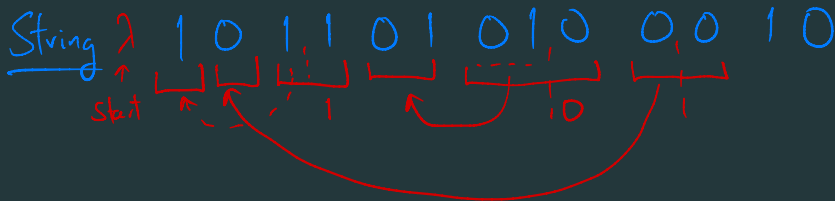
Lempel-Ziv codes (dictionary codes)

Eg - 000 | 00000 |

Eg - 000 0000000000

Lempel-Ziv-Welch coding

source substrings	λ	1	0	11	01	010	00	10
$s(n)$	0	1	2	3	4	5	6	7
$s(n)$ _{binary}	000	001	010	011	100	101	110	111
(pointer, bit)		(, 1)	(0, 0)	(01, 1)	(10, 1)	(100, 0)	(010, 0)	(001, 0)



• As $n \rightarrow \infty$, $|L(x_1, \dots, x_n, c)| \approx H(x_1, \dots, x_n)$