

problems with Huffman codes

changing ensembles

Huffman assumes iid, real world data is non-iid

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

If one common + many uncommon symbols, the +1 bit is very bad

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

$$H(X) = 0.114$$

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

the guessing game

14 1 11 1 1 1 1 1 1 2 3 1 1 1 1 1

MAJORITY OF PEOPLE -

26 5 ~6 1

FOUR IN TEN -HATE MATH

R

how to model data sources

- iid source, known distrⁿ - opt code is Huffman
- iid source, unknown distrⁿ - need to learn distrⁿ (inference)
- non iid source, known distrⁿ - separate probabilistic model from encoding/learning/commⁿ...
(eg. arithmetic coding) \uparrow
 - easy to learn
 - ⋮
 - hard to learn
- non iid source, unknown distrⁿ - 'agnostic' setting
 - combine modeling + task
 - universal codes, online learning, 'adversarial' methods, etc.

two approaches to stream coding (for comparison - see Mackey)

Arithmetic Coding - needs to know model
(known model, non iid)

(djuu, pp 3)

$$L(D) \leq H(D) + 2 \text{ bits}$$

Dictionary Coding (LZW) ('universal', i.e. unknown model)

(gzip)

$\lim_{D \rightarrow \infty} L(D) \approx H(D)$ without knowing model

arithmetic coding

idea - 'represent every database as a single real number'

$D \equiv$ 'to be or to be \square ' \xrightarrow{AC} 0.3141592652
{ decode }

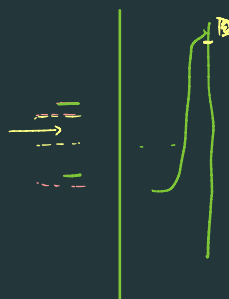
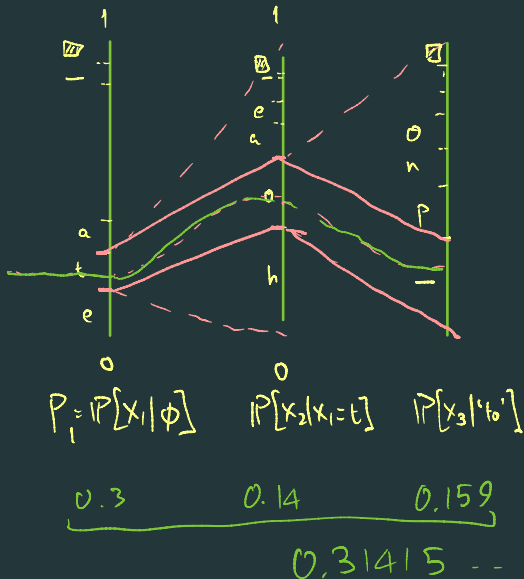
• If $D = X_1, X_2, \dots, X_n$ \square
 $\in X \cup \{\square\}$

Need - $P_t \equiv \underbrace{IP[X_t = x \mid X_1, X_2, \dots, X_{t-1}]}_{\text{probabilistic model}}$

Eg - Markovian model - $P_t \equiv IP[X_t = x \mid X_{t-1}]$ (bigram)

arithmetic coding

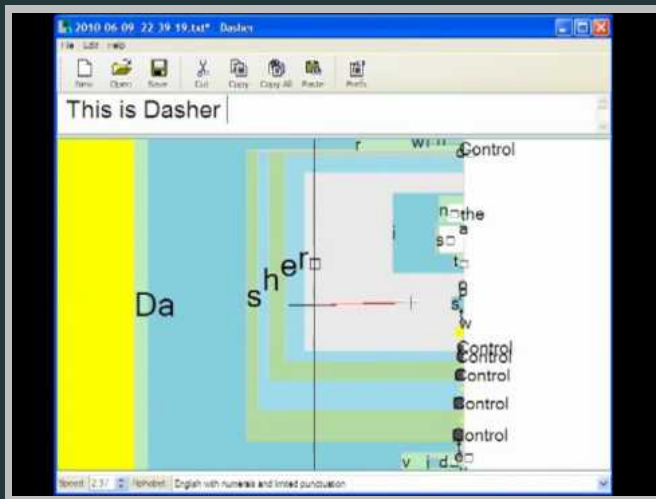
to_be_o



arithmetic coding

- for decoding - decoder knows the model
- runs the encoder 'in reverse'

application of arithmetic coding beyond compression

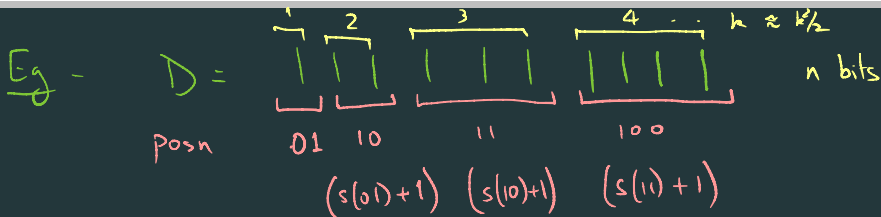


<https://www.youtube.com/watch?v=nr3s4613DX8>

Lempel-Ziv-Welch coding

source data - 1011010100010

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



$$\approx \sqrt{n} \log n \text{ bits}$$

Formal guarantee $\lim_{n \rightarrow \infty} \frac{L(D_n)}{n} \approx H(x)$

btill now - Chs 1,2 (intro) + Ch 4,5,6 (source coding)

Plan

- today - information theory for dependant r.v. (Ch 8)
 - Channel coding (Ch 9)
 - next class - Shannon's channel coding theorem (Ch 10)
-

In 2 classes - Intro to Bayesian stats
(Ch 2 of Mackay, Ch 1 of Bishop)