

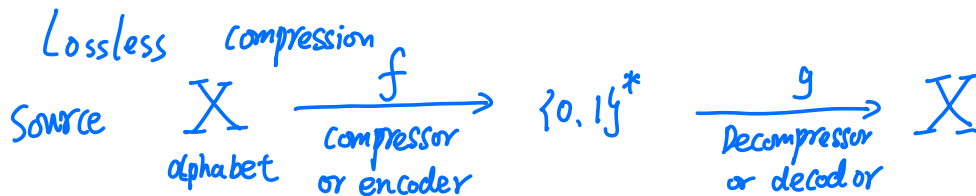
# Lossless Data Compression

"To day is Wednesday" a message as a sequence of letter  
 a letter  $\longleftrightarrow$  bytes  $\{1 \dots 256\}$   $\longleftrightarrow$  8 bits  $\{0,1\}^8$   
 "a"  $\xrightarrow{97}$   $\xrightarrow{28}$  110 0001

See table below

Binary	Oct	Dec	Hex	Glyph	Binary	Oct	Dec	Hex	Glyph	Binary	Oct	Dec	Hex	Glyph
010 0000	040	32	20	(space)	100 0000	100	64	40	@	110 0000	140	96	60	.
010 0001	041	33	21	!	100 0001	101	65	41	A	110 0001	141	97	61	a
010 0010	042	34	22	"	100 0010	102	66	42	B	110 0010	142	98	62	b
010 0011	043	35	23	#	100 0011	103	67	43	C	110 0011	143	99	63	c
010 0100	044	36	24	\$	100 0100	104	68	44	D	110 0100	144	100	64	d
010 0101	045	37	25	%	100 0101	105	69	45	E	110 0101	145	101	65	e
010 0110	046	38	26	&	100 0110	106	70	46	F	110 0110	146	102	66	f
010 0111	047	39	27	'	100 0111	107	71	47	G	110 0111	147	103	67	g
010 1000	050	40	28	(	100 1000	110	72	48	H	110 1000	150	104	68	h

Is this optimal in # of bits? No, if only english words we only need  
 $2^5 = 32 < 26 \times 2 = 52 < 64 = 2^6$   
 one letter - 6 bits



- $\{0,1\}^* = \{\emptyset, 0, 1, 00, 01, 10, 11, 000, \dots\}$  bit string
- $\forall x \in X, f(x) \in \{0,1\}^*$  code word.  $\{f(x) \mid x \in X\}$  code book <sup>countable</sup>
- Lossless compression:  $g \circ f = I_X \Rightarrow f$  injective  
 need  $|f(X)| = |X|$  (code word)
- length function:  $l: \{0,1\}^* \rightarrow \mathbb{N}$ , eg.  $l(01101) = 5$   
 an alphabet  $X$  need code word with maximal length  
 $\sup l(f(x)) = \log_2 |X|$

Can we compress more? In terms of maximal codeword length, no  
expected codeword length, Yes

Example:  $X = \{a, b, c, d\} \xrightarrow{w} \{0.1\}^2$   
 $a \rightarrow 00 \quad b \rightarrow 01 \quad c \rightarrow 10 \quad d \rightarrow 11$   
 each 2 bits  $\forall x, L(w(x)) = 2$   
 length

Now given the fact

$$P(a) = \frac{1}{2} \quad P(b) = \frac{1}{8} \quad P(c) = \frac{1}{4} \quad P(d) = \frac{1}{8}$$

Consider  $a \rightarrow 0, b \rightarrow 110, c \rightarrow 10, d \rightarrow 111$

$00110101110 \longleftrightarrow a a b c d a$

← variable  
length code

Expected codeword length

$$\frac{1}{2}(1) + \frac{1}{8}(3) + \frac{1}{4}(2) + \frac{1}{8}(3) = \frac{7}{4} < 2 \quad !$$

In the long run, do better than 2 bits/letter.

What does " $P(a) = \frac{1}{2} \quad P(b) = \frac{1}{8} \quad \dots$ " mean?

frequency of letters in two different English novel are approximately the same.

English text  $\xrightarrow{\text{empirical}}$  frequency of letters  $\xrightarrow{\text{modeled by}}$  probabilistic feature  
 distribution on Alphabet  
 R.V.

— Shannon

Objective: minimize  $\rightarrow \sup L[f(X)]$   
 $\rightarrow \mathbb{E} L[f(X)]$

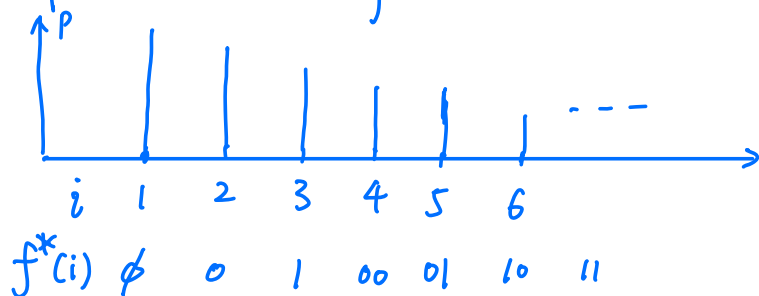
There is an optimal compressor  $f^*$  minimize both!

Main idea — Assign short code words for more probable symbols  
 longer — — — — — less probable symbols.

WLOG,  $X = \{1, 2, \dots, |X|\} \subseteq \mathbb{N}$  and reorder p.m.f  $\rightarrow$   
 $P_X(i+1) \leq P_X(i)$

Theorem (Optimal Compressor)

Define the encoder  $f^*$



Then

1. length of code word

$$L(f^*(i)) = \lfloor \log_2 i \rfloor \quad (\lfloor a \rfloor := \text{largest integer } \leq a)$$

2. Stochastically optimal:  $\forall$  encoder  $f$ ,

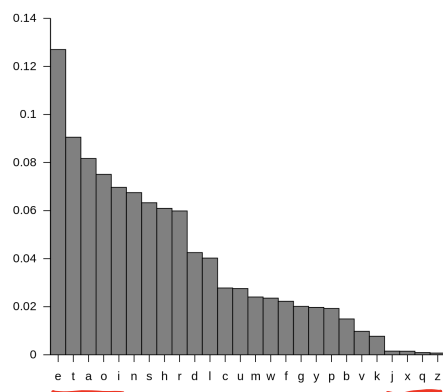
$$\forall k, P[L(f(X)) \leq k] \leq P[L(f^*(X)) \leq k] \quad (L(f^*(X)) \stackrel{\text{st.}}{\leq} L(f(X)))$$

As a consequence,

$$\sup L(f(X)) \geq \sup L(f^*(X)). \quad \mathbb{E} L(f(X)) \geq \mathbb{E} L(f^*(X))$$

# Example (Telegraph & Morse code)

## Letter frequency in English text



## Morse code

### International Morse Code

1. The length of a dot is one unit.
2. A dash is three units.
3. The space between parts of the same letter is one unit.
4. The space between letters is three units.
5. The space between words is seven units.

→ A	•••	U	•••••
B	•••••	V	•••••
C	•••••	W	•••••
D	•••••	X	•••••
→ E	•••	Y	•••••
F	•••••	Z	•••••
G	•••••		
H	•••••		
I	•••		
→ J	•••••		
K	•••••		
L	•••••		
M	•••••		
N	•••••		
→ O	•••••		
P	•••••		
Q	•••••		
R	•••••		
→ S	•••••		
T	•••		

For the other side, we recall the lemma:

Lemma. For  $Z \in \mathbb{N}$  and  $\mathbb{E}Z < \infty$ ,  $H(Z) \leq \mathbb{E}Z \log\left(\frac{1}{\mathbb{E}Z}\right)$

Entropy of Geometric distribution  $Q_p(i) = p(1-p)^i$   $H(Q_p) = \frac{h(p)}{p}$   
 $\mathbb{E}(Q_p) = \frac{1}{p}$

$\downarrow$   
 $p = \mathbb{E}Z \Rightarrow H(Z) \leq H(Q_p)$

$Z \in \mathbb{N}$

Geometric distribution has largest entropy

for integer valued RV of given mean  $p = \mathbb{E}X$

Relative entropy

How to prove? ① Lemma  $\forall P, Q$  prob.  $D(P \parallel Q) := \sum p(w) \log \frac{p(w)}{q(w)} \geq 0$

② Calculate  $D(P \parallel Q_p) = \sum p(w) \log p(w) - p(w) \log q(w)$   
 where  $p = \frac{1}{\mathbb{E}X}$ .

$$\text{Pf: } |A_k| := |\{x \mid \lfloor \log(x) \leq k \} | \leq \sum_{i=0}^k 2^i = 2^{k+1} - 1 \\ = |\{x \mid \lfloor \log^*(x) \leq k \} | = |A_k^*|$$

$$\text{Because } A_k^* = \{1, 2, \dots, 2^{k+1} - 1\}$$

$$P[\lfloor \log(x) \leq k] = \sum_{x \in A_k} P_X(x) \leq \sum_{i=1}^{2^{k+1}-1} P_X(i) = \sum_{x \in A_k^*} P_X(x) = P[\lfloor \log^*(x) \leq k]$$

□

$X, Y: \Omega \rightarrow \mathbb{R}$  real R.V.s.

Def. We say  $X$  is stochastically dominated by  $Y$ , denoted  $Y \stackrel{\text{st.}}{\leq} X$  if  
 $\forall k, \quad P_X(X \leq k) \leq P_Y(Y \leq k)$

Prop. If  $Y \stackrel{\text{st.}}{\leq} X$ , then

i)  $\sup Y \leq \sup X$  and ii)  $\mathbb{E} Y \leq \mathbb{E} X$

Pf: i)  $\sup X = \sup \{k \mid P_X(X \leq k) < 1\} \geq \sup \{k \mid P_Y(X \leq k)\} = \sup Y$

ii) special case:  $X, Y: \Omega \rightarrow \mathbb{N}$ .

$$\boxed{\mathbb{E} X = \sum_{n=1}^{\infty} P_X(X \geq n)}$$

$$\text{Indeed, } \mathbb{E} X = \sum_{n=1}^{\infty} n P_X(n) = P_X(X \geq 1) + \sum_{n=2}^{\infty} P_X(X \geq n) \\ = \sum_{n=1}^{\infty} P_X(X \geq n)$$

$$\mathbb{E} X = \sum_{n=1}^{\infty} P_X(X \geq n) \geq \sum_{n=1}^{\infty} P_Y(Y \geq n) = \mathbb{E} Y$$

Theorem (Optimal Average code length)

Given  $X \in \mathcal{N}$  and  $P_X(1) \geq P_X(2) \geq \dots$ . Then

$$\textcircled{1} \quad \mathbb{E}[L(f^*(X))] = \sum_{k=1}^{\infty} P[X \geq 2^k]$$

$$\textcircled{2} \quad H(X) - \log_2(eH(X) + 1) \leq \mathbb{E}[L(f^*(X))] \leq H(X)$$

$$\begin{aligned} \text{pf: } \textcircled{1} \quad \mathbb{E}[L(f^*(X))] &= \mathbb{E}(\lfloor \log_2 X \rfloor) = \sum_{k \geq 1} P(\lfloor \log_2 X \rfloor \geq k) \\ &= \sum_{k \geq 1} P(\log_2 X \geq k) \end{aligned}$$

$$\textcircled{2} \quad \text{Denote } L(X) = L(f^*(X))$$

$$P_X(m) \leq \frac{1}{m} \quad \text{b/c } P_X(i) \text{ decreasing}$$

$$\Rightarrow L(f^*(m)) = \lfloor \log_2 m \rfloor \leq \log_2 \frac{1}{P_X(m)}$$

$$\Rightarrow \mathbb{E}[L(X)] \leq \mathbb{E}(\log \frac{1}{P_X(m)}) = H(X)$$

For the other side, we recall the lemma:

Lemma. For  $Z \in \mathcal{N}$  and  $\mathbb{E} Z < \infty$ ,  $H(Z) \leq \mathbb{E}[Z] h(\frac{1}{\mathbb{E} Z})$   
{0, 1, 2, ..., n, ...}

$$H(X) = H(X, L) = H(X|L) + H(L)$$

$$= \sum P_L(k) H(X|L=k) + h(\frac{1}{H(L)}) (1 + \mathbb{E} L)$$

$$L \in \{0, 1, 2, \dots\}$$

$$\leq \sum P_L(k) \log \frac{2^k}{P_X(k)} + \dots$$

$$L+1 \in \mathcal{N}$$

$$= \mathbb{E} L + \log(1 + \mathbb{E} L) + (\mathbb{E} L) \log(1 + \frac{1}{\mathbb{E} L})$$

$$\leq \mathbb{E}L + \log_2(e(1+H(x)))$$

$$\left( x \log\left(1 + \frac{1}{x}\right) \leq \log e \right)_{\forall x > 0}$$

□

Cor. If  $X = S^n$  i.i.d. sequence,  $H(S^n) = n H(S)$

$$nH(S) - \log n + O(1) \leq \mathbb{E}[\ell_f^*(S^n)] \leq nH(S)$$

Hence  $\lim_{n \rightarrow \infty} \frac{\mathbb{E}[\ell_f^*(S^n)]}{n} = H(S)$  bits  
Expected length/message

$$\mathbb{E}[\ell_f^*(S^n)] = nH(S) - \frac{1}{2}\log n + O(1) \quad (\text{Szpankowski \& Verdú . 2011})$$

Remark Actually.  $\frac{\ell_f^*(S^n)}{n} \rightarrow H(S)$  in probability  
by WLLN.