

How much can we compress? - Shannon's Source Coding Theorem



On Probability and Entropy



Probability

- An ensemble X is a random variable x with a set of possible outcomes \mathbf{A}_x with probabilities \mathbf{P}_x
- Probability of a subset T of \mathbf{A}_x

$$P(T) = \sum_{a_i \in T} P(x = a_i)$$

- A **joint ensemble** XY is an ensemble for which the outcomes are ordered pairs x, y where $x \in \mathbf{A}_x$ and $y \in \mathbf{A}_y$

Three Concepts: Information '02

©Henry Tirri 2002

25

Probability continued

- Marginal probability (from the joint probability $P(x, y)$)

$$P(y) = \sum_{x \in \mathbf{A}_x} P(x, y)$$

- Conditional probability

$$P(x = a_i | y = b_j) \equiv \frac{P(x = a_i, y = b_j)}{P(y = b_j)}$$

Three Concepts: Information '02

©Henry Tirri 2002

26

Probability continued

- Product rule

$$P(x, y | H) = P(x | y, H)P(y | H)$$

- Sum rule

$$\begin{aligned} P(x | H) &= \sum_y P(x, y | H) \\ &= \sum_y P(x | y, H)P(y | H) \end{aligned}$$

Three Concepts: Information '02

©Henry Tirri 2002

27

Bayes's theorem

$$\begin{aligned} P(y | x, H) &= \frac{P(x | y, H)P(y | H)}{P(x | H)} \\ &= \frac{P(x | y, H)P(y | H)}{\sum_{y'} P(x | y', H)P(y' | H)} \end{aligned}$$



Bayesian view of probability!

Three Concepts: Information '02

©Henry Tirri 2002

28

Entropy

- The entropy of X is a measure of the information content or "uncertainty" of X

- ✓ $H(X) \geq 0$ (= iff $p_i=1$ for one i)
- ✓ $H(X) \leq \log(|X|)$ (= iff $p_i=1/|X|$ for all i)

$$H(X) \equiv \sum_{x \in A_X} P(x) \log \frac{1}{P(x)}$$



Binary entropy

$$H(X) \equiv \sum_i p_i \log_2 \frac{1}{p_i} \quad \text{Information measure?}$$

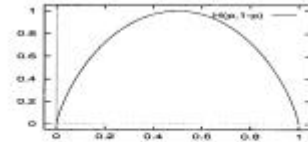


Figure 2.1. The binary entropy function $H_2(p) = H(p, 1-p) = -p \log_2 p - (1-p) \log_2 (1-p)$ as a function of p .

Information content

- First attempt: number of possible outcomes $|A_x|$

- ✓ not additive: for xy we have $|A_x| |A_y|$

- Perfect information content

- ✓ additive, but no probabilistic element

$$H_0(X) = \log_2 |A_X|$$



Shannon information

- looking for an information content of the event $x=a_i$

$$h(x) = \log_2 \frac{1}{p_i}$$

Example: letter distribution

i	a_i	P_i	$\log_2 \frac{1}{P_i}$	i	a_i	P_i	$\log_2 \frac{1}{P_i}$	i	a_i	P_i	$\log_2 \frac{1}{P_i}$
1	a	0.08	4.2	18	j	0.03	18.7	29	w	0.26	4.2
2	b	0.02	8.3	19	k	0.02	8.3	30	x	0.07	3.8
3	c	0.03	8.2	20	l	0.04	4.9	31	y	0.02	4.8
4	d	0.03	8.2	21	m	0.02	8.3	32	z	0.01	7.2
5	e	0.09	3.5	22	n	0.09	4.1	33		0.01	6.4
6	f	0.02	8.3	23	o	0.07	3.9	34		0.01	7.2
7	g	0.02	8.2	24	p	0.02	5.7	35		0.02	5.8
8	h	0.03	5.8	25	q	0.01	10.3	36		0.00	16.4
9	i	0.08	4.2	26	r	0.05	4.3	37		0.19	2.4

$$\sum_{i=1}^n p_i \log_2 \frac{1}{p_i} = 4.12$$

Figure 2.26. Probability distribution over the 27 outcomes for a randomly selected letter in an English language document (estimated from The frequency of letters in English). The picture shows the probabilities by the size of white spaces.

Entropy continued

- The joint entropy of X, Y

$$H(X, Y) \equiv \sum_{x,y \in A_X, A_Y} P(x, y) \log \frac{1}{P(x, y)}$$

- The conditional entropy of X given Y

$$H(X|Y) \equiv \sum_{y \in A_Y} P(y) \left[\sum_{x \in A_X} P(x|y) \log \frac{1}{P(x|y)} \right]$$

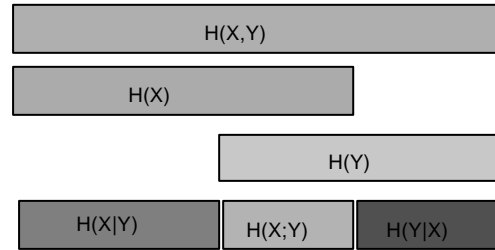
$$= \sum_{x,y \in A_X, A_Y} P(x, y) \log \frac{1}{P(x|y)}$$

"Average uncertainty that remains about x when y is known"

Entropy continued

- Chain rule for entropy
 $H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y)$
- Mutual information "Average reduction in uncertainty of x when learning the value of y"
 $H(X; Y) \equiv H(X) - H(X|Y)$
- Entropy distance
 $D_H(X, Y) \equiv H(X, Y) - H(X; Y)$

Entropy relationships



Kullback-Leibler divergence

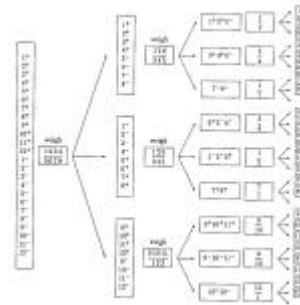
- Also known as "relative entropy"

$$D_{KL}(P \parallel Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

- Not strictly a "distance"



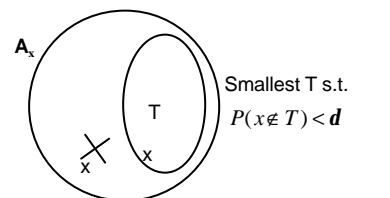
Weighting problem



Idea

- Some symbols have a smaller probability
- gamble that the rare symbols won't occur
- encode the observations in a smaller code (alphabet) C_x
- measure $\log_2 |C_x|$
- the larger the risk, the smaller the alphabet

Formalize the idea



Essential information

$$H_d(X) = \log_2 \min\{|T| : T \subseteq A_X, P(x \in T) \geq 1 - d\}$$

Example

$\mathbf{x} = (x_1, \dots, x_N)$, $x = \{0,1\}$ with probabilities $p_0 = .9$, $p_1 = .1$
Let $r(\mathbf{x})$ be the number of 1's in \mathbf{x}

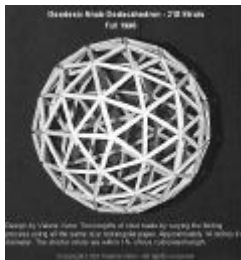
Probability of string \mathbf{x}

$$P(\mathbf{x} | p_0, p_1) = p_0^{N-r(\mathbf{x})} p_1^{r(\mathbf{x})}$$

AEP and source coding

Asymptotic Equipartition Principle: for N i.i.d. random variables $X^N = \{X_1, \dots, X_N\}$, with N sufficiently large, the outcome $\mathbf{x} = \{x_1, \dots, x_N\}$ is almost certain to belong to a subset of \mathbf{A}_x^N having only $2^{NH(X)}$ members all having probability close to $2^{-NH(X)}$

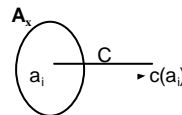
The Revenge of a Student - Symbol Codes



Symbol codes

- Notation: $\{0,1\}^+ = \{0,1,00,01,10,11,000,\dots\}$
- A symbol code C is a mapping from \mathbf{A}_x to $\{0,1\}^+$

$$c^+(x_1 x_2 x_3 \dots x_N) = c(x_1) c(x_2) c(x_3) \dots c(x_N)$$



$$l(x) = |x|$$



Decoding of symbol codes

- A code $C(X)$ is uniquely decodable if $\forall \mathbf{x}, \mathbf{y} \in A_x^+, \mathbf{x} \neq \mathbf{y} \Rightarrow c^+(\mathbf{x}) \neq c^+(\mathbf{y})$
- A code $C(X)$ is a prefix code if no codeword is a prefix of any other codeword
- The expected length $L(C, X)$ of a symbol code C for ensemble X is

$$L(C, X) = \sum_{x \in A_x} P(x) l(x)$$

Example

$$\mathbf{A}_x = \{1,2,3,4\}, P_x = \{1/2, 1/4, 1/8, 1/8\}$$

$$C: c(1) = 0, c(2) = 10, c(3) = 110, c(4) = 111$$

The entropy of X is 1.75 bits: $L(C, X)$ is also 1.75 bits

Obs!

$$l_i = \log_2(1/p_i), p_i = 2^{-l_i}$$



Kraft inequality

- Given a list of integer $\{l_i\}$, does there exist a uniquely decodable code with $\{l_i\}$?
- "Market model": total budget 1; cost per codeword of length l is 2^{-l} .

Kraft inequality: For any uniquely decodable code C over the binary alphabet $\{0,1\}$, the codeword lengths must satisfy:

$$\sum_i 2^{-l_i} \leq 1$$

Conversely, given a set of codeword lengths that satisfy this inequality, there exists a uniquely decodable prefix code with these codelengths.

Limits of unique decodeability

0	00	000	0000	Total "budget"
		001	0001	
	01	010	0010	
		011	0011	
1	10	100	0100	
		101	0101	
	11	110	0110	
		111	0111	
		1000	1000	
		1001	1001	
		1010	1010	
		1011	1011	
		1100	1100	
		1101	1101	
		1110	1110	
		1111	1111	

What can we hope for?

Lower bound on expected length: The expected length $L(C,X)$ of a uniquely decodable code is bounded below by $H(X)$.

Compression limit of symbol codes: For an ensemble X there exists a prefix code

$$H(X) \leq L(C,X) < H(X) + 1.$$



"Proof-map" of the lower bound

Define $q_i \equiv 2^{-l_i/z}$, where $z = \sum_i 2^{-l_i}$.

Thus $l_i = \log 1/q_i - \log z$.

$$L(C,X) = \sum_i p_i l_i = \sum_i p_i \log 1/q_i - \log z$$

$$\geq \sum_i p_i \log 1/p_i - \log z$$

$$\geq H(X)$$

(What happens if we use the "wrong" code?)

Assume the "true probability distribution" is $\{p_j\}$. If we use a complete code with lengths l_i , they define a probabilistic model $q_i = 2^{-l_i}$. The average length is

$$L(C,X) = H(X) + \sum_i p_i \log p_i / q_i$$

Kullback-Leibler divergence $D_{KL}(p||q)$

"Optimal" symbol code: Huffman coding

- Take two least probable symbols in the alphabet as defined by $\{p_j\}$.
- Combine these symbols into a single symbol, $p_{\text{new}} = p_1 + p_2$. Repeat (until one symbol)

