



Chapter 4: Entropy





Section 4.1. The Shannon-Game





Informal definition of Entropy

Entropy:

Amount of information contained in a message (after removing all possible redundancy)

(number of bits that the message has after compression)



Complete the Sentence



There is only one way to fill in the vowels in this sentence





Entropy of a Language: Shannon's Approach

- Show somebody the beginning of a text
- 2. Ask him/her to guess the next letter
- 3. Count the number of trials





The Shannon Game

$$T_{1}H_{1}E_{1}R_{5}E_{1}\bullet_{1}I_{2}S_{1}\bullet_{1}N_{2}O_{1}\bullet_{1}R_{15}E_{1}V_{17}E_{1}R_{1}S_{1}E_{2}\bullet_{1}O_{3}N_{2}\bullet_{1}A_{2}\bullet_{2}\\ M_{7}O_{1}T_{1}O_{1}R_{1}C_{4}Y_{1}C_{1}L_{1}E_{1}\bullet_{1}A_{3}\bullet_{1}F_{8}R_{6}I_{1}E_{3}N_{1}D_{1}\bullet_{1}O_{1}F_{1}\bullet_{1}M_{1}I_{1}\\ N_{1}E_{1}\bullet_{1}F_{6}O_{2}U_{1}N_{1}D_{1}\bullet_{1}T_{1}H_{1}I_{2}S_{1}\bullet_{1}O_{1}U_{1}T_{1}\bullet_{1}R_{4}A_{1}T_{1}H_{1}E_{1}R_{1}\bullet_{1}\\ D_{11}R_{5}A_{1}M_{1}A_{1}T_{1}I_{1}C_{1}A_{1}L_{1}L_{3}Y_{1}\bullet_{1}T_{6}H_{1}E_{1}\bullet_{1}O_{1}T_{1}H_{1}E_{1}R_{1}\bullet_{1}D_{1}A_{1}\\ Y_{1}\bullet_{1}$$

Play it yourself at:

http://www.math.ucsd.edu/~crypto/java/ENTROPY/



Cover and Kings Variation to the Shannon Game



- 1. The test person gets a certain amount of money
- 2. Show the person the beginning of the text
- 3. Ask the person to bet on the next letter
- 4. Observe how the fortune of the person evolves

IEEE TRANSACTIONS ON INFORMATION THEORY, VOL. IT-24, NO. 4, JULY 1978

A Convergent Gambling Estimate of the Entropy of English

THOMAS M. COVER, FELLOW, IEEE, AND ROGER C. KING, STUDENT MEMBER, IEEE





Example of Results

TABLE I

EXPERIMENTAL RESULTS ON ESTIMATING ENTROPY OF ENGLISH
USING SEQUENCE OF 75 SYMBOLS FROM

Jefferson the Virginian

Subject	Capital Achieved	Resultant Entropy Estimate
1	1.50×10^{78}	1.29 bits/sym
2	1.46×10^{76}	1.38
3	3.36×10^{75}	1.41
4	2.37 x 10 ⁷³	1.51
5	6.45×10^{71}	1.57
6	3.22 x 10 ⁷¹	1.59
7	2.30 x 10 ⁷⁰	1.64
8	4.00×10^{70}	1.67
9	2.21 x 10 ⁶⁹	1.68
10	9.63 x 10 ⁶⁸	1.70
11	3.88 x 10 ⁶⁷	1.76
12	3.60 x 10 ⁶⁴	1.90

Note the difference between the subjects





The Entropy of Natural Languages

From statistical analysis of text		Size of alphabet	Letter models with order:							Word		
English 26 4.70 4.14 3.56 3.3 2.61 2.14 Shannon (1951) English 26 4.70 4.12 1.65 Barnard (1955) French 26 4.70 3.98 3.02 German 26 4.70 4.10 1.08 Spanish 26 4.70 4.02 1.97 English 26+1 4.75 4.09 3.23 2.85 2.66 2.43 2.40 English 26+1 4.70 3.40 2.68 2.40 2.28 2.16 2.14 Portugese 26? 4.70? 3.92 3.51 3.15 Tamil 30 4.91 4.34 Kannada 49 5.61 4.55 Telugu 53 5.73 4.59 3.09 Chinese 4700 12.20 9.63 From experiments with subjects' best guesses English 26+1 4.75 2.9 2.6 2.8 3.0 Indian 26+1 4.75 2.9 2.6 2.8 3.0 English 26+1 4.75 3.4 3.1 3.3 3.8 From experiments with subjects using gambling English 26+1 4.75 3.4 3.1 3.3 3.8 From experiments with subjects using gambling English 26+1 4.75 3.5 2.8 2.9 3.2	Language		-1	0	i	2	3	7	11	≥100	nodel	Source
26+1 4.75 4.03 3.32 3.1 2.14		-	<u>~</u>		Fro	om stati	stical an	alysis c	of text	28%		
26+1 4.75 4.03 3.32 3.1 2.14	English	26	4.70	4.14	3.56	3.3					2,62	Shannon (1951)
Spanish 26 4.70 3.98 3.02 1.08 1.08 1.97	<u>s</u>			4.03		3.1					2.14	
French 26 4.70 3.98 3.02 German 26 4.70 4.10 1.08 Spanish 26 4.70 4.02 1.97 English 26+1 4.75 4.09 3.23 2.85 2.66 2.43 2.40 1.97 Portugese 26? 4.70? 3.92 3.51 3.15 Manfrino (1970) Tamil 30 4.91 4.34 Kannada 49 5.61 4.55 Telugu 53 5.73 4.59 3.09 Ghinese 4700 12.20 9.63 From experiments with subjects' best guesses English 26+1 4.75 upper bound (smoothed) 4.0 3.4 3.0 2.6 2.1 1.9 1.3 lower bound (smoothed) 3.2 2.5 2.1 1.8 1.2 1.1 0.6 English 26+1 4.75 Lafain 26+1 4.75 upper bound (smoothed) 3.2 2.5 2.1 1.8 1.2 1.1 0.6 English 26+1 4.75 Lafain 26+1 4.7	English	26	4.70	4.12							1.65	Barnard (1955)
Spanish 26 4.70 4.10 1.08 1.97			4.70	3.98							3.02	
Spanish 26											1.08	
Samoan 16+1 4,09 3,40 2,68 2,40 2,28 2,16 2,14 (1960) Russian 35+1 5,17 4,55 3,44 2,95 2,72 2,45 2,40 Manfrino (1970) Manfrino (1970) Siromoney (1963) Rajagopalan (1965) Rajagopalan (1965) Balasubrahmanyam Siromoney (1968)											1.97	
Samoan 16+1 4.09 3.40 2.68 2.40 2.28 2.16 2.14 (1960) Russian 35+1 5.17 4.55 3.44 2.95 2.72 2.45 2.40 (1960) Portugese 26? 4.70? 3.92 3.51 3.15 Manfrino (1970) Manfrino (1970) Tamil 30 4.91 4.34 Siromoney (1963) Rajagopalan (1965) Rajagopalan (1965) Balasubrahmanyam Siromoney (1968) Chiuse 53 5.73 4.59 3.09 Wanas et al. (1976) Wanas et al. (1976) Chinese 4700 12.20 9.63 Wong and Poon (1976) Wong and Poon (1976) Wong and Poon (1976) English 26+1 4.75 2.5 2.1 1.9 1.3 Shinnon (1951) English 26+1 4.75 2.2 1.8 1.7 Jamison and Jamison and Jamison (1960) English 26+1 4.75 3.4 3.1 3.3 3.8	English	26+1	4.75	4.09	3.23	2.85	2.66	2.43	2.40			Newman and Waugh
Russian 35+1 5.17 4.55 3.44 2.95 2.72 2.45 2.40		N-04000			2.68	2.40	2.28	2.16	2.14			(1960)
Tamil 30 4.91 4.34 Siromoney (1963) Rajagopulan (1965) Telugu 53 5.73 4.59 3.09 Balasubrahmanyam Siromoney (1968) Arabic 32 5.00 4.21 3.77 2.49 Wanas et al. (1976) Chinese 4700 12.20 9.63 Wong and Poon (1968) From experiments with subjects' best guesses English 26+1 4.75				*** T81510100			2.72	2.45	2.40			
Rajagopalan (1965) Rajagopalan (1966) Rajagop	Portugese	26?	4,70?	3.92	3.51	3.15		ä				Manfrino (1970)
Telugu 53 5.73 4.59 3.09 Balasubrahmanyam Siromoney (1968) Arabic 32 5.00 4.21 3.77 2.49 Wanas et al. (1976) Chinese 4700 12.20 9.63 Wong and Poon (1976) From experiments with subjects' best guesses English 26+1 4.75 upper bound (smoothed) 4.0 3.4 3.0 2.6 2.1 1.9 1.3 Shunnon (1951) lower bound (smoothed) 3.2 2.5 2.1 1.8 1.2 1.1 0.6 English 26+1 4.75 2.2 1.8 1.8 1.7 Jamison and Jamison Inalian 26+1 4.75 2.9 2.6 2.8 3.0 (1968) Ratian* 26+1 4.75 3.4 3.1 3.3 3.8 French* 26+1 4.75 3.5 2.8 2.9 3.2 From experiments with subjects using gambling English 26+1 4.75 1.25 Cover and King (1968)	Tamil	30	4.91	4.34								
Telugu 53 5.73 4.59 3.09 Balasubrahmanyam Siromoney (1968) Arabic 32 5.00 4.21 3.77 2.49 Wanas et al. (1976) From experiments with subjects' best guesses English 26+1 4.75 4.0 3.4 3.0 2.6 2.1 1.9 1.3 Shannon (1951) English 26+1 4.75 2.2 1.8 1.2 1.1 0.6 English 26+1 4.75 2.2 1.8 1.8 1.7 Jamison and Jamiss (1961) Italian 26+1 4.75 3.4 3.1 3.3 3.8 French* 26+1 4.75 3.5 2.8 2.9 3.2 From experiments with subjects using gambling English 26+1 4.75 2.2 2.8 2.9 3.2	Kannada	49	5.61	4.55								Rajagopalan (1965)
Siromoney (1968) Wanas et al. (1976)		53	5.73	4.59	3.09							Balasubrahmanyam and
From experiments with subjects' best guesses				13								
Chinese 4700 12.20 9.63 Wong and Poon (19	Arabic	32	5.00	4.21	3.77	2.49						Wanas et al. (1976)
English 26+1 4.75	Chinese	4700	12.20	9.63			10,000			Yes.	20	Wong and Poon (1976
upper bound (smoothed) lower bound (smoothed) lower bound (smoothed) 4.0 3.4 3.0 2.6 2.1 1.9 1.3 Shannon (1951) English 26+1 4.75 2.5 2.1 1.8 1.2 1.1 0.6 English 26+1 4.75 2.2 1.8 1.8 1.7 Jamison and Jami	pt. 15000 150		37-61	F	rom exp	eriment	s with s	ubjects	' best go	iesses		
English 26+1 4.75 2.2 1.8 1.2 1.1 0.6												EDERE MANAGE
English 26+1 4.75 2.2 1.8 1.8 1.7 Jamison and Jamison and Jamison 26+1 4.75 2.9 2.6 2.8 3.0 (196) Italian* 26+1 4.75 3.4 3.1 3.3 3.8 French* 26+1 4.75 3.5 2.8 2.9 3.2 From experiments with subjects using gambling English 26+1 4.75 1.25 Cover and King (196)												Shannon (1951)
Tablan 26+1 4.75 2.9 2.6 2.8 3.0 (196) Italian* 26+1 4.75 3.4 3.1 3.3 3.8 French* 26+1 4.75 3.5 2.8 2.9 3.2 From experiments with subjects using gambling English 26+1 4.75 1.25 Cover and King (196) Cover and King (196) 1.25 Cover	low	er bound (st	moothed)	3,2	2.5	2,1	1.8	1.2	1.1	0.6		
Halian* 26+1 4.75 3.4 3.1 3.3 3.8	English	26+1	4.75				18.00000					Jamison and Jamison
French* 26+1 4.75 3.5 2.8 2.9 3.2	Italian	26+1	4.75				2.9					(1968)
From experiments with subjects using gambling English 26+1 4.75 1.25 Cover and King (19)	Italian*	26+1	4.75									
English 26+1 4.75 1.25 Cover and King (19	French*	26+1	4.75		-		3.5	2.8	2.9	3.2		
2011 100 N				Fre	om expe	eriments	with st	ibjects (ising ga	mbling		
1.22	English	26+1	4.75							1.25		Cover and King (1978
Malay 26+1 4.75 1.32 1an (1961)	Malay	26+1	4.75							1.32		Tan (1981)

Gambling experiment gives more accurate results than simple guessing





Formal Definition of Entropy

$$H(V) = \mathsf{E}[-\log(p(V))]$$

$$= \sum_{w_i \in V} -p(w_i) \log(p(w_i))$$

Note: if you want the "unit" of the entropy to be "bit" you have to use the log to the basis 2







Vocabulary with two words:

$$V=\{a,b\}$$

$$P(a)=p$$

$$P(b)=1-p$$

$$H=-p log p - (1-p) log(1-p)$$

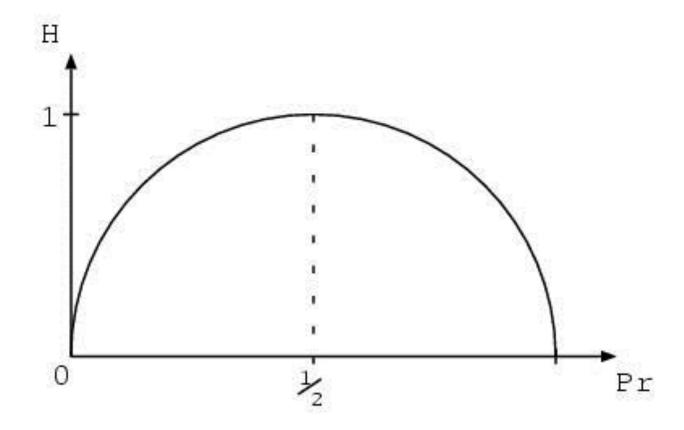
$$p=0 \mapsto H=0$$

$$p=1 \mapsto H=0$$





Entropy H=-p log p - (1-p) log(1-p)







Example 2

Vocabulary of W words w_i with uniform distribution $p(w_i)=1/W$

$$H = \sum_{i=1}^{W} -p(w_i)\log(p(w_i)) = \sum_{i=1}^{W} -\frac{1}{W}\log(\frac{1}{W})$$
$$= -W\frac{1}{W}\log(\frac{1}{W}) = -\log(\frac{1}{W}) = \log(W)$$

Entropy for uniform distribution: log of the number of symbols





Theorem: Entropy is never negative

Theorem

$$H(V) \ge 0$$

Proof:

Idea of proof: use that the log is a concave function





Proof: entropy is not negative

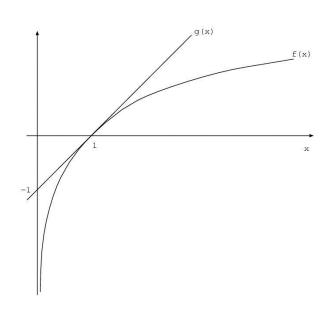
Define auxiliary functions

$$f(x)=log(x)$$
$$g(x)=x-1$$

$$\mapsto g(x) \ge f(x) \quad \forall x > 0$$

$$\mapsto x-1 \ge \log(x) \quad \forall x > 0$$

$$\mapsto -\log(x) \ge 1 - x \quad \forall x > 0$$







Proof: entropy is not negative

$$H = \sum_{i=1}^{W} -p(w_i) \log(p(w_i))$$

$$\geq \sum_{i=1}^{W} \underbrace{p(w_i)}_{\geq 0} \underbrace{(1-p(w_i))}_{\geq 0} \geq 0$$

Use

$$-\log(x) \ge 1 - x \quad \forall x > 0$$





Definition 2-Word Entropy

$$H(W,W) = -\sum_{w_1,w_2} p(w_1, w_2) \log(p(w_1, w_2))$$

Similarly higher order entropies can be defined





Example for 2-Word Entropy

Let w₁ and w₂ statistically independent

$$p(w_1, w_2) = p(w_1)p(w_2)$$

$$H(W,W) = -\sum_{w_1,w_2} p(w_1, w_2) \log(p(w_1, w_2))$$

$$= -\sum_{w_1, w_2} p(w_1) p(w_2) \log(p(w_1) p(w_2))$$

$$= -\sum_{w_1, w_2} p(w_1) p(w_2) [\log p(w_1) + \log p(w_2)]$$





Example for 2-Word Entropy

$$= -\sum_{w_1, w_2} p(w_1) p(w_2) [\log p(w_1) + \log p(w_2)]$$

$$= -\sum_{w_1, w_2} p(w_1) p(w_2) \log p(w_1) - \sum_{w_1, w_2} p(w_1) p(w_2) \log p(w_2)$$

$$= -\sum_{w_1} p(w_1) \log p(w_1) - \sum_{w_2} p(w_2) \log p(w_2)$$

$$=2H(W)$$

$$H(W,W) = 2H(W)$$

H(W,W) = 2H(W)for statistical independence



Definition: Kullback-Leibler (KL) Divergence



$$D(p \parallel q) = \sum_{i} p_{i} \log \frac{p_{i}}{q_{i}}$$

Remark:

- $-p_i$ is a short hand for $p(w_i)$
- $-q_i$ is a short hand for $q(w_i)$
- -p_i and q_i are different distributions
- -D(p||p)=0
- -D(p||q) is **not** symmetric (D(p||q) is not equal to D(q||p))



Kullback-Leibler (KL) Divergence and Entropy



$$D(p || q) = \sum_{i} p_{i} \log \frac{p_{i}}{q_{i}}$$

$$= -\sum_{i} p_{i} \log q_{i} + \sum_{i} p_{i} \log p_{i}$$

$$= E_{p}[-\log(q)] - E_{p}[-\log(p)]$$

$$= H(p,q) - H(p)$$

Cross Entropy

Entropy

We will see later: KL divergence measures how many more bits need to be transmitted if you have a mismatch in the model of the language



Kullback-Leibler (KL) Divergence: Example



$$q_A = 1/2; q_B = 1/4; q_C = 1/4;$$

 $p_A = 1/3; p_B = 1/3; p_C = 1/3;$

See white board



Theorem: KL-Divergence is nonnegative



$$D(p \parallel q) \ge 0$$

Proof:

$$D(p \parallel q) = -\sum_{i} p_{i} \log \frac{q_{i}}{p_{i}}$$

$$\geq \sum_{i} p_{i} \left(1 - \frac{q_{i}}{p_{i}} \right)$$

$$\geq \sum_{i} p_{i} - \sum_{i} q_{i} = 1 - 1 = 0$$
Use
$$-\log(x) \geq 1 - x \quad \forall x > 0$$





Theorem: H(W,W) bounded by 2H(W)

Theorem

$$H(W,W) \le 2H(W)$$

Proof: $D(p(w_1, w_2) || p(w_1) p(w_2))$

$$= \sum_{w_1, w_2} p(w_1, w_2) \log \frac{p(w_1, w_2)}{p(w_1) p(w_2)}$$

$$= \sum_{w_1,w_2} p(w_1,w_2) \log p(w_1,w_2) - \sum_{w_1,w_2} p(w_1,w_2) \log p(w_1) - \sum_{w_1,w_2} p(w_1,w_2) \log p(w_2)$$





Theorem: H(W,W) bounded by 2H(W)

$$= \sum_{w_1, w_2} p(w_1, w_2) \log p(w_1, w_2) - \sum_{w_1, w_2} p(w_1, w_2) \log p(w_1) - \sum_{w_1, w_2} p(w_1, w_2) \log p(w_2)$$

$$= \sum_{w_1, w_2} p(w_1, w_2) \log p(w_1, w_2) - \sum_{w_1} p(w_1) \log p(w_1) - \sum_{w_2} p(w_2) \log p(w_2)$$

$$=-H(W,W)+H(W)+H(W) \ge 0$$

Hence:

$$H(W,W) \le 2H(W)$$





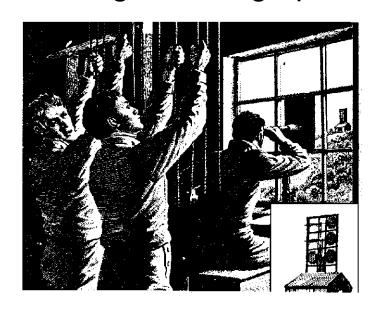
Section 4.2. Text Compression



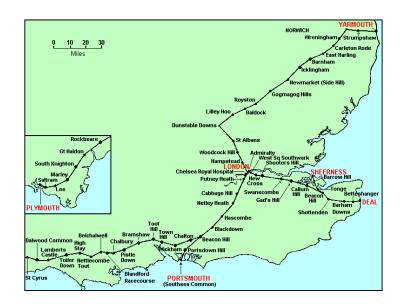
Text Compression 19th Century (Admiralty Shutter Telegraph, 1825)



Design of Telegraph



Lines in southern England



- Visual transmission of a message using 6 bit encoding
- Very limited bandwidth makes good compression necessary

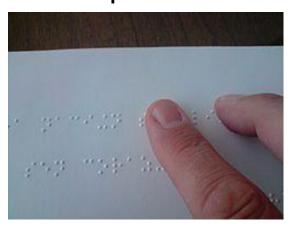
Source: http://www.douglas-self.com/MUSEUM/COMMS/telegraf/telegraf.htm



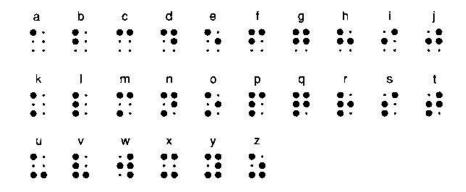
Braille-Letters: tactile writing system for the blind



Example of Text



Encoding



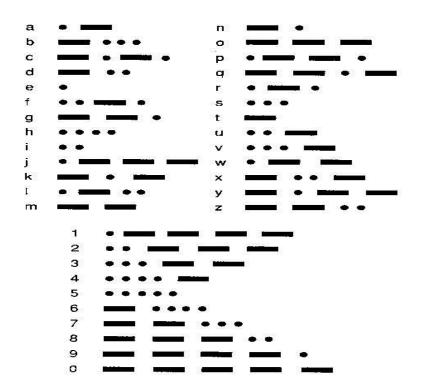
Source: Wikipedia



Morse Code



a method of transmitting text information as a series of log and short tones



Note: frequent characters Have short code words





Definition: Code

A code C(w_i) is a mapping of words w_i to finite strings of a D-nary alphabet

Example: suppose you want to code the most frequent English words "the", "and", "of" and "he" a possible binary code would be:







What is the decoded message corresponding to:

10101101110

What makes a code uniquely decodable?



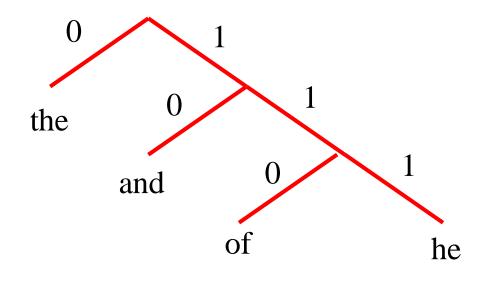




- No code word is prefix of another code word
- Organize code as a tree

Example for encoding:

Corresponding tree:









Length I_i: number of letters of D-nary alphabet

Word	Code	Length l_i
"the"	0	1
"and"	10	2
"of"	110	3
"he"	111	3

Observe:

$$\sum_{i=1}^{4} 2^{-l_i}$$

$$= \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{8}$$

$$= 1$$





Theorem: Kraft's Inequality

For every prefix code with code word length I₁ to I_m the following inequality holds true

$$\sum_{i=1}^{m} D^{-l_i} \leq 1$$

For each length distribution that satisfies this inequality there is a prefix code

D: base of the code (e.g. D=2 for binary code)

M: number of code words







The expected length (per word) of a text after coding is

$$L = \sum_{i} l_{i} p(w_{i})$$

The optimal length of code words C(w_i) is

$$l_i = -\log_D p(w_i)$$

This choice minimizes L





Proof: Optimal length distribution

The expected length of a text after coding is

$$L = \sum_{i} l_{i} p(w_{i})$$

with

$$\sum_{i=1}^{m} D^{-l_i} \le 1$$





Optimize a modified function

$$L^* = \sum_{i} l_i p(w_i) + \lambda \left(\sum_{i} D^{-l_i} - 1 \right)$$

 λ :

Lagrange multiplier (to be determined later)

Derived from constraint

$$\sum_{i=1}^m D^{-l_i} \le 1$$





Optimize a modified function

$$L^* = \sum_{i} l_i p(w_i) + \lambda \left(\sum_{i} D^{-l_i} - 1 \right)$$

Calculate first derivative

$$\frac{\partial L^*}{\partial l_j} = p(w_j) + \lambda \frac{\partial D^{-l_j}}{\partial l_j} = 0$$





$$\frac{\partial L^*}{\partial l_j} = p(w_j) + \lambda \frac{\partial D^{-l_j}}{\partial l_j} = 0$$



solve for derivative
$$\frac{\partial D^{-l_j}}{\partial l_j} = \frac{\partial}{\partial l_j} e^{-l_j \log(D)} = -\log(D) e^{-l_j \log(D)} = -\log(D) D^{-l_j}$$

insert back

$$\frac{\partial L^*}{\partial l_j} = p(w_j) - \lambda \log(D) D^{-l_j} = 0$$

solve

$$D^{-l_j} = \frac{p(w_j)}{\lambda \log(D)}$$





Goal: determine λ

Insert
$$D^{-l_i} = \frac{p(w_i)}{\lambda \log(D)}$$
 into $\sum_{i=1}^m D^{-l_i} = 1$

Equation to solve:

$$1 = \sum_{i=1}^{m} \frac{p(w_i)}{\lambda \log(D)} = \frac{1}{\lambda \log(D)} \qquad \Rightarrow \qquad \lambda = \frac{1}{\log(D)}$$

Insert λ into first equation:

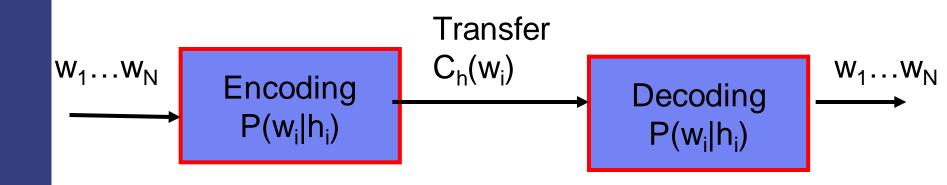
$$p(w_i) = D^{-l_i}$$

$$l_i = -\frac{\log p(w_i)}{\log D} = -\log_D p(w_i)$$





Text Compression and Perplexity



- Code words C_h(w_i) depend on history
- Length of code words $C_h(w_i) : I_h(w_i) = -log_D P(w_i|h_i)$

Question: what is the length of the compressed message?





Length of Compressed Message

$$L_{total} = \sum_{w,h} N(h, w) l_h(w) = N \sum_{w,h} f(h, w) \left(-\log_D P(w \mid h) \right)$$

$$=-N\sum_{w,h} f(h,w) \frac{\log P(w \mid h)}{\log D} = \frac{N}{\log D} \left(-\sum_{w,h} f(h,w) \log P(w \mid h)\right)$$

$$= \frac{N}{\log D} \log PP$$

$$L_{total} = \frac{N}{\log D} \log PP$$

Length of compressed text is proportional to logarithm of perplexity



Limits of the optimal Code



Theorem (entropy-bound):

Let

$$I_i = int(-log_D(p(w_i))+1)$$

then

$$H_D(W) \le L < H_D(W) + 1$$

and there is a prefix-code for li

Proof: direct calculation





Compression with a mismatched code

Theorem: The probabilities of words w_i of a text are p_i , however, it is compressed based on q_i with l_i =-log q_i

Then the length of the compressed text is

$$H(p)+D(p||q) \le L < H(p)+D(p||q)+1$$





Section 3.3. Long Range Dependencies and Conditional Entropy





Two Alternatives

- Correlation function
- Conditional Entropy







Correlation Function:

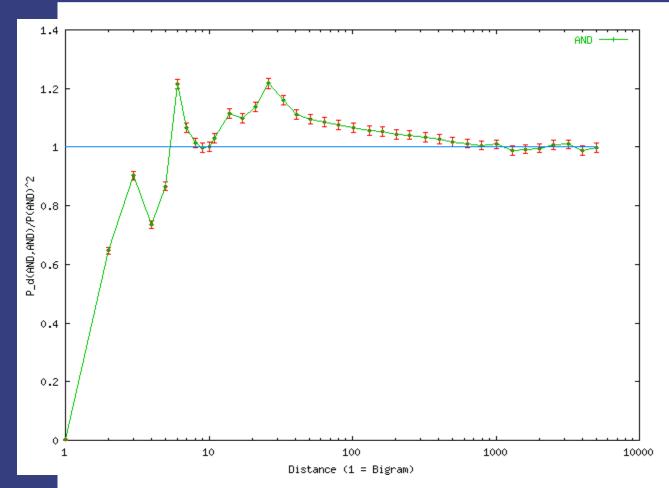
$$c_d(w) = \frac{P_d(w w)}{P(w)^2}$$

- d: distance between two observations of word w
- Statistical independence: c(w)=1
- Measures correlation





Correlation Function "and"

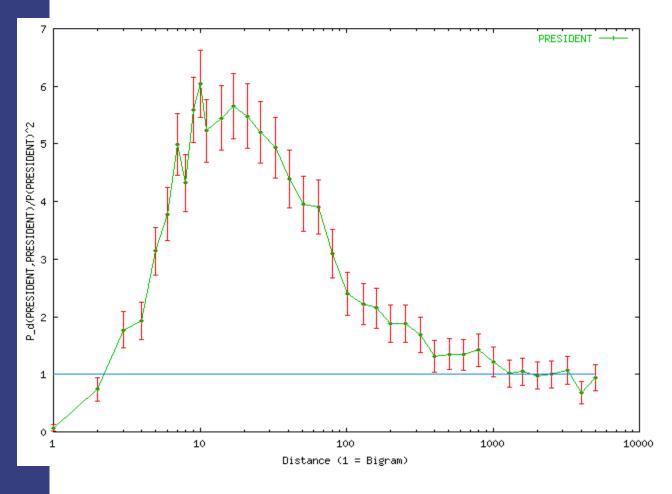


Only weak short range dependencies





Correlation Function "President"

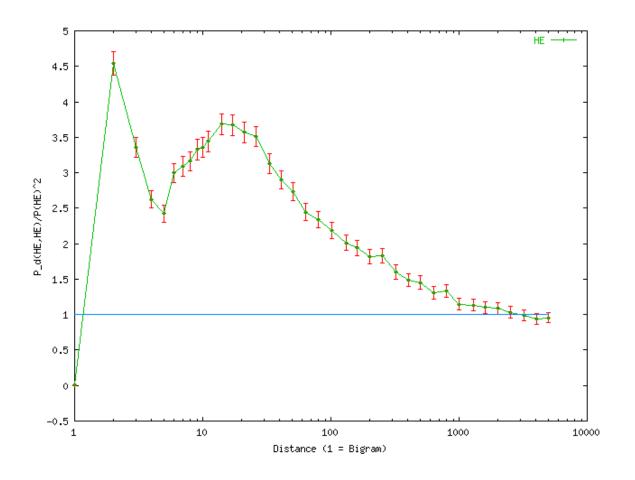


- -Long range (semantic) dependency
- -Decay of correlations after about 1000 words





Correlation Function "he"

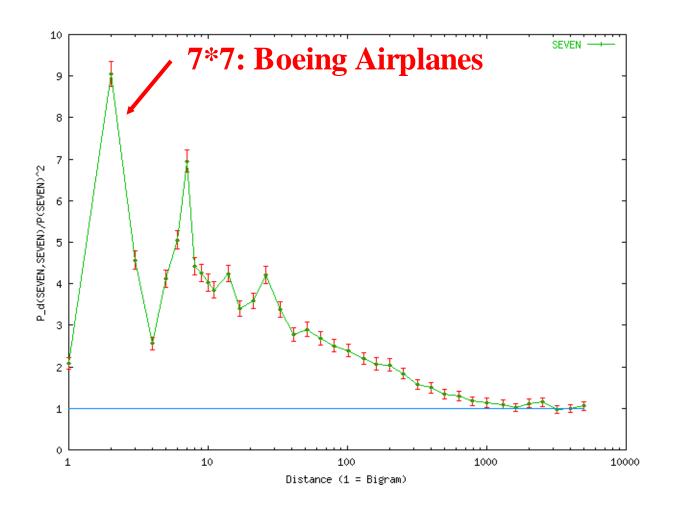


Short- and Long Range Dependencies





Correlation Function "seven"







Conditional Entropy

Definition

$$h_n = -\sum_{w_1...w_{n+1}} P(w_1...w_{n+1}) \log P(w_{n+1} | w_1...w_n)$$

Study examples to understand meaning of conditional entropy



Example I: independence of words (Bernoulli Sequence)



Assume
$$P(w_{1}...w_{n+1}) = \prod_{i=1}^{n+1} P(w_{i})$$

$$h_{n} = -\sum_{w_{1}...w_{n+1}} P(w_{1}...w_{n+1}) \log P(w_{n+1} | w_{1}...w_{n})$$

$$= -\sum_{w_{1}...w_{n+1}} \left(\prod_{i=1}^{n+1} P(w_{i}) \right) \log P(w_{n+1}) = -\sum_{w_{n+1}} P(w_{n+1}) \log P(w_{n+1})$$

=H(W)

⇒conditional entropy is independent of length ⇒conditional entropy equals one word entropy





Example II: periodic sequence

- After p symbols, the sequence repeats itself
- Example of sequence with period p=3
 - ABCABCABC...

After one period p we have $P(w_{p+1} | w_1...w_p) = 1$

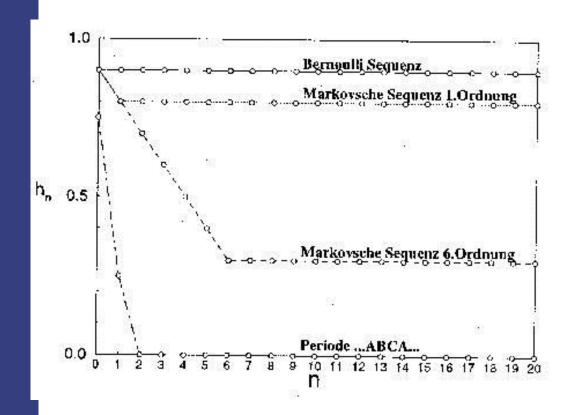
In general
$$P(w_{p+1+i} | w_1...w_{p+i}) = 1$$
 for all $i \ge 1$

$$\Rightarrow h_n = 0 \text{ for } n \ge p$$







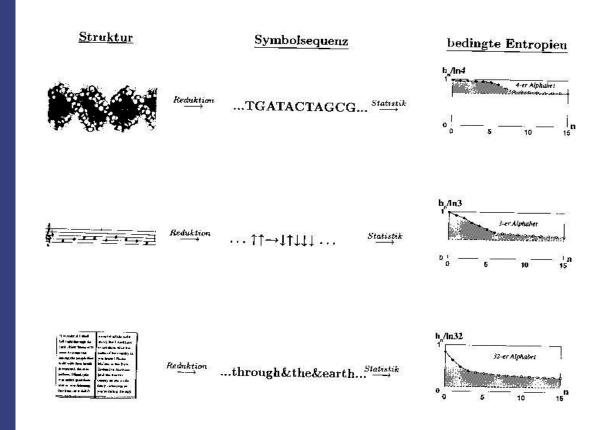


Prototypical artificial cases





Conditional Entropy: Other "Languages"



From: Ebeling





Summary

- Entropy measures the information content of a message in bits
- It is related to the compressibility of text
- Probabilities of symbols determine optimal code length
- Lower perplexity → better compression
- Long range dependencies in language:
 - Correlation function and conditional entropy