

Introduction to Information Theory

Part 3

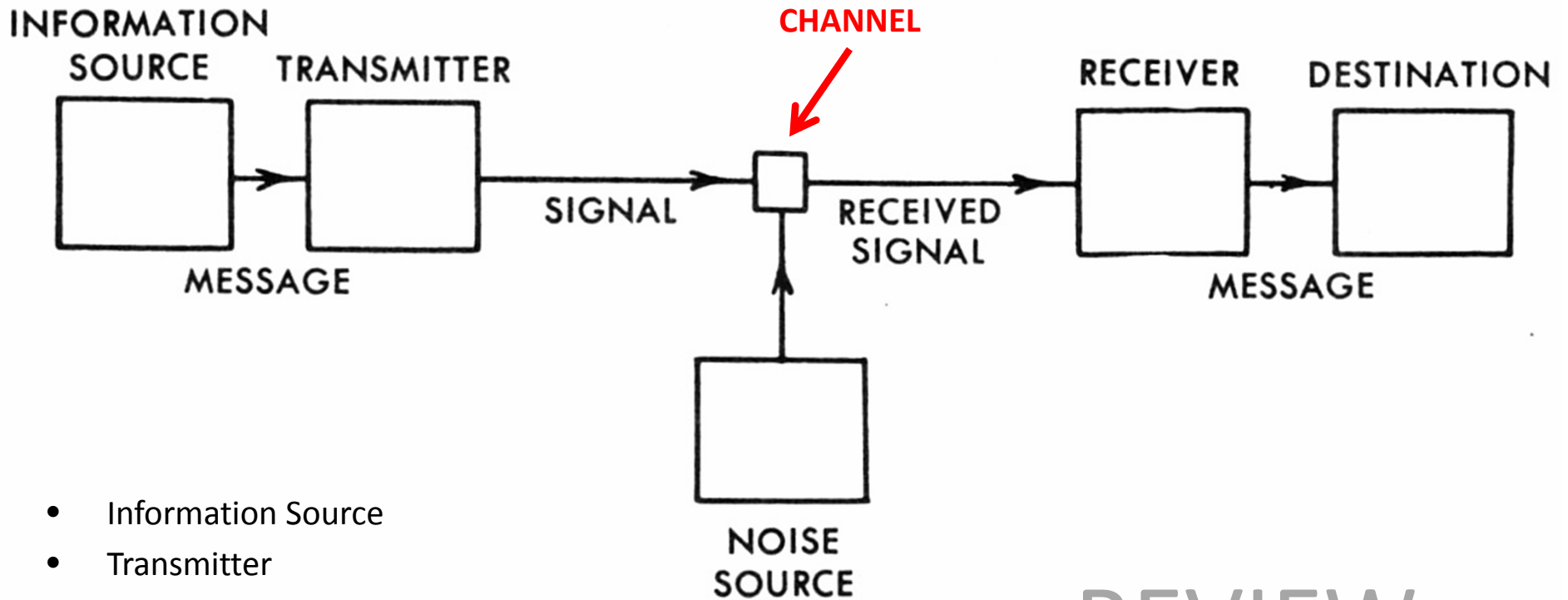
Assignment#1 Results

- List text(s) used, total # letters, computed entropy of text. Compare results.
- What is the computed average word length of 3-letter codes on

$$S = \{A, B\}, P = \{0.75, 0.25\}$$

Compare results.

A General Communication System



- Information Source
- Transmitter
- Channel
- Receiver
- Destination

9/25/2012

REVIEW...

Information Source

- Definition of Information:

$$I(p) = \log\left(\frac{1}{p}\right) = -\log(p)$$

- **Information (I)** is associated with known events/messages
- **Entropy (H)** is the average information w.r.to all possible outcomes.

Given, $P = \{p_1, p_2, \dots, p_3\}$

$$H(P) = \sum_i p_i \log\left(\frac{1}{p_i}\right)$$

REVIEW...

Source Coding: Basics

- **Block code:** When all codes have the same length. For example, ASCII (8-bits)
- **Average Word Length:**

$$L = \sum_{i=1}^m p_i l_i$$

More generally,

$$L_n = \frac{1}{n} \sum_{i=1}^m p_i l_i$$

- A code is **efficient** if it has the smallest average word length. (Turns out entropy is the benchmark...)

REVIEW...

Average Code Length & Entropy

- Average length bounds: $H \leq L < H + 1$
- Grouping n symbols together:

$$H(S^n) \leq L < H(S^n) + 1$$

$$nH(S) \leq L < nH(S) + 1$$

$$H(S) \leq \frac{L}{n} < H(S) + \frac{1}{n}$$

$$\lim_{n \rightarrow \infty} \frac{L_n}{n} = H$$

REVIEW...

Shannon's First Theorem

- By coding sequences of independent symbols (in S^n), it is possible to construct codes such that

$$\lim_{n \rightarrow \infty} \frac{L_n}{n} = H$$

The price paid for such improvement is increased coding complexity (due to increased n) and increased delay in coding.

REVIEW...

Entropy & Coding: Central Ideas

- Use short codes for highly likely events. This shortens the average length of coded messages.
- Code several events at a time. Provides greater flexibility in code design.
- Shannon's Source Coding Theorem
- Algorithms
- Applications

REVIEW...

Source Coding

- **Efficient** codes
- **Singular** codes
- **Nonsingular** codes
- **Instantaneous** codes
- **Universal Codes**

Codes that do not require knowledge of probability distribution but act in the limit as if they did have the knowledge.

Huffman Codes

- Nonsingular
- Instantaneous
- Efficient
- Non-unique
- Powers of a source lead closer to H
- Requires knowledge of symbol probabilities

REVIEW...

Entropy & Coding: Central Ideas

- Use short codes for highly likely events. This shortens the average length of coded messages.
- Code several events at a time. Provides greater flexibility in code design.
- Shannon's Source Coding Theorem
- Algorithms: Huffman Encoding, ...
- Applications: Compression... REVIEW...

xkcd(#936): Password Strength

<p>UNCOMMON (NON-GIBBERISH) BASE WORD</p> <p>ORDER UNKNOWN</p> <p>Tr0ub4dor &3</p> <p>CAPS? COMMON SUBSTITUTIONS NUMERAL PUNCTUATION</p> <p>(YOU CAN ADD A FEW MORE BITS TO ACCOUNT FOR THE FACT THAT THIS IS ONLY ONE OF A FEW COMMON FORMATS)</p>	<p>~28 BITS OF ENTROPY</p> <p>$2^{28} = 3 \text{ DAYS AT } 1000 \text{ GUESSES/SEC}$</p> <p>(PLAUSIBLE ATTACK ON A WEAK REMOTE WEB SERVICE. YES, CRACKING A STOLEN HASH IS FASTER, BUT IT'S NOT WHAT THE AVERAGE USER SHOULD WORRY ABOUT.)</p> <p>DIFFICULTY TO GUESS: EASY</p>	<p>WAS IT TROMBONE? NO, TROUBADOR. AND ONE OF THE 0s WAS A ZERO?</p> <p>AND THERE WAS SOME SYMBOL...</p> <p>DIFFICULTY TO REMEMBER: HARD</p>
<p>correct horse battery staple</p> <p>FOUR RANDOM COMMON WORDS</p>	<p>~44 BITS OF ENTROPY</p> <p>$2^{44} = 550 \text{ YEARS AT } 1000 \text{ GUESSES/SEC}$</p> <p>DIFFICULTY TO GUESS: HARD</p>	<p>THAT'S A BATTERY STAPLE.</p> <p>CORRECT!</p> <p>DIFFICULTY TO REMEMBER: YOU'VE ALREADY MEMORIZED IT</p>

THROUGH 20 YEARS OF EFFORT, WE'VE SUCCESSFULLY TRAINED EVERYONE TO USE PASSWORDS THAT ARE HARD FOR HUMANS TO REMEMBER, BUT EASY FOR COMPUTERS TO GUESS.

Some Observations

- Successive symbols from a source are not always independent. E.g. in english,

h is more likely to occur following a *t*.
- This **intersymbol dependency** must be accounted for in an accurate measure of entropy.

Lossless Compression: English Text

- How much can we compress a given text?
- What is the accurate measure of entropy of English texts?
- Zeroth-Order entropy

$$H_0 \leq \log\left(\frac{1}{27}\right) \\ \leq 4.755 \text{ bits/letter}$$

- First-Order Entropy:
- Second-Order Entropy:

$$H_1 = 3.3$$

$$H_2 = 3.1$$

Zipf's Law, 1949

$$P_n \sim 1/n^a$$

P_n is the frequency of occurrence of the n^{th} ranked item and a is close to 1.

- The most frequent word will occur approximately twice as often as the second most frequent word, three times as often as the third most frequent word, etc.
- For example, in a corpus of over 1 million words:

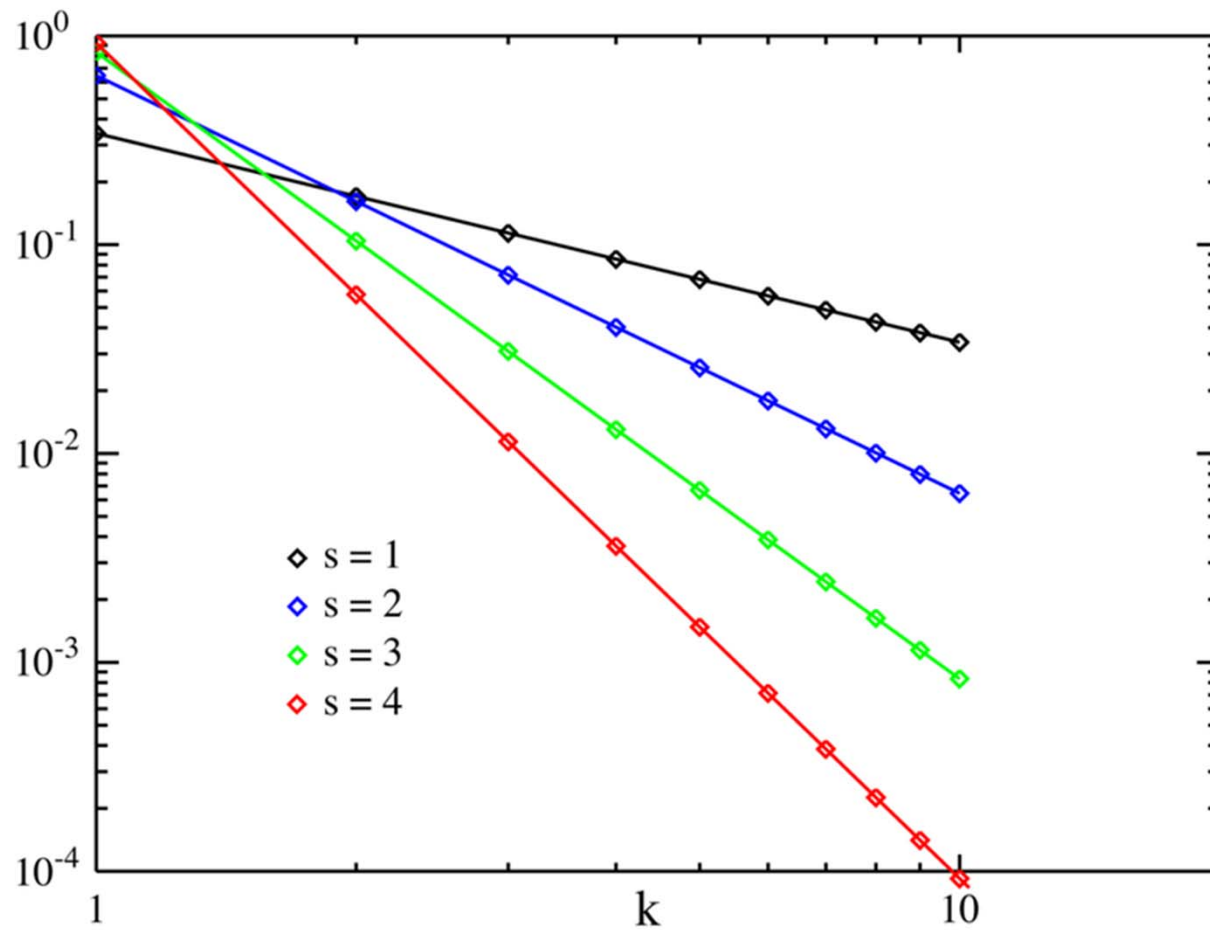
the	69,971	7%
of	36,411	3.5%
and	28,852	2.9%

- For English text:

$$p_m = \frac{A}{m}$$

where m is the word's rank, p_m is the probability of the m^{th} rank word, A is a constant that depends on the number of active words in the language.

Zipf's Law



Estimating Entropy of English Text

- For English text, W with M words:

$$p_m = \frac{A}{m}$$

where m is the word's rank, p_m is the probability of the m^{th} rank word, A is a constant that depends on the number of active words in the language.

- If $A = 0.1$, so that

$$\sum_{m=1}^M p_m = 1$$

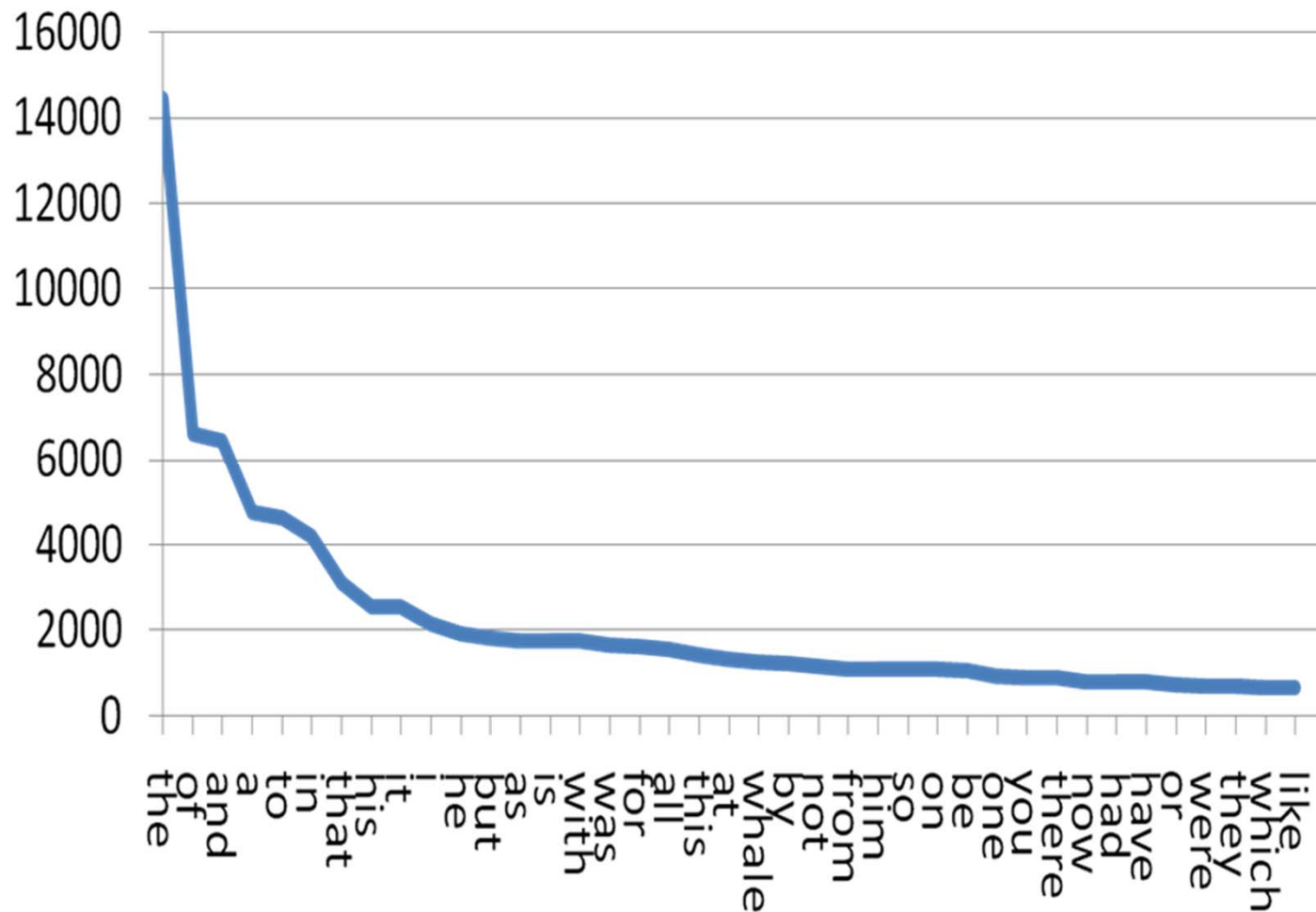
We need, $M = 12,366$.

$$H_W = \sum_{m=1}^{m=12,366} \frac{.01}{m} \log\left(\frac{m}{.01}\right) = 9.72 \text{ bits/word}$$

If $\bar{w}=4.5$ letters/word

$$H = \frac{9.72}{4.5} = 2.16 \text{ bits/word}$$

The “Long Tail” of Moby Dick



Shannon Redundancy

$$R = 1 - \frac{H}{\log M}$$

Where H is the per letter entropy, M is the size of the source alphabet. Thus redundancy of English is

$$1 - \frac{2.16}{\log 27} = 54.6\%$$

With an entropy of 1.5 we get 67% redundancy.

I.e. Huffman coding (even with an entropy of 3.3 or 3.1) will not get close to the theoretical limit.

Can we achieve compression rates close to 33%???

Lempel-Ziv Coding

- Sequences of text repeat patterns (words, phrases, etc)
- Construct a dictionary of common patterns
- Send references to patterns as triples (x, y, z)

Lempel-Ziv Coding (LZ77)



Lempel-Ziv Coding (LZ77)

Message				Search Buffer	Look-Ahead Buffer
					THIS-THESIS-IS-THE-THESIS.
0	0	T	T		THIS-THESIS-IS-THE-THESIS.

Lempel-Ziv Coding (LZ77)

Message				Search Buffer	Look-Ahead Buffer
					THIS-THESIS-IS-THE-THESIS.
0	0	T	T		THIS-THESIS-IS-THE-THESIS.
0	0	H	H		THIS-THESIS-IS-THE-THESIS.

Lempel-Ziv Coding (LZ77)

Message				Search Buffer	Look-Ahead Buffer
					THIS-THESIS-IS-THE-THESIS.
0	0	T	T		THIS-THESIS-IS-THE-THESIS.
0	0	H	H		THIS-THESIS-IS-THE-THESIS.
0	0	I	I		THIS-THESIS-IS-THE-THESIS.

Lempel-Ziv Coding (LZ77)

Message				Search Buffer	Look-Ahead Buffer
					THIS-THESIS-IS-THE-THESIS.
0	0	T	T		THIS-THESIS-IS-THE-THESIS.
0	0	H	H		THIS-THESIS-IS-THE-THESIS.
0	0	I	I		THIS-THESIS-IS-THE-THESIS.
0	0	S	S		THIS-THESIS-IS-THE-THESIS.

Lempel-Ziv Coding (LZ77)

Message				Search Buffer	Look-Ahead Buffer
					THIS-THESIS-IS-THE-THESIS.
0	0	T	T		THIS-THESIS-IS-THE-THESIS.
0	0	H	H		THIS-THESIS-IS-THE-THESIS.
0	0	I	I		THIS-THESIS-IS-THE-THESIS.
0	0	S	S		THIS-THESIS-IS-THE-THESIS.
0	0	-	-		THIS-THESIS-IS-THE-THESIS.

Lempel-Ziv Coding (LZ77)

Message				Search Buffer	Look-Ahead Buffer
					THIS-THESIS-IS-THE-THESIS.
0	0	T	T		THIS-THESIS-IS-THE-THESIS.
0	0	H	H		THIS-THESIS-IS-THE-THESIS.
0	0	I	I		THIS-THESIS-IS-THE-THESIS.
0	0	S	S		THIS-THESIS-IS-THE-THESIS.
0	0	-	-		THIS-THESIS-IS-THE-THESIS.
5	2	E	E		THIS-THE SIS-IS-THE-THESIS.

Lempel-Ziv Coding (LZ77)

Message				Search Buffer	Look-Ahead Buffer
					THIS-THESIS-IS-THE-THESIS.
0	0	T	T		THIS-THESIS-IS-THE-THESIS.
0	0	H	H		THIS-THESIS-IS-THE-THESIS.
0	0	I	I		THIS-THESIS-IS-THE-THESIS.
0	0	S	S		THIS-THESIS-IS-THE-THESIS.
0	0	-	-		THIS-THESIS-IS-THE-THESIS.
5	2	E	THE		THIS-THE SIS-IS-THE-THESIS.
5	1	I	SI		THIS-THESIS-IS-THE-THESIS.

Lempel-Ziv Coding (LZ77)

Message				Search Buffer	Look-Ahead Buffer
					THIS-THESIS-IS-THE-THESIS.
0	0	T	T		THIS-THESIS-IS-THE-THESIS.
0	0	H	H		THIS-THESIS-IS-THE-THESIS.
0	0	I	I		THIS-THESIS-IS-THE-THESIS.
0	0	S	S		THIS-THESIS-IS-THE-THESIS.
0	0	-	-		THIS-THESIS-IS-THE-THESIS.
5	2	E	THE		THIS-THE SIS-IS-THE-THESIS.
5	1	I	SI		THIS-THESIS-IS-THE-THESIS.
7	2	I	SI		THIS-THESIS-IS-THE-THESIS.

Lempel-Ziv Coding (LZ77)

Message				Search Buffer	Look-Ahead Buffer
					THIS-THESIS-IS-THE-THESIS.
0	0	T	T		THIS-THESIS-IS-THE-THESIS.
0	0	H	H		THIS-THESIS-IS-THE-THESIS.
0	0	I	I		THIS-THESIS-IS-THE-THESIS.
0	0	S	S		THIS-THESIS-IS-THE-THESIS.
0	0	-	-		THIS-THESIS-IS-THE-THESIS.
5	2	E	THE		THIS-THE SIS-IS-THE-THESIS.
5	1	I	SI		THIS-THESIS-IS-THE-THESIS.
7	2	I	SI		THIS-THESIS-IS-THE-THESIS.
10	5	-	S-THE-		THIS-THESIS-IS-THE-THESIS.

Lempel-Ziv Coding (LZ77)

Message				Search Buffer	Look-Ahead Buffer
					THIS-THESIS-IS-THE-THESIS.
0	0	T	T		THIS-THESIS-IS-THE-THESIS.
0	0	H	H		THIS-THESIS-IS-THE-THESIS.
0	0	I	I		THIS-THESIS-IS-THE-THESIS.
0	0	S	S		THIS-THESIS-IS-THE-THESIS.
0	0	-	-		THIS-THESIS-IS-THE-THESIS.
5	2	E	THE		THIS-THE SIS-IS-THE-THESIS.
5	1	I	SI		THIS-THESIS-IS-IS-THE-THESIS.
7	2	I	SI		THIS-THESIS-IS-THE-THESIS.
10	5	-	S-THE-		THIS-THESIS-IS-THE-THESIS.
14	6	.	THESIS.	THIS-THESIS-IS-THE-THESIS.	

Lempel-Ziv Coding

- Sequences of text repeat patterns (words, phrases, etc)
- Construct a dictionary of common patterns
- Send references to patterns as triples (x, y, z)
e.g. $(5, 3, F)$
go back 5 received chars
take the next 3 from there
add F to the end
- Size of Search Buffer and Look-Ahead Buffer is finite.
- Used by ZIP, PKZip, Lharc, PNG, gzip, ARJ
- Extended to LZ78 (uses dictionary), LZW (+Terry Welch)
- Achieves optimal rate of transmission in the long run w/o using probability dist.

Decode

Message

0	0	I	
0	0	-	
0	0	M	
3	1	S	
1	1	-	
5	5	L	
5	3	Y	

Decode

Message

0	0	I	I
0	0	-	
0	0	M	
3	1	S	
1	1	-	
5	5	L	
5	3	Y	

Decode

Message

0	0	I	I
0	0	-	I-
0	0	M	
3	1	S	
1	1	-	
5	5	L	
5	3	Y	

Decode

Message

0	0	I	I
0	0	-	I -
0	0	M	I - M
3	1	S	
1	1	-	
5	5	L	
5	3	Y	

Decode

Message

0	0	I	I
0	0	-	I -
0	0	M	I - M
3	1	S	I - M I S
1	1	-	
5	5	L	
5	3	Y	

Decode

Message

0	0	I	I
0	0	-	I -
0	0	M	I - M
3	1	S	I - M I S
1	1	-	I - M I S S -
5	5	L	
5	3	Y	

Decode

Message

0	0	I	I
0	0	-	I -
0	0	M	I - M
3	1	S	I - MIS
1	1	-	I - MISS -
5	5	L	I - MISS - MISS - L
5	3	Y	

Decode

Message

0	0	I	I
0	0	-	I -
0	0	M	I - M
3	1	S	I - MIS
1	1	-	I - MISS -
5	5	L	I - MISS - MISS - L
5	3	Y	I - MISS - MISS - LISSY

References

- Eugene Chiu, Jocelyn Lin, Brok McFerron, Noshirwan Petigara, Satwiksai Seshasai: *Mathematical Theory of Claude Shannon: A study of the style and context of his work up to the genesis of information theory.* MIT 6.933J / STS.420J The Structure of Engineering Revolutions
- Luciano Floridi, 2010: *Information: A Very Short Introduction*, Oxford University Press, 2011.
- Luciano Floridi, 2011: *The Philosophy of Information*, Oxford University Press, 2011.
- James Gleick, 2011: *The Information: A History, A Theory, A Flood*, Pantheon Books, 2011.
- Zhandong Liu , Santosh S Venkatesh and Carlo C Maley, 2008: *Sequence space coverage, entropy of genomes and the potential to detect non-human DNA in human samples*, *BMC Genomics* 2008, **9**:509
- David Luenberger, 2006: *Information Science*, Princeton University Press, 2006.
- David J.C. MacKay, 2003: *Information Theory, Inference, and Learning Algorithms*, Cambridge University Press, 2003.
- Claude Shannon & Warren Weaver, 1949: *The Mathematical Theory of Communication*, University of Illinois Press, 1949.
- W. N. Francis and H. Kucera: *Brown University Standard Corpus of Present-Day American English*, Brown University, 1967.
- Edward L. Glaeser: A Tale of Many Cities, New York Times, April 10, 2010. Available at: <http://economix.blogs.nytimes.com/2010/04/20/a-tale-of-many-cities/>
- Alan Rimm-Kaufman, The Long Tail of Search. Search Engine Land Website, September 18, 2007. Available at: <http://searchengineland.com/the-long-tail-of-search-12198>