



FOUNDATIONS OF CODING

**Theory and Applications
of Error-Correcting Codes
with an Introduction to
Cryptography and Information Theory**

Jiří Adámek

Czech Technical University in Prague



A Wiley-Interscience Publication
JOHN WILEY & SONS, INC.
Chichester • New York • Brisbane • Toronto • Singapore

This page intentionally left blank

FOUNDATIONS OF CODING

This page intentionally left blank

FOUNDATIONS OF CODING

**Theory and Applications
of Error-Correcting Codes
with an Introduction to
Cryptography and Information Theory**

Jiří Adámek

Czech Technical University in Prague



A Wiley-Interscience Publication

JOHN WILEY & SONS, INC.

Chichester • New York • Brisbane • Toronto • Singapore

A NOTE TO THE READER

This book has been electronically reproduced from digital information stored at John Wiley & Sons, Inc. We are pleased that the use of this new technology will enable us to keep works of enduring scholarly value in print as long as there is a reasonable demand for them. The content of this book is identical to previous printings.

In recognition of the importance of preserving what has been written, it is a policy of John Wiley & Sons, Inc., to have books of enduring value published in the United States printed on acid-free paper, and we exert our best efforts to that end.

Copyright © 1991 by John Wiley & Sons, Inc.

All rights reserved. Published simultaneously in Canada.

Reproduction or translation of any part of this work beyond that permitted by Section 107 or 108 of the 1976 United States Copyright Act without the permission of the copyright owner is unlawful. Requests for permission or further information should be addressed to the Permissions Department, John Wiley & Sons, Inc.

Library of Congress Cataloging in Publication Data:

Adámek, Jiří, ing.

Foundations of coding: theory and applications of error-correcting codes, with an introduction to cryptography and information theory / Jiří Adámek.

p. cm.

"A Wiley-Interscience publication."

Includes bibliographical references and index.

ISBN 0-471-62187-0

I. Coding theory. I. Title.

QA268.A36 1991

003:54—dc20

90-20905

CIP

Rev.

To Honza, Kuba, Jirka, and Ondřej

This page intentionally left blank

Preface

Coding theory is a fascinating field combining elegant mathematical theories with constructions of a major practical impact.

This book is devoted to constructions of

- (1) error-correcting codes,
- (2) secrecy codes, and
- (3) codes used in data compression.

The stress is on the first direction: we introduce a number of important classes of error-detecting and error-correcting codes, and we present their decoding methods. Some of these constructions require a deep background in modern algebra, and we carefully provide such background. Secret codes are treated only briefly; we mainly explain the role of error-correcting codes in modern cryptography. Data compression and other topics related to information theory are briefly discussed in the first part of the book.

The material is presented in a way making it possible to appreciate both the beauty of the theory and the scope of practical applications. We use the definition-theorem-proof style usual in mathematical texts (since the reader can thus skip a proof to keep continuity of the text and return to it later), but formalism is avoided as much as possible.

The book evolved from a series of lectures I held at the Czech Technical University in Prague in 1985–1990. They were based primarily on the following excellent textbooks which the reader may use for further reading: *Information Theory and Coding*, Abramson (1963),* *Theory and Practise of Error Control Codes*, Blahut (1983), *The Theory of Error-Correcting Codes*, MacWilliams and Sloane (1981), and *An Introduction to Cryptology*, van Tilborg (1988).

Jiří Adámek

*A name followed by a year in parentheses refers to the list of references at the end of the book.

This page intentionally left blank

Contents

Introduction	1
Part I Coding and Information Theory	3
1 Coding and Decoding	5
1.1 Coding	5
1.2 Unique Decoding	6
1.3 Block Codes and Instantaneous Codes	7
1.4 Some Important Block Codes	9
1.5 Construction of Instantaneous Codes	11
1.6 Kraft's Inequality	12
1.7 McMillan's Theorem	13
Exercises	14
Notes	16
2 Huffman Codes	17
2.1 Information Source	17
2.2 Huffman Codes	17
2.3 Construction of Binary Huffman Codes	18
2.4 Example	21
2.5 Construction of General Huffman Codes	22
Exercises	24
Notes	24
3 Data Compression and Entropy	25
3.1 An Example of Data Compression	25
3.2 The Idea of Entropy	26
3.3 The Definition of Entropy	28
3.4 An Example	29

3.5	Maximum and Minimum Entropy	30
3.6	Extensions of a Source	32
3.7	Entropy and Average Length	33
3.8	Shannon's Noiseless Coding Theorem	34
3.9	Concluding Remarks	36
	Exercises	36
	Notes	38
4	Reliable Communication Through Unreliable Channels	39
4.1	Binary Symmetric Channels	40
4.2	Information Rate	42
4.3	An Example of Increased Reliability	44
4.4	Hamming Distance	46
4.5	Detection of Errors	48
4.6	Correction of Errors	49
4.7	Channel Capacity	50
4.8	Shannon's Fundamental Theorem	56
	Exercises	58
	Notes	60
	Part II Error-Correcting Codes	61
5	Binary Linear Codes	63
5.1	Binary Addition and Multiplication	63
5.2	Codes Described by Equations	64
5.3	Binary Linear Codes	65
5.4	Parity Check Matrix	67
5.5	Hamming Codes—Perfect Codes for Single Errors	69
5.6	The Probability of Undetected Errors	75
	Exercises	77
	Notes	78
6	Groups and Standard Arrays	79
6.1	Commutative Groups	79
6.2	Subgroups and Cosets	81
6.3	Decoding by Standard Arrays	84
	Exercises	87
	Notes	89

7	Linear Algebra	91
7.1	Fields and Rings	91
7.2	The Fields \mathbf{Z}_p	93
7.3	Linear Spaces	95
7.4	Finite-Dimensional Spaces	98
7.5	Matrices	101
7.6	Operations on Matrices	105
7.7	Orthogonal Complement	108
	Exercises	111
	Notes	114
8	Linear Codes	115
8.1	Generator Matrix	115
8.2	Parity Check Matrix	119
8.3	Syndrome	121
8.4	Detection and Correction of Errors	122
8.5	Extended Codes and Other Modifications	125
8.6	Simultaneous Correction and Detection of Errors	128
8.7	MacWilliams Identity	130
	Exercises	133
	Notes	135
9	Reed-Muller Codes: Weak Codes with Easy Decoding	137
9.1	Boolean Functions	138
9.2	Boolean Polynomials	140
9.3	Reed-Muller Codes	144
9.4	Geometric Interpretation: Three-Dimensional Case	147
9.5	Geometric Interpretation: General Case	151
9.6	Decoding Reed-Muller Codes	154
	Exercises	159
	Notes	160
10	Cyclic Codes	161
10.1	Generator Polynomial	161
10.2	Encoding Cyclic Codes	167
10.3	Parity Check Polynomial	171
10.4	Decoding Cyclic Codes	175
10.5	Error-Trapping Decoding	180
10.6	Golay Code: A Perfect Code for Triple Errors	182
10.7	Burst Errors	185
10.8	Fire Codes: High-Rate Codes for Burst Errors	188

Exercises 192

Notes 194

11 Polynomials and Finite Fields 197

 11.1 Zeros of Polynomials 197

 11.2 Algebraic Extensions of a Field 201

 11.3 Galois Fields 206

 11.4 Primitive Elements 207

 11.5 The Characteristic of a Field 211

 11.6 Minimal Polynomial 213

 11.7 Order 216

 11.8 The Structure of Finite Fields 219

 11.9 Existence of Galois Fields 221

 Exercises 223

 Notes 227

12 BCH Codes: Strong Codes Correcting Multiple Errors 229

 12.1 Hamming Codes as Cyclic Codes 230

 12.2 Double-Error-Correcting BCH Codes 232

 12.3 BCH Codes 239

 12.4 Reed-Solomon Codes and Derived Burst-Error-Correcting Codes 245

 12.5 Generalized Reed-Muller Codes 246

 12.6 Goppa Codes: Asymptotically Good Codes 248

 Exercises 255

 Notes 256

13 Fast Decoding of BCH Codes 257

 13.1 Error Location and Error Evaluation 258

 13.2 Euclidean Algorithm 260

 13.3 The Decoding Algorithm 263

 Exercises 266

 Notes 267

14 Convolutional Codes 269

 14.1 Linear Codes and Convolutional Codes 269

 14.2 Generator Polynomials and Generator Matrices 274

 14.3 Maximum-Likelihood Decoding of Convolutional Codes 279

 14.4 The Viterbi Decoding Algorithm 283

 Exercises 288

 Notes 290

Part III Cryptography	291
15 Cryptography	293
15.1 A Noisy Wiretap	294
15.2 Secret-Key Encryption	296
15.3 Public-Key Encryption	303
15.4 Encryption Based on Large Prime Numbers	305
15.5 Encryption Based on Knapsack Problems	307
15.6 Data Encryption Standard	309
Exercises	316
Notes	317
Appendixes	319
A Galois Fields	321
B BCH Codes and Reed-Muller Codes	325
Bibliography	327
List of Symbols	331
Index	333

This page intentionally left blank

FOUNDATIONS OF CODING

This page intentionally left blank

Introduction

Data transmission and data storage suffer from errors created by noise. Techniques for combatting noise have been used for a long time. They range from simple ones, e.g., adding a parity check symbol to every byte, to modern complex error-correcting techniques described in this book. The basic idea of error correction by a block code (i.e., a code in which all code words have the same length) is simple: code words must be “wide apart” from each other. That is, two distinct code words have a large Hamming distance, which means the number of symbols in which the words differ. Then the code corrects errors as follows: the word received is corrected to the nearest code word (in the sense of the Hamming distance). If the number of errors created by noise is smaller than one-half of the minimum Hamming distance of code words, then the correction is well done. Thus, the theory of error-correcting block codes is concerned with a construction of “good” codes with large Hamming distances. “Good” means that (1) the number of code words is as high as possible (to keep the redundancy low) and (2) an efficient technique for error correction is known (to make the search for the nearest code word fast).

Besides block codes, there is another class of error-correcting codes, called convolutional codes, in which memory plays a role: the message is again divided into blocks, but each block sent depends on a certain number of preceding blocks. The theory of convolutional codes is less rich than that of block codes: whereas good convolutional codes have been found by computer search, good block codes result from the algebraic theory presented in this book. However, the importance of convolutional codes in practical applications is ever increasing.

The theory of error-correcting codes is closely related to the theory of information, and the first part of this book is devoted to the foundations of information theory. Both of these theories were initiated by the pioneering paper of Claude Shannon (1948) in which he introduced entropy as a measure of information contained in an average symbol of a message. Shannon proved, *inter alia*, that entropy gives a precise estimate of how

much can be achieved by data compression. Combined with the famous Huffmann construction of the shortest code, this result of Shannon leads to a simple technique of data compression, presented in Chapters 2 and 3. (However, data compression is restricted to the case of information sources without memory.) The fourth chapter discusses the Fundamental Theorem of Shannon, which states that for every channel there exist error-correcting codes which remove noise while keeping the redundancy within the channel capacity. This result is purely theoretical: no algorithm for finding such codes has ever been found. The theory of error-correcting codes today has a lesser goal, viz., constructing codes with a reasonable redundancy and a fast decoder.

Constructions of efficient error-correcting and error-detecting codes with fast decoders are presented in the second part of the book. Some of the constructions require a deeper background in modern algebra and geometry, and we provide a thorough presentation of the relevant topics. The most important classes of error-correcting codes are the following:

Hamming codes (Chapter 5), perfect codes for single errors;

Reed-Muller codes (Chapter 9), multiple-error-correcting codes with a particularly efficient and easily implemented decoder;

Golay code (Chapter 10), the unique perfect code for triple errors;

BCH codes (Chapters 12 and 13), strong multiple-error-correcting codes with a fast decoder;

Convolutional codes (Chapter 14), multiple-error-correcting codes with memory.

The last part of the book is a short introduction to modern cryptography, stressing the role which error-correcting codes play here. Some of the well-known secret codes used in cryptography are based on constructions of error-correcting codes (e.g. the cryptosystem of McEliece, see 15.3). However, the main relation between cryptography and error-correcting codes is that, since noise is fatal for decryption, secret codes are usually combined with error-correcting codes. Furthermore, since encryption is costly, secret codes are usually combined with data compression.

The book is organized in chapters numbered consecutively throughout the three parts. Each chapter is divided into sections, and cross-references are always related to the number of section. For example, Theorem 3.2 means (the only) theorem in Section 3.2 of Chapter 3.

Part I

**Coding and Information
Theory**

This page intentionally left blank

Chapter 1

Coding and Decoding

We are often faced with the task of converting a message, i.e., a sequence of symbols from a finite set (called a *source alphabet*), into a binary message, i.e., a sequence of 0's and 1's. The most common method is to translate each source symbol into a binary word. [*Word* means precisely a finite sequence; we often write $a_1a_2 \dots a_n$ instead of the more precise (a_1, a_2, \dots, a_n) . Thus, 00101 is an example of a binary word.] Then the message is encoded symbol by symbol: we simply concatenate the words corresponding to the first, second, etc., symbol of the source message. The question of *how* to encode the source symbols is very important. There are two major criteria: we want to compress data, i.e., we want the resulting binary message to be as concise as possible, and we want to protect information against noise. These two requirements are rather contradictory, since by compressing data in the presence of noise, we are apt to increase, rather than decrease, the loss of information. In the present part, we therefore disregard noise: assuming that no errors occur, we try to find a concise code.

In the present chapter we introduce the important class of instantaneous codes, i.e., codes which can be decoded letter by letter, and we show how to construct such codes. The construction of the shortest code will be presented in Chapter 2.

1.1 Coding

Definition. *Given finite sets A (source alphabet) and B (code alphabet), a coding is a rule assigning to each source symbol exactly one word in the code alphabet, i.e., a function from A to the set of all words in B . We speak about binary coding if the code alphabet B has two symbols.*

Source Symbol	Code Word
1	11000
2	10100
3	01100
4	10010
5	01010
6	00110
7	10001
8	01001
9	00101
0	00011

Decoding	01247
----------	-------

Figure 1: 2-out-of-5 code

Example of Binary Coding. Observe that among binary words of length 5, the number of those having two 1's is $\binom{5}{2} = 10$. This can be used to the *2-out-of-5 code* of decimal digits, see Figure 1.

The message “173” has the following code: 110001000101100. Observe that no space is left between the code words since “space” is a code symbol too. Thus, for example, the famous Morse code has code alphabet $B = \{ \cdot, - , \text{space} \}$.

How do we decode the 2-out-of-5 code? Of course, the first five binary digits correspond to the first decimal one, and after decoding them, we proceed to the second group of five binary digits, etc. A helpful mnemonic rule: use 01247 as a “weight” of the five columns, and add the weights of all 1's in your words. Examples: $11000 \mapsto 0+1 = 1$ and $01100 \mapsto 1+2 = 3$. Unfortunately, 0 is an exception.

Remark. Given a coding (i.e., a function K from $A = \{a_1, \dots, a_n\}$ to the set of all words in B), the words $K(a_1), \dots, K(a_n)$ are called *code words*, and the set of all code words is called a *code*. When the concrete symbol in A is not important, “code” and “coding” are usually identified.

1.2 Unique Decoding

Definition. For each coding K (of source symbols), we define the coding of source messages as the rule K^* , which to each word $x_1x_2\dots x_m$ in the

source alphabet assigns the word $K^*(x_1x_2\dots x_m) = K(x_1)K(x_2)\dots K(x_m)$ obtained by concatenation of the code words $K(x_i)$, $i = 1, \dots, m$.

The coding K is said to be uniquely decodable provided that arbitrary two distinct source messages have distinct codes. In other words, provided that K^* is one-to-one.

For example, the 2-out-of-5 code is uniquely decodable. The assignment of a binary word to "173" is a sample of coding source messages.

In contrast, the following coding

$$a \mapsto 00 \quad b \mapsto 10 \quad c \mapsto 101 \quad d \mapsto 110 \quad e \mapsto 1001$$

is not uniquely decodable: try to decode 10110.

We now introduce two important types of uniquely decodable codes.

1.3 Block Codes and Instantaneous Codes

We now introduce two important types of codes: instantaneous codes, which are codes of variable word lengths decodable symbol per symbol, and block codes, which are the special case of instantaneous codes with constant word length:

Definition. (1) A coding using only pairwise distinct code words of a certain length n is called a block coding of length n .

(2) A coding is called instantaneous provided that no code word is a prefix of another code word; i.e., if a source symbol has a code $b_1b_2\dots b_n$ then no other source symbol has a code $b_1b_2\dots b_nb_{n+1}\dots b_m$.

Remark. Block codes (e.g., the 2-out-of-5 code above) are very convenient for decoding since we know in advance which code symbols correspond to the first (second, third, etc.) source symbol. And they are certainly efficient whenever all source symbols appear with equal frequency. However, if the frequencies of various source symbols differ substantially, then block codes become clumsy, and it is preferable to use instantaneous codes of variable lengths of words.

Examples

- (1) The famous *Morse code* is exhibited in Figure 2. This is an instantaneous code with the code alphabet $\{ \cdot, - , \text{space} \}$. Since "space" is only used at the end of each code word, the decoding procedure is simple: we always look for the first "space". There is an obvious reason for not using a block code: the frequency of, say, "E" in the English language is much higher than that of "F".

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

Figure 2: Morse code

(2) An important example of a block code is the *octal code*:

0	000	4	100
1	001	5	101
2	010	6	110
3	011	7	111

(3) Suppose that we are to find a binary coding for the alphabet $\{0, 1, 2, 3\}$, and we observe that 0 appears much more often in source messages than any other symbols. Then the following coding seems reasonable:

$$0 \mapsto 0 \quad 1 \mapsto 01 \quad 2 \mapsto 011 \quad 3 \mapsto 111.$$

We can decode quite easily: count the number of 1's among the last three symbols of the encoded message. If the number is i , then the last source symbol is i . However, the above coding is not instantaneous. Indeed, when receiving a long message

$$01111111111 \dots,$$

we will not know whether the first source symbol is 0, 1, or 2 until the message stops.

1.4 Some Important Block Codes

Long binary codes are difficult to handle. It is thus often suitable to group the binary symbols: we get shorter codes in more complex alphabets. For example, by forming groups of three symbols, we obtain the octal code, see 1.3. Representation by the octal code is usually indicated by the subscript 8. Example:

$$(01)_8 = 000001.$$

By forming groups of four binary symbols, we obtain the *hexadecimal code* in Figure 3.

Binary	Hexadecimal	Binary	Hexadecimal
0000	0	1000	8
0001	1	1001	9
0010	2	1010	A
0011	3	1011	B
0100	4	1100	C
0101	5	1101	D
0110	6	1110	E
0111	7	1111	F

Figure 3: Hexadecimal code

A very important code used for a standard binary representation of alphabetic and numeric symbols is the ASCII code* (Figure 4). It has $2^7 = 128$ source symbols encoded into binary words of length 8: the first seven symbols carry the information, and the eighth one is set in such a way that the *parity* is even (i.e., each code word contains an even number of 1's). The role of this eighth symbol is to enable detection of single errors—we explain this in detail in Chapter 5. For example, the letter A has code

$$A \longleftrightarrow \underbrace{10000010}_{\substack{\text{inform. check} \\ \text{symbols symbol}}}$$

which in Figure 4 is represented by its seven information bits: $1(01)_8 = 1000001$.

*American Standard Code for Information Interchange

Source Symbol	Code	Source Symbol	Code	Source Symbol	Code	Source Symbol	Code
@	1(00) ₈	'	1(40) ₈	NUL	0(00) ₈	SP	0(40) ₈
A	1(01) ₈	a	1(41) ₈	SOH	0(01) ₈	!	0(41) ₈
B	1(02) ₈	b	1(42) ₈	STX	0(02) ₈	"	0(42) ₈
C	1(03) ₈	c	1(43) ₈	ETX	0(03) ₈	#	0(43) ₈
D	1(04) ₈	d	1(44) ₈	EOT	0(04) ₈	\$	0(44) ₈
E	1(05) ₈	e	1(45) ₈	ENQ	0(05) ₈	%	0(45) ₈
F	1(06) ₈	f	1(46) ₈	ACK	0(06) ₈	&	0(46) ₈
G	1(07) ₈	g	1(47) ₈	BEL	0(07) ₈	'	0(47) ₈
H	1(10) ₈	h	1(50) ₈	BS	0(10) ₈	(0(50) ₈
I	1(11) ₈	i	1(51) ₈	HT	0(11) ₈)	0(51) ₈
J	1(12) ₈	j	1(52) ₈	LF	0(12) ₈	*	0(52) ₈
K	1(13) ₈	k	1(53) ₈	VT	0(13) ₈	+	0(53) ₈
L	1(14) ₈	l	1(54) ₈	FF	0(14) ₈	,	0(54) ₈
M	1(15) ₈	m	1(55) ₈	CR	0(15) ₈	-	0(55) ₈
N	1(16) ₈	n	1(56) ₈	SO	0(16) ₈	.	0(56) ₈
O	1(17) ₈	o	1(57) ₈	SI	0(17) ₈	/	0(57) ₈
P	1(20) ₈	p	1(60) ₈	DLE	0(20) ₈	0	0(60) ₈
Q	1(21) ₈	q	1(61) ₈	DC1	0(21) ₈	1	0(61) ₈
R	1(22) ₈	r	1(62) ₈	DC2	0(22) ₈	2	0(62) ₈
S	1(23) ₈	s	1(63) ₈	DC3	0(23) ₈	3	0(63) ₈
T	1(24) ₈	t	1(64) ₈	DC4	0(24) ₈	4	0(64) ₈
U	1(25) ₈	u	1(65) ₈	NAK	0(25) ₈	5	0(65) ₈
V	1(26) ₈	v	1(66) ₈	SYN	0(26) ₈	6	0(66) ₈
W	1(27) ₈	w	1(67) ₈	ETB	0(27) ₈	7	0(67) ₈
X	1(30) ₈	x	1(70) ₈	CAN	0(30) ₈	8	0(70) ₈
Y	1(31) ₈	y	1(71) ₈	EM	0(31) ₈	9	0(71) ₈
Z	1(32) ₈	z	1(72) ₈	SUB	0(32) ₈	:	0(72) ₈
{	1(33) ₈	{	1(73) ₈	ESC	0(33) ₈	;	0(73) ₈
\	1(34) ₈		1(74) ₈	FS	0(34) ₈	<	0(74) ₈
]	1(35) ₈	}	1(75) ₈	GS	0(35) ₈	=	0(75) ₈
-	1(36) ₈	-	1(76) ₈	RS	0(36) ₈	>	0(76) ₈
_	1(37) ₈	DEL	1(77) ₈	US	0(37) ₈	?	0(77) ₈

Figure 4: ASCII code (7 information bits)

Let us finally mention an “everyday” code we meet in most textbooks. It is called the *international standard book number*, ISBN, and it is a block code of length 10. (Various hyphens are often inserted between the symbols, but we can ignore them here since they are used just for optical orientation.) The code alphabet has 11 symbols: 0, 1, . . . , 9 and X (read: ten). For example, the book of Lin and Costello (1983) has

ISBN 0-13-283796-X.

The first number 0 denotes the country (USA), 13 denotes the publisher (Prentice-Hall), and the next six digits are assigned by the publisher as an identification number of the book. The last symbol is a check symbol (analogously as in the ASCII code above). It is set in such a way that for each ISBN code word $a_1a_2a_3 \dots a_9a_{10}$, the sum

$$\sum_{i=1}^{10} ia_{11-i} = 10a_1 + 9a_2 + 8a_3 + \dots + 2a_9 + a_{10}$$

be divisible by 11. For example, in the ISBN above:

$$\begin{aligned} 10 \times 0 + 9 \times 1 + 8 \times 3 + 7 \times 2 + 6 \times 8 + \\ 5 \times 3 + 4 \times 7 + 3 \times 9 + 2 \times 6 + 1 \times 10 &= 132 = 11 \times 12. \end{aligned}$$

Some publishers have a three-digit identification (e.g., Wiley-Interscience has 471) and then they assign a five-digit number to each publication.

1.5 Construction of Instantaneous Codes

Suppose you want to construct a binary instantaneous code of the source alphabet $\{a_1, \dots, a_n\}$. It is sufficient to specify the lengths d_1, \dots, d_n of the expected code words. In fact, we can certainly assume that $d_1 \leq d_2 \leq \dots \leq d_n$. Then we choose an arbitrary binary word $K(a_1)$ of length d_1 . Next we choose a binary word $K(a_2)$ of length d_2 , but we avoid all those which have the prefix $K(a_1)$. This is possible: the number of all binary words of length d_2 is 2^{d_2} . The number of those having the prefix $K(a_1)$ is $2^{d_2-d_1}$ (because you can choose the d_2-d_1 digits remaining after the prefix $K(a_1)$ arbitrarily). Since $2^{d_2} \geq 2^{d_2-d_1} + 1$, we have at least one choice of $K(a_2)$.

Next, we want to choose a word of length d_3 which has neither prefix $K(a_1)$ nor $K(a_2)$. Thus, from the 2^{d_3} possible words, we must avoid all the $2^{d_3-d_1}$ words with the prefix $K(a_1)$ and all the $2^{d_3-d_2}$ words with the prefix $K(a_2)$. This is possible if (and only if)

$$2^{d_3} \geq 2^{d_3-d_1} + 2^{d_3-d_2} + 1.$$

Dividing the last inequality by 2^{d_3} , we obtain

$$1 \geq 2^{-d_1} + 2^{-d_2} + 2^{-d_3}.$$

Analogously, we can see that the following inequality

$$1 \geq 2^{-d_1} + 2^{-d_2} + \dots + 2^{-d_n}$$

makes it possible to fulfil our task. It turns out that this inequality is both necessary and sufficient for a construction of instantaneous codes.

Example. We want to find an instantaneous code with the same lengths of code words as that in Example 1.3(3) above. This is possible since

$$1 \geq 2^{-1} + 2^{-2} + 2^{-3} + 2^{-3}.$$

Here is such a code:

$$0 \mapsto 0 \quad 1 \mapsto 10 \quad 2 \mapsto 110 \quad 3 \mapsto 111.$$

1.6 Kraft's Inequality

Theorem. *Given a source alphabet of n symbols and a code alphabet of k symbols, then an instantaneous code with given lengths d_1, d_2, \dots, d_n of code words exists, whenever the following Kraft's inequality*

$$k^{-d_1} + k^{-d_2} + \dots + k^{-d_n} \leq 1$$

is fulfilled.

PROOF. We can assume that the source alphabet $\{a_1, a_2, \dots, a_n\}$ is presented in the order imposed by the lengths of the expected code words; i.e., $d_1 \leq d_2 \leq \dots \leq d_n$. We define instantaneous coding K by the following induction:

(1) Choose an arbitrary word $K(a_1)$ of length d_1 .

(2) Suppose $K(a_1), K(a_2), \dots, K(a_{s-1})$ have been chosen. Then choose an arbitrary word $K(a_s)$ of length d_s with no prefix among $K(a_1), \dots, K(a_{s-1})$. This is possible: the number of words with the prefix $K(a_i)$ is $k^{d_s-d_i}$, and thus we have a choice of

$$k^{d_s} - \sum_{i=1}^{s-1} k^{d_s-d_i}$$

words. From Kraft's inequality, we get

$$1 - \sum_{i=1}^{s-1} k^{-d_i} \geq k^{-d_s}$$

and multiplying by k^{d_s} , we conclude

$$k^{d_s} - \sum_{i=1}^{s-1} k^{d_s-d_i} \geq 1.$$

□

1.7 McMillan's Theorem

McMillan's Theorem. *Every uniquely decodable coding satisfies Kraft's inequality.*

Remark. We see that Kraft's inequality is not only sufficient, but also necessary for the construction of an instantaneous code. However, McMillan's Theorem says much more: instantaneous codes are just as efficient as uniquely decodable codes. More precisely, for every uniquely decodable code there exists an instantaneous code with the same lengths of code words.

PROOF. Let K be a uniquely decodable coding. Denote by d_i the length of the code word $K(a_i)$, $i = 1, 2, \dots, n$. Observe that for each number $j = 1, 2, 3, \dots$, we can form exactly k^j words of length j in the (k -symbol) code alphabet. By unique decodability, the number of source messages $a_{i_1} a_{i_2} \dots a_{i_r}$ whose code has length j cannot exceed k^j . The length of the code is $d_{i_1} + d_{i_2} + \dots + d_{i_r}$. Thus, we observe that the number of all sums of the form

$$d_{i_1} + d_{i_2} + \dots + d_{i_r} = j \tag{1.7.1}$$

is smaller or equal to k^j .

It is our task to prove that the number

$$c = \sum_{i=1}^n k^{-d_i}$$

is smaller or equal to 1. For this, we will verify that the numbers $\frac{c^r}{r}$ are bounded for all $r = 1, 2, 3, \dots$. In fact, each number $c > 1$ clearly fulfils

$\lim_{r \rightarrow \infty} \frac{c^r}{r} = \infty$, and, therefore, the theorem will then be proved. Let us compute the powers of c :

$$c^2 = \left(\sum_{i=1}^n k^{-d_i} \right) \left(\sum_{j=1}^n k^{-d_j} \right) = \sum_{i,j=1}^n k^{-(d_i+d_j)}$$

and, in general,

$$c^r = \sum_{i_1, \dots, i_r=1}^n k^{-(d_{i_1}+d_{i_2}+\dots+d_{i_r})}. \quad (1.7.2)$$

We can reorder the last sum by collecting all the summands k^{-j} , where j satisfies (1.7.1). The largest possible j is $j = d + d + \dots + d = rd$, where $d = \max(d_1, \dots, d_n)$. As observed above, the number of all summands k^{-j} in sum (1.7.2) is smaller or equal to k^j . Thus,

$$c^r \leq \sum_{j=1}^{rd} k^j \cdot k^{-j} = \sum_{j=1}^{rd} 1 = rd.$$

Consequently, $\frac{c^r}{r} \leq d$, which proves that $c \leq 1$. □

Exercises

1A What is the smallest length of a block code with the same source alphabet $\{A, B, \dots, Z\}$ and the same code alphabet $\{ \cdot, - , \text{space} \}$ as the Morse code?

1	...	01
2	...	011
3	...	10
4	...	1000
5	...	1100
6	...	0111

Figure 5

A	...	1010
B	...	001
C	...	101
D	...	0001
E	...	1101
F	...	1011

Figure 6

1B Is the code in Figure 5 uniquely decodable? Is it instantaneous? Can you find an instantaneous code with the same lengths of code words?

1C Is the code in Figure 6 uniquely decodable? If not, exhibit two source messages with the same code.

1D Is the code in Figure 7 uniquely decodable?

0	...	AA
1	...	AABAB
2	...	ABBBBB
3	...	ABABA
4	...	ABBAA
5	...	BABBA
6	...	BBBAB
7	...	AAAABB
8	...	AAAABA
9	...	AAAAAB

Figure 7

A	...	001	A	...	00
B	...	1001	B	...	10
C	...	0010	C	...	011
D	...	1110	D	...	101
E	...	1010	E	...	111
F	...	01110	F	...	110
G	...	0101	G	...	010

Figure 8

1E Can you decide unique decodability of the two codes in Figure 8 by using Kraft's inequality?

1F Construct a binary instantaneous code for the following source alphabet with the prescribed lengths of code words:

Symbol	A	B	C	D	E	F	G	H	I	J	K	L
Length	2	4	7	7	3	4	7	7	3	4	7	7

1G Construct a ternary (three code symbols) instantaneous code for the following source alphabet with the prescribed lengths of code words:

Symbol	1	2	3	4	5	6	7	8	9	0
Length	1	3	3	3	3	2	2	2	2	2

1H How many code symbols are needed if the following source alphabet is to be encoded into an instantaneous code with the prescribed lengths of code words:

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	2	2	2	1	2	2	2	1	2	2	2	2	2	1	2

1I Prove that for each instantaneous code for which Kraft's inequality is not an equality, it is possible to add a new source symbol and extend the given code to an instantaneous code (with the same code alphabet). Demonstrate this on the code found in Exercise 1.6.

Notes

Coding theory and information theory have their origin in the fundamental work of Claude E. Shannon (1948), reprinted in Slepian (1974), although many of the ideas were understood and used before. The small interesting book of Richard N. Hamming (1980), one of the founders of coding theory, provides further historical remarks.

Kraft's inequality was formulated by Kraft (1949). The source of McMillan's Theorem is McMillan (1956), the present proof is from Karush (1961).

Chapter 2

Huffman Codes

We have mentioned that if the frequencies of source symbols vary, then instantaneous codes can be preferable to block codes: the most frequent symbols will be encoded into the shortest code words. We now make these considerations more precise. If the frequencies of source symbols are known exactly (i.e., if the probability distribution of source symbols in messages has been determined), we want to find the most efficient coding.

2.1 Information Source

Definition. An information source is a source alphabet together with a probability distribution; i.e., a set $\{a_1, \dots, a_n\}$ together with numbers $P(a_1), \dots, P(a_n)$ satisfying $\sum_{i=1}^n P(a_i) = 1$ and $0 \leq P(a_i) \leq 1$.

More precisely, we should speak about a discrete, zero-memory information source: discrete because we have discrete source symbols, and zero-memory because we assume that the source symbols appear independently. In other words, the probabilities $P(a_i)$ fully describe the statistics of the source messages, and the probability of a message $a_{i_1} a_{i_2} \dots a_{i_r}$ can be computed from the probabilities of the individual symbols:

$$P(a_{i_1} a_{i_2} \dots a_{i_r}) = P(a_{i_1}) P(a_{i_2}) \dots P(a_{i_r}).$$

2.2 Huffman Codes

Let K be a coding of an information source. That is, for each source symbol a_i , we have a code word $K(a_i)$ and we know the probability $P(a_i)$