

Information Theory

Information Theory

Jan C.A. VAN DER LUBBE

Associate Professor

Information Theory Group

Department of Electrical Engineering

Delft University of Technology

Translated by Hendrik Jan Hoes and Steve Gee

VSSD

© 1997 VSSD and Cambridge University Press

Published by:

VSSD

Leeghwaterstraat 42, 2628 CA Delft, The Netherlands

tel. +31 15 27 82124, e-mail: dap@vssd.nl

internet: <http://www.delftacademicpress.nl/e002EN.php>

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photo-copying, recording, or otherwise, without the prior written permission of the publisher.

Printed in The Netherlands.

ISBN 90-407-1255-7

NUGI 965

Keywords: information theory.

Contents

<i>Preface</i>	<i>page xi</i>
1 Discrete information	1
1.1 The origin of information theory	1
1.2 The concept of probability	4
1.3 Shannon's information measure	8
1.4 Conditional, joint and mutual information	16
1.5 Axiomatic foundations	22
1.6 The communication model	24
1.7 Exercises	27
1.8 Solutions	29
2 The discrete memoryless information source	39
2.1 The discrete information source	39
2.2 Source coding	43
2.3 Coding strategies	49
2.4 Most probable messages	56
2.5 Exercises	60
2.6 Solutions	63
3 The discrete information source with memory	79
3.1 Markov processes	79
3.2 The information of a discrete source with memory	85
3.3 Coding aspects	91
3.4 Exercises	95
3.5 Solutions	97
4 The discrete communication channel	109
4.1 Capacity of noiseless channels	109
4.2 Capacity of noisy channels	116
4.3 Error probability and equivocation	126

4.4	Coding theorem for discrete memoryless channels	130
4.5	Cascading of channels	133
4.6	Channels with memory	136
4.7	Exercises	138
4.8	Solutions	142
5	The continuous information source	155
5.1	Probability density functions	155
5.2	Stochastic signals	164
5.3	The continuous information measure	171
5.4	Information measures and sources with memory	176
5.5	Information power	186
5.6	Exercises	190
5.7	Solutions	194
6	The continuous communication channel	209
6.1	The capacity of continuous communication channels	209
6.2	The capacity in the case of additive gaussian white noise	214
6.3	Capacity bounds in the case of non-gaussian white noise	215
6.4	Channel coding theorem	218
6.5	The capacity of a gaussian channel with memory	222
6.6	Exercises	227
6.7	Solutions	229
7	Rate distortion theory	238
7.1	The discrete rate distortion function	238
7.2	Properties of the $R(D)$ function	243
7.3	The binary case	250
7.4	Source coding and information transmission theorems	253
7.5	The continuous rate distortion function	259
7.6	Exercises	263
7.7	Solutions	265
8	Network information theory	268
8.1	Introduction	268
8.2	Multi-access communication channel	269
8.3	Broadcast channels	281
8.4	Two-way channels	292
8.5	Exercises	298
8.6	Solutions	299
9	Error-correcting codes	305
9.1	Introduction	305

9.2	Linear block codes	307
9.3	Syndrome coding	312
9.4	Hamming codes	316
9.5	Exercises	318
9.6	Solutions	319
10	Cryptology	324
10.1	Cryptography and cryptanalysis	324
10.2	The general scheme of cipher systems	325
10.3	Cipher systems	327
10.4	Amount of information and security	334
10.5	The unicity distance	337
10.6	Exercises	340
10.7	Solutions	341
	Bibliography	345
	Index	347

Preface

On all levels of society systems have been introduced that deal with the transmission, storage and processing of information. We live in what is usually called the information society. Information has become a key word in our society. It is not surprising therefore that from all sorts of quarters interest has been shown in what information really is and consequently in acquiring a better knowledge as to how information can be dealt with as efficiently as possible.

Information theory is characterized by a quantitative approach to the notion of information. By means of the introduction of measures for information answers will be sought to such questions as: How to transmit and store information as compactly as possible? What is the maximum quantity of information that can be transmitted through a channel? How can security best be arranged? Etcetera. Crucial questions that enable us to enhance the performance and to grasp the limits of our information systems.

This book has the purpose of introducing a number of basic notions of information theory and clarifying them by showing their significance in present applications. Matters that will be described are, among others: Shannon's information measure, discrete and continuous information sources and information channels with or without memory, source and channel decoding, rate distortion theory, error-correcting codes and the information theoretical approach to cryptology. Special attention has been paid to multiterminal or network information theory; an area with still lots of unanswered questions, but which is of great significance because most of our information is transmitted by networks.

All chapters are concluded with questions and worked solutions. That makes the book suitable for self study.

The content of the book has been largely based on the present lectures by the author for students in Electrical Engineering, Technical Mathematics and Informatics, Applied Physics and Mechanical Engineering at the Delft University of Technology, as well as on former lecture notes by Profs. Ysbrand Boxma, Dick Boekee and Jan Biemond. The questions have been derived from recent exams.

The author wishes to express his gratitude to the colleagues mentioned above as well as the other colleagues who in one way or other contributed to this textbook. Especially I wish to thank E. Prof. Ysbrand Boxma, who lectured on information theory at the Delft University of Technology when I was a student and who introduced me to information theory. Under his inspiring guidance I received my M.Sc. in Electrical Engineering and my Ph.D. in the technical sciences. In writing this book his old lecture notes were still very helpful to me. His influence has been a determining factor in my later career.

Delft, December 1996

Jan C.A. van der Lubbe

1

Discrete information

1.1 The origin of information theory

Information theory is the science which deals with the concept ‘information’, its measurement and its applications. In its broadest sense distinction can be made between the American and British traditions in information theory.

In general there are three types of information:

- *syntactic information*, related to the symbols from which messages are built up and to their interrelations,
- *semantic information*, related to the meaning of messages, their referential aspect,
- *pragmatic information*, related to the usage and effect of messages.

This being so, syntactic information mainly considers the form of information, whereas semantic and pragmatic information are related to the information content.

Consider the following sentences:

- (i) John was brought to the railway station by taxi.
- (ii) The taxi brought John to the railway station.
- (iii) There is a traffic jam on highway A3, between Nuremberg and Munich in Germany.
- (iv) There is a traffic jam on highway A3 in Germany.

The sentences (i) and (ii) are syntactically different. However, semantically and pragmatically they are identical. They have the same meaning and are both equally informative.

The sentences (iii) and (iv) do not differ only with respect to their syntax, but also with respect to their semantics. Sentence (iii) gives more precise information than sentence (iv).

The pragmatic aspect of information mainly depends on the context. The information contained in the sentences (iii) and (iv) for example is relevant for someone in Germany, but not for someone in the USA.

The semantic and pragmatic aspects of information are studied in the British tradition of information theory. This being so, the British tradition is closely related to philosophy, psychology and biology. The British tradition is influenced mainly by scientists like MacKay, Carnap, Bar-Hillel, Ackoff and Hintikka.

The American tradition deals with the syntactic aspects of information. In this approach there is full abstraction from the meaning aspects of information. There, basic questions are the measurement of syntactic information, the fundamental limits on the amount of information which can be transmitted, the fundamental limits on the compression of information which can be achieved and how to build information processing systems approaching these limits. A rather technical approach to information remains.

The American tradition in information theory is sometimes referred to as communication theory, mathematical information theory or in short as information theory. Well-known scientists of the American tradition are Shannon, Renyi, Gallager and Csizsár among others.

However, Claude E. Shannon, who published his article “A mathematical theory of communication” in 1948, is generally considered to be the founder of the American tradition in information theory. There are, nevertheless, a number of forerunners to Shannon who attempted to formalise the efficient use of communication systems.

In 1924 H. Nyquist published an article wherein he raised the matter of how messages (or characters, to use his own words) could be sent over a telegraph channel with maximum possible speed, but without distortion. The term information however was not yet used by him as such.

It was R.V.L. Hartley (1928) who first tried to define *a measure of information*. He went about it in the following manner.

Assume that for every symbol of a message one has a choice of s possibilities. By now considering messages of l symbols, one can distinguish s^l messages. Hartley now defined the amount of information as the logarithm of the number of distinguishable messages. In the case of messages of length l one therefore finds

$$H_H(s^l) = \log\{s^l\} = l \log\{s\}. \quad (1.1)$$

For messages of length 1 one would find

$$H_H(s^1) = \log\{s\}$$

and thus

$$H_H(s^l) = l H_H(s^1).$$

This corresponds with the intuitive idea that a message consisting of l symbols, by doing so, contains l times as much information as a message consisting of only one symbol. This also accounts for the appearance of the logarithm in Hartley's definition.

It can readily be shown that the only function that satisfies the equation

$$f\{s^l\} = l f\{s\}$$

is given by

$$f\{s\} = \log\{s\}, \quad (1.2)$$

which yields Hartley's measure for the amount of information. Note that the logarithm also guarantees that the amount of information increases as the number of symbols s increases, which is in agreement with our intuition.

The choice of the base of the logarithm is arbitrary and is more a matter of normalisation. If the natural logarithm is used, the unit of information is called the *nat* (natural unit). Usually 2 is chosen as the base. The amount of information is then expressed in *bits* (derived from binary unit, i.e. two-valued unit). In the case of a choice of two possibilities, the amount of information obtained when one of the two possibilities occurs is then equal to 1 bit. It is easy to see that the relationship between bit and nat is given by

$$1 \text{ nat} = 1.44 \text{ bits}.$$

In Hartley's approach as given above, no allowance is made for the fact that the s symbols may have unequal chances of occurring or that there could be a possible dependence between the l successive symbols.

Shannon's great achievement is that he extended the theories of Nyquist and Hartley, and laid the foundation of present-day information theory by associating information with uncertainty using the concept of chance or probability. With regard to Hartley's measure, Shannon proposed that it could indeed be interpreted as a measure for the amount of information, with the assumption that all symbols have an equal probability of occurring. For the general case, Shannon introduced an information measure based on the concept of probability, which includes Hartley's measure as a special

case. Some attention will first be paid to probability theory, during which some useful notations will be introduced, before introducing Shannon's definition of information.

1.2 The concept of probability

Probability theory is the domain dealing with the concept of probability. The starting point of probability theory is that experiments are carried out which then yield certain outcomes. One can also think in terms of an information source which generates symbols. Every occurrence of a symbol can then be regarded as an event. It is assumed that one is able to specify which possible outcomes or events can occur. The collection of all possible outcomes or events is called the *sample space*. It is now possible to speak of the probability that an experiment has a certain outcome, or of the probability that an information source will generate a certain symbol or message. Each event or outcome has a number between 0 and 1 assigned to it, which indicates how large the probability is that this outcome or event occurs. For simplicity it is assumed that the sample space has a finite number of outcomes.

Consider a so-called *probabilistic experiment* X with possible outcomes/events x_i , with $x_i \in X$ and X the probability space as defined by

$$X = \{x_1, \dots, x_i, \dots, x_n\}. \quad (1.3)$$

If we think of throwing a die, then x_1 could be interpreted as the event that "1" is thrown, x_2 the event that "2" is thrown, etc. In the case of the die it is obvious that $n = 6$.

Each event will have a certain probability of occurring. We denote the probability related to x_i by $p(x_i)$ or simply p_i . The collection of probabilities with regard to X is denoted by

$$P = \{p_1, \dots, p_i, \dots, p_n\}, \quad (1.4)$$

and is called the *probability distribution*. The probability distribution satisfies two fundamental requirements:

(i) $p_i \geq 0$, for all i .

(ii) $\sum_{i=1}^n p_i = 1$.

That is, no probability can take on a negative value and the sum of all the probabilities is equal to 1.

Sometimes we can discern two types of outcomes in one experiment, such that we have a combination of two subexperiments or subevents. When testing IC's for example, one can pay attention to how far certain requirements are met (well, moderately or badly for example), but also to the IC's type number. We are then in actual fact dealing with two sample spaces, say X and Y , where the sample space Y , relating to experiment Y , is in general terms defined by

$$Y = \{y_1, \dots, y_j, \dots, y_m\}, \quad (1.5)$$

and where the accompanying probability distribution is given by

$$Q = \{q_1, \dots, q_j, \dots, q_m\}, \quad (1.6)$$

where $q(y_j) = q_j$ is the probability of event y_j . We can now regard (X, Y) as a probabilistic experiment with pairs of outcomes (x_i, y_j) , with $x_i \in X$ and $y_j \in Y$. The probability $r(x_i, y_j)$, also denoted by r_{ij} or $p(x_i, y_j)$, is the probability that experiment (X, Y) will yield (x_i, y_j) as outcome and is called the *joint probability*. If the joint probability is known, one can derive the probabilities p_i and q_j , which are then called the *marginal probabilities*. It can be verified that for all i

$$p_i = \sum_{j=1}^m r_{ij}, \quad (1.7)$$

and for all j

$$q_j = \sum_{i=1}^n r_{ij}. \quad (1.8)$$

Since the sum of all the probabilities p_i must be equal to 1 (and likewise the sum of the probabilities q_j), it follows that the sum of the joint probabilities must also be equal to 1:

$$\sum_{i=1}^n \sum_{j=1}^m r_{ij} = 1.$$

Besides the joint probability and the related marginal probability, there is a third type, namely the *conditional probability*. This type arises when a probabilistic experiment Y is conditional for experiment X . That is, if the

6 Discrete Information

probabilities of the outcomes of X are influenced by the outcomes of Y . We are then interested in the probability of an event, x_i for example, given that another event, y_j for example, has already occurred.

Considering the words in a piece of English text, one may ask oneself, for example, what the probability is of the letter “n” appearing if one has already received the sequence “informatio”. The appearances of letters in words often depend on the letters that have already appeared. It is very unlikely, for example, that the letter “q” will be followed by the letter “t”, but much more likely that the letter “u” follows.

The conditional probability of x_i given y_j is defined as

$$p(x_i/y_j) = \frac{r(x_i, y_j)}{q(y_j)}, \text{ provided } q(y_j) > 0,$$

or in shortened notation

$$p_{ij} = \frac{r_{ij}}{q_j}, \quad \text{provided } q_j > 0. \quad (1.9)$$

The conditional probability of y_j given x_i can be defined in an analogous manner as

$$q(y_j/x_i) = \frac{r(x_i, y_j)}{p(x_i)}, \text{ provided } p(x_i) > 0,$$

or simply

$$q_{ji} = \frac{r_{ij}}{p_i}, \quad \text{provided } p_i > 0. \quad (1.10)$$

From the definitions given it follows that the joint probability can be written as the product of the conditional and marginal probabilities:

$$r(x_i, y_j) = q(y_j) p(x_i/y_j) = p(x_i) q(y_j/x_i). \quad (1.11)$$

The definition of conditional probability can be simply extended to more than two events. Consider x_i , y_j and z_k for example:

$$\begin{aligned} p(x_i, y_j, z_k) &= r(y_j, z_k) p(x_i/y_j, z_k) \\ &= p(z_k) p(y_j/z_k) p(x_i/y_j, z_k), \end{aligned}$$

hence

$$p(x_i/y_j, z_k) = p(x_i, y_j, z_k) / r(y_j, z_k).$$

Returning to the conditional probability, summation over the index i with y_j given yields

$$\sum_{i=1}^n p(x_i/y_j) = 1. \quad (1.12)$$

Whenever an event y_j has occurred, one of the events in X must also occur. Thus, summation will yield 1. Note that the converse is not true. It is generally true that

$$\sum_{j=1}^m p(x_i/y_j) \neq 1. \quad (1.13)$$

A handy aid which will be of use in the following is *Bayes' theorem*. It is often the case that the conditional probability $q(y_j/x_i)$ is known, but that we want to determine the conditional probability $p(x_i/y_j)$. One can do this by making use of the following relations:

$$r(x_i, y_j) = p(x_i) q(y_j/x_i) = q(y_j) p(x_i/y_j).$$

Hence, if $q(y_j) > 0$,

$$p(x_i/y_j) = \frac{p(x_i) q(y_j/x_i)}{q(y_j)},$$

or also

$$p(x_i/y_j) = \frac{p(x_i) q(y_j/x_i)}{\sum_{i=1}^n p(x_i) q(y_j/x_i)}. \quad (1.14)$$

We are thus able to calculate $p(x_i/y_j)$ with the help of $q(y_j/x_i)$.

Finally, a comment about the concept of independence. The situation can arise that

$$p(x_i/y_j) = p(x_i).$$

That is, the occurrence of y_j has no influence on the occurrence of x_i . But it then also follows that

$$r(x_i, y_j) = p(x_i) q(y_j)$$

and

$$q(y_j/x_i) = q(y_j).$$

In this case one says that the events are independent of each other. The reverse is also true, from $r(x_i, y_j) = p(x_i) q(y_j)$ it follows that $q(y_j/x_i) = q(y_j)$ and $p(x_i/y_j) = p(x_i)$. Two experiments X and Y are called *statistically independent* if for all i and j

$$r(x_i, y_j) = p(x_i) q(y_j). \quad (1.15)$$

An experiment X is called *completely dependent* on another, Y , if for all j , there is a unique i , say k , such that

$$p(x_k/y_j) = 1, \quad (1.16)$$

or

$$p(x_k, y_j) = p(y_j). \quad (1.17)$$

1.3 Shannon's information measure

As we saw in Section 1.1, Hartley's definition of information did not take the various probabilities of occurrence of the symbols or events into account. It was Shannon who first associated information with the concept of probability.

This association is in actual fact not illogical. If we consider a sample space where all events have an equal probability of occurring, there is great uncertainty about which of the events will occur. That is, when one of these events occurs it will provide much more information than in the cases where the sample space is structured in such a way that one event has a large probability of occurring. Information is linked to the concept of chance via uncertainty.

Before considering to what extent *Shannon's information measure* satisfies the properties one would in general expect of an information measure, we first give his definition.

Definition 1.1

Let X be a probabilistic experiment with sample space X and probability distribution P , where $p(x_i)$ or p_i is the probability of outcome $x_i \in X$. Then the average amount of information is given by

$$H(X) = - \sum_{i=1}^n p(x_i) \log p(x_i) = - \sum_{i=1}^n p_i \log p_i. \quad (1.18)$$

○

Other notations for Shannon's information measure are $H(X)$, $H(P)$ and $H(p_1, \dots, p_n)$. All of these notations will be used interchangeably in this text.

Because this measure for the amount of information is attended with the choice (selection) from n possibilities, one sometimes also speaks of the measure for the amount of *selective information*.

Because 2 is usually chosen as the base, the unit of information thereby becoming the bit, this will not be stated separately in future, but left out.

In the case of two outcomes with probabilities $p_1 = p$ and $p_2 = 1 - p$ we find

$$H(P) = -p \log p - (1 - p) \log (1 - p). \quad (1.19)$$

Figure 1.1 shows how $H(P)$ behaves as a function of p . It can be concluded that if an outcome is certain, that is, occurs with a probability of 1, the information measure gives 0. This is in agreement with the intuitive idea that certain events provide no information. The same is true for $p = 0$; in that case the other outcome has a probability of 1.

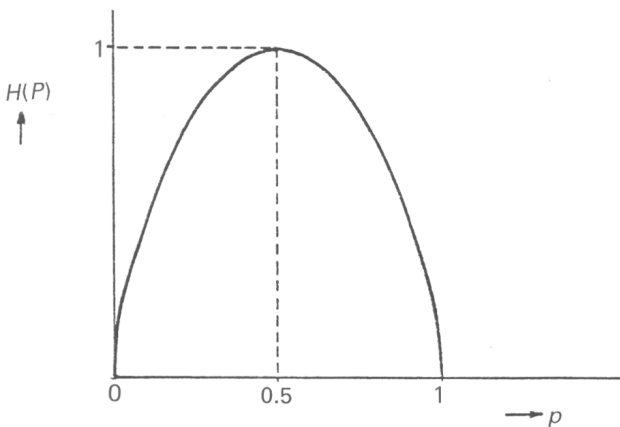
When $p = 0.5$, $H(P)$ reaches its maximum value, which is equal to 1 bit. For $p = 0.5$, both outcomes are just as probable, and one is completely uncertain about the outcome. The occurrence of one of the events provides the maximum amount of information in this case.

As an aside, note that by definition $0 \cdot \log(0) = 0$.

Returning to the general case, we can posit that the information measure satisfies four intuitive requirements :

- I $H(P)$ is *continuous* in p
- II $H(P)$ is *symmetric*. That is, the ordering of the probabilities p_1, \dots, p_n

Figure 1.1. $H(P) = H(p, 1 - p)$ as a function of p .



does not influence the value of $H(P)$.

- III $H(P)$ is *additive*. If X and Y are two sample spaces, where outcomes in X are independent of those in Y , then we find for the information relating to joint events (x_i, y_j)

$$\begin{aligned} &H(p_1 q_1, \dots, p_1 q_m, \dots, p_n q_1, \dots, p_n q_m) \\ &= H(p_1, \dots, p_n) + H(q_1, \dots, q_m). \end{aligned} \tag{1.20}$$

- IV $H(P)$ is maximum if all probabilities are equal. This corresponds with the situation where maximum uncertainty exists. $H(P)$ is minimum if one outcome has a probability equal to 1.

A short explanation of a number of the above requirements follows.

Ad II. That Shannon's information measure is symmetric means that changing the sequence in which one substitutes the probabilities does not change the amount of information. A consequence of this is that different sample spaces with probability distributions that have been obtained from the permutations of a common probability distribution will result in the same amount of information.

Example 1.1

Consider the experiments X and Y with the following sample spaces:

$$\begin{aligned} X &= \{\text{it will rain tomorrow, it will be dry tomorrow}\} \\ &\text{where } P = \{0.8, 0.2\} \end{aligned}$$

and

$$\begin{aligned} Y &= \{\text{John is younger than 30, John is at least 30}\} \\ &\text{where } Q = \{0.2, 0.8\}. \end{aligned}$$

The amount of information with relation to X is

$$H(X) = -0.8 \log 0.8 - 0.2 \log 0.2 = 0.72 \text{ bit}$$

and with relation to Y

$$H(Y) = -0.2 \log 0.2 - 0.8 \log 0.8 = 0.72 \text{ bit}$$

and thus

$$H(X) = H(Y). \quad \triangle$$

From this example, it can be concluded that Shannon's information measure is not concerned with the contents of information. The probabilities with which events occur are of importance and not the events themselves.