

# Information sources with different cost scales and the principle of conservation of entropy

by

I. Csiszár, G. Katona and G. Tusnády

Math. Inst. of the Hung. Acad. of Sci, Budapest, Hungary

## INTRODUCTION

Despite of the vast progress information theory has made in the last decade, some problems important from the point of view of the very foundations - to the authors' knowledge - still lack a rigorous and sufficiently general exposition. In this paper we attempt to fill some of these gaps, concerning problems of the following type (precise definitions will be given in Section 2):

(i) If the message symbols produced by an information source are of different cost, the entropy per unit cost can be defined either as the limit of the entropy of the message sequence of cumulative cost  $t$  divided by  $t$ , or as the entropy per symbol (limit of the entropy of the first  $n$  symbols divided by  $n$ ) divided by the average symbol cost. (Most frequently, the cost of a symbol is its duration and entropy per unit cost is entropy per second.) Dating back to Shannon's fundamental paper [16], in general the second definition is adopted (see also [12], [15], etc.) but in heuristic reasonings it is often implicitly assumed to be equivalent to the first one. In the literature consideration is usually restricted to the simplest case that each symbol of the source alphabet has some fixed cost (duration), but no proof of the equivalence of the two possible definitions of entropy per second seems to have been published even for that case. It should be noted that also the general case has considerable interest, in particular if one looks at sources producing message symbols at random times - according to some point process - and "cost" is interpreted as the length of the time interval between two subsequent symbols.

(ii) For the interpretation of entropy as the measure of the amount of information the so-called noiseless coding theorem is of basic importance. It asserts, intuitively, that the greatest lower bound of the average number  $L$  of code characters per symbol needed to encode in a uniquely decipherable way the output of a source of entropy rate  $H$  equals

$$\frac{H}{\log_2 s'},$$

where  $s'$  is the size of the coding alphabet. The "noiseless coding

theorem" is usually stated and proved, however, for rather special codes only, namely for those defined by a fixed assignment of sequences of code symbols to the letters of the source alphabet, or to sequences ("blocks") of fixed length of letters of the source alphabet (see e.g. Feinstein [9], Ash [1], etc.). On the other hand, the theorem is expected to be true for "all conceivable" codes and in order that it be really valuable from the "foundations" point of view it should be proved for "arbitrary" codes, including blockwise encodings with variable block length and the code mapping varying from block to block, in dependence on the previously encoded message symbols. The strongest results known in this direction are apparently those of Billingsley [2]. Moreover, if the message symbols or the code characters or both are of different cost, if  $H$  is interpreted as the entropy per unit cost,  $L$  as the average cost dilation due to the encoding, and  $\log_2 s'$  is replaced by the "capacity of the noiseless channel" as defined by Shannon [16], one may infer that the statement still remains valid. In fact - for the case of fixed symbol costs - this statement occurs already in [16], but, to the authors' knowledge, no exact proof has been published so far, except for special (Markovian) sources.

(iii) It is "intuitively clear" that uniquely decipherable coding gives rise to a message of entropy (per symbol)  $\frac{H}{L}$  where  $H$  is the entropy (per symbol) of the original message and  $L$  is the average number of code characters per message symbol. A rigorous proof of this assertion for finite-state Markovian sources and encodings performed by finite-state transducers has been given by Sidelnikov [17], and for arbitrary sources and simple letter codes by two of the present authors [10]. In case of symbols of different cost, a similar relation is to be expected for the entropies per unit cost. In this direction there seems nothing to have been published.

The problems listed under (i), (ii) and (iii) are very closely related to each other. As a main tool for dealing with them we introduce the concept of entropy rate with respect to a stochastic cost scale and establish a theorem on the relation of entropy rates with respect to different cost scales under general conditions. Applying this result, we obtain an apparently satisfactory solution of problems (i)-(iii), for sources with finite alphabet; in particular we prove the "principle of conservation of entropy" for a very wide class of encoding procedures. Our method is straightforward and follows

closely intuition. Our aim was to make familiar heuristic reasonings rigorous rather than to replace them by ad hoc non-information-theoretic arguments, in this respect, even for the particular cases of our results that have been proved earlier, our proofs seem preferable to the existing ones.

## § 1. PRELIMINARIES

(A) Throughout this paper, the terms "random variable", "discrete random variable" (= random variable with finite or countable state space), "integer valued random variable", "almost surely" (= with probability one), "uniformly integrable" and "if and only if" will be abbreviated as RV, DRV, IRV, a. s., u. i. and iff, respectively.

All RV's will be assumed to be defined on the same probability space  $(\Omega, \mathcal{F}, P)$ . RV's will be denoted by greek letters, omitting, as a rule, the argument  $\omega$ . Except  $\Omega$  and  $\omega$  (typical element of  $\Omega$ ), all greek letters occurring in this paper denote RV's. In case of families of RV's (= stochastic processes) we shall write the parameter as an argument rather than as an index; thus a typical element of a sequence of RV's will be denoted by  $\xi(n)$  rather than by  $\xi_n$  (of course,  $\xi(n)$  means really  $\xi(n, \omega)$ ).

If  $A \in \mathcal{F}$ ,  $P(A) > 0$ , symbols with subscript A will refer to the probability measure  $P_A(\cdot) = P(\cdot | A)$ ; if  $P(A) = 0$ , such symbols will be meant to be 0. E.g.,  $E_{\{\alpha \leq \xi \leq b\}}(\eta)$  means  $E(\eta | \alpha \leq \xi \leq b)$  if  $P(\alpha \leq \xi \leq b) > 0$  and 0 otherwise.

(B) By entropy (conditional entropy) of DRV's we shall always mean entropy in the sense of Shannon:

$$(1.1) \quad H(\xi) = - \sum_x P(\xi = x) \log_2 P(\xi = x)$$

$$(1.2) \quad \begin{aligned} H(\xi | \eta) &= \sum_y P(\eta = y) H_{\{\eta = y\}}(\xi) = \\ &= - \sum_{x,y} P(\eta = y) P(\xi = x | \eta = y) \log_2 P(\xi = x | \eta = y) \end{aligned}$$

where  $x$  and  $y$  range over the state space of  $\xi$  and  $\eta$ , respectively (eventual undefined terms in (1.1) and (1.2) are considered as zeros).

We shall need also the concept of information distance (of DRV's)

$$(1.3) \quad d(\xi, \eta) = H(\xi|\eta) + H(\eta|\xi)$$

and the mutual information (of DRV's with  $d(\xi, \eta) < \infty$ )

$$(1.4) \quad I(\xi, \eta) = H(\xi) - H(\xi|\eta) = H(\eta) - H(\eta|\xi) .$$

The equality  $H(\xi) - H(\xi|\eta) = H(\eta) - H(\eta|\xi)$  follows from (1.5) below, if  $d(\xi, \eta) < \infty$ . For the purposes of this paper, we need not define  $I(\xi, \eta)$  if  $d(\xi, \eta) = \infty$ .

The well-known basic identities and inequalities concerning entropies and conditional entropies (due essentially to Shannon [16], see also e.g. [1], [9]) such as

$$(1.5) \quad H(\xi, \eta) = H(\xi) + H(\eta|\xi) = H(\eta) + H(\xi|\eta)$$

$$(1.5') \quad H(\xi, \eta|\zeta) = H(\xi|\zeta) + H(\eta|\xi, \zeta) = H(\eta|\zeta) + H(\xi|\eta, \zeta)$$

$$(1.6) \quad 0 \leq H(\xi|\eta, \zeta) \leq H(\xi|\eta) \leq H(\xi)$$

$$(1.7) \quad 0 \leq I(\xi, \eta) \leq \max(H(\xi), H(\eta))$$

$$(1.8) \quad 0 \leq H(\xi) \leq \log_2 \{ \text{number of possible values of } \xi \}$$

$$(1.8') \quad 0 \leq H(\xi|\eta) \leq E \log_2 \{ \text{number of possible values of } \xi \text{ given } \eta \}$$

will be used freely, without any further reference. We shall need also some other simple but somewhat less standard inequalities summarized in the following lemmas:

LEMMA 1. 1. We have for arbitrary DRV' s

$$(1. 9) \quad d(\xi, \eta) + d(\eta, \zeta) \geq d(\xi, \zeta)$$

$$(1. 10) \quad |H(\xi) - H(\eta)| \leq d(\xi, \eta)$$

$$(1. 11) \quad |H(\xi_1 | \eta_1) - H(\xi_2 | \eta_2)| \leq d(\xi_1, \xi_2) + d(\eta_1, \eta_2)$$

$$(1. 12) \quad |I(\xi_1, \eta_1) - I(\xi_2, \eta_2)| \leq d(\xi_1, \xi_2) + d(\eta_1, \eta_2)$$

provided (in (1. 10)-(1. 12)) that the left hand side is meaningful.

REMARK 1. 1. This lemma means that the information-distance is a metric in the space of DRV' s with finite entropy and the different information quantities are (uniformly) continuous functions with respect to it.

LEMMA 1. 2. If  $\xi$  is an IRV with finite expectation then

$$(1. 13) \quad H(\xi) \leq E|\xi| + \log_2 3$$

and

$$(1. 14) \quad H(\xi) \leq E \log_2(|\xi| + 1) + \log_2 \left( \frac{\pi^2}{3} - 1 \right) .$$

PROOF OF LEMMA 1. 1. The triangle inequality (1. 9) follows from

$$(1. 15) \quad H(\xi | \eta) + H(\eta | \xi) \geq H(\xi | \eta, \zeta) + H(\eta | \xi) = H(\xi, \eta | \zeta) \geq H(\xi | \zeta)$$

and the corresponding inequality obtained by changing the role of  $\xi$  and  $\zeta$ . (1. 10) is an immediate consequence of (1. 5) and (1. 3). By obvious substitutions, (1. 15) gives rise also to

$$(1. 16) \quad \begin{aligned} |H(\xi_1 | \eta) - H(\xi_2 | \eta)| &\leq d(\xi_1, \xi_2) \\ |H(\xi | \eta_1) - H(\xi | \eta_2)| &\leq d(\eta_1, \eta_2) \end{aligned}$$

whence (1.11) directly follows:

$$|H(\xi_1|\eta_1) - H(\xi_2|\eta_2)| \leq |H(\xi_1|\eta_1) - H(\xi_2|\eta_1)| + \\ + |H(\xi_2|\eta_1) - H(\xi_2|\eta_2)| \leq d(\xi_1, \xi_2) + d(\eta_1, \eta_2)$$

at last, from (1.14), (1.5) and the second inequality of (1.16) we get

$$|I(\xi_1, \eta_1) - I(\xi_2, \eta_2)| = |H(\xi_1) - H(\eta_2) - H(\xi_1|\eta_1) + H(\eta_2|\xi_2)| = \\ = |H(\xi_1|\eta_2) - H(\eta_2|\xi_1) - H(\xi_1|\eta_1) + H(\eta_2|\xi_2)| \leq d(\eta_1, \eta_2) + d(\xi_1, \xi_2)$$

i. e. (1.12).

PROOF OF LEMMA 1.2. Set  $p_k = P(\xi=k)$  and  $q_k = \frac{1}{3} 2^{-|k|}$  ( $k=0, \pm 1, \pm 2, \dots$ ). Then  $\{q_k\}$  is a probability distribution and the well-known inequality

$$\sum p_k \log_2 \frac{p_k}{q_k} \geq 0$$

gives rise to (1.13).

(1.14) can be proved in the same way, with the choice

$$q_k = \frac{c}{(|k|+1)^2} \quad (k = 0, \pm 1, \pm 2, \dots) \quad c = \left(\frac{\pi^2}{3} - 1\right)^{-1}$$

(C) We shall have to do with three types of convergence of RV's: convergence in probability (or stochastic convergence), almost sure convergence (convergence with probability one) and convergence in  $L_1$ -norm.

They will be denoted by  $\xrightarrow{P}$ ,  $\xrightarrow{a.s.}$  and  $\xrightarrow{L_1}$ , respectively.

LEMMA 1.3. Let  $\xi(t)$ ,  $t \geq 0$  be a family of RV's and let  $\xi(t) \xrightarrow{P} \xi$ , ( $t \rightarrow +\infty$ ). Then the conditions

a)  $\xi(t)$  is uniformly integrable (u. i.) for  $t \rightarrow \infty$  and

b)  $E|\xi(t)| \rightarrow E|\xi| < \infty$  ( $t \rightarrow \infty$ )

are equivalent and imply  $\xi(t) \xrightarrow{L_1} \xi$  ( $t \rightarrow \infty$ ); and conversely,  $\xi(t) \xrightarrow{L_1} \xi \in L_1$  implies  $\xi(t) \xrightarrow{P} \xi$  and both (a) and (b).

Here condition (a) means that  $\overline{\lim}_{t \rightarrow \infty} \int_{|\xi(t)| \geq k} |\xi(t)| P(d\omega) \rightarrow 0$  as  $k \rightarrow +\infty$ .

REMARK 1.2. If there exists  $t_0 \geq 0$  such that for every finite  $t_1 > t_0$  the RV's  $\xi(t)$ ,  $t_0 \leq t \leq t_1$  are u.i. then, obviously, condition (a) is equivalent to saying that  $\xi(t)$  is u.i. for  $t \geq t_0$ .

PROOF. From  $\xi(t) \xrightarrow{L_1} \xi \in L_1$  obviously follows both  $\xi(t) \xrightarrow{P} \xi$  and (b). These, in turn, imply, on account of

$$E|\xi(t)| - E|\xi| =$$

$$= \int_{A(t,k)} (|\xi(t)| - |\xi|) P(d\omega) + \int_{\bar{A}(t,k)} |\xi(t)| P(d\omega) - \int_{\bar{A}(t,k)} |\xi| P(d\omega)$$

(with  $A(t,k) = \{\omega : |\xi(t) - \xi| < k, |\xi| < k\}$ ,  $\bar{A}(t,k) = \Omega \setminus A(t,k)$ ,  $k > 0$  fixed) the relation

$$\overline{\lim}_{t \rightarrow \infty} \int_{|\xi(t)| \geq 2k} |\xi(t)| P(d\omega) \leq \overline{\lim}_{t \rightarrow \infty} \int_{\bar{A}(t,k)} |\xi| P(d\omega) = \int_{|\xi| > k} |\xi| P(d\omega)$$

(we used that  $\{\omega : |\xi(t)| \geq 2k\} \subset \bar{A}(t,k)$ ) whence (a) directly follows.

Finally,  $\xi(t) \xrightarrow{P} \xi$  and (a) obviously imply  $\xi(t) \xrightarrow{L_1} \xi$  completing the proof.

(D) If  $X = \{x_1, \dots, x_s\}$  is an arbitrary finite set, we denote by  $U(X)$  and  $\tilde{U}(X)$  the set of all finite and infinite sequences, respectively, of elements of  $X$ . Here "sequence" means simply juxtaposition of elements without commas. The void sequence  $u_0$  will be also considered to belong to  $U(X)$ . The set  $X$  will be called an alphabet iff identical sequences from  $U(X) \cup \tilde{U}(X)$  are elementwise identical, too, i. e.  $x_{i_1} x_{i_2} \dots x_{i_m} = x_{j_1} x_{j_2} \dots x_{j_n}$  implies  $n = m$  and  $i_k = j_k$  ( $k = 1, 2, \dots, n$ ),  $x_{i_1} x_{i_2} \dots = x_{j_1} x_{j_2} \dots$  implies  $i_k = j_k$  ( $k = 1, 2, \dots$ ) and  $U(X)$  and  $\tilde{U}(X)$  are disjoint. This condition

means only that the elements of  $X$  are really "elementary", excluding e.g. the possibility of  $x_1 = a$ ,  $x_2 = b$ ,  $x_3 = ab$  or  $x_1 = a$ ,  $x_2 = bc$ ,  $x_3 = ab$ ,  $x_4 = c$  etc.

The elements of an alphabet will be referred to as letters. If  $X$  is an alphabet and  $u = x_{i_1} x_{i_2} \dots x_{i_m} \in U(X)$  then  $m$ , the length of the sequence  $u$ , is uniquely determined by  $u$ ; it will be denoted by  $\|u\|$ . Of course, we set  $\|u_0\| = 0$  and for  $u \in \tilde{U}(X)$  we set  $\|u\| = +\infty$ . For

$u, v \in U(X) \cup \tilde{U}(X)$  we shall write  $u \prec v$  iff  $\|u\| \leq \|v\|$  and the sequence of the first  $m = \|u\|$  letters of  $v$  is identical with  $u$ . Obviously,  $\prec$  is a partial order on  $U(X) \cup \tilde{U}(X)$ . Subsets of  $\tilde{U}(X)$  of form  $c(u) = \{\tilde{u} : \tilde{u} \succ u\}$  ( $u \in U(X)$ ) will be referred to as cylinder sets; the smallest  $\sigma$ -algebra of subsets of  $\tilde{U}(X)$  containing all cylinder sets will be denoted by  $\mathcal{B}$ .

(E) For the rest of this paper, it will be convenient to restrict the use of certain letters, attaching them some specific meanings in a consistent way. Our notational conventions will be the following ones:

$X$  : finite alphabet

$\xi(n)$  ( $n=1,2,\dots$ ) sequence of DRV's with common state space  $X$   
 $\mathfrak{X}$  abbreviation for the above sequence

$$\xi(k,n) = \begin{cases} \xi(k), \xi(k+1), \dots, \xi(n) & \text{if } k \leq n \\ \text{the void sequence } u_0 & \text{if } k > n \end{cases}$$

$\zeta(n)$  ( $n=1,2,\dots$ ) sequence of nonnegative real RV's, such that

$$\sum_{k=1}^n \zeta(k) \xrightarrow{\text{a.s.}} +\infty \quad \text{as } n \rightarrow \infty.$$

$\mathfrak{Z}$  abbreviation for the above sequence

$$\tau(n) \quad (n=1,2,\dots) : \tau(n) = \sum_{k=1}^n \zeta(k) \quad (\tau(0) = 0)$$

$\nu(t)$  ( $0 \leq t < +\infty$ ) : number of  $n$ 's with  $\tau(n) \leq t$

$$\eta(t) \quad (0 \leq t < +\infty) : \eta(t) = \xi(1, \nu(t))$$

( $\nu(t)$  and  $\eta(t)$  are well-defined a.s., due to the assumption  $\tau(n) \xrightarrow{\text{a.s.}} +\infty$ ).



Processes  $\mathfrak{Z}$  will be thought of as associated to processes (cf. Definition 2.1 below). Different  $\mathfrak{X}$  processes will be distinguished by dashes; different  $\mathfrak{Z}$  processes associated with the same  $\mathfrak{X}$  will be dashed in the same way as  $\mathfrak{X}$ , and they will be distinguished by indices. Given  $\mathfrak{X}$  and  $\mathfrak{Z}$ , the corresponding  $\xi$ 's,  $\zeta$ 's,  $\tau$ 's,  $\nu$ 's and  $\eta$ 's will be given the same dashes and (or) indices as  $\mathfrak{X}$  and  $\mathfrak{Z}$ .

Instead of  $\xi(n)$ ,  $\zeta(n)$ ,  $\tau(n)$ ,  $\nu(t)$  and  $\eta(t)$  we shall often write simply  $\xi, \zeta, \tau, \nu$  and  $\eta$ , if omitting the argument does not cause ambiguity.

Observe that  $\mathfrak{Z}$  is uniquely defined both by the RV's  $\tau$  and  $\nu$ ; each non-decreasing sequence of nonnegative RV's  $\tau(n)$  ( $n=1,2,\dots$ ) with

$\tau(n) \xrightarrow{a.s.} +\infty$  defines a sequence and so does each family of IRV's  $\nu(t) \geq 0$  ( $0 \leq t < +\infty$ ) with right-continuous sample functions tending to infinity as  $t \rightarrow \infty$ .

## § 2. INFORMATION SOURCES WITH DIFFERENT COST SCALES; COMPARISON OF THE CORRESPONDING ENTROPY RATES

**DEFINITION 2.1.** An information source  $\mathfrak{X}$  with finite alphabet  $X$  is a sequence of DRV's  $\xi(n)$  ( $n=1,2,\dots$ ) having the finite alphabet  $X$  as common state space. A cost scale  $\mathfrak{Z}$  is a point process on  $[0, \infty)$  described in terms of  $\zeta$ 's,  $\tau$ 's and  $\nu$ 's, see Section 1 (E). A cost scale  $\mathfrak{Z}$  will be called regular if  $\frac{\nu(t)}{t}$  is u. i. for  $t \rightarrow \infty$ .

Intuitively,  $\xi(n)$  represents the  $n$ 'th message symbol emitted by the source,  $\zeta(n)$  its cost,  $\tau(n)$  the cumulative cost of the first  $n$  message symbols and  $\nu(t)$  the number of message symbols with cumulative cost just not exceeding  $t$ . E.G. the "cost" may be time as in examples 2.2 and 2.3 below; then  $\tau(n)$  represents the epoch at which the emission of the  $n$ 'th message symbol terminates and  $\nu(t)$  is the number of message symbols emitted up to the epoch  $t$ .

**EXAMPLE 2.1.** The simplest cost scale is defined by

$$(2.1) \quad \zeta(n) = 1, \quad \tau(n) = n \quad (n=1,2,\dots), \quad \nu(t) = [t] \quad (0 \leq t < \infty).$$

This cost scale will be referred to as the counting scale  $\mathcal{C}$ .

EXAMPLE 2.2. Let  $l(u)$  be a nonnegative valued function on  $U(X)$  such that  $l(u_0)=0$  and  $u < v$  implies  $l(u) \leq l(v)$ . Then, for any source  $\mathfrak{X}$  with alphabet  $X$ ,

$$(2.2) \quad \zeta(n) = l(\xi(1,n)) - l(\xi(1,n-1)), \quad \tau(n) = l(\xi(1,n)) \quad (n=1,2,\dots)$$

defines a cost scale  $\mathfrak{Z}$ . Cost scales of this type will be called strictly intrinsic.

In particular, if

$$(2.3) \quad l(u) = \sum_{j=1}^n l(x_{i_j}) \quad (u = x_{i_1} x_{i_2} \dots x_{i_n})$$

the corresponding cost scale  $\mathfrak{Z}$  defined by

$$(2.4) \quad \zeta(n) = l(\xi(n)), \quad \tau(n) = \sum_{k=1}^n l(\xi(k)) \quad (n=1,2,\dots)$$

may be called a memoryless intrinsic cost scale; observe that  $l(x) \equiv 1$  gives rise to the counting scale  $\mathcal{C}$ . E.g.  $l(x)$  may be the length or duration of the symbol  $x \in X$ . Then, if the symbols are emitted consecutively, without

intervals,  $\tau(n) = \sum_{k=1}^n l(\xi(k))$  is the epoch at which the emission of the  $n$ 'th message symbol terminates.

EXAMPLE 2.3. The cost of transmission may depend on random outer disturbances independent of the symbols to be transmitted; in our model this means that  $\mathfrak{X}$  and  $\mathfrak{Z}$  are independent stochastic processes. The same holds if the symbols are emitted at random epochs, independent of the symbols themselves and "cost" means time.

EXAMPLE 2.4. A cost scale may be defined by letting  $\tau(n)$  denote the number of binary digits needed to encode the first  $n$  message symbols when a particular method of encoding is used.

REMARK 2.1. A cost scale  $\mathfrak{Z}$  is trivially regular if  $\frac{\nu(t)}{t}$  is uniformly bounded; e.g. a strictly intrinsic cost scale (cf. example 2.2) is surely regular if  $\frac{l(u)}{\|u\|}$  is bounded away from 0 ( $u \neq u_0$ ).

If  $\mathfrak{X}$  is a source, sequences of type  $\xi(1, n)$  will be called finite messages of  $\mathfrak{X}$  (for the notations cf. § 1, (E)). In particular,  $\eta(t) = \xi(1, \nu(t))$  is the message of cumulative cost just not exceeding  $t$  (with respect to the cost scale  $\mathfrak{Z}$ ). E.g. if "cost" is time then  $\eta(t)$  is the message emitted in the time interval  $[0, t]$ .

Obviously,  $\eta(t)$  is DRV; its possible "values" are finite sequences belonging to  $\mathcal{U}(\mathfrak{X})$ , including, possibly, the void sequence  $\alpha_0$ .

The entropy of  $\eta(t) = \xi(1, \nu(t))$  can be considered as the average information content of a message of cost  $t$  (with respect to the given cost scale). This suggests the following.

DEFINITION 2.2. The entropy rate of the source  $\mathfrak{X}$  with respect to the cost scale  $\mathfrak{Z}$  is the limit

$$(2.5) \quad H(\mathfrak{X} \parallel \mathfrak{Z}) = \lim_{t \rightarrow \infty} \frac{1}{t} H(\eta(t))$$

provided that it exists. If the limit does not exist, we shall denote the limsup and liminf of  $\frac{1}{t} H(\eta(t))$  by  $\bar{H}(\mathfrak{X} \parallel \mathfrak{Z})$  and  $\underline{H}(\mathfrak{X} \parallel \mathfrak{Z})$ , respectively. (The double bar is used in order to avoid confusion with conditional entropy.)

If the cost of each symbol is unity i. e.  $\mathfrak{Z} = \mathcal{E}$  (cf. example 2.1) then  $\eta(t) = \xi(1, [t])$  reduces to the usual definition of entropy per symbol

$$(2.6) \quad H(\mathfrak{X}) = H(\mathfrak{X} \parallel \mathcal{E}) = \lim_{n \rightarrow \infty} \frac{1}{n} H(\xi(1, n)) .$$

The idea underlying definition 2.2 is that the relevant information is carried by the message symbols i. e. by the process  $\mathfrak{X}$  and not by the process  $\mathfrak{Z}$ .

In the sequel we shall omit the arguments  $t$  where doing so does not cause ambiguity.

LEMMA 2.1. For every regular cost scale

$$(2.7) \quad H(\nu(t)) = \sigma(t) \quad (t \rightarrow \infty)$$

holds. Furthermore, for two arbitrary cost scales  $\mathfrak{Z}_1$  and  $\mathfrak{Z}_2$

$$(2.8) \quad d(\nu_1, \nu_2) \leq 2E|\nu_1 - \nu_2| + 2 \log_2 3$$

thus if  $\frac{v_1(t) - v_2(t)}{t} \xrightarrow{L_1} 0$ , then (2.7) holds or does not hold simultaneously for  $\mathfrak{Z}_1$  and  $\mathfrak{Z}_2$ .

PROOF. If  $\mathfrak{Z}$  is regular i. e. if  $\frac{v(t)}{t}$  is u. i. for  $t \rightarrow \infty$  then  $E(v(t)) = O(t)$  and  $E \log_2(v+1) \leq \log_2 E(v+1) = \log_2 O(t) = o(t)$  holds. Furthermore, as

$$d(v_1, v_2) = H(v_1|v_2) + H(v_2|v_1) = H(v_1 - v_2|v_2) + H(v_2 - v_1|v_1) \leq 2H(v_1 - v_2)$$

the inequality (2.8) is a consequence of lemma 1.2. The last assertion follows from (2.8) and lemma 1.1 ((1.10)).

In this section we shall be interested in the relation of entropy rates with respect to different cost scales of the same information source. Let us remark that for memoryless intrinsic cost scales (2.4), when interpreting  $\ell(x)$  as the length or duration of the symbol  $x \in X$ , "Entropy per second" is commonly defined (cf. [16], [12], etc.) as entropy per symbol (2.6) divided by the "average symbol length"  $L$  and not as in definition 2.2; in some reasonings, however, this  $\frac{H(\mathfrak{X})}{L}$  is implicitly replaced by our  $H(\mathfrak{X}||\mathfrak{Z})$ . Our results will, in particular, provide a justification for such reasonings under general conditions (cf. theorem 2.4).

For an arbitrary real number  $r$  and  $k > 0$  we set

$$(2.9) \quad |r|^+ = \max(0, r), \quad |r|^- = -\min(0, r), \quad |r| = |r|^+ + |r|^-$$

$$|r|_k^+ = \min(|r|^+, k), \quad |r|_k^- = \min(|r|^-, k), \quad |r|_k = \min(|r|, k).$$

The following estimates will play fundamental role in the sequel.

THEOREM 2.1. Let  $\mathfrak{X}$  be a source with alphabet  $X$  of size  $s$  and let  $\mathfrak{Z}_1$  and  $\mathfrak{Z}_2$  be two different cost scales for  $\mathfrak{X}$ . Then

$$(2.10) \quad H(\eta_1|\eta_2) \leq H(v_1|v_2) + \log_2 s E|v_1 - v_2|^+$$

and also, if  $A \in \mathcal{F}$  is such that on  $\bar{A} = \Omega \setminus A$  we have  $v_1 - v_2 \leq kt$  (where  $k > 0$  is arbitrary)

$$(2.11) \quad H(\eta_1 | \eta_2) \leq 1 + P(\bar{A}) H_{\bar{A}}(v_1 | v_2) + \log_2 \delta E |v_1 - v_2|_{kt} + P(A) H_A(\eta_1) .$$

PROOF. As  $v_i = \|\eta_i\|$  is uniquely determined by  $\eta_i = \xi(1, v_i)$  ( $i = 1, 2$ ), we have

$$(2.12) \quad H(\eta_1 | \eta_2) = H(v_1, \eta_1 | v_2, \eta_2) \leq H(v_1 | v_2) + H(\eta_1 | v_1, v_2, \eta_2) .$$

As for given  $v_1, v_2$  and  $\eta_2 = \xi(1, v_2)$  the number of possible "values" of  $\eta_1 = \xi(1, v_1)$  is at most  $\delta^{|v_1 - v_2|^+}$ , the last term in (2.12) is  $\leq E(\log_2 \delta^{|v_1 - v_2|^+}) = \log_2 \delta E |v_1 - v_2|^+$  proving (2.10). We may also write (setting  $\alpha = 1$  if  $\omega \in A$  and  $\alpha = 0$  otherwise),  $H(\eta_1 | \eta_2) \leq H(\alpha, \eta_1 | \eta_2) = H(\alpha) + P(\bar{A}) H_{\bar{A}}(\eta_1 | \eta_2) + P(A) H_A(\eta_1 | \eta_2)$  hence, applying (2.10) to  $H_{\bar{A}}(\eta_1 | \eta_2)$  and taking into account

$$(2.13) \quad P(\bar{A}) E_{\bar{A}} |v_1 - v_2|^+ \leq E |v_1 - v_2|_{kt}^+$$

and using the obvious inequalities  $H(\alpha) \leq 1$ ,  $H_A(\eta_1 | \eta_2) \leq H_A(\eta_1)$  we obtain (2.11).

In order that the implications of theorem 2.1 can be formulated concisely, we introduce some definitions concerning cost scales.

If  $\mathcal{Z}$  is an arbitrary cost scale and  $c > 0$ , we may define the cost scale  $c\mathcal{Z}$  as the sequence of RV's  $cZ(n)$  ( $n = 1, 2, \dots$ ). Let us denote the  $\tau$ 's and  $v$ 's corresponding to the cost scale  $c\mathcal{Z}$  by  $\tau^c$  and  $v^c$ , respectively:  $\tau^c(n) = c\tau(n)$ ,  $n = 1, 2, \dots$ ;  $v^c(t) = v(\frac{t}{c})$ ,  $0 \leq t < \infty$ .

DEFINITION 2.3. For two cost scales  $\mathcal{Z}_1$  and  $\mathcal{Z}_2$  we write

$$(2.14) \quad \mathcal{Z}_1 \sim \mathcal{Z}_2 \quad \text{iff} \quad \frac{1}{t} (v_1(t) - v_2(t)) \xrightarrow{L_1} 0 .$$

If  $\mathcal{Z}_1 \sim \mathcal{Z}_2$ , we say that  $\mathcal{Z}_1$  is equivalent to  $\mathcal{Z}_2$ .

The cost scales  $\mathcal{Z}_1$  and  $\mathcal{Z}_2$  will be said to be quasi-equivalent with quotient  $c_{12} = c > 0$  if  $\mathcal{Z}_1 \sim c\mathcal{Z}_2$ , i.e. if

$$(2.15) \quad \frac{1}{t} (v_1(t) - v_2^c(t)) = \frac{1}{t} (v_1(t) - v_2(\frac{t}{c})) \xrightarrow{L_1} 0 .$$

Of course, in case of regular cost scales, the replacement of  $L_1$ -convergence by stochastic convergence in the above definitions (in (2.14) and (2.15)) makes no difference.

Intuitively,  $\mathfrak{X}_1 \sim \mathfrak{X}_2$  means that the cost of one symbol is essentially the same for both scales.

The equivalence and the quasi-equivalence are equivalence relations, and the quotient  $c_{12}$  is uniquely determined by  $\mathfrak{X}_1$  and  $\mathfrak{X}_2$ , except for the trivial case  $E v_i(t) = \sigma(t)$  ( $i=1,2$ ). If  $\mathfrak{X}_1, \mathfrak{X}_2$  and  $\mathfrak{X}_3$  are quasi-equivalent cost scales, there obviously holds

$$(2.16) \quad c_{12} c_{23} = c_{13}, \quad c_{21} = \frac{1}{c_{12}} .$$

**THEOREM 2.** Let  $\mathfrak{X}_1$  and  $\mathfrak{X}_2$  be two cost scales for a source  $\mathfrak{X}$  with finite alphabet  $\mathcal{X}$ , and let one of the entropy rates  $H(\mathfrak{X} \parallel \mathfrak{X}_1)$  and  $H(\mathfrak{X} \parallel \mathfrak{X}_2)$  exist. If  $\mathfrak{X}_1$  is quasi-equivalent to  $\mathfrak{X}_2$  with quotient  $c_{12} = c > 0$  then

$$(2.17) \quad H(\mathfrak{X} \parallel \mathfrak{X}_1) = \frac{1}{c} H(\mathfrak{X} \parallel \mathfrak{X}_2) .$$

If the entropy rates in question do not exist, the assertion remains true both for the lower and upper entropy rates.

**PROOF.** (2.10) implies

$$(2.18) \quad d(\eta_1, \eta_2^c) \leq d(v_1, v_2^c) + \log_2 \Delta E |v_1 - v_2^c| .$$

Hence, using (1.10) and (2.8), we obtain

$$|H(\eta_1) - H(\eta_2^c)| \leq (2 + \log_2 \Delta) E |v_1 - v_2^c| + 2 \log_2 3$$

what results by definition the desired statement

$$\begin{aligned}
 H(\mathfrak{X} \parallel \mathfrak{Z}_1) &= \lim_{t \rightarrow \infty} \frac{1}{t} H(\eta_1(t)) = \lim_{t \rightarrow \infty} \frac{1}{t} H(\eta_2^c(t)) = \\
 &= \lim_{t \rightarrow \infty} \frac{1}{ct} H(\eta_2(t)) = \frac{1}{c} H(\mathfrak{X} \parallel \mathfrak{Z}_2)
 \end{aligned}$$

(where  $\eta^c(t) = \xi(1, v^c(t)) = \xi(1, v(\frac{t}{c})) = \eta(\frac{t}{c})$ ).

In order to apply theorem 2.2 to concrete problems it will be convenient to establish some simple sufficient conditions of the relation  $\sim$  for different cost scales.

LEMMA 2.2. Let  $\mathfrak{Z}_1$  and  $\mathfrak{Z}_2$  be two regular cost scales and let one of them have the property

$$(2.19) \quad b \leq \zeta(n) \leq B \quad \text{a.s.} \quad (0 < b < B; n = 1, 2, \dots).$$

Then each of the conditions

$$(2.20) \quad \frac{v_1(t)}{v_2(rt)} \xrightarrow{P} 1$$

$$(2.21) \quad \frac{1}{t} \tau_2(v_1(t)) \xrightarrow{P} r > 0$$

$$(2.22) \quad \frac{1}{t} \tau_1(v_2(t)) \xrightarrow{P} c > 0$$

is equivalent to  $\mathfrak{Z}_1 \sim_c \mathfrak{Z}_2$ , where  $c = \frac{1}{r}$ .

PROOF. For regular cost scales  $\mathfrak{Z}_1 \sim \frac{1}{r} \mathfrak{Z}_2$  means

$$(2.23) \quad \frac{1}{t} (v_1(t) - v_2(rt)) \xrightarrow{P} 0.$$

Without any loss of generality, let e.g.  $\mathfrak{Z}_2$  have the property (2.19), i. e.  $0 < b \leq \zeta_2(n) \leq B$ ; then  $[\frac{rt}{B}] \leq v_2(rt) \leq \frac{rt}{b}$  a.s., thus the equivalence of (2.20) and (2.23) is obvious. Furthermore, as by the definition

of the  $\tau$ 's and  $\nu$ 's the relation  $\tau_2(\nu_1(t)) \leq yt$  is equivalent to  $\nu_1(t) \leq \nu_2(yt)$  ( $0 < y < \infty$ ), the relation (2.21) is equivalent to

$$P(\nu_1(t) \leq \nu_2(yt)) \longrightarrow \begin{cases} 0 & \text{if } y < r \\ 1 & \text{if } y > r \end{cases} \quad (t \rightarrow \infty)$$

and this, in turn, is equivalent to (2.23), in view of the assumption  $0 < b \leq \zeta_2(n) \leq B$ . Similarly, (2.28) is equivalent to

$$P(\nu_2(t) \leq \nu_1(yt)) \longrightarrow \begin{cases} 0 & \text{if } y < c \\ 1 & \text{if } y > c \end{cases} \quad (t \rightarrow \infty)$$

i. e. to

$$P(\nu_2(y't) \leq \nu_1(t)) \longrightarrow \begin{cases} 0 & \text{if } y' = \frac{1}{y} > \frac{1}{c} = r \\ 1 & \text{if } y' < r \end{cases} \quad (t \rightarrow \infty)$$

which, again, is equivalent to (2.23).

The following consequence of theorem 2.2 and lemma 2.2 is worth being formulated as a new theorem.

**THEOREM 2.3.** Let  $\mathfrak{X}$  be a source with finite alphabet and let  $\mathfrak{X}_1$  and  $\mathfrak{X}_2$  be two regular cost scales for  $\mathfrak{X}$  such that one of them has the property (2.19). Then either of the three equivalent conditions (2.20)-(2.22) implies

$$(2.24) \quad H(\mathfrak{X} \parallel \mathfrak{X}_1) = r H(\mathfrak{X} \parallel \mathfrak{X}_2)$$

(in the sense that if either side exists so does also the other and they are equal). If the entropy rates in question do not exist, (2.24) still holds both for the lower and upper entropy rates.

The most important cost scales are those quasi-equivalent to the counting scale  $\mathcal{C}$  (cf. example 2.1). By definition 1.3, a cost scale  $\mathfrak{X}$  is quasi-equivalent to the counting scale  $\mathcal{C}$  iff there exists a constant  $c > 0$  such that  $\frac{\nu(t)}{t} \xrightarrow{L_1} \frac{1}{c}$ . Of course, all cost scales quasi-equivalent to the counting scale are regular (cf. Lemma 2.1).



If for a given cost scale  $\frac{v(t)}{t}$  converges in probability to a (finite) constant  $r$ , this  $r$  will be called the symbol rate of the source with respect to the given cost scale. Similarly, if  $\frac{\tau(n)}{n}$  converges in probability to a (finite) constant  $c$ , this  $c$  will be called the average symbol cost (with respect to the given cost scale).

In view of lemma 2.2 (with  $\mathcal{Z}_1 = \mathcal{X}$ ,  $\mathcal{Z}_2 = \mathcal{C}$ ) a symbol rate  $r > 0$  exists iff an average symbol cost  $c > 0$  exists ( $c = \frac{1}{r}$ ) and these are also necessary and sufficient conditions of  $\mathcal{X} \sim c\mathcal{C}$  for regular  $\mathcal{X}$ . Thus we obtain, as a particular case of theorem 2.3.

**THEOREM 2.4.** If for a source  $\mathcal{X}$  with finite alphabet  $X$  and with a regular cost scale  $\mathcal{Z}$  a positive symbol rate  $r$  or, equivalently, a positive symbol cost  $c$  exists ( $cr = 1$ ) then

$$(2.25) \quad H(\mathcal{X} \parallel \mathcal{Z}) = rH(\mathcal{X}) = \frac{1}{c} H(\mathcal{X})$$

if either of  $H(\mathcal{X})$  and  $H(\mathcal{X} \parallel \mathcal{Z})$  exists.

**REMARK 2.2.** If  $\mathcal{Z}$  is a cost scale of (2.4) (more general cost scales do not seem to have been considered in the literature) and the source  $\mathcal{X}$  is a stationary ergodic source (i. e.  $\xi(1), \xi(2), \dots$  is a stationary ergodic sequence of DRV's), then

$$(2.26) \quad \frac{1}{n} \sum_{k=1}^n \ell(\xi(k)) \xrightarrow{P} L > 0 .$$

In this case "entropy per unit cost" is often defined as the ratio  $\frac{H(\mathcal{X})}{L}$ . By theorem 2.4 this definition is equivalent to (2.5), provided that  $\mathcal{X}$  is regular. This is the case, in particular, if  $\ell(x) > 0$  for all  $x \in X$ , or if the message symbols  $\xi(n)$  are independent and identically distributed (in the latter case,  $L = E \zeta(1)$  by the law of large numbers and  $\frac{E(v(t))}{t} \rightarrow \frac{1}{E \zeta(1)}$  by the renewal theorem, and this ensures the regularity by lemma 1.3.).

**EXAMPLE 2.5.** Let us be given a finite-state noiseless channel as defined by Shannon [16]; such a channel is specified by the input alphabet  $X$ , the set of states  $A$ , an assignment of subsets  $X(\alpha)$  of  $X$  to each  $\alpha \in A$  and by a function  $G(x, \alpha)$  defined for  $\alpha \in A$ ,  $x \in X(\alpha)$  and taking its

values in  $A$  ( $X$  and  $A$  are finite sets). In each state  $\alpha$ , the channel is capable of transmitting letters from  $X(\alpha)$  only and if  $x \in X(\alpha)$  is transmitted, the new state will be  $\alpha' = G(x, \alpha)$ . Let an  $\alpha_{i_0} \in A$  be fixed as the initial state;

a sequence  $u = x_{i_1} x_{i_2} \dots x_{i_n} \in U(X)$  is transmissible iff  $x_{i_k} \in X(\alpha_{i_{k-1}})$  ( $k=1, \dots, n$ ) where  $\alpha_{i_k}$  is defined recursively by  $\alpha_{i_k} = G(x_{i_k}, \alpha_{i_{k-1}})$  ( $k=1, 2, \dots, n$ )

denote the set of all transmissible sequences by  $U_0$ . Let  $\ell(x, \alpha) \geq 0$ ,

( $\alpha \in A, x \in X(\alpha)$ ) be the cost of transmission of  $x$  at the state  $\alpha$ ;

then  $\ell(u) = \sum_{k=1}^n \ell(x_{i_k}, \alpha_{i_{k-1}})$  represents the cost of transmission of the

sequence  $u = x_{i_1} x_{i_2} \dots x_{i_n} \in U_0$ . We make the usual assumptions that for each pair

of states  $\alpha', \alpha'' \in A$  there exists  $u = x_{i_1} \dots x_{i_k} \dots x_{i_n} \in U_0$  such that

$\alpha_{i_k} = \alpha', \alpha_{i_n} = \alpha''$ , and that for all  $u \in U_0$  with  $\|u\| > m$  (say)  $\ell(u) > 0$ . Let  $\mathfrak{X}$

be a source transmissible by the channel (i. e.  $\xi(1) \dots \xi(n) \in U_0$  a. s.

$n=1, 2, \dots$ ) and let the (regular) cost scale  $\mathfrak{X}$  be defined by (2.2), with

the present  $\ell(u)$ . Let  $N(t)$  denote the number of different sequences  $u \in U_0$

with  $\ell(u) \leq t$  and such that  $\ell(ux) > t$  for some  $x \in X$ , with  $ux \in U_0$ ; we define the

channel capacity by

$$(2.27) \quad C = \overline{\lim}_{t \rightarrow \infty} \frac{1}{t} \log_2 N(t) .$$

Then  $\eta(t) = \xi(1, \nu(t))$  has at most  $N(t)$  possible values, implying

$$H(\eta(t)) = \log_2 N(t), \quad \bar{H}(\mathfrak{X} \| \mathfrak{X}) \leq C \quad ; \quad \text{thus if } \frac{\tau(n)}{n} = \frac{1}{n} \ell(\xi(1, n)) \xrightarrow{P} L$$

we have, according to (2.25),

$$(2.28) \quad \frac{1}{L} \bar{H}(\mathfrak{X}) \leq C .$$

The inequality (2.28), first appearing in Shannon's fundamental paper [16] is often considered to be "obvious". However, its familiar "justification" relies on the equivalence of the two possible definitions of entropy per unit cost and can be made rigorous only on the basis of theorem 2.4. The existing rigorous proofs of (2.28) concern stationary Markovian sources only ([7], [15]) or the case of one state ([5], [12]).

Let us now consider the inequality (2.28) in the most general form that still follows from theorem 2.3.

THEOREM 2.5. Let  $\mathfrak{X}$  be a source with finite alphabet and let  $\mathfrak{z}_1$  and  $\mathfrak{z}_2$  two regular cost scales for  $\mathfrak{X}$  such that one of them has the property (2.19) and

$$\frac{1}{t} \tau_2(v_1(t)) \xrightarrow{P} L > 0.$$

If  $N(t)$  denotes the number of different possible values of  $\eta_2(t) = \xi(1, v_2(t))$  and

$$C = \overline{\lim}_{t \rightarrow \infty} \frac{1}{t} \log_2 N(t)$$

then

$$(2.29) \quad \frac{1}{L} \bar{H}(\mathfrak{X} \parallel \mathfrak{z}_1) \leq C.$$

In the original paper (cf. footnote 1) we proved (2.29) under more general conditions which did not ensure the validity of (2.24).

### § 3. THE PRINCIPLE OF CONSERVATION OF ENTROPY

Let  $\mathfrak{X}$  and  $\mathfrak{X}'$  be two sources with finite alphabets  $X$  and  $X'$ , respectively. Let  $\mathfrak{z}$  be a cost scale for  $\mathfrak{X}$  and  $\mathfrak{z}'$  a cost scale for  $\mathfrak{X}'$ ; then we define the rates

$$(3.1) \quad H(\mathfrak{X}, \mathfrak{X}' \parallel \mathfrak{z}, \mathfrak{z}') = \lim_{t \rightarrow \infty} \frac{1}{t} H(\eta(t), \eta'(t))$$

$$(3.2) \quad H(\mathfrak{X} \mid \mathfrak{X}' \parallel \mathfrak{z}, \mathfrak{z}') = \lim_{t \rightarrow \infty} \frac{1}{t} H(\eta(t) \mid \eta'(t))$$

$$(3.3) \quad I(\mathfrak{X}, \mathfrak{X}' \parallel \mathfrak{z}, \mathfrak{z}') = \lim_{t \rightarrow \infty} \frac{1}{t} I(\eta(t), \eta'(t))$$

$$(3.4) \quad d(\mathfrak{X}, \mathfrak{X}' \parallel \mathfrak{z}, \mathfrak{z}') = \lim_{t \rightarrow \infty} \frac{1}{t} d(\eta(t), \eta'(t))$$

provided that the limits exist ( $\eta(t)$  and  $\eta'(t)$  are defined as in § 1, (E)); the argument  $t$  will be often omitted). If the limits do not exist, one may consider the corresponding upper and lower rates (upper and lower limits).

As an immediate consequence of theorem 2.1, lemma 2.1 and lemma 1.1 we have

**THEOREM 3.1.** All the rates (3.1)-(3.4), as well as the corresponding upper and lower rates, remains unchanged if either of  $\mathfrak{X}$  and  $\mathfrak{X}'$  is replaced by an equivalent cost scale.

In this section we shall apply the results of Section 2 to the case that  $\mathfrak{X}'$  is obtained from  $\mathfrak{X}$  by encoding (or conversely). For this purpose we define the codes in a very general sense.

**DEFINITION 3.1.** Let  $\mathfrak{X}$  and  $\mathfrak{X}'$  be finite alphabets. An arbitrary mapping  $f$  of  $\mathfrak{U}(\mathfrak{X})$  into  $\mathfrak{U}(\mathfrak{X}')$  such that  $f(u_0) = u_0$  and

$$(3.5) \quad u \prec v \quad \text{implies} \quad f(u) \prec f(v)$$

will be called a code from  $\mathfrak{X}$  to  $\mathfrak{X}'$ .

The code of a sequence  $u = x_{i_1} \dots x_{i_n}$  we can write in the form

$$(3.6) \quad f(u) = g(x_{i_1} | u_0) g(x_{i_2} | x_{i_1}) \dots g(x_{i_n} | x_{i_1} \dots x_{i_{n-1}}).$$

Each code  $f$  from  $\mathfrak{X}$  to  $\mathfrak{X}'$  defines a mapping of infinite sequences, too; in fact, for  $\tilde{u} = x_{i_1} x_{i_2} \dots \in \tilde{\mathfrak{U}}(\mathfrak{X})$  we may write

$$f(\tilde{u}) = g(x_{i_1} | u_0) g(x_{i_2} | x_{i_1}) \dots g(x_{i_n} | x_{i_1} \dots x_{i_{n-1}}) \dots \quad . \text{ Observe}$$

that the definition 3.1 does not exclude that the "code"  $f(\tilde{u})$  of some infinite sequence  $\tilde{u}$  be finite. The set of those  $\tilde{u} \in \tilde{\mathfrak{U}}(\mathfrak{X})$  for which  $f(\tilde{u})$  is infinite, will be denoted by  $\tilde{\mathfrak{D}}_f$ .

**REMARK 3.1.** It is easy to see that  $\tilde{\mathfrak{D}}_f \in \mathfrak{B}$  for any code  $f$  and also, the mapping  $f: \tilde{\mathfrak{D}}_f \rightarrow \tilde{\mathfrak{U}}(\mathfrak{X}')$  is measurable with respect to the  $\sigma$ -algebras  $\mathfrak{B}$  and  $\mathfrak{B}'$  spanned by the cylinder sets. If  $f$  is a code, the mapping  $f: \tilde{\mathfrak{D}}_f \rightarrow \tilde{\mathfrak{U}}(\mathfrak{X}')$  may be called an infinite-code; clearly, the concept of an infinite-code is much more restrictive than that of an arbitrary measurable mapping  $\tilde{\mathfrak{U}}(\mathfrak{X}) \rightarrow \tilde{\mathfrak{U}}(\mathfrak{X}')$  (though all practically realizable mappings seem to belong to this class). Two different codes  $f_1$  and  $f_2$  may give rise to the same infinite-code. In this case we shall say that  $f_1$  and  $f_2$  are equivalent and write  $f_1 \sim f_2$ . Of course, each infinite-code may be identified with the corresponding equivalence class of codes and vice-versa.

If  $f$  is a code from  $X$  to  $X'$  and  $u = x_{i_1} x_{i_2} \dots x_{i_p} \in U(X)$  let  $k_n = k_n(u)$  denote the  $n$ 'th point of increase of the sequence

$$(3.7) \quad 0, \|f(x_{i_1})\|, \dots, \|f(x_{i_1}, \dots, x_{i_m})\|, \dots, \|f(x_{i_1}, \dots, x_{i_p})\|$$

and  $k_0 = 0$ . For  $\tilde{u} \in \tilde{U}(X)$  we define  $k_n(\tilde{u})$  in a similar way; in particular, if  $\tilde{u} \in \tilde{D}_f$ , then  $k_n(\tilde{u})$  is well defined for  $n = 0, 1, 2, \dots$

EXAMPLE 3.1. If in (3.6)  $g(x|u) = g(x)$ , where  $g(x)$  is some mapping of  $X$  into  $U(X') - \{u_0\}$ , the resulting code will be called a simple letter code. In this case we have

$$(3.8) \quad \begin{aligned} f(u) &= g(x_{i_1}) g(x_{i_2}) \dots g(x_{i_n}) \quad (u = x_{i_1} \dots x_{i_n} \in U(X)) \\ f(\tilde{u}) &= g(x_{i_1}) g(x_{i_2}) \dots \quad (\tilde{u} = x_{i_1} x_{i_2} \dots \in \tilde{U}(X)). \end{aligned}$$

EXAMPLE 3.2. Let  $X$  and  $X'$  be finite alphabets, let  $A$  be a finite set to be called the set of states, and let us be given two functions  $F$  and  $G$  mapping the Cartesian product  $X \times A$  into  $U(X')$  and  $A$ , respectively. Let an initial state  $a_0 \in A$  be specified, set  $f(u_0) = u_0$  and for  $u = x_{i_1} x_{i_2} \dots x_{i_n}, n > 0$  set

$$(3.9) \quad f(u) = F(x_{i_1}, a_0) F(x_{i_2}, a_1) \dots F(x_{i_n}, a_{n-1})$$

where the states  $a_k$  are defined, recursively by  $a_k = G(x_{i_k}, a_{k-1})$

Then  $f$  is a code in the sense of definition 3.1; the encoder  $(X, X', A, a_0, F, G)$  will be called, following Shannon [16], a finite-state transducer and  $f$  will be referred to as the code generated by this finite-state transducer.

DEFINITION 3.2. Let  $\mathfrak{X}$  be a source with finite alphabet  $X$  and let  $f$  be a code from  $X$  to  $X'$ . We say that  $\mathfrak{X}$  is encodable by  $f$  if the IRV's

$$(3.10) \quad x(n) = k_n(\xi(1) \xi(2) \dots)$$

( $k_n$  has been defined in connection with (3.7)) are well defined a. s. for  $n = 1, 2, \dots$ .)

The encoding results in a new source  $\mathfrak{X}' = f(\mathfrak{X})$  defined by

$$\xi'(1), \xi'(2), \dots; \xi'(1) \xi'(2) \dots = f(\xi(1) \xi(2) \dots).$$

If there is given a cost scale  $\mathfrak{z}_1$  for  $\mathfrak{X}$  then we also define the mapped cost scale  $\mathfrak{z}'_1 = f(\mathfrak{z}_1)$  by

$$(3.11) \quad \nu'_1(t) = \|f(\eta_1(t))\|$$

or, equivalently, by

$$(3.12) \quad \tau'_1(n) = t \quad \text{iff} \quad \|f(\eta_1(t'))\| < n < \|f(\eta_1(t))\| \quad \text{for all } t' < t.$$

REMARK 3.2. If a source  $\mathfrak{X}$  is encodable by a code  $f$ , the mapped cost scale  $f(\mathfrak{z}_1)$  need not be regular, in general. In order that the regularity of  $\mathfrak{z}_1$  imply that of  $f(\mathfrak{z}_1)$ , a simple sufficient condition consists in the boundedness of  $\frac{\|f(u)\|}{\|u\|}$  ( $u \neq u_0$ ). In particular, for codes generated by finite-state transducers (cf. example 3.2), this condition is trivially fulfilled.

In order that a code be of any practical value, it ought to be possible, in some sense, to recover the original message from its encoded form. We adopt the following

DEFINITION 3.3. Let  $\mathfrak{X}$  be a source with finite alphabet  $\mathfrak{X}$  and let  $f$  be a code such that  $\mathfrak{X}$  encodable by  $f$ . The encoding  $f: \mathfrak{X} \rightarrow \mathfrak{X}'$  will be said to be finite decodable if there exists a natural number  $d$  such that to any  $u'$  there are at most  $d$   $u$ 's satisfying  $f(u) = u'$ .

THEOREM 3.2. Let  $\mathfrak{X}$  be a source with a cost scale  $\mathfrak{z}_1$  and let  $\mathfrak{X}$  be encodable by a code  $f$ . Let further the encoding result in the source  $\mathfrak{X}' = f(\mathfrak{X})$  and let  $\mathfrak{z}'_1$  be a cost scale for  $\mathfrak{X}'$ . Then if the code  $f: \mathfrak{X} \rightarrow \mathfrak{X}'$  is finite decodable and if the mapped cost scale  $\mathfrak{z}'_1 = f(\mathfrak{z}_1)$  is quasi-equivalent to  $\mathfrak{z}'_1$ , i. e.  $\mathfrak{z}'_1 \sim c' \mathfrak{z}'_1$ ;  $c' > 0$ , we have

$$(3.13) \quad H(\mathfrak{X}' \| \mathfrak{z}'_1) = c' H(\mathfrak{X} \| \mathfrak{z}_1).$$

PROOF. (3.11) implies  $\eta'_1(t) = \xi'(1, \nu'_1(t)) = f(\eta_1(t))$  and thus  $H(\eta'_1(t)) \leq H(\eta_1(t))$  i. e.

$$(3.14) \quad H(\mathfrak{X}' \| \mathfrak{z}'_1) \leq H(\mathfrak{X} \| \mathfrak{z}_1).$$

On the other hand, since  $f$  is finite decodable,

$$H(\eta_1 | \eta'_1) \leq \log_2 d = o(t)$$

holds. Hence, using (1.5), the inequality

$$(3.15) \quad H(\mathfrak{X}' \parallel \mathfrak{Z}'_1) \geq H(\mathfrak{X} \parallel \mathfrak{Z}_1)$$

follows. (3.14) and (3.15) give

$$(3.16) \quad H(\mathfrak{X}' \parallel \mathfrak{Z}'_1) = H(\mathfrak{X} \parallel \mathfrak{Z}_1) .$$

However,  $\mathfrak{Z}'_1 \sim c' \mathfrak{Z}'_1$ , thus by theorem 2.2 we obtain

$$(3.17) \quad H(\mathfrak{X}' \parallel \mathfrak{Z}'_1) = \frac{H(\mathfrak{X}' \parallel \mathfrak{Z}'_1)}{c'}$$

Finally, (3.16) and (3.17) give the desired (3.13).

The intuitive meaning of theorem 3.2 is clear. (3.13) represents the "principle of conservation of entropy", i. e. that finite decodable encoding does not change the entropy rate apart from a factor representing the quotient of average costs after and before encoding. If the coding is not decodable from (3.14) and (3.17) we can obtain

$$H(\mathfrak{X}' \parallel \mathfrak{Z}'_1) \leq c' H(\mathfrak{X} \parallel \mathfrak{Z}_1)$$

which means that in the non-decodable case some information may get lost.

The following corollary of theorem 3.2 is worth formulating as a new theorem.

**THEOREM 3.3.** Let  $\mathfrak{X}$  be a source with a given regular cost scale  $\mathfrak{Z}_1$ ; let  $\mathfrak{X}$  be encodable by a finite decodable code  $f$  and let  $\mathfrak{Z}'$  be a cost scale for  $\mathfrak{X}' = f(\mathfrak{X})$  such that

$$(3.18) \quad b \leq \zeta'(n) \leq B \quad \text{a. s.} \quad (n = 1, 2, \dots; 0 < b < B) .$$

Then if

$$(3.19) \quad \frac{1}{t} \tau'(\|f(\eta_1(t))\|) \xrightarrow{P} r > 0$$

and the code  $f$  has the property that  $\frac{\|f(u)\|}{\|u\|}$  is bounded, we have

$$(3.20) \quad H(\mathfrak{X}' \parallel \mathfrak{Z}') = \frac{H(\mathfrak{X} \parallel \mathfrak{Z}_1)}{r} .$$

PROOF. By remark 3.2 the obtained cost scale  $\mathfrak{X}' = f(\mathfrak{X}_1)$  is regular, similarly,  $\mathfrak{X}_1$  also is regular trivially. Using lemma 2.2 we obtain that (3.19) is equivalent to  $\mathfrak{X}' \sim \frac{1}{r} \mathfrak{X}'$ . Thus, theorem 3.2 leads to (3.20).

E.g. if  $\mathfrak{X}'$  has a memoryless intrinsic cost scale (cf. example 2.2) defined by fixed symbol costs  $\ell(x')$  ( $x' \in \mathfrak{X}'$ ) such that  $Z'_i(n) = \ell(\xi'_i(n))$  or somewhat more generally, if  $\mathfrak{X}'$  is to be transmitted by a finite-state noiseless channel (cf. example 2.5), the symbol costs may depend on the "state of the channel", i. e.  $Z'_i(n) = \ell(\xi'_i(n), \alpha(n-1))$  where  $\alpha(k)$  represents the state of the channel after the transmission of the  $k$ 'th message symbol. In these cases the condition (3.18) is trivially fulfilled provided that  $\ell(x')$  (or  $\ell(x', \alpha)$ ) is strictly positive.

REMARK 3.3. Theorem 3.3 is perhaps the most impressive form of the "principle of conservation of entropy". As  $\tau'(\|f(\eta_1(t))\|)$  is the cumulative cost of the code of a message of cumulative cost  $t$  (i. e. of  $\eta_1(t) = \xi_1(1, \nu_1(t))$ ) the condition (3.19) requires the existence of an average code cost  $r$  per unit message cost, in the sense of convergence in probability. If both  $\mathfrak{X}_1$  and  $\mathfrak{X}'$  are the counting scale  $\mathcal{C}$  (3.19) reduces to

$$(3.19') \quad \frac{1}{n} \|f(\xi(1, n))\| \xrightarrow{P} r$$

and the identity (3.20) becomes

$$(3.20') \quad H(\mathfrak{X}') = \frac{H(\mathfrak{X})}{r}$$

This relation, dating back to Shannon [16], has often been regarded as "obvious" but, to the authors' knowledge, it has never been proved in a rigorous way, for arbitrary sources and codes. For the case that  $f$  is a simple letter code, a proof of (3.20') appears in [10]. The more general case of codes generated by finite-state transducers (cf. example 3.2) has been considered by Sidel'nikov [17]; he however, restricted attention to Markovian sources (though, as he has remarked, some of his results hold in more general cases, too).

We conclude this section by exhibiting a general form of the "noiseless coding theorem" as a consequence of theorem 3.3. (In our original paper we prove a more general form which is not such simple consequence of theorem 3.3.)

THEOREM 3.4. Let  $\mathfrak{X}$  be a source with a given regular cost



scale  $\mathfrak{z}_1$ ; let  $\mathfrak{x}$  be encodable by a finite decodable code  $f$  and let  $\mathfrak{z}'$  be a cost scale for  $\mathfrak{x}' = f(\mathfrak{x})$  such that

$$b \leq \zeta'(n) \leq B \quad \text{a. s.} \quad (n = 1, 2, \dots; 0 < b < B).$$

Denote by  $N(t)$  the number of different possible values of  $\eta'(t) = \xi(1, \nu'(t))$  and put

$$C = \overline{\lim}_{t \rightarrow \infty} \frac{\log_2 N(t)}{t}.$$

Then if

$$\frac{1}{t} \tau'(\|f(\eta_1(t))\|) \xrightarrow{P} L > 0$$

and the code has the property that  $\frac{\|f(u)\|}{\|u\|}$  is bounded, we have

$$(3.21) \quad \frac{\overline{H}(\mathfrak{x} \| \mathfrak{z})}{C} \leq L.$$

PROOF. We have to observe only that by definition of  $C$

$$H(\eta'(t)) \leq \log_2 N(t) \quad \text{and} \quad H(\mathfrak{x}' \| \mathfrak{z}') \leq C$$

hold, thus the statement follows from (3.20).

REMARK 3.4. (3.21) has been proved for the special case if  $\mathfrak{z}_1 = \mathcal{C}$ ,  $\mathfrak{z}'$  is a memoryless intrinsic cost scale and the code is a simple letter code. For the case of stationary ergodic sources with  $\mathfrak{z}_1 = \mathcal{C}$  for  $\mathfrak{z}' = \mathcal{C}$  and for general codes defined by him, Billingsley [2] has proved even a stronger theorem than the corresponding particular case of theorem 3.4. For more general cases, to the authors' knowledge, the assertion is, though "intuitively obvious", as a mathematical theorem new.

#### § 4. CONCLUDING REMARKS

In this paper, we restricted ourselves to information sources with finite alphabet; the finiteness assumption has been essential for our basic estimations (2.12) and (2.13) (theorem 2.1) and it remains an open problem, under what conditions does the "entropy rate comparison theorem" (theorems 2.2 and 3.2) hold for countable alphabets, too. Theorem 2.5, however, remains unchanged also for countable  $X$  (except for the bound  $C \leq \frac{\log_2 s}{b}$ ). Observe, too, that the theorems involving coding (theorems 3.2, 3.3, 3.4) can obviously be extended to countable  $X$  (provided, in the case of the first two that  $X'$  remains finite). A possible approach to problems concerning countable alphabets in general would be to reduce them to the finite-alphabet case by an appropriate encoding, using the above remark.

As to the generality of the concept of coding used in this paper (definitions 3.1 and 3.2) one might make the objection that in some cases the code sequence assigned to a (finite) message sequence may conceivably depend not only on this sequence but on some subsequent letters, too. In all practical cases, however, the encoding is "of finite delay", i. e. the code sequence assigned to the first  $n$  letters of the message to be encoded is uniquely determined by these letters and  $m$  subsequent ones, where  $m$  is fixed. Then one may consider that the code sequence obtained in this way is actually assigned to the first  $n+m$  letters of the message sequence (rather to the first  $n$  cases); as  $m$  is constant, this change of viewpoint does not cause any change in the results.

It would be very desirable, both from the theoretical and practical points of view, to extend our results in order to include the case of information transmission in the presence of noise. E. g., our results may conceivably be useful in the theory of channels with error synchronisation, investigated by Dobrusin [8]. A closer study of noisy channels with arbitrary cost scales, however, is beyond the scope of the present paper.

In connection with the "noiseless coding theorem" (theorem 3.5), we did not tackle the problem whether the lower bound

$\frac{H(X||Z)}{C}$  of  $L$  can be attained (or approximated to any specified degree) by an appropriate encoding. In practically important cases, this question may be answered in the affirmative using familiar methods, though in the most general case there may arise some difficulties. Another problem we did not enter is that of generalizing McMillan's theorem for sources with cost scales; this problem, though of considerable interest, apparently requires different methods than those used in this paper.

## REFERENCES

- [1] R. ASH: Information Theory: New York, Interscience Publishers 1965.
- [2] P. BILLINGSLEY: On the coding theorem for the noiseless channel. Ann. Math. Statist. 32, 594-601 (1961)
- [3] N. M. BLACHMAN: Minimum-cost encoding of information. IRE Trans. Information Theory PGIT. 3, 139-149 (1954)
- [4] ---: Minimum-cost transmission of information. Information and Control 7, 508-511 (1964)
- [5] E. L. BLOH: Generalization of an inequality of information theory to the case of symbols of unequal duration. (In Russian) Problemy Peredaci Informacii 5, 95-99 (1960)
- [6] ---: Construction of optimal code constituted of elementary symbols of unequal duration (in Russian). Problemy Peredaci Informacii 5, 100-111 (1960)
- [7] I. CSISZÁR: Two remarks on noiseless coding. Information and Control. To appear.
- [8] R. L. DOBRUSIN: Capacity of channels with error synchronisation. Paper presented at the Colloquium of Information Theory, Debrecen, 1967.
- [9] A. FEINSTEIN: Foundations of Information Theory. New York: McGraw-Hill Book Co. 1958.
- [10] G. KATONA and G. TUSNÁDY: The principle of conservation of entropy in a noiseless channel. Studia Sci. Math. Acad. Sci. Hung. 2, 29-35 (1966)
- [11] J. R. KINNEY: Singular functions associated with Markov chains. Proc. Amer. Math. Soc. 9, 603-608 (1958)
- [12] R. M. KRAUSE: Channels which transmit letters of unequal duration. Information and Control 5, 13-24 (1962)
- [13] Ju. I. LJUVIC: Remark on the capacity of the discrete noiseless channel (in Russian). Uspehi Mat. Nauk. 17, 191-198 (1962)

- [14] W. PARRY: Intrinsic Markov chains. Trans. Amer. Math. Soc. 112, 55-66 (1964)
- [15] C. E. RADKE: Necessary and sufficient conditions on conditional probabilities to maximize entropy. Information and Control 9, 279-284 (1966)
- [16] C. E. SHANNON: A mathematical theory of communication. Bell System Techn. J. 27, 379-432, 623-656 (1948)
- [17] V. M. SIDEL'NIKOV: On statistical properties of transformations induced by finite automata (in Russian). Kibernetika (Kiev) 6, 1-14 (1965)
- [18] R. A. ZAIDMAN: On the asymptotics of certain sequences encountered in problems of non-Markov random walks and of information theory (in Russian). Vestnik Leningrad. Univ. 1, 23-33 (1965).