# Studia Scientiarum Mathematicarum Hungarica

G. KATONA and G. TUSNÁDY

## THE PRINCIPLE OF CONSERVATION OF ENTROPY IN A NOISELESS CHANNEL

1967

# THE PRINCIPLE OF CONSERVATION OF ENTROPY
# IN A NOISELESS CHANNEL

by

G. KATONA and G. TUSNÁDY

## Introduction

The main aim of this paper is to formulate precisely and to prove the following statement:

If we have an information source, more precisely, a sequence of random variables $\xi_1, \xi_2, \ldots$ with entropy $H(\mathscr{X})$ and we code this sequence in a uniquely decodable manner, the obtained sequence $\mathscr{Y}$ has the entropy

$$(1) \qquad\qquad H(\mathscr{Y}) = \frac{H(\mathscr{X})}{L},$$

where $L$ is the average length of the codes.

The intuitive meaning of (1) is clear: it expresses the principle of conservation of information, when the coding is uniquely decodable (and no noise is present). In spite of this according to our best knowledge (1) has not been proved in full generality up to now.

If we have finite number of code signals $y_1, \ldots, y_m$, the maximum of $H(\mathscr{Y})$ is $\log m$, where log denotes the logarithm with respect to the base 2. It follows from (1) that

$$\frac{H(\mathscr{X})}{\log m} \leq L.$$

This is a well known theorem of SHANNON, and according to our best knowledge only this consequence of (1) was proved (e. g. [1]).

In the case when the coding is not necessarily uniquely decodable instead of (1) we prove the inequality

$$H(\mathscr{Y}) \leq \frac{H(\mathscr{X})}{L},$$

which has also an intuitive meaning.

## Precise Formulation

Let $X = \{x_1, \ldots, x_n\}$ be the set of possible signals (the alphabet) of the information source, and let $X^\infty$ be the set of all infinite sequences formed from the letters $x_1, \ldots, x_n$. If $1 \leq i_1 \leq n, \ldots, 1 \leq i_k \leq n$, we denote by $[x_{i_1}, \ldots, x_{i_k}]$ the set of all sequences having $x_{i_1}, \ldots, x_{i_k}$ on the first $k$ places. We call such subsets of $X^\infty$ cylinder sets. Let $\mathfrak{A}_X$ denote the $\sigma$-field generated by the cylinder sets. The measure space $\mathscr{X} = (X^\infty, \mathfrak{A}_X, p_X)$

is called an information source, if $p_X$ is a probability measure on $\mathfrak{A}_X$. This space defines an other space on the sequences of length $k$: $\mathcal{X}^k = (X^k, \mathfrak{A}_X^k, p_X^k)$, where $X^k$ is the space of the sequences $(x_{i_1}, ..., x_{i_k})$. The average information contained in the first $k$ signals of the information source is

$$H(\mathcal{X}^k) = - \sum_{\substack{1 \leq i_1 \leq n \\ \vdots \\ 1 \leq i_k \leq n}} p_X^k(x_{i_1}, ..., x_{i_k}) \log p_X^k(x_{i_1}, ..., x_{i_k}) =$$

$$= - \sum_{\substack{1 \leq i_1 \leq n \\ \vdots \\ 1 \leq i_k \leq n}} p_X[x_{i_1}, ..., x_{i_k}] \log p_X[x_{i_1}, ..., x_{i_k}].$$

Finally, the definition of the entropy of $\mathcal{X}$ is

$$H(\mathcal{X}) = \lim_{k \to \infty} \frac{H(\mathcal{X}^k)}{k}$$

(the average information content of one signal), if this limit exists.

Consider now the definition of the coding. Let $Y = \{y_1, ..., y_m\}$ be the set of possible code signals. Let $c(x_i)$ $(1 \leq i \leq n)$ be a finite, non empty sequence formed from elements of $Y$. We call this function coding. Thus we may associate to every $(x_{i_1}, ..., x_{i_k})$, resp. $(x_{j_1}, x_{j_2}, ...)$ a sequence of $y_i$'s. Let us write successively the codes $c(x_{i_1}), ..., ..., c(x_{i_k})$, resp. $c(x_{j_1}), c(x_{j_2}), ...$. Denote by $c(x_{i_1}, ..., x_{i_k})$, resp. $d(x_{j_1}, x_{j_2}, ...)$ the resulting sequence. Let $Y^\infty$ be the set of all infinite $y$-sequences. Thus the function $d(x_{j_1}, x_{j_2}, ...)$ transforms the set $X^\infty$ into a subset $Y^*$ of $Y^\infty$. Let $\mathfrak{A}_Y$ be the $\sigma$-field in $Y^*$ generated by the mapping $d(x_{j_1}, x_{j_2}, ...)$ of $\mathfrak{A}_X$ and let us define the measure $P_Y$ on $\mathfrak{A}_Y$ by putting:

$$p_Y(A) = p_X(d^{-1}(A)) \quad A \in \mathfrak{A}_Y$$

where $d^{-1}(A)$ denotes the inverse image of $A$. $\mathcal{Y} = (Y^*, \mathfrak{A}_Y, p_Y)$ is the space of coded sequences. As above $[y_{i_1}, ..., y_{i_k}]$ denotes the cylinder set consisting of all sequences in $Y^\infty$ of which the first $k$ terms are $y_{i_1}, ..., y_{i_k}$.

LEMMA 1. $[y_{i_1}, ..., y_{i_k}] \cap Y^* \in \mathfrak{A}_Y$.

PROOF. We have to prove that the set of all sequences $(x_{j_1}, x_{j_2}, ...)$ having the image in $[y_{i_1}, ..., y_{i_k}] \cap Y^*$ is in $\mathfrak{A}_X$. We say that $(y_{i_1}, ..., y_{i_k})$ is a segment of $(y_{l_1}, ..., y_{l_s})$ if $k \leq s$, and $y_{i_1} = y_{l_1}, ..., y_{i_k} = y_{l_k}$. Obviously, the image of $(x_{j_1}, x_{j_2}, ...)$ is in $[y_{i_1}, ..., y_{i_k}] \cap Y^*$ if and only if $(y_{i_1}, ..., y_{i_k})$ is a segment of $c(x_{j_1}, ..., x_{j_k})$. Thus

(2) $$d^{-1}([y_{i_1}, ..., y_{i_k}] \cap Y^*) = \cup [x_{j_1}, ..., x_{j_k}],$$

where union runs over sequences $j_1, ..., j_k$ for which $(y_{i_1}, ..., y_{i_k})$ is a segment of $c(x_{j_1}, ..., x_{j_k})$. However the right side of (2) is a union of cylinder sets in $\mathfrak{A}_X$ which proves the Lemma.

Denote by $Y^k$ the set of sequences $(y_{i_1}, ..., y_{i_k})$ for which the cylinder set $[y_{i_1}, ..., y_{i_k}]$ is in $Y^*$. Let $\mathfrak{A}_Y^k$ be the $\sigma$-field of all subsets of $Y^k$. By Lemma 1 we can define

$$p_Y^k(y_{i_1}, ..., y_{i_k}) = p_Y([y_{i_1}, ..., y_{i_k}] \cap Y^*).$$

The average information content of the first $k$ signals of the coded sequence is

$$H(\mathcal{Y}^k) = - \sum_{\substack{1 \le y_{i_1} \le m \\ \vdots \\ 1 \le y_{i_k} \le m}} p_Y^k(y_{i_1}, \ldots, y_{i_k}) \log p_Y^k(y_{i_1}, \ldots, y_{i_k}),$$

where $\mathcal{Y}^k = (Y^k, \mathfrak{A}_Y^k, p_Y^k)$. Finally the definition of the entropy of $\mathcal{Y}$ is

$$H(\mathcal{Y}) = \lim_{k \to \infty} \frac{H(\mathcal{Y}^k)}{k}$$

(the average information content of one signal of the coded sequence), if this limit exists.

If $l_i$ denotes the length of the sequence $c(x_i)$ $(1 \le i \le n)$ then the length of $c(x_{i_1}, \ldots, x_{i_k})$ is $\sum_{j=1}^{k} l_{i_j}$. Let $L_k$ be a random variable in the probability space $\mathcal{X}^k$, which takes on the value $\sum_{j=1}^{k} l_{i_j}$ if we have the sequence $(x_{i_1}, \ldots, x_{i_k})$. We say that the average code length is $L$, if

$$\frac{L_k}{k} \Rightarrow L$$

for $k \to \infty$, where $\Rightarrow$ denotes convergence in probability.

A coding $c(x_i)$ $(1 \le i \le n)$ is called uniquely decodable, if

$$c(x_{i_1}, \ldots, x_{i_k}) = c(x_{j_1}, \ldots, x_{j_s})$$

holds only in the case $k = s$, $x_{i_1} = x_{j_1}, \ldots, x_{i_k} = x_{j_k}$.

We would like to point out that in the above sequence of definitions only the definitions of entropies, average code length and uniquely decodable coding are important, and the other ones are technical.

Finally, one more definition is necessary to the proof. Denote by $Z^N$ the set of the sequences $(x_{i_1}, \ldots, x_{i_s})$ satisfying the conditions

$$\sum_{j=1}^{s} l_{i_j} \geqq N, \quad \sum_{j=1}^{s-1} l_{i_j} < N.$$

Let $\mathfrak{A}_Z^N$ be the $\sigma$-field of all subsets of $Z^N$ and put

$$p_Z^N(x_{i_1}, \ldots, x_{i_s}) = p_X[x_{i_1}, \ldots, x_{i_s}].$$

It is easy to see, that

$$\sum_{(x_{i_1}, \ldots, x_{i_s}) \in Z^N} p_Z^N(x_{i_1}, \ldots, x_{i_s}) = 1$$

thus $\mathcal{Z}^N = (Z^N, \mathfrak{A}_Z^N, p_Z^N)$ is a probability space, and

$$H(\mathcal{Z}^N) = - \sum_{(x_{i_1}, \ldots, x_{i_s}) \in Z^N} p_X(x_{i_1}, \ldots, x_{i_s}) \log p_X(x_{i_1}, \ldots, x_{i_s}).$$

### Theorems and Proofs

To prove our main theorem we need a lemma.

LEMMA 2. *If L exists,*

$$\left(\frac{H(\mathcal{X}^k)}{L \cdot k} - \frac{H(\mathcal{L}^N)}{N}\right) \to 0$$

*provided that* $k = \left[\dfrac{N}{L}\right]$ *and* $N \to \infty$.

PROOF. Let $H(\mathcal{L}^N | \mathcal{X}^k)$ be the conditional entropy

$$H(\mathcal{L}^N | \mathcal{X}^k) = - \sum_{\substack{u \in X^k \\ z \in Z^N}} p_X(u, z) \log p_X(z|u),$$

where $p_X(u, z)$ denotes the probability of the subsets of elements in $X^\infty$, for which the first $k$ elements are $x_{i_1}, \ldots, x_{i_k}$ and the first $s$ elements are $x_{j_1}, \ldots, x_{j_s}$, if $u = (x_{i_1}, \ldots, x_{i_k})$, $z = (x_{j_1}, \ldots, x_{j_s})$. Further

$$p_X(z|u) = \frac{p_X(u, z)}{p_X(u)}.$$

Obviously, $p_X(u, z) \neq 0$ $(p_X(z|u) \neq 0)$ if and only of if one of the $u$ and $z$ is a segment of the other.

It is well known that (e. g. [1])

(3)                      $H(\mathcal{X}^k) + H(\mathcal{L}^N | \mathcal{X}^k) = H(\mathcal{L}^N) + H(\mathcal{X}^k | \mathcal{L}^N),$

so it is sufficient to show, that

$$H(\mathcal{L}^N | \mathcal{X}^k) = o(N)$$

and

$$H(\mathcal{X}^k | \mathcal{L}^N) = o(N)$$

if $k = \left[\dfrac{N}{L}\right]$. Namely, in this case

$$\frac{H(\mathcal{X}^k)}{N} - \frac{H(\mathcal{L}^N)}{N} = \frac{H(\mathcal{X}^k | \mathcal{L}^N)}{N} - \frac{H(\mathcal{L}^N | \mathcal{X}^k)}{N} \to 0$$

follows from (3).

However

(4)                      $H(\mathcal{L}^N | \mathcal{X}^k) = \sum_{u \in X^k} p_X(u) H(\mathcal{L}^N | u)$

where $H(\mathcal{L}^N | u)$ denotes the entropy $- \sum_{z \in Z^N} p_X(z|u) \log p_X(z|u)$. Let $L_k(u) = l(u)$ be the number $\sum_{j=1}^{k} l_{i_j}$ if $u = (x_{i_1}, \ldots, x_{i_k})$.

If $l(u) \geqq N$, $p_X(z|u) \neq 0$ can hold only if $z$ is a segment of $u$ (as $u$ cannot be a proper segment of $z$ only if $u = z$), but because of definition of $Z^N$ there can be only one segment of $u$ in $Z^N$. Thus $p_X(z|u) = 1$ for a certain $z$, and

(5)                      $H(\mathcal{L}^N | u) = 0.$

In the case $l(u) < N$, $p_X(z|u) \neq 0$ if and only if $u$ is a segment of $z$. On the other hand $s - k \leqq N - l(u)$ since in the case $s - k = N - l(u)$ for the length $l(z)$ the inequality $l(z) \geqq l(u) + s - k = N$ holds. Obviously, we have only $n^{N - l(u)}$ such sequences $z$, that is,

(6)
$$H(\mathcal{L}^N | u) \leqq \log n^{N - l(u)} = (N - l(u)) \log n$$

(as the maximum of the entropy of a distribution on a set of $M$ elements is $\log M$).

Applying (4), (5) and (6) we have

$$H(\mathcal{L}^N | \mathcal{X}^k) = \sum_{\substack{u \in X^k \\ l(u) < N}} p_X(u) H(\mathcal{L}^N | u) \leqq \sum_{\substack{u \in X^k \\ l(u) < N}} p_X(u)(N - l(u)) \log n =$$

$$= \sum_{\substack{u \in X^k \\ N(1-\varepsilon) \leqq l(u) < N}} p_X(u)(N - l(u)) \log n + \sum_{\substack{u \in X^k \\ l(u) < (1-\varepsilon)N}} p_X(u)(N - l(u)) \log n \leqq$$

$$\leqq N \varepsilon \log n + p_X(u \in X^k, \ l(u) < (1-\varepsilon)N) N \log n.$$

Thus we have the inequality

(7)
$$\frac{H(\mathcal{L}^N | \mathcal{X}^k)}{N} \leqq \varepsilon \log n + p_X(u \in X^k, l(u) < (1 - \varepsilon)N) \log n.$$

Since $\dfrac{l(u)}{k} = \dfrac{L_k(u)}{k}$ converges stochastically to $L$, $p_X\left(u \in X^k, \left|\dfrac{l(u)}{k} - L\right| > \varepsilon\right)$ converges to zero if $k \to \infty$. It follows that on the right side of (7)

$$p_X(u \in X^k, l(u) < (1 - \varepsilon)N) = p_X\left(u \in X^k, \frac{l(u)}{k} - L < (1 - \varepsilon)\frac{N}{k} - L\right) \leqq$$

$$\leqq p_X\left(u \in X^k, \frac{l(u)}{k} - L < -L \cdot \frac{\varepsilon}{2}\right)$$

tends to zero for $N \to \infty$, because of $k = \left[\dfrac{N}{L}\right]$. Thus, if $N$ is sufficiently large,

$$p_X(u \in X^k, l(u) < (1 - \varepsilon)N) < \varepsilon$$

that is

$$\frac{H(\mathcal{L}^N | \mathcal{X}^k)}{N} \leqq 2\varepsilon \log n$$

consequently

$$\lim_{N \to \infty} \frac{H(\mathcal{L}^N | \mathcal{X}^k)}{N} = 0 \qquad \left(k = \left[\frac{N}{L}\right]\right).$$

We prove in similar way that

$$\lim_{N \to \infty} \frac{H(\mathcal{X}^k | \mathcal{L}^N)}{N} = 0 \qquad \left(k = \left[\frac{N}{L}\right]\right).$$

Obviously

$$H(\mathcal{X}^k | \mathcal{L}^N) = \sum_{z \in L^N} p_X(z) H(\mathcal{X}^k | z),$$

where $H(\mathscr{X}^k|z)=0$ in the case $z=(x_{i_1}, \ldots, x_{i_s})$, $s \geq k$. Further, if $s < k$, we have only $n^{k-s}$ different $u$'s with $p_X(u|z) \neq 0$, that is, $H(\mathscr{X}^k|z) \leq (k-s) \log n$. Finally, as above

$$H(\mathscr{X}^k|\mathscr{L}^N) \leq \sum_{\substack{z \in Z^N \\ k(1-\varepsilon) \leq s < k}} p_X(z)(k-s) \log n + \sum_{\substack{z \in Z^N \\ s < (1-\varepsilon)k}} p_X(z)(k-s) \log n \leq$$

$$\leq \varepsilon k \log n + k \log n \, p_X(z \in Z^N, s < (1-\varepsilon)k).$$

Here on the right side

$$p_X(z \in Z^N, s < (1-\varepsilon)k) = p_X(u \in X^{[(1-\varepsilon)k]}, l(u) \geq N) =$$

$$= p_X\left(u \in X^M, \frac{l(u)}{M} - L \geq \frac{N}{(1-\varepsilon)k} - L\right) \leq p_X\left(u \in X^M, \frac{l(u)}{M} - L \geq \varepsilon L\right)$$

which converges to zero if $M = [(1-\varepsilon)k] \to \infty$. Thus $H(\mathscr{X}^k|\mathscr{L}^N) \leq o(N)$, indeed, which proves the Lemma.

THEOREM 1. *If the entropy $H(\mathscr{X})$ and the average code length $L$ exist, and the coding is uniquely decodable, then $H(\mathscr{Y})$ exists, and*

$$H(\mathscr{Y}) = \frac{H(\mathscr{X})}{L}.$$

PROOF. If $H(\mathscr{X})$ exists, in Lemma 2 $\dfrac{H(\mathscr{X}^k)}{L \cdot k}$ tends to $\dfrac{H(\mathscr{X})}{L}$, and so $\dfrac{H(\mathscr{L}^N)}{N}$ does, too. Thus, it is sufficient to show that

$$\lim_{N \to \infty} \frac{H(\mathscr{L}^N)}{N} = \lim_{N \to \infty} \frac{H(\mathscr{Y}^N)}{N}.$$

We can write

(8) $$H(\mathscr{Y}^N) + H(\mathscr{L}^N|\mathscr{Y}^N) = H(\mathscr{L}^N) + H(\mathscr{Y}^N|\mathscr{L}^N)$$

where

$$H(\mathscr{L}^N|\mathscr{Y}^N) = - \sum_{\substack{z \in Z^N \\ v \in Y^N}} p_X(v, z) \log p_X(z|v),$$

$$H(\mathscr{Y}^N|\mathscr{L}^N) = - \sum_{\substack{z \in Z^N \\ v \in Y^N}} p_X(v, z) \log p_X(v|z),$$

$$p_X(z|v) = \frac{p_X(v, z)}{p_Y(v)}, \quad p_X(v|z) = \frac{p_X(v, z)}{p_X(z)},$$

and $p_X(v, z)$ is the probability of the set of sequences in $X^\infty$, for which the first $s$ elements are $x_{i_1}, \ldots, x_{i_s}$ and the first $N$ elements of its code are $y_{j_1}, \ldots, y_{j_N}$, if $z = (x_{i_1}, \ldots, x_{i_s})$, $v = (y_{j_1}, \ldots, y_{j_N})$.

Obviously $p_X(v, z) \neq 0$ only if $v$ is a segment of $c(z)$. Thus for given $z$ there is only one $v$ satisfying $p_X(v, z) \neq 0$, that is, $p_X(v, z) = 1$. Applying this result we obtain

(9) $$H(\mathscr{Y}^N|\mathscr{L}^N) = \sum_{z \in Z^N} p_X(z) H(\mathscr{Y}^N|z) = 0$$

because of $H(\mathscr{Y}^N|z) = 0$. On the other hand

(10) $$H(\mathscr{L}^N|\mathscr{Y}^N) = \sum_{v \in Y^N} p_Y(v) H(\mathscr{L}^N|v).$$

Let $l$ be the maximum of the numbers $l_1, ..., l_n$. Because of definition of $Z^N$, $N \leqq$ $\leqq l(z) < N + l$ holds. Thus for a fixed $v$, the sequences $z$ satisfying $p_X(z|v) \neq 0$ are such that the first $N$ elements of $c(z)$ are equal to $v$, and the other elements are arbitrary. The number of such $c(z)$'s is at most $m^{l-1}$. Since the coding is uniquely decodable i. e. to a given $c(z)$ there exists only one $z$, the number of different $z$ is also at most $m^{l-1}$. From (10) we obtain

$$H(\mathscr{Z}^N | \mathscr{Y}^N) \leqq \sum_{v \in Y^N} p_Y(v) \log m^{l-1} = \log m^{l-1},$$

which tends to zero divided by $N$, if $N \to \infty$. The proof is finished by (8).

If the coding is not necessarily uniquely decodable, we can not prove the existence of $H(\mathscr{Y})$. In this case let us put

$$\overline{H}(\mathscr{Y}) = \varlimsup_{k \to \infty} \frac{H(\mathscr{Y}^k)}{k}.$$

In this case from the above proof we get only $\overline{H}(\mathscr{Y}^N) \leqq H(\mathscr{Z}^N)$ that is, the following theorem holds.

THEOREM 2. *If the entropy* $H(\mathscr{X})$ *and the average code length* $L$ *exist, then*

(11) $$\overline{H}(\mathscr{Y}) \leqq \frac{H(\mathscr{X})}{L}.$$

### Further Questions

**1.** A natural question is th e following: under which assumption does the limit $\lim_{k \to \infty} \dfrac{H(\mathscr{Y}^k)}{k}$ exists in the not uniquely decodable case? Probably it is not difficult to answer this question if $\mathscr{X}$ is an information source, which produces independent signals.

**2.** It is easy to see, that for independent $\mathscr{X}$, and not uniquely decodable coding the strict inequality

$$\overline{H}(\mathscr{Y}) < \frac{H(\mathscr{X})}{L}$$

holds. In other words, in the independent case equality holds in (11) if and only if the coding is uniquely decodable. What is the necessary and sufficient condition, in general, of the equality in (11)?

We are greatly indebted to A. RÉNYI and I. CSISZÁR for several helpful comments and ideas.

### REFERENCE

[1] FEINSTEIN, A.: *Foundation of information theory*, McGraw-Hill, New York, 1958.

MATHEMATICAL INSTITUTE OF THE HUNGARIAN ACADEMY OF SCIENCES, BUDAPEST