

## EPSILON ENTROPY AND DATA COMPRESSION<sup>1</sup>

EDWARD C. POSNER AND EUGENE R. RODEMICH

*Jet Propulsion Laboratory, California Institute of Technology*

**0. Summary.** This article studies efficient data transmission, or “data compression”, from the standpoint of the theory of epsilon entropy. The notion of the entropy of a “data source” is defined. This quantity gives a precise measure of the amount of channel capacity necessary to describe a data source to within a given fidelity, epsilon, with probability one, when each separate “experiment” must be transmitted without storage from experiment to experiment. We also define the absolute epsilon entropy of a source, which is the amount of capacity needed when storage of experiments is allowed before transmission. The absolute epsilon entropy is shown to be equal to Shannon’s rate distortion function evaluated for zero distortion, when suitable identifications are made. The main result is that the absolute epsilon entropy and the epsilon entropy have ratio close to one if either is large. Thus, very little can be saved by storing the results of independent experiments before transmission.

**1. Introduction.** An area of information theory that is becoming more and more important is now known as “Data Compression”. This is the theory of efficient handling of data for transmission or storage. As the space program goes to more and more distant missions, such as those to the outer planets, the importance of using the spacecraft-to-Earth communication link efficiently becomes more and more important. And here on Earth, as our technology gets more and more computer based, more and more data is being transmitted between different locations, so that the economic importance of data compression is increasing rapidly.

A need has been felt by workers in data compression to have a mathematical theory of what they are trying to do. It is the purpose of this paper to provide such a theory. Thus, we start off by abstracting the notion of the source of the data to be transmitted as a *probabilistic metric space*. Such a space is a triple  $(X, d, \mu)$  with the following two properties:

(1)  $(X, d)$  is a complete separable metric space, of points of the set  $X$  under the metric  $d$ .

(2)  $(X, B, \mu)$  is a probability space, where  $B$  is the Borel field generated by the open sets of  $X$ ; i.e.,  $B$  is the field of *Borel sets* of  $X$ . (However, we will speak of subsets of  $X$  as measurable if they belong to the completion of  $B$  under  $\mu$ .)

---

Received February 13, 1970; revised May 14, 1971.

<sup>1</sup> This paper presents the results of one phase of research carried out at the Jet Propulsion Laboratory, California Institute of Technology, under Contract No. NAS 7-100, sponsored by the National Aeronautics and Space Administration.

The points of  $X$  represent the possible experimental outcomes, i.e., the data; the metric  $d$  is a “fidelity criterion”, such that  $d(x, y)$  is the loss of fidelity when one outcome  $x$  occurs, but another outcome  $y$  is thought to have occurred.

We were given an  $\varepsilon \geq 0$ , as the allowed loss of fidelity. But we are also given a  $\delta \geq 0$ , such that the loss of fidelity can exceed  $\varepsilon$ , but only with probability  $\delta$  or less. Any physical data transmission system has a probability  $\delta > 0$  associated with it, such that, with probability  $\delta$ , the system does not work as specified. For example, an “analog-to-digital converter” has as its  $\delta$  the “off-scale” probability.

Given a probabilistic metric space  $(X, d, \mu)$ , which we abbreviate as  $X$ , we wish to transmit outcomes within  $\varepsilon$  or less with probability  $1 - \delta$  or more. What can this mean? One reasonable interpretation is the following. If, given a received message, we know the actual outcome within  $\varepsilon$ , then we know that the actual outcome falls within an  $\varepsilon$ -set, i.e., within a (measurable) set of diameter at most  $\varepsilon$ . However, with probability  $\delta$  or less, there is no transmission and no such set of diameter at most  $\varepsilon$ . We shall sometimes also use a definition involving spheres of radius  $\varepsilon/2$  or less.

Now it might appear that “mixed strategies” should be considered, in which the set of diameter  $\varepsilon$  depends on other “things”, or random variables, not depending on past or future samples from  $X$ , in addition to depending on the actual outcome. However, it turns out, and is not hard to show, that there is a “pure strategy” at least as good as any mixed strategy, insofar as reducing the load on the communication channel. A “pure strategy” assigns a given outcome to a fixed  $\varepsilon$ -set. Thus, in a pure strategy, these  $\varepsilon$ -sets are to be disjoint. That is, we have a collection of disjoint  $\varepsilon$ -sets, such that with probability  $1 - \delta$  or more, a point of the space  $X$  lies in one of the  $\varepsilon$ -sets. This leads to the definition of  $\varepsilon; \delta$  partition.

**DEFINITION.** Given  $\varepsilon \geq 0$ ,  $\delta \geq 0$ , an  $\varepsilon; \delta$  partition of the probabilistic metric space  $X$  is a finite or denumerably infinite partition of part of  $X$  by (disjoint)  $\varepsilon$ -sets, such that the union of the sets in the partition has probability  $1 - \delta$  or more.

A data compression system based on a given  $\varepsilon; \delta$  partition  $U$  of  $X$  works as follows: observe an outcome, see into which set of the partition  $U$  the outcome falls, and merely transmit information which uniquely determines that set. In fact, any (“pure strategic”) data transmission system which yields fidelities of  $\varepsilon$  or better with probability  $1 - \delta$  or more is of this type.

(Other authors [10, 4] usually consider restricting some average distance to be  $\varepsilon$  or better, instead of the actual distance between actual outcome and transmitted outcome. The actual distance is felt to be the appropriate criterion in judging most data systems.)

So we now have an  $\varepsilon; \delta$  partition  $U$  of  $X$ ; how much of a load does this imply for our communication system? A measure of the load is the entropy of the partition,  $H(U)$ . This entropy is the number of bits necessary to transmit information as to which set of  $U$  the outcome fell, given that the outcome was in some set of  $U$ . Thus, using Shannon’s formula [11], we define  $H(U)$  as follows.

$$\text{Let } U = \{U_i\}, \text{ with } \mu(U_i) = p_i, \sum p_i = p \geq 1 - \delta.$$

Let  $q_i = \frac{p_i}{p}$ , so that  $\{q_i\}$  is a probability distribution.

Then the entropy  $H(U)$  is defined as the entropy of  $\{q_j\}$ :

$$(1) \quad H(U) = \sum q_i \log \frac{1}{q_i}.$$

Thus,  $H(U)$  is nonnegative, and can be equal to  $+\infty$ .

The interpretation of  $H(U)$  is the number of bits per sample (we use base  $e$  logarithms, even though we speak of bits instead of “nats”) necessary to describe into which set of  $U$  the outcome falls, when minimum expected-length binary encoding is used for this purpose, but storage between experiments is not allowed.

However, the partition  $U$  may not have been especially well chosen with a view to minimizing  $H(U)$ . Thus, let  $U_{\varepsilon;\delta}$  denote the class of  $\varepsilon;\delta$  partitions of  $X$  (when  $\varepsilon > 0$ , this class is non-empty). Then we define  $H_{\varepsilon;\delta}(X)$ , the *epsilon;delta entropy* of  $X$ , as

$$(2) \quad H_{\varepsilon;\delta}(X) = \inf_{U \in U_{\varepsilon;\delta}} H(U).$$

(When  $\delta = 0$ , we write  $H_\varepsilon(X)$ , the  $\varepsilon$ -entropy of  $X$ . If  $\varepsilon = 0 < \delta$  and  $U_{0;\delta}$  is empty,  $H_{0;\delta}(X)$  is infinite.)

Thus, the  $\varepsilon;\delta$  entropy of  $X$ , which may equal  $+\infty$  even when  $\varepsilon > 0$ , provided that  $\delta = 0$ , measures the minimum number of bits necessary to describe at least  $1 - \delta$  of  $X$  within accuracy  $\varepsilon$ . When  $\delta = 0$ , we can say that the  $\varepsilon$ -entropy of  $X$  is the number of bits necessary to describe elements of  $X$  within  $\varepsilon$  (without storage, but we will say more about this later), with probability 1. We are mainly concerned in this paper with  $\varepsilon$ -entropy;  $\varepsilon;\delta$  entropy enters chiefly as a device in proving theorems.

We have chosen to use sets of diameter  $\varepsilon$  in our definition of epsilon entropy to agree with the definition used in the theory of the epsilon entropy of compact metric spaces [13]. However, other workers in information theory use in effect sets of radius  $\varepsilon/2$  in the definition. We will treat both cases; when it is necessary to distinguish the two definitions, we will use the terms “diametric entropy” and “radial entropy”, respectively. The original paper on the subject of data compression and epsilon entropy [7] used diametric entropy, although Kolmogorov [2] uses radial entropy.

The purpose of the remainder of this paper is to make the statements of this section more precise. Various techniques will be given for obtaining upper and lower bounds for  $\varepsilon$ -entropy. The  $\varepsilon$ -entropy defined here will be related to the epsilon entropy of compact metric spaces and to Shannon’s rate distortion theory. Finally, a very precise “channel coding theorem” and its converse will be stated and proved for  $\varepsilon$ -entropy, which makes the relevance of  $\varepsilon$ -entropy to data compression clear.

**2. Connection with  $\varepsilon$ -entropy for compact metric spaces.** A definition for the  $\varepsilon$ -entropy of a compact metric space finds wide use in various kinds of approximation theory; see [13] for a good expository treatment and many references.

DEFINITION. An  $\varepsilon$ -partition ( $\varepsilon > 0$ ) of a compact metric space  $X$  with metric  $d$  is a partition of all of  $X$  by disjoint Borel sets of diameters at most  $\varepsilon$  (caution: other workers use  $2\varepsilon$  instead of  $\varepsilon$ ).

DEFINITION. The  $\varepsilon$ -entropy  $K_\varepsilon(X)$  of a compact metric space  $X$  is the logarithm of the minimum number of sets in any  $\varepsilon$ -partition of  $X$  (since  $X$  is compact,  $K_\varepsilon(X)$  is finite).

An information-theoretic interpretation of  $K_\varepsilon(X)$  can be given as follows. Suppose data is generated as elements of a compact metric space  $X$ . We wish to transmit outcomes of  $X$ , not allowing storage from one experiment to the next, outputting binary words of fixed length, using as short a length as possible. What is this shortest length? The  $\varepsilon$ -entropy  $K_\varepsilon(X)$  is the answer, except for round-off in the logarithm. If, however, any probabilistic information is known about  $X$ , that is, if the data is selected according to a Borel distribution  $\mu$  on  $X$ , the average word length may be able to be shortened by taking advantage of  $\mu$  in assigning word lengths. That is,

$$H_\varepsilon(X) \leq K_\varepsilon(X)$$

since the entropy of an  $\varepsilon$ -partition of  $X$  by  $n$  sets is bounded from above by  $\log n$ . In general, one does not consider  $K_\varepsilon(X)$  for non-compact metric spaces, since this number can be infinite. (In fact, if  $K_\varepsilon(X) < \infty$ , all  $\varepsilon > 0$ , then  $X$  is totally bounded.) In the general (non-compact) case, then, words of fixed length cannot be used if one wishes to communicate with probability 1.

Examples can be given (a pentagon with center not connected to one vertex) of a complete separable metric space in which  $K_\varepsilon(X)$  is not the supremum of  $H_\varepsilon(X)$  taken over all Borel probability measures  $\mu$  on  $X$ , although equality is often obtained.

We remark that the construction of  $\varepsilon$ -entropy for certain product spaces, done in the next section for probabilistic metric spaces, can also be done for compact metric spaces. Thus, we can speak of the *absolute epsilon entropy* of a compact metric space, which is essentially the infimum of the number of bits necessary to describe all of  $X$  to within  $\varepsilon$ , when words of fixed length are used, but when a large number  $n$  of experiments from  $X$  are performed before transmission. We shall not go into this matter further in this paper, but will elsewhere [3].

**3. Epsilon entropy of product spaces.** In this section, we begin the main subject of this paper, which is the relation of  $\varepsilon$ -entropy to channel coding theorems and their converse. We first define the kind of finite products of probabilistic metric spaces we consider. This definition is motivated by the desire to know each outcome of a sequence of experiments to within  $\varepsilon$ .

DEFINITION. Let  $X = (X, d, \mu)$  and  $Y = (Y, e, \nu)$  be probabilistic metric spaces. Then the *product (probabilistic metric) space*  $X \times Y = (X \times Y, \max(d, e), \mu \times \nu)$  is the probabilistic metric space whose point set is the Cartesian product  $X \times Y$  of

$X$  and  $Y$ , whose metric  $f$  is defined by  $f[(x, y), (x', y')] = \max [d(x, x'), e'(y, y')]$ , and whose measure is the product measure of  $\mu$  and  $\nu$ .

DEFINITION. Let  $X$  be a probabilistic metric space. Define  $X^{(1)} = X, X^{(n+1)} = X^{(n)} \times X, n \geq 1$ . Thus  $X^{(n)}$  is the  $n$ -fold product of  $X$  with itself.

Since each  $X^{(n)}$  is a probabilistic metric space, we can study the sequence  $H_\epsilon(X^{(n)})$ . We have the following lemma.

LEMMA 1. *The sequence  $\{H_\epsilon(X^{(n)})\}$  is subadditive in  $n$ . That is,*

$$(3) \quad H_\epsilon(X^{(m+n)}) \leq H_\epsilon(X^{(m)}) + H_\epsilon(X^{(n)}).$$

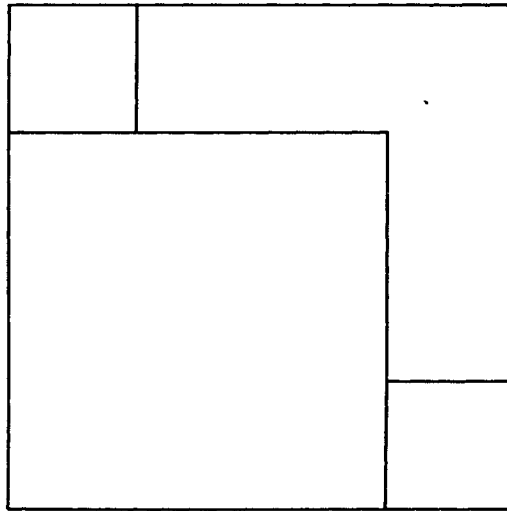


FIG. 1.  $\frac{3}{4}$ -partition of  $X \times X$ .

PROOF. By [8a], Lemma 8, if we have an  $\epsilon$ -partition  $U$  of  $X^{(m)}$  and  $U'$  of  $X^{(n)}$ , then the partition  $U \times U'$  of  $X^{(m+n)}$ , defined as the partition consisting of all products of sets in  $U$  with sets in  $U'$ , has the property that

$$(4) \quad H(U \times U') = H(U) + H(U').$$

By the definition of the metric on  $X^{(m+n)}$ ,  $U \times U'$  is an  $\epsilon$ -partition of  $X^{(m+n)}$ . Since we can demand  $H(U) = H_\epsilon(X^{(m)})$ ,  $H(U') = H_\epsilon(X^{(n)})$  ([7], Theorem 2), we have an  $\epsilon$  partition  $V$  of  $X^{(m+n)}$  with

$$(5) \quad H(V) = H_\epsilon(X^{(m)}) + H_\epsilon(X^{(n)}).$$

The lemma follows.

We observe that  $H_\epsilon(X^{(m+n)})$  can be less than  $H_\epsilon(X^{(m)}) + H_\epsilon(X^{(n)})$ . For example, let  $m = n = 1$ , and let  $X$  be the unit interval with Lebesgue measure and linear metric; let  $\epsilon = \frac{3}{4}$ . Then  $2H_\epsilon(X) = 2((\frac{3}{4}) \log(\frac{4}{3}) + (\frac{1}{4}) \log 4)$ , whereas  $H_\epsilon(X \times X) = (\frac{9}{16}) \log(\frac{16}{9}) + (\frac{5}{16}) \log(\frac{16}{5}) + 2(\frac{1}{16}) \log 16 < 2H_\epsilon(X)$ ; in fact, it can be proved that the  $\frac{3}{4}$  partition of  $X \times X$  in Figure 1 achieves  $H_\epsilon(X \times X)$ ; diametric entropy or radial entropy can be used here.

Another inequality which will be useful is

$$(6) \quad H_\varepsilon(X^{(m)}) \leq H_\varepsilon(X^{(m+n)}).$$

This is a special case of the following lemma:

LEMMA 2. *Let  $X$  and  $Y$  be separable metric spaces with metrics  $d_1$  and  $d_2$ , and let  $Z = X \times Y$  have the metric  $d_3$ , such that if  $z_j = (x_j, y_j)$ , then*

$$d_3(z_1, z_2) \geq \max [d_1(x_1, x_2), d_2(y_1, y_2)].$$

*Let  $X, Y, Z$  be probabilistic metric spaces under some Borel probability distribution on  $Z$  and under the induced marginal distributions on  $X$  and  $Y$ . Thus  $H_\varepsilon(X) \leq H_\varepsilon(Z)$ .*

PROOF. Let  $\mu_3$  be measure on  $Z$  and  $\mu_1$  measure on  $X$ . Take an  $\varepsilon$ -partition  $U = \{U_j\}$  of  $Z$  for which

$$H_\varepsilon(Z) = \sum_j \mu_3(U_j) \log \frac{1}{\mu_3(U_j)}.$$

Let the  $\{U_j\}$  be arranged in order of decreasing probability. If  $V_j$  is the closure of the projection of  $U_j$  on  $X$ ,  $V_j$  is a Borel set of diameter  $\leq \varepsilon$ , by hypothesis, and  $\mu_1(V_j) \geq \mu_3(U_j)$ , since  $U_j \subset V_j \times Y$ . Similarly,

$$\mu_1(V_1 UV_2 U \cdots UV_j) \geq \mu_3(U_1 U U_2 U \cdots U U_j).$$

Hence, if we define  $W_1 = V_1$ , and

$$W_j = V_j - (V_1 UV_2 U \cdots UV_{j-1}), \quad j \geq 2,$$

$W = \{W_j\}$  is an  $\varepsilon$ -partition of  $X$  with

$$\sum_{k=1}^j \mu_1(W_k) \geq \sum_{k=1}^j \mu_3(U_k).$$

It follows [7, Lemma 2] that  $H(W) \leq H(U)$ . Therefore  $H_\varepsilon(X) \leq H_\varepsilon(Z)$ . The lemma follows.

We define the *absolute  $\varepsilon$ -entropy*  $I_\varepsilon(X)$  as

$$(7) \quad I_\varepsilon(X) = \lim_{n \rightarrow \infty} \frac{1}{n} H_\varepsilon(X^{(n)}).$$

After Lemma 1, a well-known result on subadditive sequences [3] implies that this limit exists if the  $H_\varepsilon(X^{(n)})$  are all finite, and

$$(8) \quad I_\varepsilon(X) \leq H_\varepsilon(X), \quad n = 1, 2, \dots$$

Furthermore, Lemma 1 and equation (6) ensure that if any  $H_\varepsilon(X^{(n)})$  is finite, they are all finite. If these entropies are all infinite, we interpret (7) to mean  $I_\varepsilon(X) = \infty$ .

$I_\varepsilon(X)$  is clearly nonnegative, and the following result holds.

LEMMA 3.  *$I_\varepsilon(X)$  is zero if and only if  $H_\varepsilon(X)$  is zero, and infinite if and only if  $H_\varepsilon(X)$  is itself infinite.*

PROOF. After (8) and the remarks above, the only thing to show is that if  $H_\varepsilon(X) > 0$ , then  $I_\varepsilon(X) > 0$ .

When  $H_\varepsilon(X) > 0$ , there is a number  $p < 1$  such that any  $\varepsilon$ -set in  $X$  has measure  $\leq p$  [7, Theorem 3]. Any  $\varepsilon$ -set in  $X^{(n)}$  lies in the product of  $n$   $\varepsilon$ -sets of the component spaces (its  $n$  projections), hence has measure at most  $p^n$ . It follows that if  $U = \{U_j\}$  is any  $\varepsilon$ -partition of  $X^{(n)}$ ,

$$H(U) = \sum \mu^{(n)}(U_j) \log \frac{1}{\mu^{(n)}(U_j)} \geq \sum \mu^{(n)}(U_j) \log \frac{1}{p^n} = n \log \frac{1}{p}.$$

Thus  $1/n H_\varepsilon(X^{(n)}) \geq \log 1/p$  for all  $n$ , which implies  $I_\varepsilon(X) \geq \log 1/p$ . Lemma 3 is proved.

The following definition and lemma are needed before we can get down to channel coding theorems.

DEFINITION.  $I_{\varepsilon;[\delta]}(X) = \lim_{n \rightarrow \infty} \inf n^{-1} H_{\varepsilon;\delta}(X^{(n)})$ .

LEMMA 4. If  $H_\varepsilon(X) < \infty$ ,  $I_\varepsilon(X) = \lim_{\delta \rightarrow 0^+} I_{\varepsilon;[\delta]}(X)$ .

PROOF. The existence of an  $\varepsilon$ -partition  $U = \{U_j\}$  of  $X$  with finite entropy will be used to estimate the difference between  $H_{\varepsilon;\delta}(X^{(n)})$  and  $H_\varepsilon(X^{(n)})$ .

Let  $V$  be an  $\varepsilon; \delta$  partition of  $X^{(n)}$  with  $H(V) = H_{\varepsilon;\delta}(X^{(n)})$ , and let  $B$  be the part of  $X^{(n)}$  not covered by  $V$ . Let  $W = \{W_j\}$  be the restriction to  $B$  of the product partition  $U^{(n)} = \{U_j^{(n)}\}$  of  $X^{(n)}$ . If  $B$  has measure  $\mu_B$ , the partition  $WUV$  of  $X^n$  has entropy

$$H(WuV) = \mu_B H(W) + (1 - \mu_B) H(V) + \mu_B \log \frac{1}{\mu_B} + (1 - \mu_B) \log \frac{1}{1 - \mu_B}.$$

We have

$$H(W) = \sum_j \frac{\mu^{(n)}(W_j)}{\mu_B} \log \frac{\mu_B}{\mu^{(n)}(W_j)}.$$

Hence

$$(9) \quad H_\varepsilon(X^{(n)}) \leq (1 - \mu_B) H_{\varepsilon;\delta}(X^{(n)}) + (1 - \mu_B) \log \frac{1}{1 - \mu_B} + \sum_j \mu^{(n)}(W_j) \log \frac{1}{\mu^{(n)}(W_j)}.$$

We group the  $W_j$  into two classes:

- (I)  $\mu^{(n)}(W_j) < \delta^{\frac{1}{2}} \mu^{(n)}(U_j^{(n)})$ ,
- (II)  $\mu^{(n)}(W_j) \geq \delta^{\frac{1}{2}} \mu^{(n)}(U_j^{(n)})$ .

If  $n$  is sufficiently large, every  $U_j^{(n)}$  has measure less than  $1/e$ . Then

$$(10) \quad \sum_{(I)} \mu^{(n)}(W_j) \log \frac{1}{\mu^{(n)}(W_j)} \leq \sum_{(I)} \delta^{\frac{1}{2}} \mu^{(n)}(U_j^{(n)}) \log \frac{1}{\delta^{\frac{1}{2}} \mu^{(n)}(U_j^{(n)})} \\ \leq \delta^{\frac{1}{2}} \log \frac{1}{\delta^{\frac{1}{2}}} + n \delta^{\frac{1}{2}} H(U),$$

since  $H(U^{(n)}) = nH(U)$ . For the second class we have

$$\sum_{(II)} \mu^{(n)}(W_j) \log \frac{1}{\mu^{(n)}(W_j)} \leq \sum_{(II)} \mu^{(n)}(U_j^{(n)}) \log \frac{1}{\mu^{(n)}(U_j^{(n)})}.$$

Each  $U_j^{(n)}$  is the product of a certain  $n$  sets of  $U$ :

$$U_j^{(n)} = U_{k_1} \times U_{k_2} \times \dots \times U_{k_n},$$

so we have

$$\sum_{(II)} \mu^{(n)}(W_j) \log \frac{1}{\mu^{(n)}(W_j)} \leq \sum_{i=1}^n \sum_{(II)} \mu(U_{k_1}) \dots \mu(U_{k_n}) \log \frac{1}{\mu(U_{k_i})}.$$

Since

$$\sum_{(II)} \mu^{(n)}(W_j) \leq \mu_B \leq \delta^{\frac{1}{2}},$$

we must have

$$\sum_{(II)} \mu^{(n)}(U_j^{(n)}) \leq \delta^{\frac{1}{2}}.$$

Also,

$$\sum_{k_i \text{ fixed}} \mu^{(n)}(U_j^{(n)}) = \mu(U_{k_i}).$$

Hence

$$\sum_{(II)} \mu^{(n)}(W_j) \log \frac{1}{\mu^{(n)}(W_j)} \leq \sum_{i=1}^n \sum_{k_i} \min [\delta^{\frac{1}{2}}, \mu(U_{k_i})] \log \frac{1}{\mu(U_{k_i})},$$

or

$$\sum_{(II)} \mu^{(n)}(W_j) \log \frac{1}{\mu^{(n)}(W_j)} \leq n \sum_k \min [\delta^{\frac{1}{2}}, \mu(U_k)] \log \frac{1}{\mu(U_k)}.$$

Combining the last inequality with (9) and (10), and replacing the first terms on the right in (10) by upper bounds,

$$\begin{aligned} \frac{1}{n} H_\epsilon(X^{(n)}) &\leq \frac{1}{n} H_{\epsilon;\delta}(X^{(n)}) + \frac{1}{ne} + \sum_k \min [\delta^{\frac{1}{2}}, \mu(U_k)] \log \frac{1}{\mu(U_k)} \\ &\quad + \delta^{\frac{1}{2}} H(U) + \frac{1}{n} \delta^{\frac{1}{2}} \log \frac{1}{\delta^{\frac{1}{2}}}. \end{aligned}$$

Take the lower limit as  $n \rightarrow \infty$ :

$$I_\epsilon(X) \leq I_{\epsilon;[\delta]}(X) + \sum_k \min [\delta^{\frac{1}{2}}, \mu(U_k)] \log \frac{1}{\mu(U_k)} + \delta^{\frac{1}{2}} H(U).$$

The series on the right approaches zero as  $\delta \rightarrow 0$ , for it is bounded, term by term, by the series for  $H(U)$ , and each term approaches zero. Hence

$$\liminf_{\delta \rightarrow 0^+} I_{\epsilon;[\delta]}(X) \geq I_\epsilon(X).$$

On the other hand, it follows directly from the definitions and the inequality  $H_{\epsilon;\delta}(X^{(n)}) \leq H_\epsilon(X^{(n)})$  that  $I_{\epsilon;[\delta]}(X) \leq I_\epsilon(X)$ . Hence,  $I_{\epsilon;[\delta]}(X) \rightarrow I_\epsilon(X)$  as  $\delta \rightarrow 0$ . Lemma 4 is proved.

Because of needs in the last section, we have the following slight strengthening of Lemma 4.

LEMMA 5. Let  $H_\epsilon(X) < \infty$ . Then

$$(11) \quad I_\epsilon(X) = \inf_{[\delta_n] \rightarrow 0} \liminf_{n \rightarrow \infty} \frac{1}{n} H_{\epsilon;\delta_n}(X^{(n)}).$$



PROOF. Let  $\{\delta_n\} \rightarrow 0$ , and  $\delta > 0$  be given. Then for  $n$  sufficiently large,  $\delta_n \leq \delta$ . Thus,

$$H_{\varepsilon;\delta}(X^{(n)}) \leq H_{\varepsilon;\delta_n}(X^{(n)}),$$

if  $n$  is sufficiently large. Hence, for the given  $\delta$ ,

$$\liminf_{n \rightarrow \infty} \frac{1}{n} H_{\varepsilon;\delta}(X^{(n)}) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} H_{\varepsilon;\delta_n}(X^{(n)}),$$

so, by Lemma 4, we conclude

$$I_\varepsilon(X) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} H_{\varepsilon;\delta_n}(X^{(n)}).$$

Consequently

$$(12) \quad I_\varepsilon(X) \leq \inf_{[\delta_n] \rightarrow 0} \liminf_{n \rightarrow \infty} \frac{1}{n} H_{\varepsilon;\delta_n}(X^{(n)}).$$

On the other hand,

$$(13) \quad H_{\varepsilon;\delta_n}(X^{(n)}) \leq H_\varepsilon(X^{(n)}),$$

so that

$$(14) \quad \liminf_{n \rightarrow \infty} \frac{1}{n} H_{\varepsilon;\delta_n}(X^{(n)}) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} H_\varepsilon(X^{(n)}).$$

(7) then implies

$$\liminf_{n \rightarrow \infty} \frac{1}{n} H_{\varepsilon;\delta_n}(X^{(n)}) \leq I_\varepsilon(X)$$

for any sequence  $\{\delta_n\}$  of nonnegative numbers. In particular,

$$(15) \quad \inf_{[\delta_n] \rightarrow 0} \liminf_{n \rightarrow \infty} \frac{1}{n} H_{\varepsilon;\delta_n}(X^{(n)}) \leq I_\varepsilon(X).$$

Equations (15) and (12) taken together complete the proof of Lemma 5.

**4. Noisy channel coding theorems.** We can now relate our concept of epsilon entropy to channel coding. First note that (11) states that if  $H_\varepsilon(X)$  is finite, then  $I_\varepsilon(X)$  is the number of bits per sample necessary to describe  $X$  to within  $\varepsilon$  over a noiseless channel with probability approaching 1 when arbitrarily long storage is allowed.

The careful reader will at this juncture observe that storage is necessary, however, to achieve even  $H_\varepsilon(X)$ . For the "single-use" number of bits necessary is actually not  $H_\varepsilon(X)$ , but can be arbitrarily close to  $\log 2$  less [1, page 71]. However, since we are interested chiefly in the case  $H_\varepsilon(X)$  large, we shall ignore this difference, and refer to  $H_\varepsilon(X)$  itself as the number of bits necessary when storage between experiments is not allowed.

A noisy channel coding theorem for  $I_e(X)$  can be stated and proved as in Shannon [12]. Before we can think of stating a precise result, we must review some definitions [5, Chapter 2]. Let  $X$  and  $Y$  be probability spaces with measures  $\mu$  and  $\nu$ , and let  $\rho$  be a probability distribution on  $X \times Y$  defined on the product  $\sigma$ -field of  $\mu$  and  $\nu$ . Further, let

$$\rho(A \times Y) = \mu(A),$$

for  $A$   $\mu$ -measurable, and let

$$\rho(X \times B) = \nu(B) \quad (\text{or sometimes } \nu_\rho(B))$$

if  $B$   $\nu$ -measurable. Then the *mutual information* of  $\rho$  (more precisely, of  $\mu$  and  $\nu$  with respect to  $\rho$ ),  $I(\rho)$ , is defined as follows. Let

$$\frac{d\rho}{d\mu d\nu}$$

denote the Radon–Nikodym derivative of  $\rho$  with respect to product measure  $\mu \times \nu$ . Then

$$I(\rho) = \infty$$

if  $\rho$  is not absolutely continuous with respect to  $\mu \times \nu$ , that is, if  $d\rho/d\mu d\nu$  is infinite on a set of positive  $\mu \times \nu$ -measure. If, on the other hand,  $\rho$  is absolutely continuous with respect to  $\mu \times \nu$ , then

$$\begin{aligned} (16) \quad I(\rho) &= E_\rho \left[ \log \left( \frac{d\rho}{d\mu d\nu} \right) \right] = \int_{X \times Y} \left( \log \frac{d\rho}{d\mu d\nu} \right) d\rho \\ &= \int_{X \times Y} \left( \log \frac{d\rho}{d\mu d\nu} \right) \left( \frac{d\rho}{d\mu d\nu} \right) d\mu d\nu, \end{aligned}$$

which may still be infinite. However, since we have

$$(17) \quad -\frac{1}{e} \leq t \log t, \quad \text{for } t \geq 0,$$

$I(\rho)$  is well defined as a real number or  $+\infty$ . And since the function  $t \log t$  is convex in  $t \geq 0$ , we have, from (16),

$$I(\rho) \geq \left( \log \left[ \int_{X \times Y} \left( \frac{d\rho}{d\mu d\nu} \right) d\mu d\nu \right] \right) \left( \int_{X \times Y} \left( \frac{d\rho}{d\mu d\nu} \right) d\mu d\nu \right).$$

Since

$$\int_{X \times Y} d\rho = 1,$$

we conclude

$$(18) \quad I(\rho) \geq 0.$$

It turns out to be easier to first consider the case of radial entropy. The diametric case will be done afterwards with a more complicated definition which will be related to the radial definition. Let  $R_\epsilon(X)$  be defined as that class of probability distributions  $\rho$  on  $X \times X$  such that:

- (i)  $\rho(A \times X) = \mu(A)$  for  $A$  a Borel set in  $X$ ;
- (ii)  $\rho(\{(x, y) | d(x, y) \leq \epsilon/2\}) = 1$ .

Then the noisy channel coding theorem in question turns out to be the statement that

$$(19) \quad I_\epsilon(X) = \inf_{\rho \in R_\epsilon(X)} I(\rho) = I'_\epsilon(X) \quad \text{say.}$$

This result will now be shown to follow from Shannon's proof in [12], coupled with our definitions and results. We will now explain what is going on, under the added restriction that  $H_\epsilon(X)$  be finite. This restriction will be removed in the next section.

In his Rate Distortion paper [12], Shannon studies coding for a source  $X$  with a fidelity criterion  $F$  which, for our purposes, is a nonnegative Borel function of the distance  $d$  between two points of  $X$ , which is a non-decreasing function of distance. He then asks the question of how much channel capacity is needed to send independent outcomes of  $X$  over a given noisy channel, if it is desired to keep the average distortion  $D$  between the actual outcome  $x$  and the decoded outcome  $y$  bounded by the constant  $A$ . The distortion  $D$  is defined as the expectation of  $F(x, y)$ , the expected loss of fidelity between actual outcome and decoded value. In our case,

$$(20) \quad \begin{aligned} F(x, y) &= 0, & d(x, y) &\leq \epsilon/2; \\ F(x, y) &= 1, & d(x, y) &> \epsilon/2. \end{aligned}$$

It is desired to know if there is a value  $C_0$  of channel capacity such that, over any channel with capacity  $C > C_0$ , outcomes of  $X$  can be transmitted over the channel in such a way that the distortion is zero, with probability approaching 1, whereas, if  $C < C_0$ , outcomes cannot be transmitted over the channel in such a way that the distortion is zero, with probability approaching 1, as the block length becomes infinite. Shannon's Theorems 1 and 2 of [12] show that this result is indeed true, where

$$(21) \quad C_0 = \inf_{\rho \in R_\epsilon(X)} I(\rho).$$

We state the desired result in a lemma. We mention that although the proof is adapted from Shannon's Theorem 1, the result does not follow from Shannon's statement.

LEMMA 6. *Let  $H_\epsilon(X) < \infty$ . Then*

$$I'_\epsilon(X) = I_\epsilon(X).$$

*In any case,*

$$I'_\epsilon(X) \leq I_\epsilon(X).$$

PROOF. First we shall prove that  $I_\varepsilon(X) \leq I'_\varepsilon(X)$ . Suppose that  $I'_\varepsilon(X) < \infty$ . Let  $\delta, \eta$  be positive numbers. We shall show that (39) is valid for  $\eta$  sufficiently large (so that (24), (31), (35) are (37) are satisfied) and use this to show  $I_\varepsilon(X) \leq I'_\varepsilon(X)$ .

By the definition of  $I'_\varepsilon(X)$ , we can choose  $\rho_\varepsilon R_\varepsilon(X)$  such that

$$(22) \quad I(\rho) \leq I'_\varepsilon(X) + \eta.$$

Let  $\nu = \nu_\rho$  be the second marginal distribution of  $\rho$  on  $X \times X$ . For a positive  $n$ , we consider the space  $Z_n$  which is the product of  $n$  independent copies of  $X \times X$ ,

$$Z_n = (X \times X)^{(n)}$$

with measure  $\rho^{(n)}$ , the product of the measure  $\rho$  on each  $X \times X$ .  $Z_n$  can also be written as

$$(23) \quad Z_n = X^{(n)} \times X^{(n)},$$

where the first factor has product measure  $\mu^{(n)} = \mu \times \mu \times \dots \times \mu$ , and the second factor has product measure  $\nu^{(n)}$ . Also,  $\mu^{(n)}$  and  $\nu^{(n)}$  are the marginal distributions of  $\rho^{(n)}$  in this factorization.

Let  $\bar{x} = (x_1, \dots, x_n), \bar{y} = (y_1, \dots, y_n)$  be the coordinates of a point of  $Z_n$  in the representation (23). Consider the function

$$J(\bar{x}, \bar{y}) = \frac{1}{n} \sum_{i=1}^n \log \frac{d\rho(x_i, y_i)}{d\mu(x_i)d\nu(y_i)}.$$

By the weak law of large numbers, for the given  $\delta$  there is an  $n_0$  such that

$$(24) \quad n \geq n_0$$

implies that on a set  $\Gamma$  in  $Z_n$  of probability  $\geq 1 - \delta^2$  we have

$$(25) \quad J(\bar{x}, \bar{y}) \leq I(\rho) + \eta.$$

Because of condition (ii) in the definition of  $R_\varepsilon(X)$ , we can demand that

$$d^{(n)}(\bar{x}, \bar{y}) \leq \varepsilon/2, \quad (\bar{x}, \bar{y}) \in \Gamma.$$

From (25), we have

$$(26) \quad \frac{d\rho^{(n)}(\bar{x}, \bar{y})}{d\mu^{(n)}(\bar{x})d\nu^{(n)}(\bar{y})} \leq \exp [n(I(\rho) + \eta)], \quad (x, y) \in \Gamma.$$

For any  $\bar{x} \in X^{(n)}$ , define

$$I_{\bar{x}} = \{\bar{y} \mid (\bar{x}, \bar{y}) \in \Gamma\}.$$

The quantity

$$(27) \quad \rho_{\bar{x}}^{(n)}(L_{\bar{x}}) = \int_{L_{\bar{x}}} \frac{d\rho^{(n)}(\bar{x}, \bar{y})}{d\mu^{(n)}(\bar{x})d\nu^{(n)}(\bar{y})} d\nu^{(n)}(\bar{y})$$

is the probability of  $L_{\bar{x}}$  given  $\bar{x}$ , and

$$(28) \quad \int_{X^{(n)}} \rho_{\bar{x}}^{(n)}(I_{\bar{x}}) d\mu^{(n)}(\bar{x}) = \rho^{(n)}(\Gamma) \geq 1 - \delta^2.$$

Let  $A$  be the set of  $\bar{x}$  such that

$$(29) \quad \rho_{\bar{x}}^{(n)}(L_{\bar{x}}) \geq 1 - \delta, \quad \bar{x} \in A.$$

It follows easily from (28) that

$$(30) \quad \mu^{(n)}(A) \geq 1 - \delta.$$

Using (26) to estimate the integral in (27) we get

$$\rho_{\bar{x}}^{(n)}(L_{\bar{x}}) \leq v^{(n)}(L_{\bar{x}}) \exp [n(I(\rho) + \eta)].$$

Combining with (29)

$$v^{(n)}(L_{\bar{x}}) \geq (1 - \delta) \exp [-n(I(\rho) + \eta)], \quad \bar{x} \in A.$$

Suppose that  $n$  is so large that

$$(31) \quad e^{n\eta}(1 - \delta) \geq 1.$$

Then

$$(32) \quad v^{(n)}(L_{\bar{x}}) \geq \exp [-n(I(\rho) + 2\eta)], \quad \bar{x} \in A.$$

Let  $N$  be a positive integer. Choose  $N$  independent points  $\bar{y}_1, \dots, \bar{y}_N$  in  $X^{(n)}$  according to the distribution  $v^{(n)}$ . Let  $P_N$  be the expected  $\mu^{(n)}$ -probability of the part of  $X^{(n)}$  not covered by the union of the spheres  $S_{\epsilon/2}(\bar{y}_j), j = 1, \dots, N$ . We have

$$\begin{aligned} P_N &= E_{\bar{y}_1, \dots, \bar{y}_N} \left\{ \int_{X^{(n)}} [1 - \chi_{\bar{x}} (\bigcup_{i=1}^N S_{\epsilon/2}(\bar{y}_i))] d\mu^{(n)}(\bar{x}) \right\} \\ &= \int_{X^{(n)}} E_{\bar{y}_1, \dots, \bar{y}_N} \left\{ \prod_{i=1}^N [1 - \chi_{\bar{x}}(S_{\epsilon/2}(\bar{y}_i))] \right\} d\mu^{(n)}(\bar{x}). \end{aligned}$$

By the independence, the integrand is equal to

$$\prod_{i=1}^N E_{\bar{y}_i} [1 - \chi_{\bar{x}}(S_{\epsilon/2}(\bar{y}_i))] = \prod_{i=1}^N [1 - v^{(n)}(S_{\epsilon/2}(\bar{x}))].$$

Thus,

$$P_N = \int_{X^{(n)}} [1 - v^{(n)}(S_{\epsilon/2}(\bar{x}))]^N d\mu^{(n)}(\bar{x}).$$

Since  $L_{\bar{x}} \in S_{\epsilon/2}(\bar{x})$ , (30) and (32) imply

$$P_N \leq \delta + (1 - \delta) \{1 - \exp [-n(I(\rho) + 2\eta)]\}^N.$$

Applying the inequality  $1 - t \leq e^{-t} (0 \leq t \leq 1)$ , we get

$$(33) \quad P_N \leq \delta + \exp \{-N \exp [-n(I(\rho) + 2\eta)]\}.$$

Let  $n$  be the integer for which

$$(34) \quad N - 1 < \exp [n(I(\rho) + 3\eta)] \leq N.$$

Then from (33) we have  $P_N \leq \delta + \exp (-e^{n\eta})$ .

For  $n$  sufficiently large,

$$(35) \quad \exp (-e^{n\eta}) < \delta.$$

Then

$$(36) \quad P_N < 2\delta.$$

Suppose that

$$(37) \quad n\eta > \log 2.$$

Then by virtue of the inequality

$$\log(1 + e^t) = t + \log(1 + e^{-t}) \leq t + \log 2, \quad t \geq 0,$$

(34) implies

$$(38) \quad \log N \leq n[I(\rho) + 4\eta].$$

From the definition of  $P_N$ , it follows that  $\bar{y}_1, \dots, \bar{y}_N$  can be selected so that the set

$$B = \bigcup_{i=1}^N S_{\varepsilon/2}(\bar{y}_i)$$

has  $\mu^{(n)}$ -measure at least  $1 - P_N$ . The set  $B$  can be partitioned into  $N$  sets of radius  $\varepsilon/2$  (constructed from the  $S_{\varepsilon/2}(\bar{y}_i)$ ). Hence, by (36) and (38)

$$H_{\varepsilon; 2\delta}(X^{(n)}) \leq n[I_\varepsilon(X) + 4\eta],$$

where  $\mu^{(n)}$ -measure is understood, or, from (22),

$$(39) \quad H_{\varepsilon; 2\delta}(X^{(n)}) \leq n[I'_\varepsilon(X) + 5\eta].$$

From the definition of  $I_{\varepsilon; [\delta]}$ ( $X$ ) and Lemma 4, (39) implies that if  $H_\varepsilon(X) < \infty$ ,

$$I_\varepsilon(X) \leq I'_\varepsilon(X) + 5\eta.$$

Now let  $\eta \rightarrow 0$ . This shows that

$$I_\varepsilon(X) \leq I'_\varepsilon(X),$$

as promised.

The reverse inequality is easier: given a partition  $U$  of  $X^{(n)}$  by sets  $U_i$ , of radius  $\varepsilon/2$  and centers  $\bar{y}_i$ , of entropy  $H(U)$ , we proceed as follows. The partition  $U$  induces a joint distribution  $\sigma$  on  $(X \times X)^{(n)}$  by associating  $n$ -tuples  $\bar{x}$  with the centers  $\bar{y}$  of the sets  $U_i$  containing  $\bar{x}$ . In the obvious notation, we have measures  $\rho_1, \rho_2, \dots, \rho_n$  on  $X \times X$  where  $\rho_i$  corresponds to the  $i$ th coordinate, and each  $\rho_i$  is in  $R_\varepsilon(X)$ . We have, by a simple calculation,

$$I(\bar{x}, \bar{y}) = H(U).$$

We can demand, given  $\eta > 0$ , that

$$\frac{1}{n} H(U) \leq I_\varepsilon(X) + \eta,$$

so that

$$(40) \quad \frac{1}{n} I(\bar{x}, \bar{y}) \leq I_\varepsilon(X) + \eta.$$

On the other hand, simple properties of  $I$  yield

$$I(\bar{x}, \bar{y}) \geq \sum_{i=1}^n I(x_i, \bar{y}) \text{ (since the } \{x_i\} \text{ are independent [5, Section 2.2])}$$

$$\geq \sum_{i=1}^n I(x_i, y_i) = \sum_{i=1}^n I(\rho_i) \text{ (since } I(A, (B, C)) \geq I(A, B)\text{).}$$

We conclude that for some  $i$ ,  $1 \leq i \leq n$ , we have

$$I(\rho_i) \leq \frac{1}{n} I(\bar{x}, \bar{y}),$$

and recall  $\rho_i \in R_\epsilon(X)$ . (37) then yields

$$I(\rho_i) \leq I_\epsilon(X) + \eta.$$

Definition (23) yields

$$I'_\epsilon(X) \leq I_\epsilon(X) + \eta,$$

for all  $\eta > 0$ . We finally conclude

$$(41) \quad I'_\epsilon(X) \leq I_\epsilon(X)$$

which completes the proof of Lemma 6.

The restriction that  $H_\epsilon(X) < \infty$  will be removed when we prove that  $H_\epsilon(X)$  is finite if  $I'_\epsilon(X)$  is finite.

To prove the channel coding result from Lemma 6, first we show that if  $K$  is a memoryless channel of capacity  $C$ , with

$$C > I_\epsilon(X),$$

then we can communicate over  $K$  keeping the distance between every received outcome and its corresponding transmitted outcome at most  $\epsilon/2$  with probability approaching 1 as the block length increases. To do this, merely take a partition  $U$  of  $X^{(n)}$  with

$$C > H(U) \geq I_\epsilon(X).$$

Use the centers of the  $U_i$  as a new source of entropy  $H(U)$  and transmit over  $K$  with probability of (word) error approaching zero.

To prove the converse, suppose, as in [12], Theorem 1, that we have a block code which encodes all of the  $n$ -tuples  $\bar{x}$  into  $m$ -tuples  $\bar{w}$ , and suppose that  $\bar{w}$  is received as say the  $m$ -tuple  $\bar{z}$ . Suppose also that we have a decoding procedure which decodes  $m$ -tuples  $\bar{z}$  into  $n$ -tuples  $\bar{y}$ , such that  $d^{(n)}(\bar{x}, \bar{y}) \leq \epsilon/2$  with probability at least  $1 - \eta$ . Shannon's result uses the condition

$$\frac{1}{n} \sum_{i=1}^n d(x_i, y_i) \leq \frac{\epsilon}{2} (1 + \eta), \quad \text{with } \eta > 0,$$

but we want the stronger condition with  $\eta = 0$ . Simple properties of  $I$  yield

$$I(\bar{w}, \bar{z}) \geq I(\bar{x}, \bar{y}).$$

Also, if  $\bar{w}_k, \bar{z}_k$  denote the first  $k$  components of  $\bar{w}, \bar{z}$  respectively, we have

$$\begin{aligned}\Delta_k(I(\bar{w}, \bar{z})) &= I(\bar{w}, \bar{z}_{k+1}) - I(\bar{w}, \bar{z}_k) \\ &= I(\bar{w}, (\bar{z}_k, z_{k+1})) - I(\bar{w}, \bar{z}_k) \\ &= I(z_{k+1}, (\bar{w}, \bar{z}_k)) - I(z_{k+1}, \bar{z}_k).\end{aligned}$$

Since the channel  $K$  is memoryless,

$$I(z_{k+1}, (\bar{w}, \bar{z}_k)) = I(z_{k+1}, w_{k+1}).$$

Thus

$$\begin{aligned}\Delta_k I(\bar{w}, \bar{z}) &= I(z_{k+1}, w_{k+1}) - I(z_{k+1}, \bar{z}_k) \\ &\leq I(z_{k+1}, w_{k+1}) = I(w_{k+1}, z_{k+1}).\end{aligned}$$

But the channel capacity  $C$  per channel letter is defined as  $\sup I(w, v)$  over all encodings  $w$ , decodings  $v$ . Hence

$$\Delta_k I(\bar{w}, \bar{z}) \leq C', \quad 1 \leq k \leq m,$$

where  $C'$  is the capacity of  $K$  per coded letter. Therefore

$$I(\bar{w}, \bar{z}) \leq mC', \quad \frac{1}{n} I(\bar{x}, \bar{y}) \leq \frac{m}{n} C' = C,$$

where  $C$  is the capacity per source letter.

Finally, by Lemma 5, since  $\eta$  is arbitrary  $> 0$ ,  $C \geq I_\varepsilon(X)$ , as promised. We shall not go into further detail.

To do the diametric case involves a slight complication in the definition of  $R_\varepsilon(X)$ . This happens because we can no longer choose a point  $y$  to send when we observe  $x$ , but instead must send an  $\varepsilon$ -set  $Y$  containing  $x$  when  $x$  occurs.

So we define matters as follows. Consider the complete separable metric space  $Y(X)$  of closed subsets of  $X$  under the Hausdorff metric [9, Section 1.4], or rather its closed subspace  $Y_\varepsilon(X)$  of sets of diameters at most  $\varepsilon$ . Then the diametric  $R_\varepsilon(X)$  is the class of Borel probability distributions  $\rho$  on  $X \times Y_\varepsilon(X)$  such that

- (i)  $\rho(W \times Y_\varepsilon(X)) = \mu(W)$  for Borel  $W \subset X$ ;
- (ii)  $\rho(\{x, Y\} \mid x \in Y) = 1$ .

In other words,  $\rho$  sends points of  $X$  into closed  $\varepsilon$ -sets containing the point, with  $\nu_\rho$ -probability 1.

We then define

$$I'_\varepsilon(X) = \inf_{\rho \in R_\varepsilon(X)} I(\rho),$$

and the proof goes through as before, with

$$\nu_\rho(B) = \rho(X \times B) \quad \text{for Borel } B \subset Y_\varepsilon(X).$$

The question can be raised as to what happens in the radial case if we define  $R_\varepsilon(X)$  as above, using the closed subset  $Z_\varepsilon(X)$  of  $Y_\varepsilon(X)$  of closed sets of radius  $\varepsilon/2$ .



The answer is that  $I'_\epsilon(X)$  is unchanged. The main point is that unique centers of the  $\epsilon/2$  spheres called out in the definition can be chosen in a measurable fashion. We omit details.

**5. Inequalities for  $I_\epsilon$  and  $H_\epsilon$ .** This section gives some inequalities on  $I_\epsilon$  and  $H_\epsilon$ , and, in particular, proves that  $H_\epsilon(X) < \infty$  if  $I_\epsilon(X) < \infty$ , closing a gap in the last section. We remark that all previous lower bounds for  $H_\epsilon(X)$  were lower bounds for  $I_\epsilon(X)$ , and, conversely, all previous upper bounds for  $I_\epsilon(X)$  were actually upper bounds for  $H_\epsilon(X)$  ([7], [8b]). The reason it is so hard to get bounds for one and not the other is that, as we shall see,  $H_\epsilon(X)$  and  $I_\epsilon(X)$  are close in ratio if either is large. We shall have to distinguish between radial and diametric entropy some of the time; when necessary, we use  ${}_D I_\epsilon(X)$ ,  ${}_D H_\epsilon(X)$  for the diametric case, and  ${}_R I_\epsilon(X)$ ,  ${}_R H_\epsilon(X)$  for the radial case. Observe the trivial inequalities

$$(42) \quad \begin{aligned} {}_R H_{2\epsilon}(X) &\leq {}_D H_\epsilon(X) \leq {}_R H_\epsilon(X), \\ {}_R I_{2\epsilon}(X) &\leq I_\epsilon(X) \leq {}_R I_\epsilon(X). \end{aligned}$$

When a result is true for either definition, we shall omit the presubscript. Our next lemma is quite useful, but will be strengthened later. It provides an upper bound on  $I_\epsilon(X)$ , but not on  $H_\epsilon(X)$ .

LEMMA 7.

$$I'_\epsilon(X) \leq E_\mu \left[ \log \frac{1}{\mu(S_{\epsilon/2}(x))} \right].$$

PROOF. It is enough to prove the result for  ${}_R I'_\epsilon(X)$ , in view of (42). We shall find a  $\rho \in R_\epsilon(X)$  such that

$$(43) \quad I(\rho) \leq E_\mu \left[ \log \frac{1}{\mu(S_{\epsilon/2}(x))} \right]$$

to prove the lemma. This  $\rho$  is, roughly speaking, the measure which takes an outcome  $x$  in  $X$ , and spreads its probability throughout  $S_{\epsilon/2}(x)$  according to  $\mu$  restricted to  $S_{\epsilon/2}(x)$ . That is, for  $A$  a Borel subset of  $X \times X$ ,

$$(44) \quad \rho(A) = \int \frac{\mu(S_{\epsilon/2}(x) \cap A_x)}{\mu(S_{\epsilon/2}(x))} d\mu(x),$$

where

$$A_x = \{y \mid (x, y) \in A\}.$$

The definition makes sense, since

$$\mu(\{x \mid \mu(S_{\epsilon/2}(x)) = 0\}) = 0.$$

Note that for  $B$  Borel in  $X$ ,

$$(45) \quad \nu_\rho(B) = \rho(X \times B) = \int_X \frac{\mu(S_{\epsilon/2}(x) \cap B)}{\mu(S_{\epsilon/2}(x))} d\mu(x).$$

We now evaluate

$$I(\rho) = \int_{X \times X} \log \left( \frac{d\rho}{d\mu dv_\rho} \right) d\rho.$$

Observe that

$$(46) \quad \frac{dv_\rho(y)}{d\mu(y)} = \int_{S_{\varepsilon/2}(y)} \frac{d\mu(x)}{\mu(S_{\varepsilon/2}(x))} = q(y), \quad \text{say.}$$

A simple calculation shows

$$(47) \quad \frac{d\rho(x, y)}{d\mu(x)dv_\rho(y)} = \frac{1}{q(y)} \frac{\chi_x(S_{\varepsilon/2}(y))}{\mu(S_{\varepsilon/2}(x))}.$$

Now the support of  $\rho$  is contained in

$$\{(x, y) \mid \chi_x(S_{\varepsilon/2}(y)) = 1\}.$$

Hence, continuing from (47), we have

$$\begin{aligned} I(\rho) &= \int_{X \times X} \log \frac{1}{\mu(S_{\varepsilon/2}(x))} d\rho(x, y) + \int_{X \times X} \log \frac{1}{q(y)} d\rho(x, y) \\ &= \int_X \log \frac{1}{\mu(S_{\varepsilon/2}(x))} d\mu(x) + \int_X \log \frac{1}{q(y)} dv_\rho(y). \end{aligned}$$

So, from (46), we have

$$(48) \quad I(\rho) + E_\mu \left[ \log \frac{1}{\mu(S_{\varepsilon/2}(x))} \right] + \int_X q(y) \log \frac{1}{q(y)} d\mu(y).$$

Now from (46) again

$$\begin{aligned} \int q(y) d\mu(y) &= \iint \frac{\chi_x(S_{\varepsilon/2}(y))}{\mu(S_{\varepsilon/2}(x))} d\mu(x) d\mu(y) \\ &= \int \frac{1}{\mu(S_{\varepsilon/2}(x))} \left[ \int \chi_y(S_{\varepsilon/2}(x)) d\mu(y) \right] d\mu(x), \end{aligned}$$

so that

$$(49) \quad \int q(y) d\mu(y) = 1.$$

Convexity of  $t \log 1/t$  then shows

$$(50) \quad \int q(y) \log \frac{1}{q(y)} d\mu(y) \leq 0.$$

This, coupled with (48), proves the lemma.

The prime on  $I'_\varepsilon(X)$  will be removed after the main theorem. The next lemma is sometimes a useful lower bound on  $I_\varepsilon(X)$ , and, of course, on  $H_\varepsilon(X)$ .

LEMMA 8. *Let*

$$\alpha = \sup_{S \text{ an } \varepsilon\text{-set}} \mu(S),$$

$$\beta = \sup_{X \in X} \mu(S_{\varepsilon/2}(x)),$$

so that  $\beta \leq \alpha$  (by [7], Theorem 3, these sup's are max's). Then

$$(51) \quad {}_D I_\varepsilon(X) \geq \log \frac{1}{\alpha},$$

$$(52) \quad {}_R I_\varepsilon(X) \geq \log \frac{1}{\beta},$$

PROOF. We shall do the  ${}_D I_\varepsilon(X)$  case, the other being similar. If  ${}_D I_\varepsilon(X) = \infty$ , there is nothing to prove. If  ${}_D H_\varepsilon(X) < \infty$ , let  $n$  be a positive integer, and let  $U$  be an  $\varepsilon$ -partition of  $X^{(n)}$ . Every set  $U_i$  in  $U$  has probability

$$\mu^{(n)}(U_i) \leq \alpha^n,$$

by hypothesis. Hence

$$H(U) \geq \log \frac{1}{\alpha^n},$$

$$H_\varepsilon(X^{(n)}) \geq \log \frac{1}{\alpha^n},$$

$$\frac{1}{n} H_\varepsilon(X^{(n)}) \geq \log \frac{1}{\alpha}.$$

Reference to (7) completes the proof of the lemma.

The next lemma is useful to bound the entropy of a union of spaces; it was implicitly used in Lemma 4. By abuse of language, we shall speak of the epsilon entropy of subsets of  $X$  even if they are merely measurable and not closed. The interpretation is quite natural: to send  $X$  say whether you are in  $Y$  or  $Z$ , then send  $Y$  or  $Z$ .

LEMMA 9. *Let  $X = \bigcup_k Y_k$ ,  $Y_k$  disjoint and measurable. Regard each  $Y_k$  as a probabilistic metric space in its own right, with measure and metric inherited from  $X$ . The renormalised probability measure on  $Y_k$  is  $\mu_k = p_k^{-1}\mu$ , where  $\mu(Y_k) = p_k$ . Define*

$$\mathcal{P} = \{p_k\}, \quad H(\mathcal{P}) = \sum p_k \log \frac{1}{p_k},$$

the entropy of  $\mathcal{P}$ . Then

$$(53) \quad H_\varepsilon(X) \leq \sum p_k H_\varepsilon(Y_k) + H(\mathcal{P}),$$

$$I'_\varepsilon(X) \leq \sum p_k I'_\varepsilon(Y_k) + H(\mathcal{P}).$$

Furthermore, these inequalities become equalities if, for each  $k, l$  with  $k \neq l$ ,

$$(54) \quad \inf_{y_k \in Y_k, y_l \in Y_l} d(y_k, y_l) > \varepsilon.$$

PROOF. If  $V_k = \{V_{ik}\}$  is an  $\varepsilon$ -partition of  $Y_k$  with  $H(V_k) = H_\varepsilon(Y_k)$ , let  $U$  be the  $\varepsilon$ -partition of  $X$  given by

$$(55) \quad U = \{V_{ik}, \text{ all } i, k\}.$$

Then

$$\begin{aligned} H_\varepsilon(X) &\leq \sum_{i,k} \mu(V_{ik}) \log \frac{1}{\mu(V_{ik})} \\ &= \sum_k p_k \sum_i \frac{\mu(V_{ik})}{p_k} \log \frac{p_k}{\mu(V_{ik})} + \sum_k p_k \log \frac{1}{p_k} \\ &= \sum_k p_k H_\varepsilon(Y_k) + H(\mathcal{P}), \end{aligned}$$

as required. When (54) is satisfied, every  $\varepsilon$ -partition of  $X$  is of the form (55), and so the assertion about equality holds.

To prove the result for  $I'_\varepsilon(X)$ , let  $\rho_k \in R_\varepsilon(Y_k)$  with

$$I(\rho_k) \leq I'_\varepsilon(Y_k) + \eta/2^k,$$

$\eta$  arbitrary  $> 0$ . Let  $\rho$  be the measure on  $X \times X$  such that

$$(56) \quad \rho(A) = \sum_k p_k \rho_k(A \cap (Y_k \times Y_k)).$$

Then  $\rho \in R_\varepsilon(X)$ . Also,

$$\begin{aligned} \nu_\rho(B) = \rho(X \times B) &= \sum p_k \rho_k(Y_k \times (B \cap Y_k)) \\ &= \sum p_k \nu_{\rho_k}(B \cap Y_k). \end{aligned}$$

Thus,

$$(57) \quad \begin{aligned} \frac{d\rho}{d\mu d\nu_\rho} &= \frac{1}{p_k} \frac{1}{d\mu_k d\nu_{\rho_k}} \quad \text{on } Y_k \times Y_k, \quad k = 1, 2, \dots, \\ &= 0 \quad \text{outside } \bigcup_k (Y_k \times Y_k). \end{aligned}$$

So

$$\begin{aligned} I(\rho) &= E_\rho \left( \log \frac{d\rho}{d\mu d\nu_\rho} \right) = \sum_k p_k E_{\rho_k} \left[ \log \left( \frac{1}{p_k} \left( \frac{d\rho_k}{d\mu_k d\nu_{\rho_k}} \right) \right) \right] \\ &\leq \sum_k p_k I'_\varepsilon(Y_k) + H(\mathcal{P}) + \eta. \end{aligned}$$

Since  $\eta$  is arbitrary, the result follows. When (54) holds, every  $\rho$  in  $R_\varepsilon(X)$  is of the form (56). This proves the lemma.

REMARK. The primes will be removed from  $I'_\varepsilon(X)$  in (53) after we prove the main theorem.

**6. Packings of sets.** This section prepares for the main theorem by introducing a new kind of random coding argument and using it to prove a result on “asymptotic close packing”. The argument was first used in Lemma 6.

LEMMA 10. *Let  $V$  be the random  $\varepsilon$ -partition of  $X$  defined in the following way. Choose a sequence  $\{x_i\}$  of points of  $X$  randomly  $d\mu$  and independently. Define*

$$W_i = S_{\varepsilon/2}(x_i), \quad i \geq 1.$$

Define the  $\varepsilon$ -partition  $V = \{V_n\}$  by

$$V_n = W_n - \bigcup_{i=1}^{n-1} W_i, \quad n \geq 1.$$

Then

$$(58) \quad \mu\left(\bigcup_{n=1}^{\infty} V_n\right) = 1 \quad \text{with probability } 1.$$

Also,

$$(59) \quad H_\varepsilon(X) \leq E(H(V)) \leq E\left(\log \frac{1}{\mu(S_{\varepsilon/2}(x))}\right) + E\left(\log \frac{1}{1 - \mu(S_{\varepsilon/2}(x))}\right) \cdot E\left(\frac{1 - \mu(S_{\varepsilon/2}(x))}{\mu(S_{\varepsilon/2}(x))}\right),$$

and

$$(60) \quad E(\mu(V_n)) = \int \mu(S_{\varepsilon/2}(x))(1 - \mu(S_{\varepsilon/2}(x)))^{n-1} d\mu(x).$$

PROOF. We have

$$(61) \quad E(H(V)) = \sum_{n \geq 1} E\left(\mu(V_n) \log \frac{1}{\mu(V_n)}\right);$$

by convexity,

$$(62) \quad E(H(V)) \leq \sum_{n \geq 1} E(\mu(V_n)) \log \frac{1}{E(\mu(V_n))}.$$

But

$$\begin{aligned} \sum_{i=1}^n E(\mu(V_i)) &= E\left(\sum_{i=1}^n \mu(V_i)\right) = E\mu\left(\bigcup_{i=1}^n W_i\right) \\ &= E\left(\int [1 - \prod_{i=1}^n (1 - \chi_x(W_i))] d\mu(x)\right) \end{aligned}$$

which, by the independence of the  $x_i$ , can be written as

$$\begin{aligned} &\int [1 - \prod_{i=1}^n (1 - E_{x_i}(\chi_x(S_{\varepsilon/2}(x_i))))] d\mu(x) \\ &= \int [1 - \prod_{i=1}^n (1 - E_{x_i}(\chi_{x_i}(S_{\varepsilon/2}(x))))] d\mu(x) \\ &= \int [1 - \prod_{i=1}^n (1 - \mu(S_{\varepsilon/2}(x)))] d\mu(x). \end{aligned}$$

That is,

$$(63) \quad E(\mu(\bigcup_{i=1}^n W_i)) = \sum_{i=1}^n E(\mu(V_i)) = 1 - \int (1 - \mu(S_{\varepsilon/2}(x)))^n d\mu(x).$$

Observe that the random variable  $1 - \mu(S_{\varepsilon/2}(x))$  takes the value 1 with probability zero; therefore

$$(64) \quad \int (1 - \mu(S_{\varepsilon/2}(x)))^n d\mu(x) \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

so

$$E(\mu(\bigcup_{n=1}^{\infty} V_n)) = 1,$$

and hence

$$\mu(\bigcup_{n=1}^{\infty} V_n) = 1$$

with probability 1, which proves (58). So we do indeed have a random family of  $\varepsilon$ -partitions  $V$  of  $X$ , whose entropies have expectation  $E(H(V))$ . Then there is at least one  $V$  with entropy at most  $E(H(V))$ , which proves the first inequality of (59).

Finally, (63) shows that (60) holds. The entropy  $H(\{E(\mu(V_n))\})$  of the sequence of probabilities  $\{E(\mu(V_n))\}$  is therefore given by

$$\begin{aligned} & H(\{E(\mu(V_n))\}) \\ &= \sum_{n \geq 1} \left[ \int \mu(S_{\varepsilon/2}(x))(1 - \mu(S_{\varepsilon/2}(x)))^{n-1} d\mu(x) \right] \log \\ & \quad \times \frac{1}{\int \mu(S_{\varepsilon/2}(x))(1 - \mu(S_{\varepsilon/2}(x)))^{n-1} d\mu(x)}. \end{aligned}$$

By convexity of the function  $\log(1/t)$ , we conclude

$$\begin{aligned} (65) \quad & H(\{E(\mu(V_n))\}) \\ & \leq \sum_{n \geq 1} \left[ \int \mu(S_{\varepsilon/2}(x))(1 - \mu(S_{\varepsilon/2}(x)))^{n-1} d\mu(x) \right] \\ & \quad \cdot \int \log \frac{1}{\mu(S_{\varepsilon/2}(x))(1 - \mu(S_{\varepsilon/2}(x)))^{n-1}} d\mu(x). \end{aligned}$$

(65) yields

$$\begin{aligned} (66) \quad & H(\{E(\mu(V_n))\}) = \left[ \int \log \frac{1}{\mu(S_{\varepsilon/2}(x))} d\mu \right] \cdot \int \sum_{n \geq 1} \mu(S_{\varepsilon/2}(x))(1 - \mu(S_{\varepsilon/2}(x)))^{n-1} d\mu \\ & \quad + \left[ \int \log \frac{1}{1 - \mu(S_{\varepsilon/2}(x))} d\mu \right] \\ & \quad \cdot \int \sum_{n \geq 1} (n-1)\mu(S_{\varepsilon/2}(x))(1 - \mu(S_{\varepsilon/2}(x)))^{n-1} d\mu. \end{aligned}$$

Evaluating the sums in (66) and using (62) completes the proof of (59). Lemma 10 is proved.

As a consequence of Lemma 10, we have the following result. Note that Lemma 7 yields

$$I_{\varepsilon}(X) \leq \log \frac{1}{\alpha} < \infty$$

if  $\mu(S_{\varepsilon/2}(x)) \geq \alpha > 0$  with probability 1.

LEMMA 11. *Let*

$$\mu(S_{\varepsilon/2}(x)) \geq \alpha > 0 \quad \text{with probability } 1.$$

*Then  $H_\varepsilon(X) < \infty$ . In fact,*

$$(67) \quad H_\varepsilon(X) \leq \log \frac{1}{\alpha} + \frac{1-\alpha}{\alpha} \log \frac{1}{1-\alpha} < \log \frac{1}{\alpha} + 1.$$

PROOF. (63) yields

$$(68) \quad \sum_{i=1}^n E(\mu(V_i)) \geq 1 - (1-\alpha)^n,$$

so that [7, Lemma 2] the entropy of the  $\{E(\mu(V_n))\}$  sequence is bounded by the entropy of the sequence  $\alpha(1-\alpha)^{n-1}$ . Consequently,

$$(69) \quad H_\varepsilon(X) \leq \sum_{n \geq 1} \alpha(1-\alpha)^{n-1} \log \frac{1}{\alpha(1-\alpha)^{n-1}},$$

which yields (67) and proves the lemma.

We are now ready for the asymptotic close packing result.

THEOREM 1. *Let  $X$  be a probabilistic metric space such that*

$$(70) \quad \mu(S_{\varepsilon/2}(x)) = \alpha$$

*with probability 1. Then*

$$(71) \quad {}_R I_\varepsilon(X) = \log \frac{1}{\alpha}$$

*and*

$$(72) \quad 0 \leq {}_R H_\varepsilon(X) - {}_R I_\varepsilon(X) < 1.$$

*If in addition*

$$(73) \quad \mu(A) \leq \alpha$$

*for every measurable set  $A$  of diameter at most  $\varepsilon$ , then*

$$(74) \quad {}_D I_\varepsilon(X) = \log \frac{1}{\alpha}$$

*and*

$$(75) \quad 0 \leq {}_D H_\varepsilon(X) - {}_D I_\varepsilon(X) < 1.$$

REMARK. The probability  $\alpha$  must be positive since  $\mu(X) = 1$ , so that  $H_\varepsilon(X) < \infty$  by Lemma 11. The assumption that  $\mu(S_{\varepsilon/2}(x)) = \alpha$  with probability 1 and that no measurable set of diameter at most  $\varepsilon$  has probability exceeding  $\alpha$  is satisfied by spheres, toruses, etc., under the natural measure given by surface hyper-area normalized to 1.

PROOF OF THEOREM. First assume only (70). By Lemma 8,

$${}_R I_\varepsilon(X) \geq \log \frac{1}{\alpha}.$$

By Lemma 7,

$${}_R I_\varepsilon(X) \leq \log \frac{1}{\alpha}.$$

Hence (71) holds. Lemma 11 then implies that (72) holds.

Now let both (70) and (73) hold. Lemma 8 gives

$${}_D I_\varepsilon(X) \geq \log \frac{1}{\alpha}.$$

(42) coupled with (71) gives

$${}_D I_\varepsilon(X) \leq \log \frac{1}{\alpha}.$$

Hence (74) holds. (42), coupled with (71), (74), and (72), yields (75). Theorem 1 is proved.

REMARK. We have seen that  $H_\varepsilon(X) - I_\varepsilon(X) < 1$  if  $\mu(S_{\varepsilon/2}(x))$  is constant with probability 1. It is not known whether this difference can be arbitrarily large, but we shall prove in the next section

$$H_\varepsilon(X) - I_\varepsilon(X) = o(H_\varepsilon(X)),$$

for  $H_\varepsilon(X)$  large.

*Discussion of Theorem 1.* Theorem 1 is called an “asymptotic close packing” result for the class of spaces to which it applies, for it says that the space can be partitioned by  $\varepsilon$ -sets in such a way that the resulting partition has entropy less than one more than the entropy of a partition consisting of  $1/\alpha$  sets each of probability  $\alpha$ . When  $\alpha \rightarrow 0$ , then, the space can be partitioned by a partition practically as “nice” as a close-packing one. Such nice partitions use up “most” of the probability with sets of probability “almost” equal to  $\alpha$ . One can use Theorem 1 to state and prove apparently otherwise difficult theorems about partitions of  $n$ -spheres, regular graphs, etc. ([6]).

**7. The main theorem.** In order to get the strongest possible form of the main theorem, we first need the following lemma, a strengthening of the known result ([7], page 54) that if  $\{a_n\}$  is a nonnegative sequence summing to 1, such that

$$\sum a_n \log n < \infty,$$

then

$$\sum a_n \log \frac{1}{a_n} < \infty$$

also (the converse is true, as well, but we do not need it here).



LEMMA 12. For  $\alpha, B > 0$ , consider all nonnegative sequences  $a_1, a_2, \dots$  with

$$(76) \quad \sum_{n=1}^{\infty} a_n \leq 1,$$

$$(77) \quad \sum_{n=1}^{\infty} a_n(\alpha + \log n) \leq B.$$

Let

$$(78) \quad M(B) = \sup \sum_{n=1}^{\infty} a_n \log \frac{1}{a_n},$$

over all such sequences.

Let  $\lambda$  be the unique solution on  $(1, \infty)$  of the equation

$$(79) \quad e^{-\lambda\alpha} [-\zeta'(\lambda) + \alpha\zeta(\lambda)] = Be,$$

where  $\zeta(s)$  is the Riemann Zeta-function. For  $B > \alpha$ , let  $\sigma$  be the unique solution on  $(1, \infty)$  of

$$(80) \quad -\zeta'(\sigma)/\zeta(\sigma) = B - \alpha.$$

Let  $B_0$  be the value of  $B$  at which  $e^{-\lambda\alpha}\zeta(\lambda) = e$ . Then

$$M(B) = \lambda B + e^{-1-\lambda\alpha}\zeta(\lambda), \quad B \leq B_0,$$

$$= \sigma(B - \alpha) + \log \zeta(\sigma), \quad B > B_0.$$

As  $B \rightarrow \infty$ ,

$$M(B) = B + \log B + O(1);$$

as  $B \rightarrow 0$ ,

$$M(B) = \frac{B}{\alpha} \log \frac{\alpha}{B} + o(B).$$

PROOF. For  $N$  a positive integer, define

$$(78') \quad M_N(B) = \sup \sum_{n=1}^N a_n \log \frac{1}{a_n},$$

under the conditions (76) and (77). Clearly

$$M_N(B) \leq M_{N+1}(B) \leq M(B), \quad N \geq 1.$$

For any  $\varepsilon > 0$ , if  $\{a_n\}$  is a sequence for which the sum in (78) is greater than  $M(B) - \varepsilon$ , then for  $N$  sufficiently large

$$M_N(B) \geq \sum_{n=1}^N a_n \log \frac{1}{a_n} > M(B) - 2\varepsilon.$$

Hence

$$M(B) = \lim_{N \rightarrow \infty} M_N(B).$$

The region of admissible values of  $(a_1, \dots, a_N)$  in  $N$ -space is a compact region. Hence there is a maximizing sequence for the sum in (78'):

$$M_N(B) = \max \sum_{n=1}^N a_n \log \frac{1}{a_n}.$$

Under the single condition

$$(76') \quad \sum_{n=1}^N a_n \leq 1,$$

this sum has a unique local maximum, at  $a_n = N^{-1}, n = 1, \dots, N$ . Here

$$\sum_{n=1}^N a_n(\alpha + \log n) > B,$$

if  $N$  is sufficiently large. Considering only such large values of  $N$ , we see that the maximizing sequence for  $M_N(B)$  must have

$$(77') \quad \sum_{n=1}^N a_n(\alpha + \log n) = B.$$

Suppose first that  $\sum_{n=1}^N a_n < 1$  for the maximizing sequence. Then  $M_N(B)$  is a local maximum of the series in (78') when  $a_1, \dots, a_N$  vary over the region in the hyperplane (77') where  $a_n \geq 0, n = 1, \dots, N$ . This maximum clearly does not occur on the boundary of the region, since the function  $x \log(1/x)$  has an infinite derivative at  $0+$ . Hence, by the method of Lagrange multipliers, at the maximum we have

$$\log \frac{1}{a_n} - 1 = \lambda_n(\alpha + \log n), \quad n = 1, 2, \dots, N,$$

or

$$(81) \quad a_n = e^{-1 - \alpha \lambda_n} n^{-\lambda_n}.$$

Define

$$\zeta_N(s) = \sum_{n=1}^N n^{-s}.$$

Then if  $\{a_n\}$  is given by (81), the condition (76') is

$$(82) \quad e^{-\alpha \lambda_N} \zeta_N(\lambda_N) \leq e,$$

and (77') becomes

$$(83) \quad e^{-\alpha \lambda_N} [\alpha \zeta_N(\lambda_N) - \zeta_N'(\lambda_N)] = Be.$$

In this case, the value of  $M_N(B)$  is

$$(84) \quad \begin{aligned} M_{N1}(B) &= \sum_{n=1}^N e^{-1 - \alpha \lambda_n} n^{-\lambda_n} [1 + \alpha \lambda_n + \lambda_n \log n] \\ &= \lambda_N B + e^{-1 - \alpha \lambda_N} \zeta_N(\lambda_N). \end{aligned}$$

Now suppose that equality holds in (76'). Then  $M_N(B)$  is the local maximum of the series in (78'), when  $a_1, \dots, a_N$  vary over a region in an  $(N-2)$ -dimensional hyperplane. Again this maximum cannot occur at a boundary point. Hence, at the maximum,

$$\log \frac{1}{a_n} - 1 = \mu_N + \sigma_N(\alpha + \log n), \quad n = 1, 2, \dots, N,$$

or

$$a_n = \exp(-1 - \mu_N - \alpha \sigma_N) n^{-\sigma_N}.$$

From (76') and (77')

$$\exp(-1 - \mu_N - \alpha\sigma_N)\zeta_N(\sigma_N) = 1.$$

Eliminating  $\mu_N$ , we get

$$\begin{aligned} \exp(-1 - \mu_N - \alpha\sigma_N)[\alpha\zeta_N(\sigma_N) - \zeta_N'(\sigma_N)] &= B. \\ -\zeta_N'(\sigma_N)/\zeta_N(\sigma_N) &= B - \alpha, \end{aligned} \tag{85}$$

and

$$a_n = n^{-\sigma_N}/\zeta_N(\sigma_N). \tag{86}$$

Now the value of  $M_N(B)$  is

$$\begin{aligned} M_{N2}(B) &= \zeta_N(\sigma_N)^{-1} \sum_{n=1}^N n^{-\sigma_N} [\sigma_N \log n + \log \zeta_N(\sigma_N)] \\ &= \sigma_N(B - \alpha) + \log \zeta_N(\sigma_N). \end{aligned}$$

The functions of  $\lambda_N, \sigma_N$  in (83) and (85) are strictly decreasing on  $(-\infty, \infty)$  for  $N \geq 2$ , taking all positive values. Thus there is a unique solution  $\lambda_N$  of (83) for any  $B > 0$ , and a unique  $\sigma_N$  satisfying (85) for any  $B > \alpha$ .

If (82) is violated, we must have  $M_N(B) = M_{N2}(B)$ . We shall show that  $M_N(B) = M_{N1}$  when (82) is satisfied. Let  $B_N$  be the value of  $B$  for which

$$e^{-\alpha\lambda_N} \zeta_N(\lambda_N) = e.$$

Then for  $B \leq B_N$ , (82) is satisfied, and (83) gives us

$$e^{-\alpha\lambda_N} [\alpha\zeta_N(\lambda_N) - \zeta_N'(\lambda_N)] \geq B e^{-\alpha\lambda_N} \zeta_N(\lambda_N),$$

or

$$-\zeta_N'(\lambda_N)/\zeta_N(\lambda_N) \geq B - \alpha.$$

If  $B > \alpha$ , then  $\sigma_N$  is defined by (85), and we see that

$$\sigma_N \geq \lambda_N, \tag{88}$$

with equality for  $B = B_N$ .

Expressing  $M_{N1}(B)$  and  $M_{N2}(B)$  by (84) and (87), and differentiating, one has

$$\begin{aligned} \frac{d}{dB} [M_{N1}(B) - M_{N2}(B)] &= \lambda_N - \sigma_N + [-B + \alpha - \zeta_N'(\sigma_N)/\zeta_N(\sigma_N)] \frac{d\sigma_N}{dB} \\ &\quad + [B + e^{-1 - \alpha\lambda_N} \{\zeta_N'(\lambda_N) - \alpha\zeta_N(\lambda_N)\}] \frac{d\lambda_N}{dB}. \end{aligned}$$

Applying (83) and (85), this reduces to

$$\frac{d}{dB} [M_{N1}(B) - M_{N2}(B)] = \lambda_N - \sigma_N,$$

and by (88), this is nonpositive for  $\alpha < B \leq B_N$ . We have  $M_{N1}(B_N) = M_{N2}(B_N)$ , since at  $B = B_N, \sigma_N = \lambda_N$  and

$$e^{-\alpha\lambda_N} \zeta_N(\lambda_N) = e^{-\alpha\sigma_N} \zeta_N(\sigma_N) = e.$$

Hence  $M_{N_1}(B) \geq M_{N_2}(B)$  for  $\alpha < B \leq B_N$ .  
Now we know that

$$\begin{aligned} M_N(B) &= M_{N_1}(B), & B \leq B_N, \\ &= M_{N_2}(B), & B > B_N. \end{aligned}$$

As  $N \rightarrow \infty$ , the functions  $\alpha\zeta_N(s) - \zeta_N'(s)$  and  $\zeta_N'(s)/\zeta_N(s)$  tend monotonically to the corresponding non-subscripted functions on  $(1, \infty)$ , while for  $s \leq 1$  they approach  $+\infty$  uniformly. Since the limit functions are strictly decreasing on  $(1, \infty)$ , we must have  $\lambda_N \rightarrow \lambda$  and  $\sigma_N \rightarrow \sigma$ , the solutions of (79) and (80). Also  $B_N \rightarrow B_0$ .

If  $B < B_0$ , or  $B > B_0$ , then for  $N$  sufficiently large the same inequality is valid with  $B_0$  replaced by  $B_N$ . Hence

$$\begin{aligned} M(B) &= \lim_{n \rightarrow \infty} M_{N_1}(B) = \lambda B + e^{-1-\lambda\alpha} \zeta(\lambda), & B < B_0, \\ &= \lim_{n \rightarrow \infty} M_{N_2}(B) = \sigma(B-\alpha) + \log \zeta(\sigma), & B > B_0. \end{aligned}$$

The value of  $M(B_0)$  is forced by the obvious monotonicity of the function  $M(B)$ .

For the asymptotic form of  $M(B)$  as  $B \rightarrow \infty$ , note that  $\sigma \rightarrow 1+$  as  $B \rightarrow \infty$ ; hence

$$\zeta(\sigma) = \frac{1}{\sigma-1} + O(1), \quad -\zeta'(\sigma) = \frac{1}{(\sigma-1)^2} + O(1).$$

From (80)

$$\frac{1}{\sigma-1} = B + O(1), \quad \sigma = 1 + B^{-1} + O(B^{-2}).$$

Then we have  $\zeta(\sigma) = B + O(1)$ , and

$$\begin{aligned} M(B) &= (B-\alpha)[1 + B^{-1} + O(B^{-2})] + \log B + O(B^{-1}) \\ &= B + \log B + O(1). \end{aligned}$$

As  $B \rightarrow 0$ ,  $\lambda \rightarrow \infty$ . Then

$$\zeta(\lambda) = 1 + O(2^{-\lambda}), \quad -\zeta'(\lambda) = 2^{-\lambda} \log 2 + O(3^{-\lambda}),$$

and by (79),

$$\begin{aligned} e^{-\lambda\alpha}[\alpha + O(2^{-\lambda})] &= Be, \\ e^{-\lambda\alpha} &= \frac{Be}{\alpha} [1 + O(2^{-\lambda})], \\ \lambda &= \frac{1}{\alpha} \log \frac{\alpha}{Be} + O(2^{-\lambda}). \end{aligned}$$

Hence

$$\begin{aligned} M(B) &= B \left[ \frac{1}{\alpha} \log \frac{\alpha}{Be} + O(2^{-\lambda}) \right] + \frac{B}{\alpha} [1 + O(2^{-\lambda})] \\ &= \frac{B}{\alpha} \log \frac{\alpha}{B} + o(B). \end{aligned}$$

This completes the proof of Lemma 12.

The next lemma gives a lower bound to  $I(\rho)$  for every  $\rho$  in  $R_\epsilon(X)$ .

LEMMA 13. *Let  $\rho \in {}_R R_\epsilon(X)$ . Then*

$$(89) \quad I(\rho) \geq E_\mu \left( \log \frac{1}{v_\rho(S_{\epsilon/2}(x))} \right) = {}_R K(\rho) \quad \text{say.}$$

*Let  $\rho \in {}_D R_\epsilon(X)$ . Then*

$$(90) \quad I(\rho) \geq E_\mu \left( \log \frac{1}{v_\rho(A_x)} \right) = {}_D K(\rho) \quad \text{say,}$$

where

$$A_x = \{Y \in Y_\epsilon(X) \mid x \in Y\}.$$

PROOF. Let  $\rho \in {}_R R_\epsilon(X)$ . We have

$$(91) \quad I(\rho) = \iint \left[ \log \left( \frac{d\rho}{d\mu dv_\rho} \right) \frac{d\rho}{d\mu dv_\rho} dv_\rho \right] d\mu.$$

Note that

$$\begin{aligned} \int \left( \log \frac{d\rho(x, y)}{d\mu(x) dv_\rho(y)} \right) \frac{d\rho(x, y)}{d\mu(x) dv_\rho(y)} dv_\rho(y) &= \int_{S_{\epsilon/2}(x)} \left( \log \frac{d\rho(x, y)}{d\mu(x) dv_\rho(y)} \right) \\ &\quad \cdot \frac{d\rho(x, y)}{d\mu(x) dv_\rho(y)} dv_\rho(y) \end{aligned}$$

by property (ii) of  ${}_R R_\epsilon(X)$ .

Hence, by convexity of  $t \log t$

$$(92) \quad \begin{aligned} &\int \left( \log \frac{d\rho(x, y)}{d\mu(x) dv_\rho(y)} \right) \frac{d\rho(x, y)}{d\mu(x) dv_\rho(y)} dv_\rho(y) \\ &\geq v_\rho(S_{\epsilon/2}(x)) \left[ \log \int \frac{d\rho(x, y)}{d\mu(x) dv_\rho(y)} \frac{dv_\rho(y)}{v_\rho(S_{\epsilon/2}(x))} \right] \int_{S_{\epsilon/2}(x)} \frac{d\rho(x, y)}{d\mu(x) dv_\rho(y)} \frac{dv_\rho(y)}{v_\rho(S_{\epsilon/2}(x))}. \end{aligned}$$

But if

$$\int_{S_{\epsilon/2}(x)} \frac{d\rho(x, y)}{d\mu(x) dv_\rho(y)} dv_\rho(y) = \int_X \frac{d\rho(x, y)}{d\mu(x) dv_\rho(y)} dv_\rho(y) = f(x),$$

then

$$\int_B f(x) d\mu(x) = \int \int_{B \times X} d\rho(x, y) = \mu(B),$$

so that  $f(x) = 1$  with  $\mu$ -probability 1. Thus, (92) yields

$$(93) \quad \int \left( \log \frac{d\rho(x, y)}{d\mu(x)dv_\rho(y)} \right) \frac{d\rho(x, y)}{d\mu(x)dv_\rho(y)} dv_\rho \geq \log \frac{1}{v_\rho(S_{\varepsilon/2}(x))},$$

with  $\mu$ -probability 1. (91) then proves (89). The proof of (90) is similar and omitted.

The next lemma upper bounds  $H_\varepsilon(X)$  in terms of the right-hand side of (89) or (90).

LEMMA 14. *Let  $\rho \in {}_R R_\varepsilon(X)$ . Then there exists an absolute constant  $C$  such that*

$$(94) \quad {}_R H_\varepsilon(X) \leq {}_R K(\rho) + \log^+ {}_R K(\rho) + C.$$

*Let  $\rho \in {}_D R_\varepsilon(X)$ . Then, for the same absolute constant  $C$ ,*

$$(95) \quad {}_D H_\varepsilon(X) \leq {}_D K(\rho) + \log^+ {}_D K(\rho) + C.$$

PROOF. As in Lemma 10, let  $V$  be the random  $\varepsilon$ -partition of  $X$  defined by choosing a sequence  $\{X_i\}$  of points of  $X$  randomly  $dv_\rho$ , defining

$$W_i = S_{\varepsilon/2}(x_i), \quad i \geq 1,$$

and then

$$V_n = W_n - \bigcup_{i=1}^{n-1} W_i, \quad n \geq 1.$$

The expected  $(d\rho^{(n)})$   $\mu$ -probability not covered by  $\bigcup_{i=1}^n V_i$  is given by

$$(96) \quad E \left[ \int \prod_{i=1}^n (1 - \chi_x(S_{\varepsilon/2}(x_i))) d\mu(x) \right] = \int (1 - v_\rho(S_{\varepsilon/2}(x)))^n d\mu(x),$$

as in Lemma 10. Let the set  $E$  be defined by

$$E = \{x \mid v_\rho(S_{\varepsilon/2}(x_i)) = 0\}.$$

We claim

$$\mu(E) = 0.$$

For

$$\mu(E) = \rho(E \times X) = \int_E \int_{S_{\varepsilon/2}(x)} \frac{d\mu}{d\mu dv_\rho} dv_\rho d\mu = 0,$$

since if  $d\rho$  is not absolutely continuous  $d\mu dv_\rho$ , there is nothing to prove.

We conclude that with  $dv_\rho$ -probability 1,  $V$  is an  $\varepsilon$ -partition of  $X$  under  $\mu$ -measure, as in Lemma 10. As in (62),

$$(97) \quad E(H(V)) \leq \sum_{n \geq 1} E(\mu(V_n)) \log \frac{1}{E(\mu(V_n))}.$$

And, from (96),

$$(98) \quad \sum_{i=1}^n E(\mu(V_i)) = 1 - \int (1 - v_\rho(S_{\varepsilon/2}(x)))^n d\mu.$$

Define

$$G(s) = \mu\{x \mid 1 - v_\rho(S_{\varepsilon/2}(x)) < s\},$$

and define

$$(99) \quad m_n = \int (1 - v_\rho(S_{\varepsilon/2}(x)))^n d\mu = \int_0^1 s^n dG(s),$$

the  $n$ th moment of  $s$  distributed  $dG(s)$ . As in Lemma 10, (97) yields, since  $(1 - m_1) \log 1/(1 - m_1) \leq 1/e$ ,

$$(100) \quad \begin{aligned} {}_R H_\varepsilon(X) &\leq \sum_{n=1}^\infty (m_n - m_{n+1}) \log \frac{1}{m_n - m_{n+1}} + (1 - m_1) \log \frac{1}{1 - m_1} \\ &\leq \sum_{n=1}^\infty (m_n - m_{n+1}) \log \frac{1}{m_n - m_{n+1}} + \frac{1}{e}. \end{aligned}$$

On the other hand

$${}_R K(\rho) = E_\mu \left[ \log \frac{1}{v_\rho(S_{\varepsilon/2}(x))} \right] = \int_0^1 \log \frac{1}{1-s} dG(s).$$

Expanding

$$\log \frac{1}{1-s} = s + \frac{s^2}{2} + \frac{s^3}{3} + \dots,$$

and then integrating, we conclude

$$(101) \quad {}_R K(\rho) = \sum_{n=1}^\infty \frac{m_n}{n}.$$

Let us sum the right-hand side of (101) by parts to obtain, since  $m_n \rightarrow 0$ ,

$$\begin{aligned} \sum_{n=1}^\infty \frac{m_n}{n} &= \sum_{n=1}^\infty \frac{1}{n} \sum_{i=n+1}^\infty (m_{i-1} - m_i) = \sum_{i=2}^\infty (m_{i-1} - m_i) \sum_{n=1}^{i-1} \frac{1}{n} \\ &\geq \sum_{i=1}^\infty (m_i - m_{i+1})(\log i + \gamma), \end{aligned}$$

where  $\gamma$  is Euler's constant. Hence

$$(102) \quad {}_R K(\rho) \geq \sum_{i=1}^\infty (m_i - m_{i+1})(\log i + \gamma).$$

Now

$$\sum_{i=1}^\infty (m_i - m_{i+1}) = m_1 \leq 1;$$

Lemma 12 then gives

$$(103) \quad \sum_{i=1}^\infty (m_i - m_{i+1}) \log \frac{1}{m_i - m_{i+1}} \leq {}_R K(\rho) + \log^+ {}_R K(\rho) + C',$$

for some universal constant  $C'$ . (100) next yields

$$(104) \quad {}_R H_\varepsilon(X) \leq {}_R K(\rho) + \log^+ {}_R K(\rho) + C' + 1/e.$$

Define  $C' + 1/e = C$  to prove (94) from (104). The proof of (95) is similar and omitted. This proves Lemma 14.

The next lemma is not really needed for our main purposes, but is inserted to show that  $H_\varepsilon(X)$  must approach 0 if  $I_\varepsilon(X)$  does.

LEMMA 15. *There exists a universal constant  $D$  such that, if  $K(\rho) \leq \frac{1}{4}(1 - 1/e)^2$ ,*

$$H_\epsilon(X) \leq D(K(\rho))^{\frac{1}{2}}.$$

PROOF. We shall improve the inequality (100) by getting an upper bound for

$$(1 - m_1) \log \frac{1}{1 - m_1}$$

in terms of  $K(\rho)$ . From definition (89), we find

$$\mu\{x \mid \log \frac{1}{v_\rho(S_{\epsilon/2}(x))} \leq K(\rho)^{\frac{1}{2}}\} \geq 1 - K(\rho)^{\frac{1}{2}},$$

or

$$\mu\{x \mid 1 - v_\rho(S_{\epsilon/2}(x)) \geq 1 - \exp[-K(\rho)^{\frac{1}{2}}]\} \geq K(\rho)^{\frac{1}{2}}.$$

From (99)

$$m_1 \leq 1 - \exp[-K(\rho)^{\frac{1}{2}}] + K(\rho)^{\frac{1}{2}} \leq 2K(\rho)^{\frac{1}{2}} \leq 1 - e^{-1}$$

and

$$1 - m_1 \geq 1 - 2K(\rho)^{\frac{1}{2}} \geq e^{-1}.$$

Hence

$$(1 - m_1) \log \frac{1}{1 - m_1} \leq [1 - 2K(\rho)^{\frac{1}{2}}] \log \frac{1}{1 - 2K(\rho)^{\frac{1}{2}}} \leq 2K(\rho)^{\frac{1}{2}}.$$

Thus, (102) and Lemma 12 yield a universal constant  $D_1$  such that

$$H_\epsilon(X) \leq \frac{K(\rho)}{\gamma} \log \frac{\gamma}{K(\rho)} + D_1 K(\rho) + 2(K(\rho))^{\frac{1}{2}},$$

which proves the lemma.

Note that  $D$  can be made arbitrarily close to 2 as  $K(\rho) \rightarrow 0$ .

We are now ready to state and prove Theorem 2, the main result of this paper.

THEOREM 2. *There exist universal constants  $C$  and  $D$  such that*

$$(105) \quad H_\epsilon(X) \leq I_\epsilon(X) + \log^+ I_\epsilon(X) + C,$$

and, for

$$I_\epsilon(X) \leq \frac{1}{8}(1 - 1/e)^2,$$

$$(106) \quad H_\epsilon(X) \leq D(I_\epsilon(X))^{\frac{1}{2}}.$$

Furthermore,  $I_\epsilon(X) = I'_\epsilon(X)$ .

PROOF.

Choose, given  $\eta > 0$ ,

$$(107) \quad \rho \in R_\epsilon(X), \quad \text{with } I'_\epsilon(X) \leq I(\rho) \leq I'_\epsilon(X) + \eta.$$

Then, by Lemma 13,

$$(108) \quad K(\rho) \leq I'(X) + \eta.$$



By Lemma 14, on the other hand, we found

$$H_\varepsilon(X) \leq K(\rho) + \log^+ K(\rho) + C,$$

so

$$H_\varepsilon(X) \leq I_\varepsilon'(X) + \log^+[I_\varepsilon'(X) + \eta] + C + \eta,$$

for every  $\eta > 0$ . This proves (105) for  $I_\varepsilon'(X)$  instead of  $I_\varepsilon(X)$ . By Lemma 6 and Lemma 3, if  $I_\varepsilon'(X)$  is infinite, there is nothing to prove. If  $I_\varepsilon(X)$  is finite, however, (109) forces  $H_\varepsilon(X)$  to be finite. By Lemma 6, then,

$$(110) \quad I_\varepsilon(X) = I_\varepsilon'(X),$$

and (105) is proved in the desired form (if  $I_\varepsilon'(X)$  is infinite, equality certainly holds).

Now let us prove (106). Let  $I_\varepsilon'(X) \leq \frac{1}{8}(1 - 1/e)^2$ , and choose  $\rho$  by (107), with  $\eta \leq \frac{1}{8}(1 - 1/e)^2$ . (108) holds, so that, by Lemma 15,

$$(111) \quad H_\varepsilon(X) \leq D(I_\varepsilon'(X) + \eta)^{\frac{1}{2}}.$$

Since (111) holds for all sufficiently small  $\eta > 0$ , and since (110) holds, (106) is true. This completes the proof of the main theorem.

The rest of the paper is devoted to consequences of the main theorem.

**8. Consequences of the main theorem.** This final section gives various consequences of Theorem 2. The first result is of independent interest and also needed for the second result. We recall [7, page 1008] that  $H_0(X)$  is defined as the infimum of the entropies of all partitions of  $X$  less a set of measure zero by atoms; when there is no such partition,  $H_0(X)$  is infinite. With this definition,  $H_\varepsilon(X)$  is continuous from above in  $\varepsilon$ , even at  $\varepsilon = 0$ , as was shown in the above cited reference.

**COROLLARY 1.**  $I_\varepsilon(X)$  is continuous from above in  $\varepsilon$  on  $\varepsilon \geq 0$ .

**PROOF.** Fix  $n$  a positive integer, and consider,

$$(112) \quad H_{\varepsilon'}(X^{(n)}) \leq I_{\varepsilon'}(X^{(n)}) + \log^+ I_{\varepsilon'}(X^{(n)}) + C.$$

First let  $H_\varepsilon(X) < \infty$ . We can then choose an  $\varepsilon' > \varepsilon$  so that

$$(113) \quad H_\varepsilon(X^{(n)}) < H_{\varepsilon'}(X^{(n)}) + 1,$$

since  $H(X^{(n)})$  is continuous from above in  $\varepsilon$ .

Now

$$I_\varepsilon(X^{(n)}) = nI_\varepsilon(X)$$

from (7), so

$$\begin{aligned} nI_\varepsilon(X) &\leq H_\varepsilon(X^{(n)}) < H_{\varepsilon'}(X^{(n)}) + 1 \\ &< nI_{\varepsilon'}(X) + \log^+ I_{\varepsilon'}(X) + \log n + C + 1. \end{aligned}$$

We then find

$$(114) \quad I_\varepsilon(X) < I_{\varepsilon'}(X) + \frac{1}{n} \log^+ I_{\varepsilon'}(X) + (\log n + C + 1)/n.$$

Since

$$I_{\varepsilon'}(X) \leq I_{\varepsilon}(X) < \infty,$$

given  $\eta > 0$ , we can by (114) choose an  $n$  so large that

$$(115) \quad I_{\varepsilon}(X) < I_{\varepsilon'}(X) + \eta,$$

which proves continuity from above in  $\varepsilon$  in case  $H_{\varepsilon}(X)$  is finite.

Now let  $H_{\varepsilon}(X)$ , and so  $I_{\varepsilon}(X)$ , be infinite. Given a large  $N > 0$ , we are to find an  $\varepsilon_0 > \varepsilon$  such that

$$(116) \quad I_{\varepsilon'}(X) > N \quad \text{when} \quad \varepsilon < \varepsilon' < \varepsilon_0.$$

We can assume that

$$(117) \quad H_{\varepsilon'}(X) < \infty, \quad \varepsilon' \geq \varepsilon.$$

By the continuity of  $H_{\varepsilon}(X)$  from above in  $\varepsilon$ , we can find an  $\varepsilon_0(n) > \varepsilon$  such that

$$(118) \quad H_{\varepsilon'}(X^{(n)}) > n^2, \quad \text{if} \quad \varepsilon < \varepsilon' \leq \varepsilon_0(n).$$

Thus (112) becomes

$$n^2 < H_{\varepsilon'}(X^{(n)}) < nI_{\varepsilon'}(X) + \log^+ I_{\varepsilon'}(X) + \log n + C,$$

or

$$(119) \quad n < I_{\varepsilon'}(X) + \frac{1}{n} \log^+ I_{\varepsilon'}(X) + \frac{1}{n} (\log n + C) \quad \text{for} \quad \varepsilon < \varepsilon' \leq \varepsilon_0(n).$$

But

$$\frac{1}{n} \log^+ I_{\varepsilon'}(X) \leq \log^+ I_{\varepsilon'}(X) \leq I_{\varepsilon'}(X),$$

so (119) can be written in the following form if  $n \geq n_0$ , where  $(\log n_0 + C)/n_0 \leq 1$ :

$$(120) \quad I_{\varepsilon'}(X) \geq \frac{1}{2}n - 1 \quad \text{for} \quad \varepsilon < \varepsilon' \leq \varepsilon_0(n).$$

Then (116) is proved by taking

$$n > \max(n_0, 2N + 2).$$

Corollary 1 is proved.

The proof of Corollary 1 can be modified to prove continuity of  $I_{\varepsilon}(X)$  in  $\varepsilon$  from below for certain cases in which it is known that  $H_{\varepsilon}(X)$  is continuous in  $\varepsilon$  from below. One case is that of mean-continuous Gaussian processes on the unit interval, by a modification of [8b], Theorem 1, to prove that  $H_{\varepsilon}(X^{(n)})$  is continuous from below in  $\varepsilon$ ,  $n \geq 2$ . Details are omitted.

**COROLLARY 2.** *For any probabilistic metric space  $X$ , one has*

$$(121) \quad H_{\varepsilon}(X) \sim I_{\varepsilon}(X) \quad \text{as} \quad \varepsilon \rightarrow 0.$$

PROOF. Since

$$I_\varepsilon(X) \leq H_\varepsilon(X) \leq I_\varepsilon(X) + \log^+ I_\varepsilon(X) + C$$

by Theorem 2, the result is true in case  $H_\varepsilon(X) \rightarrow \infty$  as  $\varepsilon \rightarrow 0$ . So suppose  $H_\varepsilon(X)$  remains bounded as  $\varepsilon \rightarrow 0$ . Then by the continuity of  $H_\varepsilon(X)$  from above in  $\varepsilon$  even at  $\varepsilon = 0$  ([7], page 1008), we conclude that  $H_0(X)$  is finite.

Now we always have

$$(122) \quad H_0(X) = I_0(X)$$

as (7) shows. By Corollary 1,

$$(123) \quad I_0(X) = \lim_{\varepsilon \rightarrow 0^+} I_\varepsilon(X).$$

If  $I_0(X) = 0$ , then  $H_0(X) = 0$ , and, in fact,

$$H_\varepsilon(X) = I_\varepsilon(X) = 0, \quad \text{all } \varepsilon > 0,$$

and there is nothing to prove. If  $I_0(X) > 0$ , then  $H_0(X) > 0$ , and, by the continuity in  $\varepsilon$  from above of  $H_\varepsilon(X)$

$$(124) \quad H_0(X) = \lim_{\varepsilon \rightarrow 0^+} H_\varepsilon(X).$$

Equations (123) and (124) taken in conjunction with (122) give

$$(125) \quad \lim_{\varepsilon \rightarrow 0^+} \frac{H_\varepsilon(X)}{I_\varepsilon(X)} = 1,$$

which proves (121) and Corollary 2.

The next batch of corollaries picks up some loose ends involving  $I'_\varepsilon(X)$ , defined in (19) as the inf of  $I(\rho)$  for  $\rho$  in  $R_\varepsilon(X)$ . We restate the following, although it was proved in Theorem 2.

COROLLARY 3.  $I'_\varepsilon(X) = I_\varepsilon(X)$ .

COROLLARY 4. Equation (53) holds for  $I_\varepsilon(X)$  in place of  $I'_\varepsilon(X)$ .

PROOF. Substitute  $I_\varepsilon(X)$  for  $I'_\varepsilon(X)$  in (53) to prove this corollary. We know of no other proof.

Now define for  $0 \leq \delta \leq 1$

$$(126) \quad \bar{I}_{\varepsilon;\delta}(X) = \inf_{A \text{ measurable}; \mu(A) \leq \delta} I_\varepsilon(X - A),$$

which is always finite for  $\delta > 0$ . Notice that

$$\bar{I}_{\varepsilon;0}(X) = I_\varepsilon(X).$$

Also,

$$(127) \quad \bar{I}_{\varepsilon;\delta}(X) \leq H_{\varepsilon;\delta}(X).$$

The next lemma is needed in the proof of Corollary 5.

LEMMA 16. *Let  $X$  be a probabilistic metric space with  $H_\epsilon(X)$  finite. Then for every  $\alpha > 0$  there exists a finite  $\beta$  such that, if  $Y$  is a measurable subset of  $X$ ,*

$$\mu(Y) \geq \alpha$$

*implies*

$$H_\epsilon(Y) \leq \beta.$$

PROOF. Let  $U = \{U_i\}$  be an  $\epsilon$ -partition of  $X$  with  $\{\mu(U_i)\}$  non-increasing and

$$\sum \mu(U_i) \log \frac{1}{\mu(U_i)} = H_\epsilon(X) < \infty.$$

Let

$$V = \{V_i\} = \{U_i \cap Y\},$$

an  $\epsilon$ -partition of  $Y$ . Its entropy can be bounded as follows:

$$\begin{aligned} H(V) &= \sum \frac{\mu(V_i)}{\mu(Y)} \log \frac{\mu(Y)}{\mu(V_i)} \\ &= \frac{1}{\mu(Y)} \sum \mu(V_i) \log \frac{1}{\mu(V_i)} + \frac{1}{\mu(Y)} \log \mu(Y) \\ &\leq \frac{1}{\alpha} \sum \mu(V_i) \log \frac{1}{\mu(V_i)} + \frac{1}{e}. \end{aligned}$$

Now  $\mu(U_i) \leq 1/e, i \geq 2$ , so

$$\begin{aligned} H(V) &\leq \frac{1}{\alpha} \left( \frac{2}{e} + \sum_{i \geq 2} \mu(U_i) \log \frac{1}{\mu(U_i)} \right) + \frac{1}{e} \\ &\leq \frac{1}{\alpha} \left( \frac{2}{e} + H_\epsilon(X) \right) + \frac{1}{e}. \end{aligned}$$

Thus, if we define  $\beta$  by

$$\beta = \frac{1}{\alpha} \left( \frac{2}{e} + H_\epsilon(X) \right) + \frac{1}{e},$$

then

$$H_\epsilon(Y) \leq H(V) \leq \beta,$$

and Lemma 16 is proved.

We then have the following result, proved for  $H_\epsilon(X)$  in [7], Theorem 5.

COROLLARY 5.  $I_\epsilon(X) = \lim_{\delta \rightarrow 0^+} \bar{I}_{\epsilon;\delta}(X)$ , and, in fact,  $\bar{I}_{\epsilon;\delta}(X)$  is continuous from above in  $\delta, 1 > \delta \geq 0$ . Also, given  $\eta > 0$ , there is a  $\delta_0 > 0$  such that, for  $A$  measurable of measure at most  $\delta_0$ ,

$$(128) \quad |I_\epsilon(X - A) - I_\epsilon(X)| < \eta,$$

if  $I_\epsilon(X)$  is finite. If  $I_\epsilon(X)$  is infinite, then, given  $N > 0$ , there is a  $\delta_0 > 0$  such that, for  $A$  measurable of measure at most  $\delta_0$ ,

$$(129) \quad I_\epsilon(X-A) > N.$$

Finally, if  $\delta > 0$ ,  $\bar{I}_{\epsilon;\delta}(X)$  is continuous from below in  $\delta$  if  $X$  is nonatomic.

PROOF. We first prove (128). Let  $n$  be so large that

$$(130) \quad \left| \frac{1}{n} H_\epsilon(X^{(n)}) - I_\epsilon(X) \right| < \frac{\eta}{3}.$$

Let  $\delta$  be so small that  $\delta \leq \frac{1}{2}$  and

$$(131) \quad \left| \frac{1}{n} H_\epsilon(X^{(n)}-B) - \frac{1}{n} H_\epsilon(X^{(n)}) \right| < \frac{\eta}{3},$$

provided  $\mu^{(n)}(B) \leq \delta$ . Then

$$\left| \frac{1}{n} H_\epsilon(X^{(n)}-B) - I_\epsilon(X) \right| \leq \frac{2\eta}{3},$$

whenever  $\mu^{(n)}(B) \leq \delta$ .

Define  $\delta_0$  such that

$$(132) \quad (1 - \delta_0)^n = 1 - \delta:$$

$\delta_0 \leq \frac{1}{2}$ . If  $A$  is measurable with  $\mu(A) \leq \delta_0$ , then

$$\mu^{(n)}((X-A)^{(n)}) \geq 1 - \delta, \quad \mu^{(n)}(X^{(n)} - (X-A)^{(n)}) \leq \delta,$$

and

$$(133) \quad \left| \frac{1}{n} H_\epsilon((X-A)^{(n)}) - I_\epsilon(X) \right| \leq \frac{2\eta}{3}.$$

Now by Theorem 2,

$$|H_\epsilon(X-A)^{(n)} - I_\epsilon((X-A)^{(n)})| \leq \log^+ I_\epsilon((X-A)^{(n)}) + C,$$

or

$$(134) \quad \left| \frac{1}{n} H_\epsilon((X-A)^{(n)}) - I_\epsilon((X-A)) \right| \leq \frac{1}{n} \log^+ I_\epsilon(X-A) + \frac{1}{n} (\log n + C).$$

By Lemma 16,

$$(135) \quad \sup_{\mu(A) \leq \frac{1}{2}} H_\epsilon(X-A) \leq E \quad \text{say.}$$

Choose  $n$  in advance so large that

$$(136) \quad \frac{1}{n} \log^+ E + \frac{1}{n} (\log n + C) < \frac{\eta}{3}.$$

Then (134) becomes

$$(137) \quad \left| \frac{1}{n} H_\varepsilon((X-A)^{(n)}) - I_\varepsilon(X-A) \right| < \frac{\eta}{3},$$

and (133) then yields (128).

Now to prove (129). In this case, (131) becomes

$$(138) \quad \frac{1}{n} H_\varepsilon(X^{(n)}-B) > 2N$$

whenever

$$\mu^{(n)}(B) \leq \delta.$$

Define  $\delta_0$  by (132), so that (133) becomes

$$(139) \quad \frac{1}{n} H_\varepsilon((X-A)^{(n)}) > 2N$$

if  $\mu(A) \leq \delta_0$ .

For  $H_\varepsilon(X-A)$  finite, the only case of interest, (134) becomes

$$(140) \quad \frac{1}{n} H_\varepsilon((X-A)^{(n)}) \leq I_\varepsilon(X-A) + \frac{1}{n} \log^+ I_\varepsilon(X-A) + \frac{1}{n} (\log n + C).$$

If

$$\frac{1}{n} \log^+ I_\varepsilon(X-A) \geq \frac{1}{n} \log N,$$

then (129) holds. So assume

$$(141) \quad \frac{1}{n} \log^+ I_\varepsilon(X-A) \leq \frac{1}{n} \log N.$$

Given  $N$  in advance, we choose  $n$  in advance so that

$$\frac{1}{n} (\log N + \log n + C) < N,$$

which makes (140) and (139) imply

$$2N < \frac{1}{n} H_\varepsilon((X-A)^{(n)}) < I_\varepsilon(X-A) + N,$$

which proves (129).

To prove the continuity of  $\bar{I}_{\varepsilon;\delta}(X)$  in  $\delta$  from above, we can assume  $\delta > 0$ ; the case  $\delta = 0$  is done by (128) and (129). Given  $\bar{I}_{\varepsilon;\delta}(X)$ , finite since  $\delta > 0$ , and  $\eta > 0$ , we are to prove that there is a  $\lambda_0 > 0$  with

$$(143) \quad I_\varepsilon(X-A) \geq \bar{I}_{\varepsilon;\delta}(X) - \eta$$

whenever  $\mu(A) \leq \delta + \lambda_0$ . And if  $\mu(A) \leq \delta$ , there is nothing to prove, since then (143) is satisfied even with  $\eta = 0$ .

Let us digress to prove the following result:

Given  $\delta > 0$ , there is a  $\lambda > 0$  such that if  $A$  satisfies  $\mu(A) \geq \delta$ , then, for some  $x$  in  $A$ ,

$$(144) \quad \mu(S_{\epsilon/2}(x) \cap A) \geq \lambda.$$

Here is a proof of that fact. Let  $B$  be such that  $\mu(B) \geq 1 - \delta/2$  and  $H_\epsilon(B) < \infty$ . Consider  $B \cap A$ , which has

$$\mu(B \cap A) \geq \delta/2.$$

By Lemma 16,

$$H_\epsilon(B \cap A) < E$$

for some constant  $E$ . By Lemma 8, there is an  $x$  in  $B \cap A$  with

$$\frac{2}{\delta} \mu(S_{\epsilon/2}(x) \cap B \cap A) \geq e^{-E},$$

and so, for some  $x$  in  $A$ , (144) holds, with

$$\lambda = \frac{\delta}{2} e^{-E}.$$

To prove (143), let

$$\delta \leq \mu(A) \leq \delta + \lambda.$$

Write  $W = S_{\epsilon/2}(x)$  where  $\mu(W \cap A) \geq \lambda$ . As in [7], Theorem 1, we can reduce to the case  $X$  nonatomic, in which case, given  $\lambda_0 \leq \lambda$ , there is a subset of  $V$  of  $W \cap A$  with  $\mu(V) = \lambda_0$ . Consider the set  $(X - A) \cup V$ , where

$$\mu((X - A) \cup V) \geq 1 - \delta.$$

This forces

$$(145) \quad I_\epsilon((X - A) \cup V) \geq \bar{I}_{\epsilon, \delta}(X).$$

On the other hand, Corollary 4 gives

$$(146) \quad I_\epsilon((X - A) \cup V) \leq pI_\epsilon(X - A) + (1 - p)I_\epsilon(V) + H(p),$$

with

$$p = \mu(X - A) / (\mu(X - A) + \mu(V));$$

$$H(p) = p \log \frac{1}{p} + (1 - p) \log \frac{1}{1 - p}.$$

Since

$$V \subset D_{\epsilon/2}(x),$$

however,

$$I_\epsilon(V) = 0,$$

and (146) becomes

$$(147) \quad I_\epsilon((X - A) \cup V) \leq pI_\epsilon(X - A) + H(p).$$

Together with (145), (147) gives

$$(148) \quad I_\epsilon(X - A) \leq \frac{1}{p} \bar{I}_{\epsilon;\delta}(X) - \log \frac{1}{p} - \frac{1-p}{p} \log(1-p).$$

Now

$$(149) \quad 1 < \frac{1}{p} \leq 1 + \frac{\lambda_0}{1-\delta}.$$

If we choose  $\lambda_0$  so small that

$$\log \left( 1 + \frac{\lambda_0}{1-\delta} \right) + \frac{\lambda_0}{1-\delta} \log \frac{1}{H\left(\frac{1-\delta}{\lambda_0}\right)} < \eta,$$

(143), and hence continuity from above, follows. To prove continuity from below when  $X$  is nonatomic is similar and omitted. This completes the proof of Corollary 5.

To get a stronger “strong converse” define, as in Lemma 5,

$$(150) \quad I_\epsilon''(X) = \inf_{[\delta_n] \rightarrow 0} \liminf_{n \rightarrow \infty} \frac{1}{n} H_{\epsilon;\delta_n}(X^{(n)}),$$

so that

$$(151) \quad I_\epsilon''(X) = I_\epsilon(X)$$

if  $H_\epsilon(X) < \infty$ . As in Lemma 5,  $I_\epsilon''(X)$  is the amount of capacity needed to transmit outcomes of  $X$  when it is desired that an arbitrarily large fraction of a block of outcomes be known to within  $\epsilon$ , with probability approaching 1. Alternatively,  $I_\epsilon''(X)$  is the Rate Distortion Function for zero distortion. The stronger converse is really just that  $I_\epsilon(X) = I_\epsilon''(X)$ . This means that if  $K$  is a memoryless channel of capacity  $C < I_\epsilon(X)$ , we cannot transmit outcomes of  $X$  over  $K$  such that, with probability approaching 1, an arbitrarily large fraction of a long block of outcomes is known to within  $\epsilon$ .

First recall

$$I_{\epsilon;[\delta]}(X) = \liminf_{n \rightarrow \infty} \frac{1}{n} H_{\epsilon;\delta}(X^{(n)})$$

for  $0 \leq \delta < 1$ . We then have the following corollary.

COROLLARY 6.

$$I_{\epsilon;[\delta]}(X) = \liminf_{n \rightarrow \infty} \frac{1}{n} \bar{I}_{\epsilon;\delta}(X^{(n)}).$$



PROOF. If

$$H_{\varepsilon;\delta}(X^{(n)}) \rightarrow \infty \quad \text{as } n \rightarrow \infty,$$

then (105) shows

$$H_{\varepsilon;\delta}(X^{(n)}) \sim \bar{I}_{\varepsilon;\delta}(X^{(n)}).$$

If, on the other hand,

$$(152) \quad H_{\varepsilon;\delta}(X^{(n)}) \leq F, \quad \text{say,}$$

for all  $n > 0$ , then we show that  $H_\varepsilon(X) = 0$ , and so

$$I_\varepsilon(X) = \bar{I}_{\varepsilon;\delta}(X) = H_{\varepsilon;\delta}(X^{(n)}) = \bar{I}_{\varepsilon;\delta}(X^{(n)}) = 0,$$

and the result follows. So suppose (152) holds. Then, by Lemma 8, there is a constant  $E > 0$  such that, for each  $n \geq 1$ , there is an  $\varepsilon$ -set  $A_n$  in  $X^{(n)}$  with

$$\mu^{(n)}(A_n) \geq E.$$

On the other hand,

$$\mu^{(n)}(A_n) \leq [\max_{S \text{ an } \varepsilon\text{-set}} \mu(S)]^n,$$

so

$$\max_{S \text{ an } \varepsilon\text{-set}} \mu(S) = 1,$$

and  $H_\varepsilon(X) = 0$ , which proves Corollary 6.

In order to prove the stronger converse, and strengthen Lemmas 4 and 5 by removing the requirement from them that  $H_\varepsilon(X)$  be finite, we need a lemma. We first need to define the class  $R_{\varepsilon;\eta}(X)$  of probability distributions for  $0 \leq \eta$ ; note that

$$R_{\varepsilon;0}(X) = R_\varepsilon(X).$$

DEFINITION.  $R_{\varepsilon;\eta}(X)$  is that class of Borel measures  $\rho$  on  $X \times X$  whose marginal distribution on the first coordinate is  $\mu$ , and such that

$$\rho(\{(x, y \mid d(x, y) \leq \varepsilon/2\}) \geq 1 - \eta.$$

The class  $D R_{\varepsilon;\eta}(X)$  is defined similarly for the diametric case.

Then define

$$I_{\varepsilon;\eta}(X) = \inf_{\rho \in R_{\varepsilon;\eta}(X)} I(\rho),$$

so that  $I_{\varepsilon;\eta}(X)$  is the Rate Distortion function evaluated for distortion  $\eta$ . Therefore,

$$I_{\varepsilon;\eta}(X) \leq \bar{I}_{\varepsilon;\eta}(X) + \eta \log \frac{1}{\eta},$$

and

$$I_{\varepsilon;0}(X) = I_\varepsilon(X).$$

Finally, define

$$I_\varepsilon'''(X) = \lim_{\eta \rightarrow 0^+} I_{\varepsilon;\eta}(X).$$

We have been unable to determine, except when either is 0 or  $\infty$ , in which case the answer is affirmative, whether

$$I_\varepsilon'''(X) = I_\varepsilon(X).$$

The next lemma shows that the two are however asymptotic when either is large. A strengthening of Theorem 2, it is enough for the stronger converse. The result for  $\eta > 0$  gives a result on the relation between  $H_{\varepsilon;\eta^{\frac{1}{2}}}(X)$  and  $I_{\varepsilon;\eta}(X)$ .

LEMMA 17. Let  $0 \leq \eta \leq 1/e^2$ . Define

$$L_{\varepsilon;\eta}(X) = \frac{1}{1-\eta} \left( I_{\varepsilon;\eta}(X) + \frac{2}{e} \right) + \log \frac{1}{1-\eta}.$$

Then, for  $C$  the same constant as in Theorem 2,

$$(153) \quad H_{\varepsilon;\eta^{\frac{1}{2}}}(X < (1 + 6\eta^{\frac{1}{2}})\{L_{\varepsilon;\eta}(X) + \log^+ L_{\varepsilon;\eta}(X) + C\} + 3\eta^{\frac{1}{2}}):$$

$$(154) \quad H_{\varepsilon}(X) \leq I_{\varepsilon}'''(X) + \log^+ I_{\varepsilon}'''(X) + C + \frac{4}{3}.$$

PROOF. Given  $\lambda > 0$ , choose  $\rho$  in  $R_{\varepsilon;\eta}(X)$  with

$$(155) \quad I(\rho) \leq I_{\varepsilon;\eta}(X) + \lambda.$$

Let  $S$  be that subset of  $X \times X$  on which, in the notation of Lemma 14,

$$A_x \subset S_{\varepsilon/2}(x),$$

so that

$$\rho(S) \geq 1 - \eta.$$

Define  $\rho_1$  to be the probability measure on  $X \times X$  given by

$$\rho_1(T) = \rho(S \cap T) / \rho(S),$$

with marginals  $\mu_1$  and  $\nu_1$  respectively. We observe, but do not use, that

$$\begin{aligned} \frac{d\mu_1(x)}{d\mu(x)} &= \frac{1}{\rho(S)} \int_{S_x} \frac{d\rho(x, y)}{d\mu(x)d\nu_{\rho}(y)} d\nu_{\rho}(y), \\ \frac{d\nu_1(y)}{d\nu_{\rho}(y)} &= \frac{1}{\rho(S)} \int_{S_y} \frac{d\rho(x, y)}{d\mu(x)d\nu_{\rho}(y)} d\mu(x). \end{aligned}$$

More important for our purpose is

$$(156) \quad \begin{aligned} \frac{d\rho_1(x, y)}{d\rho(x, y)} &= 0, & (x, y) \in S^c, \\ &= \frac{1}{\rho(S)}, & (x, y) \in S. \end{aligned}$$

As a consequence of (156), we have

$$(157) \quad \begin{aligned} \frac{d\rho_1(x, y)}{d\mu_1(x)d\nu_1(y)} &= 0, & (x, y) \in S^c \\ &= \frac{d\rho(x, y)}{d\mu(x)d\nu_{\rho}(y)} \frac{1}{\rho(S)} \frac{d\mu(x)}{d\mu_1(x)} \frac{d\nu_{\rho}(y)}{d\nu_1(y)}, & (x, y) \in S. \end{aligned}$$

From (157), we can conclude

$$(158) \quad I(\rho_1) = \frac{1}{\rho(S)} \int_S \log \left( \frac{d\rho}{d\mu dv_\rho} \right) d\rho + \log \frac{1}{\rho(S)} + \int \log \left( \frac{d\mu}{d\mu_1} \right) d\rho_1 + \int \log \left( \frac{dv_\rho}{dv_1} \right) d\rho_1.$$

Continuing from (158), we have the inequality

$$I(\rho_1) \leq \frac{1}{\rho(S)} \int \left| \log \frac{d\rho}{d\mu dv_\rho} \right| d\rho + \log \frac{1}{\rho(S)} + \int \left( \log \frac{d\mu}{d\mu_1} \right) \frac{d\mu_1}{d\mu} du + \int \left( \log \frac{dv_\rho}{dv_1} \right) \frac{dv_1}{dv_\rho} dv_\rho$$

which yields

$$(159) \quad I(\rho_1) \leq \frac{1}{\rho(S)} \left( I(\rho) + \frac{2}{e} \right) + \log \frac{1}{\rho(S)}.$$

The first term in (159) arises from the fact that the function  $t \log t \geq -1/e$  in  $t \geq 0$ ; there is no other term because  $-t \log t$  is convex in  $t \geq 0$ , and

$$\int \frac{d\mu_1}{d\mu} d\mu = \int \frac{dv_1}{dv_\rho} dv_\rho = 1.$$

Now form the random partition of Lemma 14, using the measure  $\rho_1$ , instead of  $\rho$ . Introduce the notation

$$H(U; \mu)$$

for  $U$  an  $\varepsilon$ -partition of  $X$  and  $\mu$  a Borel probability measure on  $X$  to denote the entropy of  $U$  when the measure in question is  $\mu$ . From Lemma 14, we conclude that there is an  $\varepsilon$ -partition  $U$  of a set of  $\mu_1$ -measure 1 in  $X$  such that

$$(160) \quad H(U; \mu_1) \leq I(\rho_1) + \log^+ I(\rho_1) + C,$$

for the  $C$  of Theorem 2. However, there are several obstacles to translating (160) into an inequality that will lead to (153).

The first of these is that  $U$  is an  $\varepsilon$ -partition of  $X$  with respect to  $\mu_1$ . What can be left uncovered by  $U$ ? If  $Z$  denotes the set of  $X$  not covered by  $U$ , then observe that

$$Z \times X \subset S^c.$$

Now let

$$U = \{U_i\}, \quad i \geq 1:$$

we have

$$\mu(\bigcup_i U_i) = \lambda \geq \rho(S) \geq 1 - \eta.$$

Let  $F$  be the set of indices  $i$  such that

$$(161) \quad \rho(U_i \times X) / \rho((U_i \times X) \cap S) \geq 1 + 2\eta^{\frac{1}{2}}$$

for  $i$  in  $F$ . If we define

$$V_F = \bigcup_{i \in F} U_i,$$

then we observe that  $V_F$  and  $Z$  are disjoint. Furthermore, we claim that

$$(162) \quad \rho(V_F \times X) = \mu(V_F) \leq \eta^{\frac{1}{2}} - \mu(Z).$$

For if not, then, since

$$\rho(V_F \times X) / \rho((V_F \times X) \cap S) \geq 1 + 2\eta^{\frac{1}{2}},$$

we have

$$\rho((V_F \times X) \cap S) \leq \rho(V_F \times X) / (1 + 2\eta^{\frac{1}{2}}),$$

or

$$\rho(V_F \times X) - \rho((V_F \times X) \cap S) \geq \frac{2\eta^{\frac{1}{2}}}{1 + 2\eta^{\frac{1}{2}}} \rho(V_F \times X).$$

Consequently,

$$\rho(V_F \times X) - \rho((V_F \times X) \cap S) \geq \frac{2\eta}{1 + 2\eta^{\frac{1}{2}}} - \frac{2\eta^{\frac{1}{2}}}{1 + 2\eta^{\frac{1}{2}}} \mu(Z) \geq \eta - \mu(Z),$$

since  $\eta < \frac{1}{4}$ , we would therefore have

$$\rho((V_F \times X) \cap S^c) + \rho((Z \times X) \cap S^c) > \eta,$$

and so

$$\rho(S^c) > \eta,$$

a contradiction. We conclude

$$\mu(V_F \cup Z) \leq \eta^{\frac{1}{2}}.$$

Now let  $U'$  be the  $\varepsilon; \eta^{\frac{1}{2}}$ -partition of  $X$  with respect to  $\mu$  consisting of those sets  $U_i$  of  $U$  with  $i$  not in  $F$ . The second obstacle to overcome is to estimate  $H(U'; \mu)$  in terms of  $H(U; \mu)$ .

We have

$$(163) \quad \begin{aligned} H(U'; \mu) &= \sum_{i \notin F} \frac{\mu(U_i)}{\lambda - \mu(V_F)} \log \frac{\lambda - \mu(V_F)}{\mu(U_i)} \\ &\leq \frac{1}{\lambda - \mu(V_F)} \sum_{i \notin F} \mu(U_i) \log \frac{1}{\mu(U_i)} \\ &= \frac{1}{\lambda - \mu(V_F)} T, \quad \text{say.} \end{aligned}$$

From the convexity of the function  $t \log (1/t)$  in  $t \geq 0$ , however, we can conclude

$$T \leq (\lambda - \mu(V_F)) \log \frac{1 - \mu_1(V_F)}{\lambda - \mu(V_F)} + \sum_{i \notin F} \mu(U_i) \log \frac{1}{\mu_1(U_i)},$$

and so

$$(164) \quad T \leq ((\lambda - \mu(V_F)) \log \frac{1 - \mu_1(V_F)}{\lambda - \mu(V_F)} + \sum_{i \notin F} \rho(U_i \times X) \log \frac{1}{\rho((U_i \times X) \cap S)}).$$

(164), coupled with the definition (161) of  $F$ , yields

$$T \leq (\lambda - \mu(V_F)) \log \frac{1 - \mu_1(V_F)}{\lambda - \mu(V_F)} + (1 + 2\eta^{\frac{1}{2}}) \sum_{\text{all } i} \rho((U_i \times X) \cap S) \log \frac{1}{\rho((U_i \times X) \cap S)},$$

and finally

$$(165) \quad T \leq (\lambda - \mu(V_F)) \log \frac{1}{\lambda - \mu(V_F)} + (1 + 2\eta^{\frac{1}{2}})H(U; \mu_1).$$

(163) coupled with (165) now gives

$$(166) \quad H_{\varepsilon; \eta^{\frac{1}{2}}}(X) < \log \frac{1}{\lambda - \mu(V_F)} + \frac{1 + 2\eta^{\frac{1}{2}}}{\lambda - \mu(V_F)} H(U; \mu_1),$$

and finally

$$(167) \quad H_{\varepsilon; \eta^{\frac{1}{2}}}(X) < 3\eta^{\frac{1}{2}} + (1 + 6\eta^{\frac{1}{2}})H(U; \mu_1).$$

Combining (167) with (160), (159), and (155) proves

$$(168) \quad H_{\varepsilon; \eta^{\frac{1}{2}}}(X) < 3\eta^{\frac{1}{2}} + (1 + 6\eta^{\frac{1}{2}}) \left\{ \frac{1}{1 - \eta} \left( I_{\varepsilon; \eta}(X) + \frac{2}{e} \right) + \log \frac{1}{1 - \eta} + \log \left[ \frac{1}{1 - \eta} \left( I_{\varepsilon; \eta}(X) + \frac{2}{e} \right) + \log \frac{1}{1 - \eta} \right] + C \right\},$$

which thus proves (153). To prove (154), let  $\eta \rightarrow 0^+$  in (168). We know from [5], Theorem 4, that

$$\lim_{\delta \rightarrow 0^+} H_{\varepsilon; \delta}(X) = H_{\varepsilon}(X),$$

and

$$\lim_{\delta \rightarrow 0^+} I_{\varepsilon; \delta}(X) = I_{\varepsilon}'''(X),$$

by definition. The simple inequality

$$\log \left( 1 + \frac{2}{e} \right) + \frac{2}{e} < \frac{4}{3}$$

completes the proof of (154) and hence of Lemma 17.

Now comes the stronger converse; define as in Lemma 5

$$I_{\varepsilon}''(X) = \inf_{[\delta_n] \rightarrow 0} \liminf_{n \rightarrow \infty} \frac{1}{n} H_{\varepsilon; \delta_n}(X^{(n)}).$$

Lemma 5 proved Corollary 7 when  $H_{\varepsilon}(X)$  is finite.

**COROLLARY 7.**  $I_{\varepsilon}''(X) = I_{\varepsilon}(X)$ .

PROOF. Since

$$I_\epsilon''(X) \leq I_\epsilon(X),$$

there is nothing to prove if  $I_\epsilon''(X)$  is infinite. So let  $I_\epsilon''(X)$  be finite. Then, as in Lemma 6, given an  $\epsilon; \delta_n$  partition  $U$  of  $X^{(n)}$ , such that

$$\frac{1}{n}H(U) \leq I_\epsilon''(X) + \eta, \quad \text{say,} \quad \text{with } \eta > 0,$$

we can obtain a  $\rho$  in  $R_{\epsilon; \delta_n}(X)$  with

$$I(\rho) \leq (1 - \delta_n) \left[ I_\epsilon''(X) + \eta + \log \frac{1}{1 - \delta_n} \right] + \delta_n \log \frac{1}{\delta_n}.$$

That is,

$$I_{\epsilon; \delta_n}(X) \leq (1 - \delta_n) \left[ I_\epsilon''(X) + \eta + \log \frac{1}{1 - \delta_n} \right] + \delta_n \log \frac{1}{\delta_n}.$$

By the definition of  $I_\epsilon''(X)$ , then,

$$I_\epsilon'''(X) \leq I_\epsilon''(X) + \eta.$$

Hence  $I_\epsilon'''(X)$  is finite, and by Lemma 17, so is  $H_\epsilon(X)$ . By Lemma 5, then,

$$I_\epsilon''(X) = I_\epsilon(X),$$

as promised. Corollary 7 is proved, and the restriction that  $H_\epsilon(X)$  be finite is removed from Lemma 5.

We can likewise show the following, which removes the assumption from Lemma 4 that  $H_\epsilon(X)$  be finite.

COROLLARY 8.  $I_\epsilon(X) = \lim_{\delta \rightarrow 0^+} I_{\epsilon; [\delta]}(X)$ .

PROOF. As in Corollary 7, we construct for every  $\delta > 0$ ,  $\eta > 0$ , a  $\rho$  in  $R_{\epsilon; \delta}(X)$  such that

$$I(\rho) \leq (1 - \delta) \left[ I_{\epsilon; [\delta]}(X) + \eta + \log \frac{1}{1 - \delta} \right] + \delta \log \frac{1}{\delta}.$$

Let the limit in the statement of the corollary be finite and equal to  $I$  say. Then we find that  $I_\epsilon'''(X)$  is also bounded by  $I$ . That is,  $I_\epsilon''(X)$ , and so  $H_\epsilon(X)$ , are finite, whenever

$$(169) \quad \lim_{\delta \rightarrow 0^+} I_{\epsilon; [\delta]}(X)$$

is finite. By Lemma 4, then, we can assume (169) infinite. But since the limit in (169) is at most  $I_\epsilon(X)$ ,  $I_\epsilon(X)$  is infinite if (160) is infinite. This proves Corollary 8.

We now summarize a combination of our principal results. If  $K$  is a memoryless channel with capacity  $\Gamma$  less than  $H_\epsilon(X) - \log^+ H_\epsilon(X) - C$ ,  $C$  a universal constant (or of finite capacity if  $H_\epsilon(X)$  is infinite) then it is not possible to transmit outcomes of  $X$  over  $K$  such that, with probability approaching 1, an arbitrarily large fraction

of a sequence of outcomes are known within  $\varepsilon$  or even a little more than  $\varepsilon$ . On the other hand, if  $\Gamma$  exceeds  $H_\varepsilon(X)$ , then outcomes of  $X$  can be transmitted over  $K$  such that, with probability approaching 1, all of a sequence of outcomes are known within  $\varepsilon$ , and such that block coding of outcomes need not be done before the channel encoder. In other words, the “one-shot” epsilon entropy tells whether the maximum error in a long block can be kept as little more than  $\varepsilon$  as we please, with probability approaching 1.

We close the paper with the outstanding open problem in this theory. Can

$$H_\varepsilon(X) - I_\varepsilon(X)$$

be arbitrarily large when either is finite? More generally, determine the function

$$f(x) = \sup_{H_\varepsilon(X) \leq x} [H_\varepsilon(X) - I_\varepsilon(X)],$$

which is finite-valued if  $x$  is finite. It is clear that  $f(x)$  is strictly increasing as  $x$  increases, that  $f(0) = 0$ , and that  $f(x)$  is continuous in  $x \geq 0$ . The main problem is whether  $f(x)$  is bounded in  $x \geq 0$ .

#### REFERENCES

- [1] FANO, R. M. (1961). *Theory of Information*. Wiley, New York.
- [2] KOLMOGOROV, ANDREI N. (1956). On the Shannon theory of information transmission in the case of continuous signals. *IEEE Trans. Information Theory* **IT-2** 102–108.
- [3] McELIECE, ROBERT J. and POSNER, EDWARD C. (1971). Hide and seek, data storage, and entropy. *Ann. Math. Statist.* **42** 1706–1716.
- [4] PINKSTON, JOHN T. III (1967). Encoding independent sample information sources. Res. Lab. Elect. Technical Report No. 462, Massachusetts Institute of Technology.
- [5] PINSKER, M. S. (1964). *Information and Information Stability of Random Variables and Processes*. (Translated from Russian.) Holden-Day, San Francisco.
- [6] POSNER, EDWARD C. and RODEMICH, EUGENE R. (1969). Differential entropy and tiling. *J. Statist. Phys.* **1** 57–69.
- [7] POSNER, EDWARD C., RODEMICH, EUGENE E. and RUMSEY, HOWARD JR. (1967). Epsilon entropy of stochastic processes. *Ann. Math. Statist.* **38** 1000–1020.
- [8a] POSNER, EDWARD C., RODEMICH, EUGENE R. and RUMSEY, HOWARD JR. (1969). Product entropy of Gaussian distributions. *Ann. Math. Statist.* **40** 870–904.
- [8b] POSNER, EDWARD C., RODEMICH, EUGENE R. and RUMSEY, HOWARD JR. (1969). Epsilon entropy of Gaussian processes. *Ann. Math. Statist.* **40** 1272–1296.
- [9] PROKHOROV, YU V. (1956). Convergence of random processes and limit theorems in probability theory. (English translation in *Theor. Probability Appl.* **1** 157–214.)
- [10] SAKRISON, DAVID J. (1968). The rate distortion function for a class of sources. Space Sciences Laboratory Report, Series No. 9, Univ. of California, Berkeley.
- [11] SHANNON, CLAUDE E. (1948). A mathematical theory of communication. *Bell System Tech. J.* **17** 623–656.
- [12] SHANNON, CLAUDE E. (1960). Coding theorems for a discrete source with a fidelity criterion. *Information and Decision Processes*. 93–126. McGraw-Hill, New York.
- [13] VITUSHKIN, A. (1962). *Complexity of the Tabulating Problem*. (Translated from Russian as *Theory of the Transmission and Processing of Information*.) Pergamon, New York.