

# Computation of Channel Capacity and Rate-Distortion Functions

RICHARD E. BLAHUT, MEMBER, IEEE

**Abstract**—By defining mutual information as a maximum over an appropriate space, channel capacities can be defined as double maxima and rate-distortion functions as double minima. This approach yields valuable new insights regarding the computation of channel capacities and rate-distortion functions. In particular, it suggests a simple algorithm for computing channel capacity that consists of a mapping from the set of channel input probability vectors into itself such that the sequence of probability vectors generated by successive applications of the mapping converges to the vector that achieves the capacity of the given channel. Analogous algorithms then are provided for computing rate-distortion functions and constrained channel capacities. The algorithms apply both to discrete and to continuous alphabet channels or sources. In addition, a formalization of the theory of channel capacity in the presence of constraints is included. Among the examples is the calculation of close upper and lower bounds to the rate-distortion function of a binary symmetric Markov source.

## I. INTRODUCTION

CHANNEL capacity, a fundamental concept in information theory, was introduced by Shannon [1] to specify the asymptotic limit on the maximum rate at which information can be conveyed reliably over a channel. The rate-distortion function, also introduced by Shannon [1], [2], serves an analogous function in the area of data compression coding for sources. These two basic concepts are discussed in detail in Gallager [3], Jelinek [4], and Berger [5].

Evaluation of a channel capacity  $C$  or a rate-distortion function  $R(D)$  involves the solution of a convex programming problem. In most cases analytic solutions cannot be found. Programmed computer search techniques have proved to be tedious even for small alphabet sizes and to be impractical for the larger alphabet sizes.

This paper reformulates the problems of computing  $C$  and  $R(D)$  from a new and slightly broader perspective, based on the observation that average mutual information  $I(p, Q)$  can be written in either of the two following forms:

$$I(p, Q) = \max_P \sum_j \sum_k p_j Q_{k|j} \log \frac{P_{j|k}}{p_j}$$

$$I(p, Q) = \min_q \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{q_k},$$

where  $P$  is an arbitrary transition matrix from the channel output alphabet to the channel input alphabet and  $q$  is an arbitrary probability distribution on the output alphabet.

Manuscript received April 22, 1971; revised July 29, 1971. This work was supported in part by the IBM Resident Study Program. It was part of the author's doctoral dissertation, Department of Electrical Engineering, Cornell University, Ithaca, N.Y.

The author is with the IBM Corporation, Owego, N.Y., and the Department of Electrical Engineering, Cornell University, Ithaca, N.Y.

Arimoto [13] used the first of the preceding expressions in an investigation of  $C$ , thereby obtaining Theorems 1 and 3 as well as Corollary 2 of this paper.<sup>1</sup>

This approach places the existing theory of  $C$  and  $R(D)$  in a more transparent setting and suggests several new results. In particular, the approach in question results in algorithms for determining  $C$  and  $R(D)$  by means of mappings from probability vectors to probability vectors. Under the first of these mappings, the sequence of average mutual informations associated with the successive channel input probability vectors increases monotonically to  $C$ . The other mapping produces a sequence of (information, distortion) pairs  $(I, D)$  that converges to a point on the  $R(D)$  curve; the convergence is monotonic in the  $(I, D)$  plane in the direction perpendicular to the slope of  $R(D)$  at the limiting point.

## II. CAPACITY OF UNCONSTRAINED DISCRETE CHANNELS

For the purposes of information theory, a discrete channel is described by a probability transition matrix  $Q = [Q_{k|j}]$  where  $Q_{k|j}$  is the probability of receiving the  $k$ th output letter given that the  $j$ th input letter was transmitted. In general,  $Q$  is not square. The capacity of the channel is defined as

$$C = \max_{p \in P^n} I(p, Q) = \max_{p \in P^n} \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}},$$

where

$$P^n = \{p \in R^n: p_j \geq 0 \forall j; \sum_j p_j = 1\}$$

is the set of all probability distributions on the channel input, and  $I(p, Q)$  is known as the mutual information between the channel input and channel output. The choice of logarithm base affects  $C$  only by a scale factor. It is usually convenient in applications to take base 2 so that  $C$  is expressed in terms of bits-per-channel use; for theoretical work, natural logs are more convenient.

The utility of the concept of capacity is widely discussed in the literature. Intuitively, the capacity of a channel expresses the maximum rate at which information can be reliably conveyed by the channel. Any coding scheme that superficially appears to operate at a rate higher than  $C$  will cause enough data to be lost because of uncorrectable channel errors so that the actual information rate is not to be greater than  $C$ .

Our concern in this section is with the calculation of capacity. The approach is to broaden the definition of

<sup>1</sup> The author is indebted to the editor for pointing out the prior existence of the Arimoto paper.

capacity to a larger maximization problem, which allows greater flexibility. This is done in the following theorem. Here, and in the sequel, maxima or minima are understood to be over the appropriate space of probability vectors or probability transition matrices (unless the domain is explicitly stated).

*Theorem 1:* Suppose the channel transition matrix  $Q$  is  $n \times m$ . For any  $m \times n$  transition matrix  $P$ , let

$$J(p, Q, P) = \sum_j \sum_k p_j Q_{k|j} \log \frac{P_{j|k}}{p_j}$$

Then the following is true.

- a)  $C = \max_p \max_P J(p, Q, P)$ .
- b) For fixed  $p$ ,  $J(p, Q, P)$  is maximized by

$$P_{j|k} = \frac{p_j Q_{k|j}}{\sum_j p_j Q_{k|j}}.$$

- c) For fixed  $P$ ,  $J(p, Q, P)$  is maximized by

$$p_j = \frac{\exp(\sum_k Q_{k|j} \log P_{j|k})}{\sum_j \exp(\sum_k Q_{k|j} \log P_{j|k})}.$$

*Proof:*

- a) It suffices to show that

$$I(p, Q) = \max_P \sum_j \sum_k p_j Q_{k|j} \log \frac{P_{j|k}}{p_j}.$$

Let

$$P_{j|k}^* = \frac{p_j Q_{k|j}}{\sum_j p_j Q_{k|j}}$$

and

$$q_k = \sum_j p_j Q_{k|j}$$

so that

$$I(p, Q) = \sum_j \sum_k q_k P_{j|k}^* \log \frac{P_{j|k}^*}{p_j}.$$

Then

$$\begin{aligned} I(p, Q) - \sum_j \sum_k p_j Q_{k|j} \log \frac{P_{j|k}}{p_j} &= \sum_j \sum_k q_k P_{j|k}^* \log \frac{P_{j|k}^*}{P_{j|k}} \\ &\geq \sum_j \sum_k q_k P_{j|k}^* - \sum_j \sum_k q_k P_{j|k} \\ &= 0 \end{aligned}$$

with equality<sup>2</sup> iff  $P_{j|k} = P_{j|k}^*$ .

b) This fact is an immediate consequence of the equality condition of part a).

c) If for some  $k$ ,  $P_{j|k} = 0$ , then  $p_j$  should be set equal to zero in order to maximize  $J$  as it is. Such a  $j$  can be deleted from the sum and dropped from further consideration.  $J(p, Q, P)$  can now be maximized over  $p$  by temporarily ignoring the constraint  $p_j \geq 0$ , and using a Lagrange multiplier to constrain

<sup>2</sup> The inequality used here is the well-known  $\log x \geq 1 - (1/x)$  with equality iff  $x = 1$ . This inequality will be used in the sequel without further comment.

$$\sum_j p_j = 1.$$

$$\begin{aligned} \frac{\partial}{\partial p_j} \left\{ \sum_j \sum_k p_j Q_{k|j} \log \frac{P_{j|k}}{p_j} + \lambda (\sum_j p_j - 1) \right\} &= 0 \\ -\log p_j - 1 + \sum_k Q_{k|j} \log P_{j|k} + \lambda &= 0. \end{aligned}$$

Hence,

$$p_j = \frac{\exp \sum_k Q_{k|j} \log P_{j|k}}{\sum_j \exp \sum_k Q_{k|j} \log P_{j|k}},$$

where  $\lambda$  is selected so that

$$\sum_j p_j = 1.$$

Notice that this  $p_j$  is always positive so that the inequality constraint  $p_j \geq 0$  is not operative.

The following corollary states a familiar condition on the solution of the basic problem. It is stated here both because it follows immediately from Theorem 1 and because the particular form that arises motivates the remainder of this section.

*Corollary 1:* If  $p$  achieves capacity, then

$$p_j = \frac{p_j \exp \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}}}{\sum_j p_j \exp \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}}}.$$

*Proof:* This is just the simultaneous satisfaction of parts b) and c) of the theorem.

The form of the equation in Corollary 1 is meant to suggest that any  $p$  can be used in the right-hand side in order to generate a new  $p$  on the left. Under appropriate conditions, this new  $p$  gives a better estimate of capacity as proved in Theorem 3.

*Corollary 2:*

$$C = \max_P \log \sum_j \exp(\sum_k Q_{k|j} \log P_{j|k}).$$

*Proof:* This follows from substituting part c) into part a).

The following specialization of the Kuhn–Tucker theorem will be used in the proof of Theorem 3.

*Theorem 2:* A vector  $p \in P^n$  achieves capacity for the channel with transition matrix  $Q$  if and only if there exists a number  $C$  such that

$$\sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} = C, \quad p_j \neq 0$$

$$\sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} \leq C, \quad p_j = 0.$$

For a proof, see Gallager or Jelinek. The conditions are sometimes called the Kuhn–Tucker conditions. The number  $C$  is then the channel capacity. It proves convenient to restate Theorem 2 as follows.

*Corollary 3:* A vector  $p \in P^n$  achieves capacity for the channel with transition matrix  $Q$  if and only if there exists a number  $C$  such that

$$\exp(-C) \exp \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} = 1, \quad p_j \neq 0$$

$$\exp(-C) \exp \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} \leq 1, \quad p_j = 0.$$

**Theorem 3:** For any  $p \in P^n$ , let

$$c_j(p) = \exp \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}}.$$

Then, if  $p^0$  is any element of  $P^n$  with all components strictly positive, the sequence of probability vectors defined by

$$p_j^{r+1} = p_j^r \frac{c_j^r}{\sum_j p_j^r c_j^r}$$

is such that  $I(p^r, Q) \rightarrow C$  as  $r \rightarrow \infty$ .

*Proof:* Given any  $p^r$ , we increase  $J(p, Q, P)$  by using Theorem 1-b) to pick  $P_{j|k}$  and then, with  $P_{j|k}$  fixed, using Theorem 1-c) to pick a new  $p$  vector. The composition of these two operations is just the operation that appears in the theorem. Hence, the algorithm in question increases mutual information. It also follows easily that the mutual information is strictly increasing unless Corollary 1 is satisfied by  $p^r$ , which in turn implies satisfaction of the first condition of Corollary 3. Thus,  $I(p, Q)$  is stable only for those  $p$  for which the first of the Kuhn-Tucker conditions is satisfied. We shall show that  $I^r$  can converge only to values of  $I(p, Q)$  that are stable in this way, and furthermore, that convergence is impossible unless the second of the Kuhn-Tucker conditions also is satisfied at the limit point.

Since  $I(p^r, Q)$  is increasing and is bounded by  $C$ ,  $I^r$  must converge to some number  $I^\infty \leq C$ . Let  $V(p^r) = I(p^{r+1}, Q) - I(p^r, Q)$ . Then  $V(p^r) \rightarrow 0$  since  $I^r$  converges. By the Bolzano-Weierstrass Theorem, the sequence  $(p^r)$  has a limit point  $p^*$  and a subsequence  $(p^r)$  converging to  $p^*$ . Therefore, by continuity of  $V$ ,  $V(p^r) \rightarrow V(p^*)$ . But  $V(p^r) \rightarrow 0$ . Therefore,  $V(p^*) = 0$  and hence  $p^*$  satisfies the first of the Kuhn-Tucker conditions.

Now suppose  $p^*$  does not achieve capacity. Then by the sufficiency condition of Corollary 3,

$$\frac{c_j^*}{\sum_j p_j^* c_j^*} > 1$$

for some  $j$ , where  $c_j^* = c_j(p^*)$ .

Since some subsequence  $\{p^{r_k}\}$  converges to  $p^*$ , then by continuity  $\{c_j^{r_k}\}$  converges to  $c_j^*$  for all  $j$ . But,

$$p_j^{r_k} = p_j^0 \prod_{n=0}^{r_k} b_j^n$$

where

$$b_j^n = \frac{c_j^n}{\sum_j p_j^n c_j^n}$$

and  $\{b_j^n\}$  has a subsequence converging to a number greater than 1. Therefore, the sequence of partial products does not converge and  $p_j^{r_k}$  does not converge, which is a contradiction.

Therefore,  $p^*$  achieves capacity and  $I^\infty = C$ . This completes the proof of the theorem.

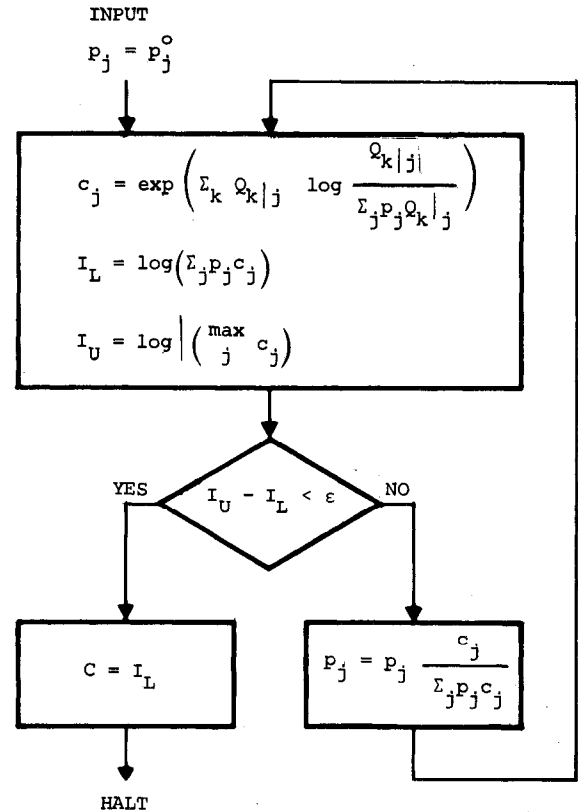


Fig. 1. Capacity algorithm.

The application of Theorem 3 to the computation of channel capacity is illustrated in Fig. 1. The termination is based on the fact that for any probability assignment  $p$  the following holds

a)

$$C \geq \log \sum_j p_j c_j$$

b)

$$C \leq \log (\max_j c_j),$$

where

$$c_j = \exp \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}}.$$

Part a) is a simple consequence of Corollary 2 and part b) appears as a problem in Gallager [3, p. 524].

### III. RATE-DISTORTION FUNCTIONS FOR DISCRETE SOURCES

A discrete-alphabet memoryless source, which produces the  $j$ th letter with probability  $p_j$ , is to be reproduced in terms of a second alphabet that need not be of the same size, although often it is identical to the source alphabet. A distortion matrix with elements  $\rho_{jk}$  specifies the distortion associated with reproducing the  $j$ th source letter by the  $k$ th reproducing letter ( $0 \leq j \leq m-1$ ,  $0 \leq k \leq n-1$ ). Without loss of generality, it can be assumed that for each source letter, there is at least one reproducing letter such that the resulting distortion equals zero.

Rate-distortion theory is concerned with the average amount of information about the source output that must be preserved by any data compression scheme such that

the reproduction can be subsequently generated from the compressed data with average distortion less than or equal to some specified  $D$ . The rate-distortion function is defined as

$$R(D) = \min_{Q \in \mathcal{Q}_D} \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} = \min_{Q \in \mathcal{Q}_D} I(p, Q),$$

where

$$\mathcal{Q}_D = \{Q \in R^m \times R^n : \sum_k Q_{k|j} = 1, Q_{k|j} \geq 0, d(Q) \leq D\}$$

and

$$d(Q) = \sum_j \sum_k p_j Q_{k|j} \rho_{jk}.$$

The definition of rate-distortion functions is justified by source compression theorems, which are widely reported [3]–[5]. Intuitively, if average distortion  $D$  is specified, then any compression must retain an average of at least  $R(D)$  bits per source letter, and conversely, compression to a level arbitrarily close to  $R(D)$  is possible by appropriate selection of the compression scheme.

The investigation of rate-distortion functions is usually carried out parametrically in terms of a parameter  $s$ , which is introduced as a Lagrange multiplier. This parameter turns out to be equal to the slope of the rate-distortion curve at the point it parameterizes [5]. These facts will be assumed in the following and the discussion will begin with the following parametric expression for  $R(D)$ .

$$R(D) = \min_Q \left[ \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} - s(\sum_j \sum_k p_j Q_{k|j} \rho_{jk} - D) \right],$$

where

$$D = \sum_j \sum_k p_j Q_{k|j}^* \rho_{jk}$$

and  $Q^*$  is the point that achieves the above minimum.

The minimization is now over all transition matrices  $Q$ . The value of  $D$ , however, is no longer an input to the computation; rather, a value of  $s$  is specified whereupon both  $D$  and  $R(D)$  are generated for the point on the  $R(D)$  curve that has slope  $s$ .

*Theorem 4:* Let

$$F(p, Q, q) = \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{q_k} - s \sum_j \sum_k p_j Q_{k|j} \rho_{jk}.$$

Then

a)

$$R(D) = sD + \min_q \min_Q F(p, Q, q),$$

where

$$D = \sum_j \sum_k p_j Q_{k|j}^* \rho_{jk}$$

and  $Q^*$  achieves the above minimum.

b) For fixed  $Q_{k|j}$ ,  $F(p, Q, q)$  is minimized by

$$q_k = \sum_j p_j Q_{k|j}.$$

c) For fixed  $q$ ,  $F(p, Q, q)$  is minimized by

$$Q_{k|j} = \frac{q_k \exp(s \rho_{jk})}{\sum_k q_k \exp(s \rho_{jk})}.$$

*Proof:*

a) It suffices to prove that  $I(p, Q) = \min_q F(p, Q, q)$ .

$$\begin{aligned} \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{q_k} - I(p, Q) &= \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{q_k} \\ &\quad - \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} \\ &= \sum_j \sum_k p_j Q_{k|j} \log \frac{\sum_j p_j Q_{k|j}}{q_k} \\ &\geq \sum_j \sum_k p_j Q_{k|j} - \sum_k q_k = 0 \end{aligned}$$

with equality if and only if

$$q_k = \sum_j p_j Q_{k|j}.$$

b) This follows immediately from the equality condition of part a).

c) Temporarily ignore the inequality constraint  $Q_{k|j} \geq 0$  and introduce a Lagrange multiplier to constrain  $\sum_k Q_{k|j} = 1$ .

$$\frac{\partial}{\partial Q_{k|j}} \left[ \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{q_k} - s \sum_j \sum_k p_j Q_{k|j} \rho_{jk} + \sum_j \lambda_j \sum_k Q_{k|j} \right] = 0$$

$$p_j \log Q_{k|j} - p_j \log q_k + p_j - s p_j \rho_{jk} + \lambda_j = 0.$$

Hence

$$Q_{k|j} = \frac{q_k \exp(s \rho_{jk})}{\sum_k q_k \exp(s \rho_{jk})},$$

where  $\lambda_j$  has been selected as that

$$\sum_k Q_{k|j} = 1.$$

Notice that this is always nonnegative so that the inequality constraint  $Q_{k|j} \geq 0$  is satisfied.

A familiar condition on the minimizing  $Q$  is the following.

*Corollary 4:* If  $Q$  achieves a point on the  $R(D)$  curve parameterized by  $s$ , then

$$Q_{k|j} = \frac{q_k \exp(s \rho_{jk})}{\sum_k q_k \exp(s \rho_{jk})},$$

where

$$q_k = \sum_j p_j Q_{k|j} = q_k \sum_j p_j \frac{\exp(s \rho_{jk})}{\sum_k q_k \exp(s \rho_{jk})}.$$

*Proof:* This is just the simultaneous satisfaction of parts b) and c). The first equation of Corollary 4 defines a transition matrix  $Q(q)$  given any  $q$ . This will form the basis for the algorithm of Theorem 6.

*Corollary 5:* In terms of the parameter  $s$ ,

$$R(D_s) = sD_s + \min_q [-\sum_j p_j \log \sum_k q_k \exp(s \rho_{jk})]$$

$$D_s = \sum_j p_j \frac{q_k^* \exp(s \rho_{jk})}{\sum_k q_k^* \exp(s \rho_{jk})} \rho_{jk},$$

where  $q_k^*$  achieves  $R(D_s)$ .

*Proof:* This follows immediately by substituting part c) of the theorem into part a).

Corollary 5 expresses the substance of a theorem by Haskell [6]. The following variation is also useful.

*Corollary 6:*

$$R(D) = \max_{s \in [-\infty, 0]} \min_q [sD - \sum_j p_j \log \sum_k q_k \exp(s\rho_{jk})].$$

*Proof:* Let  $Q$  achieve  $R(D)$  and let  $D_s$  be the average distortion value parameterized by  $s$ . Then

$$\begin{aligned} R(D) - \min_q [sD - \sum_j p_j \log \sum_k \exp(s\rho_{jk})q_k] \\ = R(D) - sD - \min_q [sD_s - \sum_j p_j \log \sum_k q_k \exp(s\rho_{jk})] \\ + sD_s \\ = R(D) - sD - R(D_s) + sD_s \\ = \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j} \exp(s\rho_{jk})}{\sum_j p_j Q_{k|j}} \\ - \min_q \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j} \exp(s\rho_{jk})}{\sum_j p_j Q_{k|j}} \\ \geq 0. \end{aligned}$$

The content of this corollary can be expressed in a pleasant form if

$$\sum_j \exp(s\rho_{jk})$$

is independent of  $k$ . We digress further to illustrate this in a special case.

*Corollary 7:* Suppose the alphabet consists of binary  $n$ -tuples and the distortion is Hamming distance. Then

$$R(D) = \max_{\rho \in [0, \frac{1}{2}]} \min_q [D \log \rho + (1 - D) \log (1 - \rho) + \sum_j p_j \log \sum_k A_{jk}(\rho) q_k],$$

where  $A_{jk}(\rho)$  is the  $n$ -tuple transition matrix of a binary symmetric channel of transition probability  $\rho$ . In particular, if

$$\sum_j A_{jk}^{-1}(D) p_j \geq 0 \quad \forall k$$

then

$$R(D) = D \log D + (1 - D) \log (1 - D) - \sum_j p_j \log p_j.$$

*Proof:* Let the superscript  $n$  denote the block length and notice that the matrix  $\exp(s\rho_{jk}^n)$  can be expressed inductively by

$$\exp(s\rho_{jk}^n) = \begin{vmatrix} \exp(s\rho_{jk}^{n-1}) & \exp(s) \exp(s\rho_{jk}^{n-1}) \\ \exp(s) \exp(s\rho_{jk}^{n-1}) & \exp(s\rho_{jk}^{n-1}) \end{vmatrix}.$$

Define  $\rho$  by  $\exp(s) = \rho/(\rho - 1)$ . The result then follows from the previous corollary.

The analog of Theorem 2 is the following.

*Theorem 5:* A necessary and sufficient condition on an output probability assignment  $q$  to yield a point on the  $R(D)$  curve via the transition matrix

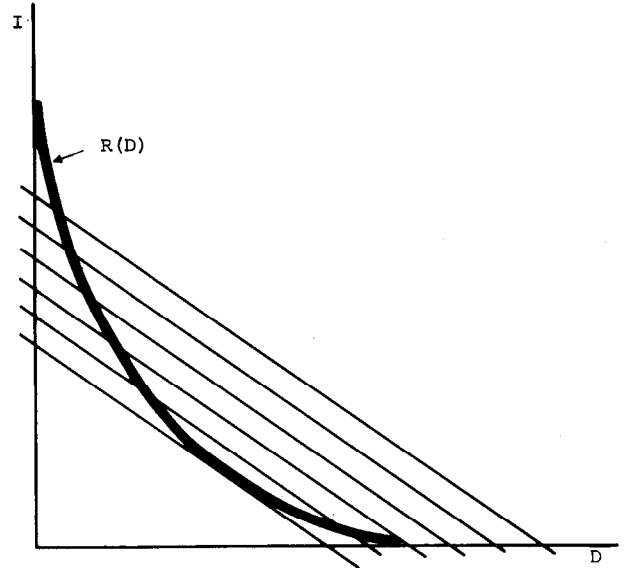


Fig. 2.  $I$ - $D$  plane.

$$Q_{k|j} = \frac{q_k \exp(s\rho_{jk})}{\sum_k q_k \exp(s\rho_{jk})}$$

is that  $q$  satisfy

$$c_k = \sum_j p_j \frac{\exp(s\rho_{jk})}{\sum_k q_k \exp(s\rho_{jk})} = 1, \quad q_k \neq 0$$

$$c_k = \sum_j p_j \frac{\exp(s\rho_{jk})}{\sum_k q_k \exp(s\rho_{jk})} \leq 1, \quad q_k = 0.$$

For a proof, see Berger or Gallager.

The major theorem of this section is the following.

*Theorem 6:* Let the parameter  $s < 0$  be given. Let  $q^0$  be any probability vector such that all components are non-zero. Let  $q^{r+1}$  be given in terms of  $q^r$  by

$$q_k^{r+1} = q_k^r \sum_j \frac{p_j A_{jk}}{\sum_k A_{jk} q_k^r},$$

where  $A_{jk} = \exp(s\rho_{jk})$ . Then,

$$D(Q(q^r)) \rightarrow D_s, \quad \text{as } r \rightarrow \infty$$

$$I(p, Q(q^r)) \rightarrow R(D_s), \quad \text{as } r \rightarrow \infty,$$

where  $(D_s, R(D_s))$  is a point on the  $R(D)$  curve parameterized by  $s$ .

*Proof:* Theorem 4 can be used to provide the first part of the proof. The following proof will, however, bring out the geometrical role of the parameter  $s$ .

For any probability vector  $q$ , recall that  $Q(q)$  is given by

$$Q_{k|j}(q) = \frac{A_{jk} q_k}{\sum_k A_{jk} q_k}.$$

Consider the  $I$ - $D$  plane of Fig. 2. For any probability vector  $q$ , let  $V(q) = I(q) - sD(q)$ , where  $I(q) = I(p, Q(q))$ . Then  $V(q)$  is the value at which a line of slope  $s$  through the point  $(I(q), D(q))$  intercepts the  $I$  axis. The point in the  $R(D)$  curve parameterized by  $s$  has a tangent that is parallel

to every such line of slope  $s$ , and lies beneath them. We will show that the sequential values  $V(q^r)$  are strictly decreasing unless  $(I(q), D(q))$  is a point on the  $R(D)$  curve in which case  $V(q)$  is stationary.

Let

$$Q^{r+1} = Q(q^r).$$

Then

$$Q_{k|j}^{r+1} = \frac{A_{jk} q_k^r}{\sum_k A_{jk} q_k^r}$$

and

$$q_k^{r+1} = \sum_j p_j Q_{k|j}^{r+1}.$$

Then

$$V(q^{r+1})$$

$$\begin{aligned} &= \sum_j \sum_k p_j Q_{k|j}^{r+1} \log \frac{Q_{k|j}^{r+1}}{q_k^{r+1}} - s \sum_j \sum_k p_j Q_{k|j}^{r+1} \rho_{jk} \\ &= \sum_j \sum_k p_j Q_{k|j}^{r+1} \log \frac{A_{jk} q_k^r}{q_k^{r+1} \sum_k A_{jk} q_k^r} - \sum_j \sum_k p_j Q_{k|j}^{r+1} \log A_{jk} \\ &= -\sum_j \sum_k p_j Q_{k|j}^{r+1} \log \sum_k A_{jk} q_k^r + \sum_j \sum_k p_j Q_{k|j}^{r+1} \log \frac{q_k^r}{q_k^{r+1}} \\ &= -\sum_j \sum_k p_j Q_{k|j}^r \log \sum_k A_{jk} q_k^r + \sum_k q_k^{r+1} \log \frac{q_k^r}{q_k^{r+1}}. \end{aligned}$$

Now let

$$W(q^r) = V(q^r) - V(q^{r+1})$$

so that

$$\begin{aligned} W(q^r) &= \sum_j \sum_k p_j Q_{k|j}^r \log \frac{Q_{k|j}^r \sum_k A_{jk} q_k^r}{q_k^r A_{jk}} + \sum_k q_k^{r+1} \log \frac{q_k^r}{q_k^{r+1}} \\ &\geq \sum_j \sum_k p_j Q_{k|j}^r \left[ 1 - \frac{q_k^r A_{jk}}{Q_{k|j}^r \sum_k A_{jk} q_k^r} \right] \\ &\quad + \sum_k q_k^{r+1} \left[ 1 - \frac{q_k^r}{q_k^{r+1}} \right] = 0 + 0 = 0 \end{aligned}$$

with strict inequality unless  $q_k^r = q_k^{r+1} \forall k$ .

Thus,  $V(q^r)$  is nonincreasing and is strictly decreasing unless

$$q_k^r = q_k^r \sum_j p_j \frac{A_{jk}}{\sum_k A_{jk} q_k^r},$$

which is just the first condition of Theorem 5. Since  $V(q^r)$  is decreasing and is bounded below by  $R(D) - sD$ , it must converge to some number  $V^\infty$ . We now argue as in the proof of Theorem 3 to show that  $V^\infty = R(D) - sD$ . That is, by the Bolzano-Weierstrass Theorem, the sequence  $q^r$  has a limit point  $q^*$  and by continuity of  $V(q)$  this limit point satisfies

$$q_k^* = q_k^* \sum_j p_j \frac{A_{jk}}{\sum_k A_{jk} q_k^*}.$$

In addition, this limit point must satisfy the second of the

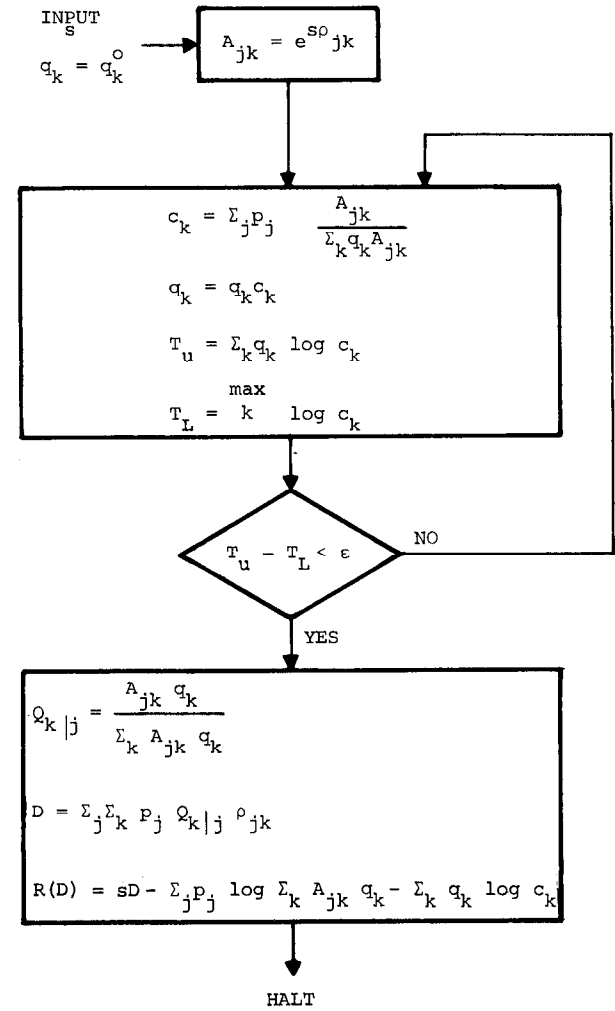


Fig. 3. Rate-distortion algorithm.

Kuhn-Tucker conditions since otherwise convergence could not occur. This completes the proof of the theorem.

The application of this theorem to the numerical computation of rate-distortion functions is illustrated in Fig. 3. In order to estimate the accuracy after any finite number of steps, the following theorem is employed.

**Theorem 7:** Let the parameter  $s \leq 0$  be given and let  $A_{jk} = \exp(s \rho_{jk})$ . Suppose  $q$  is any output probability vector and let

$$c_k = \sum_j p_j \frac{A_{jk}}{\sum_k A_{jk} q_k}.$$

Then at the point

$$D = \sum_j \sum_k p_j \frac{A_{jk} q_k}{\sum_k A_{jk} q_k} \rho_{jk}$$

we have

a)

$$R(D) \leq sD - \sum_j p_j \log \sum_k A_{jk} q_k - \sum_k q_k c_k \log c_k$$

b)

$$R(D) \geq sD - \sum_j p_j \log \sum_k A_{jk} q_k - \max_k \log c_k$$

*Proof:*

a)

$$Q_{k|j} = \frac{A_{jk}q_k}{\sum_k A_{jk}q_k}$$

is a transition matrix giving distortion  $D$ . Hence

$$\begin{aligned} R(D) &\leq I(p, Q) = \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} \\ &= \sum_j \sum_k p_j Q_{k|j} \log \frac{A_{jk}q_k}{(\sum_k A_{jk}q_k)(\sum_j p_j Q_{k|j})} \\ &= sD - \sum_j p_j \log \sum_k A_{jk}q_k - \sum_k q_k c_k \log c_k. \end{aligned}$$

b) A lower bound theorem for rate-distortion functions states that

$$R(D) \geq sD + \sum_j p_j \log \lambda_j,$$

where  $\lambda_j$  is any vector such that

$$\sum_j p_j \lambda_j A_{jk} \leq 1$$

(see Berger or Gallager). Let

$$c_{\max} = \max_k \sum_j p_j \frac{A_{jk}}{\sum_k A_{jk}q_k}$$

and let

$$\lambda_j = (c_{\max} \sum_k A_{jk}q_k)^{-1}.$$

Then

$$\sum_j p_j \lambda_j A_{jk} \leq 1$$

and

$$R(D) \geq sD - \sum_j p_j \log \sum_k A_{jk}q_k - \max_k \log c_k.$$

#### IV. CAPACITY OF CONSTRAINED DISCRETE CHANNELS

Many channels have an associated expense of using each channel letter. A common example is the power associated with each output symbol. A constrained discrete channel is a discrete channel with the requirement that the average expense be less than or equal to some specified number  $E$ .

Although capacity at an expense  $E$  has been investigated in the past, and occasionally the function  $C(E)$  has been determined, there does not seem to have been developed any formalization of the theory of  $C(E)$  functions. This formalization is straightforward and is provided in the Appendix.

A vector  $e_j$  is specified, where  $e_j$  is called the expense of using the  $j$ th input letter. The capacity at expense  $E$  is then defined as

$$C(E) = \max_{p \in P_E} \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} = \max_{p \in P_E} I(p, Q),$$

where

$$P_E = \{p \in P^n: \sum_j p_j e_j \leq E\}.$$

As discussed in the Appendix, this can be rewritten parametrically as

$$C(E) = \max_p \left[ \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} - s(\sum_j p_j e_j - E) \right],$$

where

$$E = \sum_j p_j^* e_j$$

and  $p^*$  achieves the above maximum.

The maximization is now over all input probability vectors  $p$ . The generalization of Theorem 1 is the following.

*Theorem 8:* Let

$$J(p, Q, P) = \sum_j \sum_k p_j Q_{k|j} \log \frac{P_{j|k}}{p_j} - s \sum_j p_j e_j.$$

Then

a)

$$C(E) = sE + \max_p \max_P J(p, Q, P),$$

where

$$E = \sum_j p_j^* e_j$$

and  $p^*$  achieves the above maximum.

b) For fixed  $p$ ,  $J(p, Q, P)$  is maximized by

$$P_{j|k} = \frac{p_j Q_{k|j}}{\sum_j p_j Q_{k|j}}.$$

c) For fixed  $P$ ,  $J(p, Q, P)$  is maximized by

$$p_j = \frac{\exp(\sum_k Q_{k|j} \log P_{j|k} - s e_j)}{\sum_j \exp(\sum_k Q_{k|j} \log P_{j|k} - s e_j)}.$$

*Proof:* The proof is essentially the same as that of Theorem 1.

*Corollary 8:* If  $p$  achieves capacity at expense  $E$ , then for some  $s \in [0, \infty]$

$$p_j = \frac{p_j \exp\left(\sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} - s e_j\right)}{\sum_j p_j \exp\left(\sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} - s e_j\right)}.$$

*Proof:* This is just the simultaneous satisfaction of parts b) and c).

*Corollary 9:* A parametric solution in terms of  $s$  is

$$C(E_s) = sE_s + \max_P [\log \sum_j \exp(\sum_k Q_{k|j} \log P_{j|k} - s e_j)]$$

$$E_s = \sum_j e_j \frac{\exp(\sum_k Q_{k|j} \log P_{j|k}^* - s e_j)}{\sum_j \exp(\sum_k Q_{k|j} \log P_{j|k}^* - s e_j)},$$

where  $P^*$  achieves the maximum.

*Corollary 10:*

$$C(E) = \min_{s \in [0, \infty]} \max_P [sE + \log \sum_j \exp(\sum_k Q_{k|j} \log P_{j|k} - s e_j)].$$

*Proof:* Let  $p^*$  achieve  $C(E)$  and let  $E_s$  be the expense parameterized by  $s$ . Then

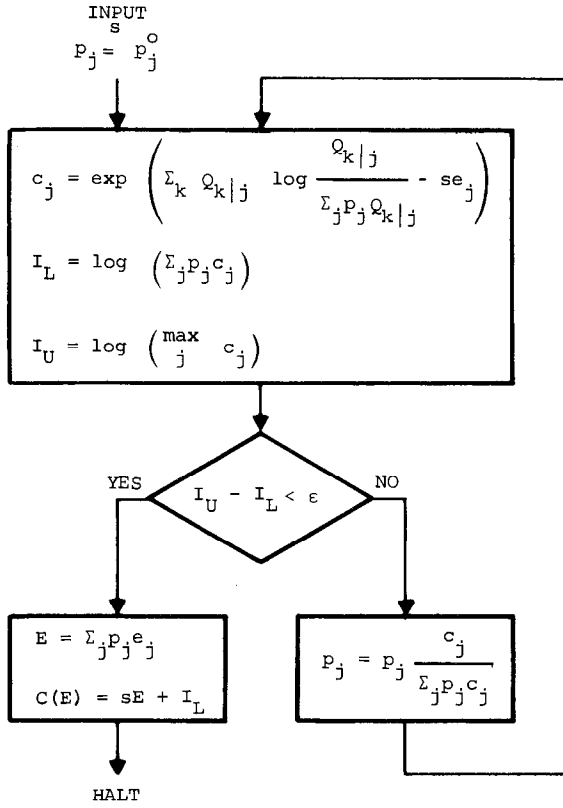


Fig. 4. Constrained capacity algorithm.

$$\begin{aligned}
 C(E) &= \max_p [sE + \log \sum_j \exp(\sum_k Q_{k|j} \log P_{j|k} - se_j)] \\
 &= C(E) - sE - \max_p [sE_s + \log \sum_j \exp(\sum_k Q_{k|j} \log P_{j|k} - se_j)] + sE_s \\
 &= C(E) - sE + C(E_s) + sE_s \\
 &= \sum_{jk} p_j^* Q_{k|j} \log \frac{Q_{k|j} \exp(-se_j)}{\sum_j p_j^* Q_{k|j}} \\
 &\quad - \max_p \sum_{jk} p_j Q_{k|j} \log \frac{Q_{k|j} \exp(-se_j)}{\sum_j p_j Q_{k|j}} \\
 &\leq 0.
 \end{aligned}$$

**Theorem 9:** A vector  $p \in \mathbf{P}^n$  achieves capacity at some expense  $E_s$  parameterized by  $s$  for the channel with transition matrix  $Q$  and expense vector  $e$  if and only if there exists a number  $V$  such that

$$\begin{aligned}
 \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} - se_j &= V, & p_j \neq 0 \\
 \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} - se_j &\leq V, & p_j = 0.
 \end{aligned}$$

*Proof:* The proof is lengthy but only trivially different from the proof of Theorem 2. The reader can readily modify any published proof of Theorem 2.

**Theorem 10:** Let  $s \in [0, \infty]$  be given, and for any  $p \in \mathbf{P}^n$  let

$$c_j(p) = \exp \left( \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} - se_j \right).$$

Then if  $p^0$  is any element of  $\mathbf{P}^n$  with all components strictly positive, the sequence of probability vectors defined by

$$p_j^{r+1} = p_j^r \frac{c_j^r}{\sum_j p_j^r c_j^r}$$

is such that

$$\begin{aligned}
 I(p^r, Q) &\rightarrow C(E_s), & \text{as } r \rightarrow \infty \\
 e(p^r) &\rightarrow E_{s^2}, & \text{as } r \rightarrow \infty,
 \end{aligned}$$

where  $E_s$  is the expense of the point parameterized by  $s$ .

*Proof:* Let

$$V(p) = I(p) - se(p) = \sum_j p_j \log c_j$$

and show that  $V(p)$  is increasing. Let

$$\begin{aligned}
 W(p) &= V(p^{r+1}) - V(p^r) \\
 &= \sum_j p_j^r \frac{c_j^r}{\sum_j p_j^r c_j^r} \log c_j^{r+1} - \sum_i p_i^r \log c_i^r \\
 &= \frac{1}{\sum_j p_j^r c_j^r} [\sum_i \sum_j p_i^r p_j^r c_j^r \log c_j^{r+1} \\
 &\quad - \sum_i \sum_j p_i^r p_j^r c_j^r \log c_i^r] \\
 &= \frac{1}{\sum_j p_j^r c_j^r} \left[ \sum_i p_i^r \sum_j p_j^r c_j^r \log \frac{c_j^{r+1}}{c_i^r} \right] \\
 &\geq 1 - \sum_j p_j^r \frac{c_j^r}{c_j^{r+1}}
 \end{aligned}$$

with equality iff

$$c_j^{r+1}/c_i^r = 1 \quad \forall i, j \ni p_i \neq 0 \neq p_j.$$

We now substitute the defining equation for  $c_j$  and apply Jensen's inequality.

$$\begin{aligned}
 W(p^r) &\geq 1 - \sum_j p_j^r \exp \sum_k Q_{k|j} \log \frac{\sum_j p_j^{r+1} Q_{k|j}}{\sum_j p_j^r Q_{k|j}} \\
 &\geq 1 - \sum_j p_j^r \sum_k Q_{k|j} \exp \log \frac{\sum_j p_j^{r+1} Q_{k|j}}{\sum_j p_j^r Q_{k|j}}.
 \end{aligned}$$

Therefore

$$W(p^r) \geq 1 - \sum_k \sum_j p_j^{r+1} Q_{k|j} = 0.$$

Thus,  $V(p)$  is increasing; moreover,  $V(p)$  is strictly increasing unless

$$c_j^{r+1} = c_i^r, \quad \forall i, j \ni p_i \neq 0 \neq p_j,$$

which condition reduces to the first condition of Theorem 9. We now argue as in the proof of Theorem 3 to show that  $V(p)$  converges to  $C(E) - sE$ . That is, by the Bolzano-Weierstrass Theorem,  $\{p^r\}$  has a limit point and by continuity it must satisfy the above Kuhn-Tucker condition. In addition, this limit point must satisfy the second of the Kuhn-Tucker conditions since otherwise convergence could not occur.

A flow diagram for the algorithm of Theorem 10 is shown in Fig. 4.



The following theorem provides a termination for this algorithm.

*Theorem 11:* Let the parameter  $s$  be given. Suppose  $p$  is any probability vector, and let

$$c_j = \exp \left( \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} - s e_j \right).$$

Then, at the point

$$E = \sum_j p_j c_j$$

a)

$$C(E) \geq sE + \log \sum_j p_j c_j$$

b)

$$C(E) \leq sE + \log \max_j c_j.$$

*Proof:*

a)  $p$  is a probability vector giving expense  $E$ . Hence

$$\begin{aligned} C(E) &\geq I(p, Q) = \sum_j p_j \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} \\ &= \sum_j p_j \log c_j + s e_j = sE + \sum_j p_j \log c_j. \end{aligned}$$

b) Suppose  $p^*$  achieves capacity at expense parameterized by  $s$ . Then by Corollary 10,

$$C(E) \leq sE + \log \sum_j p_j^* \exp \left( \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j^* Q_{k|j}} - s e_j \right).$$

Hence

$$\begin{aligned} C(E) - (sE + \max_j c_j) &\leq \log \sum_j p_j^* \exp \left( \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j^* Q_{k|j}} - s e_j - \max_j \log c_j \right) \\ &\leq \log \sum_j p_j^* \exp \left( \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j^* Q_{k|j}} - s e_j - \log c_j \right). \end{aligned}$$

We now use

$$\log c_j + s e_j = \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}}$$

so that

$$\begin{aligned} C(E) - (sE + \max_j c_j) &\leq \log \sum_j p_j^* \exp \sum_k Q_{k|j} \log \frac{\sum_j p_j Q_{k|j}}{\sum_j p_j^* Q_{k|j}} \\ &\leq \log \sum_j p_j^* \sum_k Q_{k|j} \exp \log \frac{\sum_j p_j Q_{k|j}}{\sum_j p_j^* Q_{k|j}} \\ &\quad \text{(Jensen's inequality)} \\ &= \log \sum_k \sum_j p_j Q_{k|j} = 0. \end{aligned}$$

## V. CONTINUOUS CHANNEL AND SOURCE ALPHABETS

The discussion of the preceding sections has been confined to discrete channels and sources. If we turn attention to channels or sources that are described by probability

densities, then the earlier discussion can be mimicked in order to provide the analogous theory for continuous probability distributions.

We shall not develop this continuous distribution theory in detail here, both because this would be largely a repetition of the discrete case and because a detailed treatment is available elsewhere [12]. However, several comments will be made to indicate the necessary modifications.

Suppose for any input  $x$ ,  $Q(y/x)$  is a probability density function describing the channel. Capacity is defined as

$$C(E) = \sup_{p(x) \in P_E} \iint p(x) Q(y/x) \log \frac{Q(y/x)}{\int Q(y/x) p(x) dx} dx dy,$$

where

$$P_E = \left\{ p: \mathcal{R} \rightarrow \mathcal{R} \mid \int p(x) dx = 1, p(x) \geq 0, \int p(x) e(x) dx \leq E \right\}.$$

Rate-distortion functions are similarly defined as an infimum of a mutual information over a space of conditional probability distributions.

The use of the supremum and infimum suggest that, in general, these are not actually achieved by any continuous probability distribution (e.g., convergence is to a discrete distribution) so that Kuhn-Tucker-like conditions on the extremizing probability distribution may be vacuously true. However, these conditions can nonetheless be stated and are useful for recognizing points that do not achieve the solution.

The search for extremizing probability distributions is now a problem in the calculus of variations with constraints, but otherwise closely follows the discrete case. The continuous versions of Theorems 6 and 10 can be stated. However, since the extremum might not be achieved, the proof cannot assert the existence of a limiting distribution. The proof must be modified to show that any point below the supremum (respectively above the infimum) cannot be a limit point.

## VI. MULTIPLE CONSTRAINTS

Some channels may have more than one constraint specified simultaneously. The most common example is a continuous channel that is constrained both in peak power and in average power. It is straightforward to generalize capacity-expense theory to handle this situation. The basic definition for the discrete channel is as follows

$$C(E^1, E^2) = \max_{p \in P_{E^1 E^2}} \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}},$$

where

$$P_{E^1 E^2} = \{p \in P^n: \sum_j p_j e_j^1 \leq E^1 \text{ and } \sum_j p_j e_j^2 \leq E^2\}.$$

The generalization of Theorem 10 is the following.

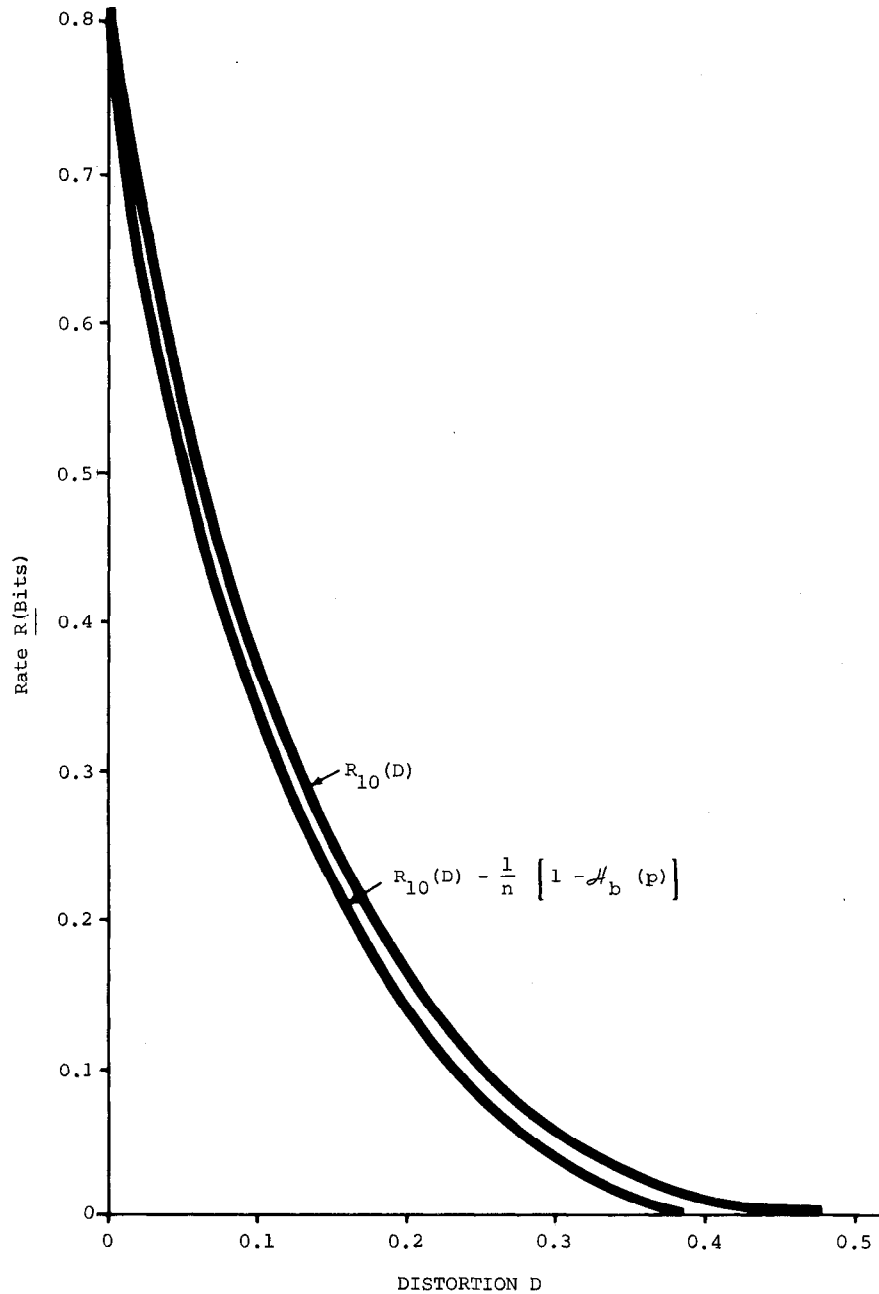


Fig. 5. Upper and lower bounds for the rate-distortion function of a binary symmetric Markov source  $p = 0.25$ .

**Theorem 13:** Let  $(s_1, s_2) \in [0, \infty] \times [0, \infty]$  be given and for any  $p \in \mathbf{P}^n$  let

$$c_j(p) = \exp \left( \sum_k Q_{k|j} \log \frac{Q_{k|j}}{\sum p_j Q_{k|j}} - s_1 e_j^1 - s_2 e_j^2 \right).$$

Then, if  $p^0$  is any element of  $\mathbf{P}^n$  with all components strictly positive, the sequence of probability vectors defined by

$$p_j^{r+1} = p_j^r \frac{c_j^r}{\sum p_j^r c_j^r}$$

is such that

$$I(p^r, Q) \rightarrow C(E_{s_1}^1, E_{s_2}^2), \quad \text{as } r \rightarrow \infty$$

$$e^1(p^r) \rightarrow E_{s_1}^1, \quad \text{as } r \rightarrow \infty$$

$$e^2(p^r) \rightarrow E_{s_2}^2, \quad \text{as } r \rightarrow \infty,$$

where  $C(E_{s_1}^1, E_{s_2}^2)$  is a point on the capacity-expense surface parameterized by  $(s_1, s_2)$ .

This theorem is offered without proof.

The analogous situation for rate-distortion functions can be considered. Thus, it may be desired that two (or more) separate definitions of distortion be satisfied [8]. One situation where this would occur is if the reproduced data is to be made available to two different users with different

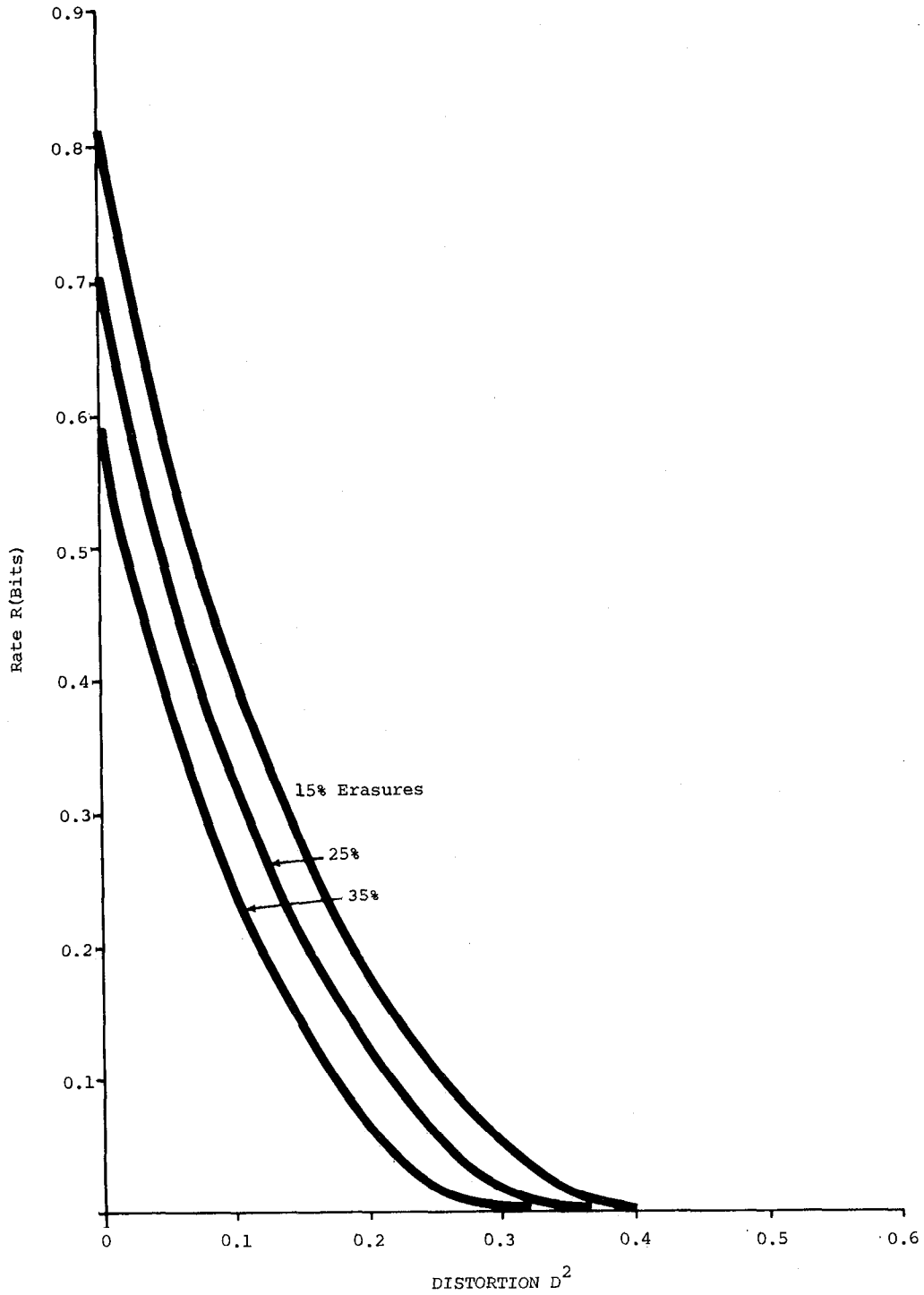


Fig. 6. Rate-distortion function of binary symmetric source with erasure option.

applications in mind. The appropriate definition is as follows

$$R(D^1, D^2) = \min_{Q \in Q_{D^1 D^2}} \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}},$$

where

$$Q_{D^1 D^2} = \{Q: d^1(Q) \leq D^1, d^2(Q) \leq D^2\}$$

$$d^1(Q) = \sum_j \sum_k p_j Q_{k|j} \rho_{jk}^1$$

$$d^2(Q) = \sum_j \sum_k p_j Q_{k|j} \rho_{jk}^2.$$

The generalization of Theorem 6 is the following.

*Theorem 14:* Let  $s_1 \leq 0, s_2 \leq 0$  be given. Let  $q^0$  be any output probability vector such that all components are non-zero. Let  $q^{r+1}$  be given in terms of  $q^r$  by

$$q_k^{r+1} = q_k^r \sum_j \frac{p_j A_{jk}}{\sum_k A_{jk} q_k^r}$$

where

$$A_{jk} = \exp(s_1 \rho_{jk}^1 + s_2 \rho_{jk}^2).$$

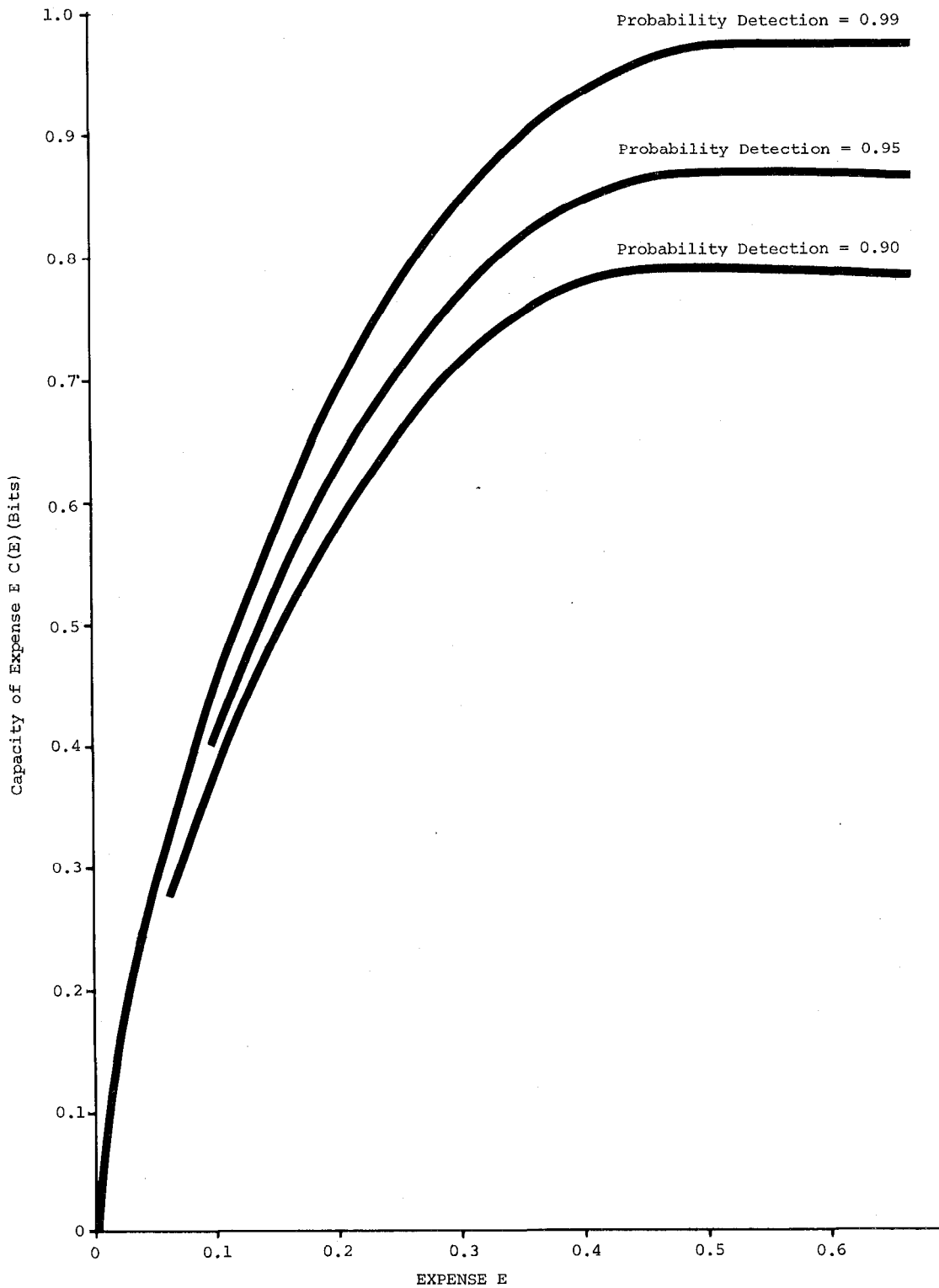


Fig. 7. Capacity-expense function for binary asymmetric channel with expense equal to percent of ones transmitted. False-alarm probability = 0.001.

Then

$$\begin{aligned}
 d^1(q^r) &\rightarrow D_{s_1}^1, & \text{as } r \rightarrow \infty \\
 d^2(q^r) &\rightarrow D_{s_2}^2, & \text{as } r \rightarrow \infty \\
 I(q^r) &\rightarrow R(D_{s_1}^1, D_{s_2}^2), & \text{as } r \rightarrow \infty.
 \end{aligned}$$

This theorem is offered without proof.

## VII. EXAMPLES

A long-standing problem in information theory is the determination of the rate-distortion function for a binary symmetric Markov source [7]–[9]. Gray [10] has recently solved this problem for a range of small  $D$ , but the problem for arbitrary  $D$  is unsolved. The rate-distortion function

for a source with memory is defined as

$$R(D) = \liminf R_n(D),$$

where  $R_n(D)$  is the rate-distortion function of a source whose alphabet is the set of words of length  $n$  with probabilities assigned to these words by the Markov source starting in an equiprobable state.

The algorithm of Theorem 6 has been used to calculate  $R_{10}(D)$  as shown in Fig. 5. Also shown is a lower bound to  $R(D)$  based on a recent theorem of Wyner and Ziv [11]. This theorem states that

$$R(D) \geq R_n(D) + H - \frac{1}{n} H(p_n),$$

where  $H$  is the source entropy rate and  $H(p_n)$  is the entropy of the set of  $n$ -words. For the binary symmetric Markov case, this becomes

$$R(D) \geq R_n(D) - \frac{1}{n} [p \log p + (1-p) \log (1-p) + 1],$$

where  $p$  is the transition probability.

Tighter bounds can be obtained by calculating  $R_n(D)$  for  $n > 10$ . However, the tightness is improving as  $1/n$  while the computations increase exponentially. Computation to an accuracy of  $10^{-3}$  bits of all  $R_n(D)$  curves from  $n = 2$  to  $n = 10$  by taking 9 points per curve required 12 min of execution time on the IBM 360 model 65.

The second example is a multiple-distortion problem. A memoryless source produces equiprobable i.i.d. outputs from a binary alphabet. In order to facilitate compression, a user agrees to allow a certain percentage  $D_1$  of erasures. Of the unerased data, he requires at most a percentage  $D_2$  be in error. Thus, the relevant distortion matrices are

$$\rho_{jk}^1 = \begin{vmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \end{vmatrix} \quad \rho_{jk}^2 = \begin{vmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \end{vmatrix}.$$

The numerical solution of the problem is shown in Fig. 6. These curves were prepared by computing  $R(D^1, D^2)$  for 1600 different values of  $(D^1, D^2)$  to an accuracy of  $10^{-3}$  bits. This required 83 s of computation time on an IBM 360 model 65.

The final example postulates the existence of a noisy binary channel, which transmits a one by the presence of a pulse and a zero by the absence of a pulse. The receiver is characterized by a probability of detection and by a probability of false alarm. The only design option available to the user is to conserve power by minimizing the percentage of ones used in a message. Fig. 7 shows the capacity-expense functions. These were computed by generating 300 points to an accuracy of  $10^{-3}$  bits, which required 5 s of computation time on an IBM 360 model 65.

#### ACKNOWLEDGMENT

The author wishes to acknowledge the advice and criticism of Prof. T. Berger of Cornell University.

#### APPENDIX

##### CAPACITY-EXPENSE FUNCTIONS

*Definition:* An *expense schedule* for a channel  $Q$  is a vector  $e$  whose  $j$ th component  $e_j$  is called the expense of using the  $j$ th channel input letter.

*Definition:* The *capacity at expense  $E$*  is

$$C(E) = \max_{p \in P_E} \sum_{jk} p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} = \max_{p \in P_E} I(p, Q),$$

where

$$P_E = \{p \in P^n : \sum_j p_j e_j \leq E\}.$$

This is well defined if  $P_E$  is nonempty since  $P_E$  is compact and hence  $I(p, Q)$  attains its maximum on  $P_E$ .

*Remark:* Without loss of generality, we can assume that  $E_{\min} = 0$  and  $C(E)$  exists for all  $E \geq E_{\min}$ . This is equivalent to assuming  $\min e_j = 0$ , which can be obtained by adding an appropriate constant to all  $e_j$ , thereby performing a simple horizontal translation of the  $C(E)$  graph.

*Remark:* If  $E' > E$  then  $P_{E'} \subset P_E$  and hence  $C(E)$  is a monotonic nondecreasing function.

*Theorem:*  $C(E)$  is a convex upward function. That is, given  $E'$ ,  $E''$ , and  $\lambda \in [0, 1]$ , then  $C(\lambda E' + (1-\lambda)E'') \geq \lambda C(E') + (1-\lambda)C(E'')$ .

*Proof:* Let  $p', p''$  achieve  $(E', C(E'))$ ,  $(E'', C(E''))$ , respectively. Let  $p^* = \lambda p' + \bar{\lambda} p''$ , where  $\bar{\lambda} = (1-\lambda)$ . Then

$$e(p^*) = \sum_j (\lambda p'_j + \bar{\lambda} p''_j) e_j = \lambda E' + \bar{\lambda} E''.$$

Hence  $p^* \in P_{\lambda E' + \bar{\lambda} E''}$  so that

$$\begin{aligned} C(\lambda E' + \bar{\lambda} E'') &\geq I(p^*, Q) \\ C(\lambda E' + \bar{\lambda} E'') - \lambda C(E') - \bar{\lambda} C(E'') &\geq I(p^*, Q) - \lambda I(p', Q) - \bar{\lambda} I(p'', Q) \\ &= \lambda \sum_{jk} p'_j Q_{k|j} \log \frac{\sum_j p_j Q_{k|j}}{\sum_j p'_j Q_{k|j}} + \bar{\lambda} \sum_{jk} p''_j Q_{k|j} \log \frac{\sum_j p_j Q_{k|j}}{\sum_j p''_j Q_{k|j}} \\ &\geq \lambda (\sum_{jk} p'_j Q_{k|j} - \sum_{jk} p_j^* Q_{k|j}) + \bar{\lambda} (\sum_{jk} p''_j Q_{k|j} - \sum_{jk} p_j^* Q_{k|j}) \\ &= 0. \end{aligned}$$

*Corollary:*  $C(E)$  is continuous except possibly at  $E = 0$ .

*Proof:*  $C(E)$  is convex and monotonic.

*Corollary:*

$$\lim_{E \rightarrow E_{\max}} C(E) = C,$$

where  $C$  is the channel capacity,

$$E_{\max} = \sum_j p_j^* e_j$$

and  $p^*$  achieves  $C$ .

*Proof:*  $C(E)$  is continuous.

*Corollary:*  $C(E)$  is strictly increasing in  $E < E_{\max}$ .

*Proof:*  $C(E)$  is convex.

*Corollary:* If  $E \leq E_{\max}$  then  $(E, C(E))$  is achieved by some  $P$  such that

$$e(P) = \sum_j p_j e_j = E.$$

*Proof:*  $C(E)$  is strictly increasing if  $E < E_{\max}$ .

*Theorem:* If  $p', p''$  both achieve the point  $(E, C(E))$ , then so

does

$$p = \lambda p' + \bar{\lambda} p'', \quad \forall \lambda \in [0,1].$$

*Proof:*

$$e(p) = \sum_j (\lambda p'_j + \bar{\lambda} p''_j) e_j = \lambda E + \bar{\lambda} E = E.$$

Therefore,  $p \in P_E$  so that

$$C(E) \geq I(p, Q) \geq \lambda I(p', Q) + \bar{\lambda} I(p'', Q) = C(E).$$

*Theorem:* Suppose  $0 \leq E \leq E_{\max}$ . Then  $C(E)$  can be expressed parametrically in terms of a parameter  $s \in [0, \infty]$  by

$$C(E_s) = sE_s + V_s$$

$$E_s = \sum_j p_j^* e_j,$$

where

$$V_s = \max_{p \in P} \sum_j p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} - s \sum_j p_j e_j$$

and  $p^*$  achieves this maximum.

*Proof:* Any such point  $(E_s, C(E_s))$  is clearly on the  $C(E)$  curve. It is only necessary to prove that every point on the  $C(E)$  curve can be so generated.

Since  $C(E)$  is concave, it has a derivative everywhere except possibly at a countable set of points and it has a left and a right derivative everywhere. Given the point  $E$ , let  $s$  be the left derivative of  $C(E)$  at  $E$ . Then for any  $E'$ , by convexity of  $C(E)$ ,

$$C(E') \leq C(E) + s(E' - E).$$

Now, the parameter  $s$  generates some point on  $C(E)$ . Let  $(E_s, C(E_s))$  be this point. Then

$$\begin{aligned} C(E_s) &= \max_{p \in P} \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} - s \sum_j p_j e_j + sE_s \\ &\geq \max_{p \in \{p: \sum_j p_j e_j = E\}} \left[ \sum_j \sum_k p_j Q_{k|j} \log \frac{Q_{k|j}}{\sum_j p_j Q_{k|j}} \right. \\ &\quad \left. - s \sum_j p_j e_j + sE_s \right] \end{aligned}$$

$$C(E_s) \geq C(E) - sE + sE_s.$$

Therefore,  $C(E_s) = C(E) + s(E_s - E)$  so that either  $E = E_s$  or they are connected by a straight line of slope  $s$ . In the latter case, the convexity of  $C(E)$  assures that every intermediate point on this straight line is also a point of  $C(E)$  and it is straightforward to verify that every point on this connecting line satisfies the parametric equation of the theorem.

*Corollary:* If  $C(E)$  is strictly concave in the neighborhood of some point, then the value of  $s$  that generates this point generates only this point.

*Corollary:* If  $s_1, s_2$  are the left and right derivatives at a point  $E$ , then  $s$  generates  $(E, C(E))$  if and only if  $s \in [s_1, s_2]$ .

## REFERENCES

- [1] C. Shannon, "A mathematical theory of communication," *Bell Syst. Tech. J.*, vol. 27, pp. 379-423, 623-656, 1948.
- [2] C. Shannon, "Coding theorems for a discrete source with a fidelity criterion," in *1959 IRE Nat. Conv. Rec.*, pt. 4, pp. 142-163.
- [3] R. G. Gallager, *Information Theory and Reliable Communication*. New York: Wiley, 1968.
- [4] F. Jelinek, *Probabilistic Information Theory*. New York: McGraw-Hill, 1968.
- [5] T. Berger, *Rate Distortion Theory: A Mathematical Basis for Data Compression*. Englewood Cliffs, N.J.: Prentice-Hall, 1971.
- [6] B. G. Haskell, "The computation and bounding of rate-distortion functions," *IEEE Trans. Inform. Theory*, vol. IT-15, pp. 525-531, Sept. 1969.
- [7] R. L. Dobrushin, "Mathematical problems in the Shannon theory of optimal coding of information," in *Proc. 4th Berkeley Symp. Mathematical Statistics and Probability*, vol. 1. Berkeley and Los Angeles: Univ. California Press, 1962, pp. 211-252.
- [8] T. J. Goblick, Jr., "Coding for a discrete information source with a distortion measure," Ph.D. dissertation, Dep. Elec. Eng., M.I.T., Cambridge, Mass, 1962.
- [9] T. Berger, "Nyquist's problem in data transmission," Ph.D. dissertation, Harvard Univ., Cambridge, Mass, 1965.
- [10] R. M. Gray, "Information rates of autoregressive sources," Ph.D. dissertation, Univ. Southern California, Los Angeles.
- [11] A. D. Wyner and J. Ziv, "Bounds on the rate-distortion function for stationary sources with memory," *IEEE Trans. Inform. Theory*, vol. IT-17, pp. 508-513, Sept. 1971.
- [12] R. E. Blahut, "Computation of information measures," Ph.D. dissertation, Cornell Univ., Ithaca, N.Y., 1972.
- [13] S. Arimoto, "An algorithm for computing the capacity of arbitrary discrete memoryless channels," *IEEE Trans. Inform. Theory*, vol. IT-18, pp. 14-20, Jan. 1972.

# Orthogonal Functionals of the Poisson Process

HISANAO OGURA

**Abstract**—In analogy to the orthogonal functionals of the Brownian-motion process developed by Wiener, Itô, and others, a theory of the orthogonal functionals of the Poisson process is presented making use of the concept of multivariate orthogonal polynomials. Following a brief discussion of Charlier polynomials of a single variable, multivariate

Charlier polynomials are introduced. An explicit representation as well as an orthogonality property are given. A multiple stochastic integral of a multivariate function with respect to the Poisson process, called the multiple Poisson-Wiener integral, is defined using the multivariate Charlier polynomials. A multiple Poisson-Wiener integral, which gives a polynomial functional of the Poisson process, is orthogonal to any other of different degree. Several explicit forms are given for the sake of application. It is shown that any nonlinear functional of the Poisson process with finite variance can be developed in terms of these orthogonal functionals, corresponding to the Cameron-Martin theorem in the case

Manuscript received April 9, 1970; revised December 10, 1971.  
The author was with the Department of Electronics, Kyoto University, Kyoto, Japan. He is now on leave at the Department of Electrical Engineering, University of Toronto, Toronto, Canada.