

A Rate-Distortion Theorem for Arbitrary Discrete Sources

Po-Ning Chen and Fady Alajaji

Abstract—A rate-distortion theorem for arbitrary (not necessarily stationary or ergodic) discrete-time finite-alphabet sources is given. This result, which provides the expression of the minimum ϵ -achievable fixed-length coding rate subject to a fidelity criterion, extends a recent data compression theorem by Steinberg and Verdú.

Index Terms—Arbitrary discrete sources, data compression, rate-distortion theory, Shannon theory.

I. INTRODUCTION

We consider the problem of source coding with a fidelity criterion for arbitrary (not necessarily stationary or ergodic) discrete-time finite-alphabet sources. We prove a general rate-distortion theorem by establishing the expression of the minimum ϵ -achievable block coding rate subject to a fidelity criterion.

In [3, Theorem 10, part a)], Steinberg and Verdú demonstrate a data compression theorem for arbitrary sources under the restriction that the probability of excessive distortion due to the achievable data compression codes is asymptotically equal to zero (cf. [3, Definitions 30 and 31]). In this work, we provide a variant of their result by relaxing the restriction on the probability of excessive distortion (cf. (3.1)).

II. PRELIMINARIES

Consider a random process \mathbf{X} defined by a sequence of finite-dimensional distributions [2]

$$\mathbf{X} \triangleq \{X^n = (X_1^{(n)}, \dots, X_n^{(n)})\}_{n=1}^{\infty}.$$

Let

$$\mathbf{Y} \triangleq \{Y^n = (Y_1^{(n)}, \dots, Y_n^{(n)})\}_{n=1}^{\infty}$$

be the corresponding output process induced by \mathbf{X} via the channel

$$\mathbf{W} \triangleq \{W_{Y^n|X^n} = P_{Y^n|X^n} : \mathcal{X}^n \rightarrow \mathcal{Y}^n\}_{n=1}^{\infty},$$

which is an arbitrary sequence of n -dimensional conditional distributions from \mathcal{X}^n to \mathcal{Y}^n , where \mathcal{X} and \mathcal{Y} are the input and output alphabets, respectively. We assume that \mathcal{X} to \mathcal{Y} are finite.

Definition 2.1 ([2]): Given a joint distribution $P_{X^n Y^n} = W_{Y^n|X^n} P_{X^n}$ on $\mathcal{X}^n \times \mathcal{Y}^n$ with marginals P_{X^n} and P_{Y^n} , the *information density* is defined by

$$i_{X^n Y^n}(a^n; b^n) = \log \frac{W_{Y^n|X^n}(b^n | a^n)}{P_{Y^n}(b^n)}.$$

Manuscript received April 15, 1996; revised January 20, 1998. The work of P.-N. Chen was supported in part by National Chiao-Tung University. The work of F. Alajaji was supported in part by Natural Sciences and Engineering Research Council of Canada (NSERC) under Grant OGP0183645. The material in this correspondence was presented in part at the International Symposium on Information Theory and Its Applications, Victoria, BC, Canada, September 1996.

P.-N. Chen is with the Department of Communication Engineering, National Chiao-Tung University, HsinChu, Taiwan, R.O.C.

F. Alajaji is with the Department of Mathematics and Statistics, Queen's University, Kingston, Ont. K7L 3N6, Canada.

Publisher Item Identifier S 0018-9448(98)03468-3.

Definition 2.2 ([2], [3]): The *sup-information rate* $\bar{I}(\mathbf{X}; \mathbf{Y})$ of the joint process $\mathbf{X}\mathbf{Y}$ is defined as the *limsup in probability*¹ of the sequence of normalized information densities $\frac{1}{n} i_{X^n Y^n}(X^n; Y^n)$.

Analogously, the *inf-information rate* $\underline{I}(\mathbf{X}; \mathbf{Y})$ between \mathbf{X} and \mathbf{Y} is defined as the *liminf in probability* of the sequence of normalized information densities $\frac{1}{n} i_{X^n Y^n}(X^n; Y^n)$.

When \mathbf{X} is equal to \mathbf{Y} , $\bar{I}(\mathbf{X}; \mathbf{X})$ (respectively, $\underline{I}(\mathbf{X}; \mathbf{X})$) is referred to as the *sup* (respectively, *inf*) *entropy rate* of \mathbf{X} and is denoted by $\bar{H}(\mathbf{X})$ (respectively, $\underline{H}(\mathbf{X})$).

Definition 2.3 ([2], [3]): Given a joint distribution $P_{X^n Y^n} = W_{Y^n|X^n} P_{X^n}$, the *conditional entropy density* is defined by

$$i_{Y^n|X^n}(b^n | a^n) = -\log W_{Y^n|X^n}(b^n | a^n).$$

The *conditional sup-entropy rate* $\bar{H}(\mathbf{Y} | \mathbf{X})$ of \mathbf{Y} given \mathbf{X} is defined as the *limsup in probability* of the sequence of normalized conditional entropy densities $\frac{1}{n} i_{Y^n|X^n}(Y^n | X^n)$.

Analogously, the *conditional inf-entropy rate* $\underline{H}(\mathbf{Y} | \mathbf{X})$ of \mathbf{Y} given \mathbf{X} is defined as the *liminf in probability* of $\frac{1}{n} i_{Y^n|X^n}(Y^n | X^n)$.

III. GENERAL DATA COMPRESSION THEOREM

Definition 3.1 (e.g., [1]): Given a finite source alphabet \mathcal{X} and a finite reproduction alphabet \mathcal{Y} , a block code for data compression of blocklength n and size M is a mapping $f_n(\cdot) : \mathcal{X}^n \rightarrow \mathcal{Y}^n$ that results in $\|f_n\| = M$ codewords of length n , where each codeword is a sequence of n reproducing letters.

Definition 3.2: A distortion measure $\rho_n(\cdot, \cdot)$ is a mapping

$$\rho_n : \mathcal{X}^n \times \mathcal{Y}^n \rightarrow \mathbb{R}^+ \triangleq [0, \infty).$$

We can view the distortion measure as the cost of representing a source n -tuple X^n by a reproduction n -tuple $f_n(X^n)$.

Definition 3.3: Let \mathbf{X} and $\{\rho_n(\cdot, \cdot)\}_{n \geq 1}$ be given. Let

$$\mathbf{f}(\mathbf{X}) \triangleq \{f_n(X^n)\}_{n=1}^{\infty}$$

denote a sequence of data compression codes for \mathbf{X} . The *distortion spectrum* $\underline{\lambda}_{\mathbf{X}\mathbf{f}(\mathbf{X})}(\theta)$ for $\mathbf{f}(\cdot)$ is defined by

$$\underline{\lambda}_{\mathbf{X}\mathbf{f}(\mathbf{X})}(\theta) \triangleq \liminf_{n \rightarrow \infty} \Pr \left\{ \frac{1}{n} \rho_n(X^n, f_n(X^n)) \leq \theta \right\}.$$

Definition 3.4: Fix $D > 0$ and $1 > \epsilon > 0$. R is an ϵ -achievable data compression rate at distortion D for a source \mathbf{X} if there exists a sequence of data compression codes $f_n(\cdot)$ with

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \|f_n\| = R$$

and

$$\sup\{\theta : \underline{\lambda}_{\mathbf{X}\mathbf{f}(\mathbf{X})}(\theta) \leq \epsilon\} \leq D. \quad (3.1)$$

¹If A_n is a sequence of random variables, then its *liminf in probability* is the largest extended real number α such that for all $\xi > 0$,

$$\lim_{n \rightarrow \infty} \Pr [A_n \leq \alpha - \xi] = 0.$$

Similarly, its *limsup in probability* is the smallest extended real number β such that for all $\xi > 0$ [2]

$$\lim_{n \rightarrow \infty} \Pr [A_n \geq \beta + \xi] = 0.$$

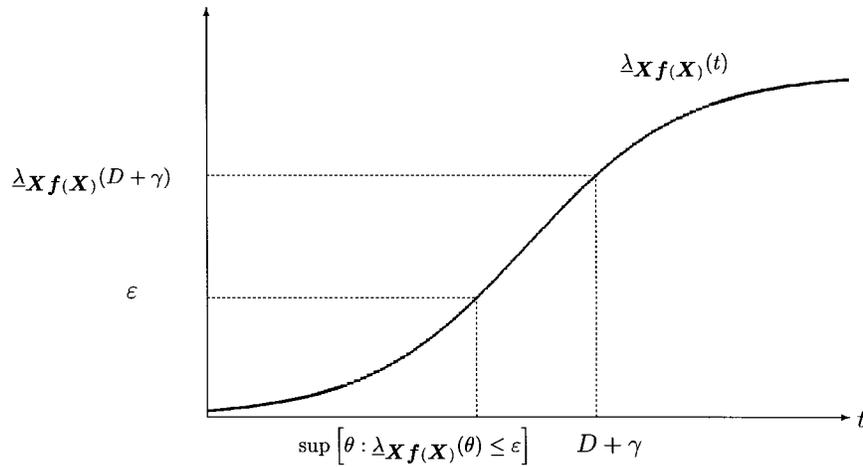


Fig. 1. $\lambda_{\mathbf{X}} f(\mathbf{X})(D + \gamma) > \varepsilon \Rightarrow \sup[\theta : \lambda_{\mathbf{X}} f(\mathbf{X})(\theta) \leq \varepsilon] \leq D + \gamma$.

Note that (3.1) is equivalent to stating that the limsup of the probability of excessive distortion (i.e., distortion larger than D) is smaller than $1 - \varepsilon$.

The infimum ε -achievable data compression rate at distortion D for \mathbf{X} is denoted by $T_\varepsilon(D, \mathbf{X})$.

Theorem 3.1 (General Data Compression Theorem): Fix $D > 0$ and $1 > \varepsilon > 0$. Let \mathbf{X} and $\{\rho_n(\cdot, \cdot)\}_{n \geq 1}$ be given. Then

$$T_\varepsilon(D, \mathbf{X}) = R_\varepsilon(D)$$

where

$$R_\varepsilon(D) \triangleq \inf_{\{\mathbf{W} : \sup_{\theta : \lambda_{\mathbf{X}\mathbf{Y}}(\theta) \leq D\}} \bar{I}(\mathbf{X}; \mathbf{Y})$$

where the infimum is taken over all conditional distributions $\mathbf{W} = \{P_{Y^n|X^n}\}_{n=1}^\infty$ for which the joint distribution $P_{\mathbf{X}\mathbf{Y}} = P_{\mathbf{X}}\mathbf{W}$ satisfies the distortion constraint.

Proof:

1) *Forward part (achievability):* Choose $\gamma > 0$. We will prove the existence of a sequence of data compression codes with

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \|f_n\| \leq R_\varepsilon(D) + 2\gamma$$

and

$$\sup[\theta : \lambda_{\mathbf{X}} f(\mathbf{X})(\theta) \leq \varepsilon] \leq D + \gamma.$$

Step 1: Let $\bar{\mathbf{W}}$ be the channel distribution achieving $R_\varepsilon(D)$, and let $P_{\bar{\mathbf{Y}}}$ be the \mathbf{Y} -marginal of $P_{\mathbf{X}}\bar{\mathbf{W}}$.

Step 2: Let $\bar{R} = R_\varepsilon(D) + 2\gamma$. Choose $M = e^{n\bar{R}}$ n -blocks independently according to $P_{\bar{\mathbf{Y}}}$, and denote the resulting random set by \mathcal{C}_n .

Step 3: For a given \mathcal{C}_n , we denote by $A(\mathcal{C}_n)$ the set of sequences $x^n \in \mathcal{X}^n$ such that there exists $y^n \in \mathcal{C}_n$ with

$$\frac{1}{n} \rho_n(x^n, y^n) \leq D + \gamma.$$

Step 4: Claim:

$$\limsup_{n \rightarrow \infty} E_{\bar{\mathbf{Y}}} [P_{X^n}(A^c(\mathcal{C}_n))] < 1 - \varepsilon.$$

The proof of this claim is provided in the Appendix. Therefore, there exists (a sequence of) \mathcal{C}_n^* such that

$$\limsup_{n \rightarrow \infty} P_{X^n}(A^c(\mathcal{C}_n^*)) < 1 - \varepsilon.$$

Step 5: Define a sequence of codes $\{f_n\}$ by

$$f_n(x^n) = \begin{cases} \arg \min_{y^n \in \mathcal{C}_n^*} \rho_n(x^n, y^n), & \text{if } x^n \in A(\mathcal{C}_n^*) \\ \underline{0}, & \text{otherwise} \end{cases}$$

where $\underline{0}$ is a fixed default n -tuple in \mathcal{Y}^n .

Then

$$\left\{ x^n \in \mathcal{X}^n : \frac{1}{n} \rho_n(x^n, f_n(x^n)) \leq D + \gamma \right\} \supset A(\mathcal{C}_n^*)$$

since $(\forall x^n \in A(\mathcal{C}_n^*))$ there exists $y^n \in \mathcal{C}_n^*$ such that $(1/n)\rho_n(x^n, y^n) \leq D + \gamma$, which by definition of f_n implies that $(1/n)\rho_n(x^n, f_n(x^n)) \leq D + \gamma$.

Step 6: Consequently,

$$\begin{aligned} & \lambda_{\mathbf{X}} f(\mathbf{X})(D + \gamma) \\ &= \liminf_{n \rightarrow \infty} P_{X^n} \left\{ x^n \in \mathcal{X}^n : \frac{1}{n} \rho_n(x^n, f(x^n)) \leq D + \gamma \right\} \\ &\geq \liminf_{n \rightarrow \infty} P_{X^n}(A(\mathcal{C}_n^*)) \\ &= 1 - \limsup_{n \rightarrow \infty} P_{X^n}(A^c(\mathcal{C}_n^*)) \\ &> \varepsilon. \end{aligned}$$

Hence

$$\sup[\theta : \lambda_{\mathbf{X}} f(\mathbf{X})(\theta) \leq \varepsilon] \leq D + \gamma$$

where the last step is clearly depicted in Fig. 1.

This proves the forward part.

2) *Converse part:* We show that for any sequence of encoders $\{f_n(\cdot)\}_{n=1}^\infty$, if

$$\sup[\theta : \lambda_{\mathbf{X}} f(\mathbf{X})(\theta) \leq \varepsilon] \leq D$$

then

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \|f_n\| \geq R_\varepsilon(D).$$

Let

$$\hat{W}^n(y^n | x^n) \triangleq \begin{cases} 1, & \text{if } y^n = f_n(x^n) \\ 0, & \text{otherwise.} \end{cases}$$

Let \hat{Y}^n denote the output corresponding to input X^n and channel \hat{W}^n . Then to evaluate the statistical properties of the random sequence $\{(1/n)\rho_n(X^n, f_n(X^n))\}_{n=1}^\infty$ under distribution

P_{X^n} is equivalent to evaluating those of the random sequence $\{(1/n)\rho_n(X^n, \hat{Y}^n)\}_{n=1}^\infty$ under distribution $P_{X^n} \hat{W}^n$. Therefore,

$$\begin{aligned} R_\varepsilon(D) &\triangleq \inf_{\{W: \sup_{\theta: \lambda_{\mathbf{X}\hat{\mathbf{Y}}}(\theta) \leq \varepsilon} D\}} \bar{I}(\mathbf{X}; \mathbf{Y}) \\ &\leq \bar{I}(\mathbf{X}; \hat{\mathbf{Y}}) \\ &\leq \bar{H}(\hat{\mathbf{Y}}) - \underline{H}(\hat{\mathbf{Y}} | \mathbf{X}) \\ &\leq \bar{H}(\hat{\mathbf{Y}}) \\ &\leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \|f_n\|, \end{aligned}$$

where the second inequality follows from [4, Theorem 8, property (d)] and the third inequality follows from the fact that $\underline{H}(\hat{\mathbf{Y}} | \mathbf{X}) \geq 0$. \square

APPENDIX

Claim (cf. Proof of Theorem 3.1):

$$\limsup_{n \rightarrow \infty} E_{\hat{\mathbf{Y}}^n} [P_{X^n}(A^c(C_n^*))] < 1 - \varepsilon.$$

Proof:

Step 1: Let

$$D^{(\varepsilon)} \triangleq \sup\{\theta : \lambda_{\mathbf{X}\hat{\mathbf{Y}}}(\theta) \leq \varepsilon\}.$$

Define

$$\begin{aligned} A_{n,\gamma}^{(\varepsilon)} &\triangleq \left\{ (x^n, y^n) : \frac{1}{n} \rho_n(x^n, y^n) \leq D^{(\varepsilon)} + \gamma \right. \\ &\quad \left. \text{and } \frac{1}{n} i_{X^n \hat{Y}^n}(x^n, y^n) \leq \bar{I}(\mathbf{X}; \hat{\mathbf{Y}}) + \gamma \right\}. \end{aligned}$$

Since

$$\liminf_{n \rightarrow \infty} \Pr \left(\mathcal{D} \triangleq \left\{ \frac{1}{n} \rho_n(X^n, \hat{Y}^n) \leq D^{(\varepsilon)} + \gamma \right\} \right) > \varepsilon$$

and

$$\begin{aligned} \liminf_{n \rightarrow \infty} \Pr \left(\mathcal{E} \triangleq \left\{ \frac{1}{n} i_{X^n \hat{Y}^n}(X^n; \hat{Y}^n) \right. \right. \\ \left. \left. \leq \bar{I}(\mathbf{X}; \hat{\mathbf{Y}}) + \gamma \right\} \right) = 1 \end{aligned}$$

we have

$$\begin{aligned} \liminf_{n \rightarrow \infty} \Pr(A_{n,\gamma}^{(\varepsilon)}) &= \liminf_{n \rightarrow \infty} \Pr(\mathcal{D} \cap \mathcal{E}) \\ &\geq \liminf_{n \rightarrow \infty} \Pr(\mathcal{D}) + \liminf_{n \rightarrow \infty} \Pr(\mathcal{E}) - 1 \\ &> \varepsilon + 1 - 1 = \varepsilon. \end{aligned}$$

Step 2: Let $K(x^n, y^n)$ be the indicator function of $A_{n,\gamma}^{(\varepsilon)}$

$$K(x^n, y^n) = \begin{cases} 1, & \text{if } (x^n, y^n) \in A_{n,\gamma}^{(\varepsilon)} \\ 0, & \text{otherwise.} \end{cases}$$

Step 3: By following a similar argument in [3, Eqs. (9)–(12)], we obtain

$$\begin{aligned} E_{\hat{\mathbf{Y}}^n} [P_{X^n}(A^c(C_n^*))] &= \sum_{C_n^*} P_{\hat{\mathbf{Y}}^n}(C_n^*) \sum_{x^n \notin A(C_n^*)} P_{X^n}(x^n) \\ &= \sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) \sum_{C_n^*: x^n \notin A(C_n^*)} P_{\hat{\mathbf{Y}}^n}(C_n^*) \end{aligned}$$

$$\begin{aligned} &= \sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) \left(1 - \sum_{y^n \in \mathcal{Y}^n} P_{\hat{\mathbf{Y}}^n}(y^n) K(x^n, y^n) \right)^M \\ &\leq \sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) \left(1 - e^{-n(\bar{I}(\mathbf{X}; \hat{\mathbf{Y}}) + \gamma)} \right. \\ &\quad \left. \times \sum_{y^n \in \mathcal{Y}^n} P_{\hat{\mathbf{Y}}^n | X^n}(y^n | x^n) K(x^n, y^n) \right)^M \\ &\leq 1 - \sum_{x^n \in \mathcal{X}^n} \sum_{y^n \in \mathcal{Y}^n} P_{X^n}(x^n) P_{\hat{\mathbf{Y}}^n | X^n}(x^n, y^n) K(x^n, y^n) \\ &\quad + \exp\{-e^{n(R - R_\varepsilon(D) - \gamma)}\}. \end{aligned}$$

Therefore,

$$\begin{aligned} \limsup_{n \rightarrow \infty} E_{\hat{\mathbf{Y}}^n} [P_{X^n}(A^c(C_n^*))] &\leq 1 - \liminf_{n \rightarrow \infty} \Pr(A_{n,\gamma}^{(\varepsilon)}) \\ &< 1 - \varepsilon. \end{aligned} \quad \square$$

ACKNOWLEDGMENT

The authors would like to thank Prof. S. Verdú for his valuable advice and encouragements.

REFERENCES

- [1] R. E. Blahut, *Principles and Practice of Information Theory*. Reading, MA: Addison Wesley, 1988.
- [2] T. S. Han and S. Verdú, "Approximation theory of output statistics," *IEEE Trans. Inform. Theory*, vol. 39, pp. 752–772, May 1993.
- [3] Y. Steinberg and S. Verdú, "Simulation of random processes and rate-distortion theory," *IEEE Trans. Inform. Theory*, vol. 42, pp. 63–86, Jan. 1996.
- [4] S. Verdú and T. S. Han, "A general formula for channel capacity," *IEEE Trans. Inform. Theory*, vol. 40, pp. 1147–1157, July 1994.

On One Useful Inequality in Testing of Hypotheses

Marat V. Burnashev

Abstract—A simple proof of one probabilistic inequality is presented.

Index Terms—Error probabilities, testing of hypotheses.

I. MAIN INEQUALITY

Let P and Q be two given probability measures on a measurable space $(\mathcal{X}, \mathcal{A})$. We consider testing of hypotheses P and Q using one observation. For an arbitrary decision rule, let α and β denote the two kinds of error probabilities. If both error probabilities have equal costs (or we want to minimize the maximum of them) then it is natural to investigate the minimal possible sum $\inf\{\alpha + \beta\}$ for the best decision rule.

Manuscript received July 10, 1997; revised November 24, 1997. This work was supported by the Russian Foundation for Fundamental Research under Grants N 95-01-00136a and INTAS-94-469.

The author is with the Institute for Problems of Information Transmission, Russian Academy of Sciences, 101447 Moscow, Russia.

Publisher Item Identifier S 0018-9448(98)03470-1.