

value. That is to say, the SNR always appears to increase until a 2.5-dB improvement limit is reached, which has been confirmed by experiments.

It is clear from the foregoing that the average estimate error in using average of logs depends on the statistics of the power estimates being averaged. The error ranges from 2.5 dB in the Gaussian noise case to 0-dB error for a noiseless sinusoid signal. A sinusoidal signal in noise falls within this 0- to 2.5-dB spread. Since the error is a function of the incoming waveform statistics it is not amenable to calibration except with difficulty. This estimation error is considered large for some applications but with many it is quite acceptable.

#### IV. CONCLUSIONS

The effect of averaging the logarithm of power, relative to power, for detection and for estimation of power has been shown. Detection with log-power leads to significant loss in performance only for a very large number of integrations. With many of today's narrowband processors, the number of integrations falls well below 20, a level where log-power can be considered to offer considerable processing advantages since only 0.2-dB detectability loss is encountered.

When estimating signals in white Gaussian noise by logarithms, the presented curves allow good estimates of both power and SNR. The errors for unknown signal and noise distributions are bounded, and with a single observation SNR of several decibels or more the estimation error may still be usable in a number of applications.

These two considerations make the log-power processing of signals very attractive and often offer advantages when compared to power processing.

#### REFERENCES

- [1] J. I. Marcum, "A statistical theory of target detection by pulsed radar," *IRE Trans. Inform. Theory*, vol. IT-6, pp. 259-267, Apr. 1960.
- [2] V. G. Hansen, "Post detection integration loss for logarithmic detectors," *IEEE Trans. Aerosp. Electron. Syst.* (Corresp.), vol. AES-8, pp. 386-388, May 1972.
- [3] B. A. Green, Jr., "Radar detection probability with logarithmic detectors," *IRE Trans. Inform. Theory*, vol. IT-4, pp. 50-52, Mar. 1958.
- [4] R. L. Hershey, "Analysis of the difference between log mean and mean log averaging," *J. Acoust. Soc. Amer.*, p. 1194, Apr. 1972.
- [5] H. Cox, "Linear versus logarithmic averaging," *J. Acoust. Soc. Amer.*, no. 39, pp. 688-690, 1966.
- [6] S. Mitchell, "Comment on linear versus logarithmic averaging," *J. Acoust. Soc. Amer.*, vol. 41, pp. 863-864, 1967.
- [7] C. H. Chen, "Signal-to-Noise ratios in logarithmic amplifiers," *Proc. IEEE (Lett.)*, vol. 57, pp. 1667-1668, Sept. 1969.
- [8] H. M. Musal, Jr., "Logarithmic compression of Rayleigh and Maxwell distributions," *Proc. IEEE (Lett.)*, vol. 57, pp. 1311-1313, July 1969.
- [9] I. Sugai and P. F. Christopher, "Comments on 'logarithmic compression of Rayleigh and Maxwell distributions,'" *Proc. IEEE (Lett.)*, vol. 58, pp. 263-264, Feb. 1970.
- [10] V. G. Hansen and H. R. Ward, "Detection performance of cell averaging LOG/CFAR receiver," *IEEE Trans. Aerosp. Electron. Syst.*, vol. AES-8, pp. 648-652, Sept. 1972.
- [11] J. Taylor and J. Mattern, "Receivers," *Radar Handbook*, M. Skolnik, Ed. New York: McGraw-Hill, 1970, ch. 5, pp. 5-29-5-36.
- [12] C. N. Pryor, "Calculation of minimum detectable signal for practical spectrum analyzers," Naval Ordnance Lab. Silver Spring, Md., Tech. Rep. NOLTR 71-92, Aug. 1971, p. 39.
- [13] H. Cramer, *Mathematical Methods of Statistics*. Princeton, N.J.: Princeton Univ. Press, 1946.
- [14] I. M. Ryzhik and I. S. Gradshteyn, *Tables of Integrals, Series, and Products*. New York: Academic Press, 1965.
- [15] J. R. Williams and G. G. Ricker, "Detection sensitivity performance of optimum Fourier receivers," *IEEE Trans. Audio Electroacoust.*, vol. UA-20, pp. 264-270, Oct. 1972.
- [16] J. V. DiFranco and W. L. Rubin, *Radar Detection*. Englewood Cliffs, N.J.: Prentice-Hall, 1968, p. 306.
- [17] L. Blake, "Prediction of radar range," *Radar Handbook*, M. I. Skolnik, Ed. New York: McGraw-Hill, 1970, ch. 2, pp. 2-26-2-29.
- [18] A. Papoulis, *Probability, Random Variables, and Stochastic Processes*. New York: McGraw-Hill, 1965, p. 250.
- [19] W. Grobner and N. Hofreiter, *Integraltafel*. Vienna and Innsbruck: Springer, 1950, no. 83, p. 81.
- [20] C. W. Helstrom, *Statistical Theory of Signal Detection*, 2nd Ed. New York: Pergamon, 1968, pp. 226-228.
- [21] J. R. Williams, "Fourier receiver sensitivity," in *Proc. 29th Navy Symp. Underwater Acoustics*, 1972, vol. II, pp. 503-514.

## Finite-Memory Algorithms for Estimating the Mean of a Gaussian Distribution

M. E. HELLMAN

**Abstract**—Let  $\{X_n\}_{n=1}^{\infty}$  be independent random variables, each having a  $\mathcal{N}(\mu, \sigma^2)$  distribution. If we try to estimate  $\mu$  with an  $m$ -state learning algorithm, then the minimum mean-squared error is bounded below by that obtained by the best  $m$ -level quantizer (which requires knowledge of  $\mu$ ). Here we show that this lower bound is tight. The results are easily extended to a number of other problems, such as estimating the mean  $\theta$  of a uniform distribution.

#### I. INTRODUCTION

The problem of estimating the mean  $\mu$  of a sequence of independent  $\mathcal{N}(\mu, \sigma^2)$  observations is a classic problem in statistics. However, until recently the effects of finite memory were never considered. Since in the real world any estimation algorithm is of necessity a finite-memory algorithm, the study of these effects is basic to the problem.

The model we shall use assumes that  $\{X_n\}_{n=1}^{\infty}$  are independent random variables, each having a  $\mathcal{N}(\mu, \sigma^2)$  distribution. The mean  $\mu$  has a known *a priori* distribution  $p(\mu)$ , but the distribution of the variance  $\sigma^2$  is unknown. Using  $m$ -state algorithms of the form

$$\begin{aligned} T_n &= f(T_{n-1}, X_n) \in \{1, 2, \dots, m\} \\ d_n &= d(T_n) \in \mathbb{R} \end{aligned} \quad (1)$$

for  $n = 1, 2, 3, \dots$ , we wish to minimize the limiting value of the mean-squared error

$$J = \lim_{L \rightarrow \infty} E \left\{ \frac{1}{L} \sum_{n=1}^L (\mu - d_n)^2 \right\}. \quad (2)$$

The crucial point in the formulation is the restriction that  $T_n$ , the value of our statistic at time  $n$ , take on one of only  $m$  possible values. For example a 10-bit memory corresponds to  $m = 2^{10} = 1024$ .

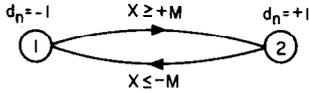
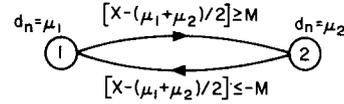
Of the recent work [9]-[22] on finite-memory learning algorithms, this problem is most closely related to Roberts and Tooley [18], Wagner [22], and Cover [23]. The first two papers deal with time-varying algorithms for estimating the mean of a distribution. In contrast, we are concerned with time-invariant algorithms (i.e., in (1) neither the state transition function  $f$  nor the decision function  $d$  depend on  $n$ ). Time-invariant rules are of greater use in practical applications because of their simpler structure. Cover applied Sagalowicz's results [24] on hypothesis testing to obtain solutions to certain time-invariant estimation problems. The reader is referred to [9]-[22], to the statistical literature [1]-[3], and to the literature on automata theory [4]-[8] for a more thorough history of the subject.

#### II. RELATION TO QUANTIZATION PROBLEM

In the quantization problem we are asked to find  $m$  real numbers  $\{\mu_i\}_{i=1}^m$ , and a mapping (quantizer)  $Q: \mathbb{R} \rightarrow \{\mu_i\}_{i=1}^m$ . Nature then generates  $\mu$  according to the *a priori* distribution  $p(\mu)$  and tells us the value of  $\mu$ . We must then represent  $\mu$  by

Manuscript received June 11, 1973; revised December 5, 1973. This work was supported by the National Science Foundation under Grant GK-33250. Portions of this correspondence were presented at the Third (Soviet) International Symposium on Information Theory, Tallinn, Estonia, USSR, June 18-23, 1973.

The author is with the Information Systems Laboratory, Department of Electrical Engineering, Stanford University, Stanford, Calif. 94305.

Fig. 1. Optimal algorithm for testing  $\mu = +1$  versus  $\mu = -1$ .Fig. 2. Optimal algorithm for testing  $\mu_1$  versus  $\mu_2$ .

$Q(\mu) \in \{\mu_i\}_{i=1}^m$ , incurring a loss  $[\mu - Q(\mu)]^2$ . The mean-squared error achieved by an  $m$ -state learning algorithm can be no smaller than the minimum mean-squared error obtainable by the optimum  $m$ -level quantizer. It will be strictly larger unless we use the optimal quantizer values  $\{\mu_i^*\}_{i=1}^m$  as our estimates (decisions) and ensure that as  $n \rightarrow \infty$  the memory is in its optimal state for the value of  $\mu$  chosen by nature. That is,  $\lim_{n \rightarrow \infty} P[d_n = Q^*(\mu)] = 1$  for almost all  $\mu$ , where  $Q^*$  denotes the optimal quantizer.

From the previous argument one might assume that this lower bound is not tight. However, as shown in the following, this intuitive guess is wrong.

As a starting point, consider the problem of testing between two hypotheses, where under  $H_0$  the observations have a  $\mathcal{N}(+1, 1)$  distribution, and under  $H_1$  they have a  $\mathcal{N}(-1, 1)$  distribution. As shown in [9], even a two-state memory can have as small an error probability as desired. The class of optimal algorithms is shown in Fig. 1.

If  $T_{n-1} = 1$  and  $X_n \geq M$ , then the memory transits to state 2. If  $T_{n-1} = 2$  and  $X_n \leq -M$ , then the memory transits to state 1. In state 1 the algorithm decides  $\mu = -1$ , while in state 2 it decides  $\mu = +1$ . By symmetry the error probability is the same under  $H_0$  as under  $H_1$ . Therefore let us assume  $H_0$  is true, so that each observation has mean  $+1$  and variance 1. If we choose  $M = 100$ , then transitions from state 1 to state 2 occur with very low probability since observations must be located at least  $99\sigma$  from the mean to cause such transitions. However, the probability of transitions from state 2 to state 1 is even lower since observations must be at least  $101\sigma$  from the mean to cause these transitions. The independence of the observations implies that the sequence of states occupied by the memory is a Markov chain. Thus it is only the ratio of the probabilities of transition which determines the steady state probabilities of being in state 1 and state 2.

Since the likelihood ratio  $l(x)$  is equal to  $\exp(2x)$ , the probability of a transition from state 1 to state 2 is approximately  $e^{200}$  times as large as the probability of a transition from state 2 to state 1. The error probability is therefore approximately  $e^{-200}$ . For other large values of  $M$ , similar reasoning shows that the error probability is approximately  $e^{-2M}$ . By choosing  $M$  sufficiently large any nonzero error rate can be achieved. Of course if the sample size  $N$  is not truly infinite, then  $M$  cannot be made too large or no transitions will occur. Recent work [19] has shown that in the large but finite sample size problem the optimal value of  $M$  is given by  $M_N^* = (2 \ln N)^{1/2} + 1 - \delta_N$ , where  $\delta_N \rightarrow 0$  as  $N \rightarrow \infty$ . The associated optimal error rate is

$$P^*(2, N) \sim \exp\{-2[(2 \ln N)^{1/2} + 1]\} \quad (3)$$

where  $\sim$  means that as  $N \rightarrow \infty$  the ratio of the two sides tends to 1. These results will prove useful when discussing the finite sample size estimation problem. The quantity  $(2 \ln N)^{1/2}$  comes from the fact that the maximum of  $N$  independent  $\mathcal{N}(0, 1)$  random variables tends to  $(2 \ln N)^{1/2}$  in probability [25].

Note that the variance of the observations does not matter when  $N = \infty$ . As  $M \rightarrow \infty$  the previous machine's error rate tends to zero no matter what the variance. Note also that if we change the problem so that we are testing between  $\mathcal{N}(\mu_1, \sigma_1^2)$

and  $\mathcal{N}(\mu_2, \sigma_2^2)$  distributions, with  $\mu_2 > \mu_1$ , then similar reasoning shows that the class of machines depicted in Fig. 2 will have a limiting error rate which can be made arbitrarily small by choosing  $M$  large enough.

At this point let us turn to a class of machines which can estimate  $\mu$  with a mean-squared error arbitrarily close to the lower bound provided by the quantization problem. As before let  $\mu_1^* < \mu_2^* < \dots < \mu_m^*$  denote the optimal quantization points for an  $m$ -level quantizer. If  $\mu$  is closest to  $\mu_i^*$ , then  $Q(\mu) = \mu_i^*$  is the optimal quantizer output when  $\mu$  is known to be the mean.

Consider the machine which transits from state  $i$  to state  $i + 1$  when  $[X - (\mu_i^* + \mu_{i+1}^*)/2] \geq M$ , and from state  $i$  to state  $i - 1$  when  $[X - (\mu_{i-1}^* + \mu_i^*)/2] \leq -M$ ; and which estimates  $\mu$  to be  $\mu_i^*$  when it is in state  $i$ . For  $M$  sufficiently large this machine has a mean-squared error which is arbitrarily close to the lower bound. Of course, if  $M$  is set equal to  $\infty$ , the mean-squared error becomes very large. This is analogous to the non-existence of optimal machines in the infinite-sample finite-memory hypothesis-testing problem [9].

Note the resemblance of the machine previously described to Fig. 2. To see that this machine is in fact  $\epsilon$ -optimal, let us consider an example in which  $m = 3$ ,  $\mu_1^* = -1$ ,  $\mu_2^* = +1$ ,  $\mu_3^* = +5$ . Then the  $\epsilon$ -optimal machine is shown in Fig. 3. Suppose nature chooses  $\mu < 0$ , so that if we knew the value of  $\mu$  we would set  $Q^*(\mu) = -1$ , the closest quantization point. The machine in Fig. 3 has a probability of transition from state 1 to state 2 which is much smaller than the probability of transition from state 2 to state 1. For large  $M$  the ratio is approximately  $\exp(2\mu M/\sigma)$ , so that as  $M \rightarrow \infty$  the ratio tends to zero. However, this implies [9, eq. (17)] that the steady-state probability of occupying state 1 is approximately  $\exp(2|\mu|M/\sigma)$  times as large as the steady-state probability of occupying state 2. Similarly the probability of transition from state 2 to state 3 is much smaller than the probability of transition from state 3 to state 2. For large  $M$ , the ratio is approximately  $\exp(-2(3 - \mu)M/\sigma)$ . Thus the steady-state probability of occupying state 3 is even smaller than the steady-state probability of occupying state 2, and as  $M \rightarrow \infty$  the steady-state probability of occupying state 1 tends to one. Since the estimate in state 1 equals  $Q^*(\mu)$ , this three-state estimator performs as well as the optimal three-level quantizer when  $\mu < 0$ .

Similarly, if  $\mu > 3$ , so that  $Q^*(\mu) = +5$ , it is seen that the steady-state probability of occupying state 3 tends to one as  $M \rightarrow \infty$ . And, with only minor changes in the argument, if  $0 < \mu < 3$ , so that  $Q^*(\mu) = +1$ , it is seen that the steady-state probability of occupying state 2 tends to one as  $M \rightarrow \infty$ . Thus, as  $M \rightarrow \infty$ , no matter what the values of  $\mu$  and  $\sigma$  are

$$\lim_{n \rightarrow \infty} P[d_n = Q^*(\mu)] = 1 \quad (4)$$

and the "optimal"  $m$ -state learning algorithm estimates  $\mu$  as well as the optimal  $m$ -level quantizer can quantize  $\mu$ . The extension to  $m > 3$  is self-evident.

### III. DISCUSSION

It is seen that there is no optimal rule for the previous infinite-sample finite-memory estimation problem. As noted, this is analogous to the nonexistence of optimal rules for the infinite-sample finite-memory hypothesis-testing problem. However, when the sample size is finite, recent work [19] has shown that optimal rules do exist for hypothesis-testing problems. The same proof is easily carried over to the estimation problem.

While discussing the finite-sample problem, let us also note that, in view of the previously mentioned behavior of the maximum of  $N$  independent  $\mathcal{N}(0,1)$  random variables, we can conjecture that when  $\sigma$  is known the optimal value of  $M$  will be approximately  $\sigma(2 \ln N)^{1/2}$ . This adds a positive term to the minimum achievable mean-squared error, and, although this term tends to zero as  $N \rightarrow \infty$ , it does so more slowly than any algebraic function (being of the form  $\exp [-(\ln N)^{1/2}]$ ).

Also note that, if the sample size is finite and  $\sigma$  is unknown, then the problem becomes entirely different. Some memory must now be devoted to estimating  $\sigma$  even though no estimate of  $\sigma$  has been asked for. This is because the optimal value of  $M$  depends on  $\sigma$ . This is to be contrasted with the infinite-sample problem, where no *a priori* knowledge was needed about  $\sigma$ , and no memory was used to estimate  $\sigma$ .

Returning to the infinite-sample problem, certain extensions are quite simple. First, another cost function or approach (e.g., minimum mean absolute error or minimax) makes no difference, provided the optimal  $m$ -level quantizer produces  $m$  ordered intervals as the partitioning of the parameter space  $\Theta$ . Second, any statistics which yield  $\gamma = \infty$  for all differing values of  $\theta$  and which are "tail monotonic-likelihood" ratio (see following definitions) are equivalent for this estimation problem.

**Definition:** The  $\gamma$  value between two probability distributions  $\mathcal{P}_{\theta_0}$  and  $\mathcal{P}_{\theta_1}$  is equal to the maximum value of the likelihood ratio divided by the minimum value of the likelihood ratio. As shown in [9],  $\gamma$  is a measure of the resolvability of finite-memory tests of  $\theta_0$  versus  $\theta_1$ . In [9] certain pathologies concerning the definition of  $\gamma$  are also dealt with (e.g., if the maximum-likelihood ratio occurs only on a set with measure zero).

**Definition:** A family of distributions  $\mathcal{P}_{\theta}$  indexed by  $\theta \in \mathbb{R}$  is tail monotonic-likelihood ratio if for all  $\theta$  and all  $B < \infty$  there exist values  $M_1$  and  $M_2$  such that  $\mathcal{P}_{\theta_1} \{X > M_1\} / \mathcal{P}_{\theta_1} \{X < M_2\} > B$  for  $\theta_1 > \theta$ , and this same ratio is less than  $B^{-1}$  for  $\theta_1 < \theta$ .

For example, if  $\mathcal{P}_{\theta}$  is the uniform distribution  $U(\theta - \frac{1}{2}, \theta + \frac{1}{2})$  and the optimal quantizer partitioning of  $\Theta = \mathbb{R}$  is  $\{(-\infty, \theta_1^*), (\theta_1^*, \theta_2^*), \dots, (\theta_{m-2}^*, \theta_{m-1}^*), (\theta_{m-1}^*, \infty)\}$  then moving from state  $i$  to state  $i + 1$  when  $X \geq \theta_i^* + \frac{1}{2}$  and from state  $i$  to state  $i - 1$  when  $X \leq \theta_{i-1}^* - \frac{1}{2}$  results in the machine's being in the "correct" state with probability one.

A word of caution is in order. The infinite and very large sample properties depend crucially on the tail behavior of the distributions  $\{\mathcal{P}_{\theta}\}$  and should not be relied upon in real world problems unless the distributions of the tails are well known. However, if memory is limited but not too small, the concepts developed in this correspondence provide qualitative insights into the design of good "suboptimal" finite-memory estimators.

### REFERENCES

- [1] H. Robbins, "A sequential decision problem with a finite memory," *Proc. Nat. Acad. Sci. U.S.A.*, vol. 42, pp. 920-933, 1956.
- [2] C. V. Smith and R. Pyke, "The Robbins-Isbell two-armed bandit

- problem with finite memory," *Ann. Math. Statist.*, vol. 36, pp. 1375-1386, 1965.
- [3] S. M. Samuels, "Randomized rules for the two-armed bandit with finite memory," *Ann. Math. Statist.*, vol. 39, pp. 2103-2107, 1968.
- [4] M. L. Tsetlin, "On the behavior of finite automata in random media," *Automat. Telemekh.*, vol. 22, pp. 1345-1354, Oct. 1961 (available in English translation).
- [5] V. Y. Krylov, "On one automaton that is asymptotically optimal in a random medium," *Avtomat. Telemekh.*, vol. 24, pp. 1226-1228, Sept. 1963 (available in English translation).
- [6] V. I. Varshavskii and I. P. Vorontsova, "On the behavior of stochastic automata with a variable structure," *Avtomat. Telemekh.*, vol. 24, Mar. 1963 (available in English translation).
- [7] K. S. Fu and T. J. Li, "On the behavior of learning automata and its applications," Purdue Univ., Lafayette, Ind., Tech. Rep. TR-EE 68-20, 1968.
- [8] B. Chandrasekaran and D. W. C. Shen, "On expediency and convergence in variable-structure automata," *IEEE Trans. Syst. Sci. Cybern.*, vol. SSC-4, pp. 52-60, Mar. 1968.
- [9] M. E. Hellman and T. M. Cover, "Learning with finite memory," *Ann. Math. Statist.*, vol. 41, pp. 765-782, June 1970.
- [10] T. M. Cover and M. E. Hellman, "The two-armed bandit problem with time-invariant finite memory," *IEEE Trans. Inform. Theory*, vol. IT-16, pp. 185-195, Mar. 1970.
- [11] M. E. Hellman and T. M. Cover, "On memory saved by randomization," *Ann. Math. Statist.*, vol. 42, pp. 1075-1078, 1971.
- [12] M. E. Hellman, "The effects of randomization on finite memory decision schemes," *IEEE Trans. Inform. Theory*, vol. IT-18, pp. 499-502, July 1972.
- [13] J. A. Horos and M. E. Hellman, "A confidence model for finite-memory learning systems," *IEEE Trans. Inform. Theory (Corresp.)*, vol. IT-18, pp. 811-813, Nov. 1972.
- [14] P. F. Lynn and R. R. Boorstyn, "Bounds on finite memory detectors," in *1972 IEEE Int. Symp. Information Theory*, Asilomar, Calif., Jan. 31-Feb. 3, 1972.
- [15] C. T. Mullis and R. A. Roberts, "Memory limitation and multistage decision processes," *IEEE Trans. Syst. Sci. Cybern.*, vol. SSC-4, pp. 307-316, Sept. 1968.
- [16] T. M. Cover, "Hypothesis testing with finite statistics," *Ann. Math. Statist.*, vol. 40, pp. 828-835, 1969.
- [17] R. W. Muise and R. R. Boorstyn, "Detection with time-varying finite-memory receivers," *1972 IEEE Int. Symp. Information Theory (Abstracts of Papers)*, pp. 35-36.
- [18] R. A. Roberts and J. R. Tooley, "Estimation with finite memory," *IEEE Trans. Inform. Theory*, vol. IT-16, pp. 685-691, Sept. 1970.
- [19] T. Cover, M. Freedman, and M. Hellman, "Optimal finite-memory learning algorithms for the finite sample problem," submitted to *Inform. Contr.*
- [20] M. Freedman, "A finite-memory, finite-time, Gaussian hypothesis testing problem," M.S. thesis, Dep. Elec. Eng., M.I.T., Cambridge, Mass., 1971.
- [21] M. Freedman and M. Hellman, "A finite-memory, finite-time, Gaussian hypothesis testing problem," in *Program and Abstracts for the 1972 IEEE Symp. Information Theory*, Asilomar, Calif., Jan. 31-Feb. 3, 1972.
- [22] T. Wagner, "Estimation of the mean with time-varying finite memory," *IEEE Trans. Inform. Theory (Corresp.)*, vol. IT-18, pp. 523-525, 1972.
- [23] T. Cover, "Some problems in estimation with finite statistics," in *Proc. 1970 IEEE Symp. Adaptive Processes (9th Decision and Control)*, Austin, Tex., Dec. 7-9.
- [24] D. Sagalowicz, "Hypothesis testing with finite memory," Ph.D. dissertation, Dep. Elec. Eng., Stanford Univ., Stanford, Calif., Sept. 1970.
- [25] W. Feller, *An Introduction to Probability Theory and Its Applications*, Vol. I, Third Ed. New York: Wiley, 1968, p. 193, problem #1.

### Finite-Memory Classification of Bernoulli Sequences Using Reference Samples

BRUNO O. SHUBERT

**Abstract**—It is shown that in two-way Bernoulli classification problems deterministic machines can perform as well as optimal randomized machines if their memory is increased by less than one bit. This is accomplished by allowing the algorithm to observe samples from both classes, thus in effect using the data source itself to provide the necessary randomization. An application to a simple communication problem is indicated.

Manuscript received June 6, 1973, revised November 30, 1973.  
The author is with the Department of Operations Research and Administration Sciences, Naval Postgraduate School, Monterey, Calif. 93940.