

# Lecture Notes in Statistics

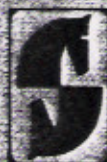
Edited by D. Brillinger, S. Fienberg, J. Gani,  
J. Hartigan, and K. Krickeberg

20

## Mathematical Learning Models— Theory and Algorithms

Proceedings of a Conference

Edited by  
Ulrich Herkenrath, Dieter Kalin,  
and Walter Vogel



Springer-Verlag  
New York Berlin Heidelberg Tokyo



ON A CLASS OF LEARNING ALGORITHMS WITH SYMMETRIC BEHAVIOR UNDER  
SUCCESS AND FAILURE

by

M. R. Meybodi and S. Lakshmivarahan  
School of Electrical Engineering and Computer Science  
University of Oklahoma, Norman, Oklahoma

1. Introduction:

Learning algorithms have been extensively studied in Mathematical Psychology [13] [14] [18] for decades and more recently in Learning Automata Theory and Mathematical Statistics [1] [2] [8] [11] [16]. In Mathematical Psychology the interest in learning stems from the desire to understand the observed animal learning and associated changes in their behavior. However, in Learning Automata Theory and Mathematical Statistics the aim is to build algorithms that exhibit prespecified behaviour. Our interests, in this paper, are in the latter approach.

We study a new class of absorbing barrier algorithms of the reward-penalty type with identical behaviour under the occurrence of success and failure. Necessary and sufficient conditions for strong absolute expediency of this class of algorithm is obtained.

The concept of absolute expediency was originally introduced by Lakshmivarahan and Thathachar [3] in 1973. Since then this concept has played a major role in the analysis and design of  $\epsilon$ -optimal absorbing barrier algorithms. [1] [2] [4] [5] [7]. The results of this paper sheds further light on this concept. It should be interesting to note that  $\epsilon$ -optimal non-absorbing learning algorithms have also been recently developed in [15] and discussed extensively in the book [21].

A number of learning algorithms of the reward-penalty type whose behaviour is asymmetrical with respect to the occurrence of success and failure have been extensively studied in the literature - Varshavskii and Vorontsova [10], Fu [8], McMurty and Fu [9], Shapiro and Narendra [20], Viswanathan and Narendra [19], Chandrasekaran and Shen [17], Lakshmivarahan and Thathachar [3] and Sawaragi and Baba [7] to mention a few. For a similar class of general learning algorithms discussed in this paper Aso and Kimura [5] derived necessary and sufficient conditions for absolute expediency [3]. Interestingly Aso and Kimura [5] call this class of algorithms as "Stochastic Vector Automaton" algorithms. Recently for the same class of algorithms considered in this paper Herkenrath, Kalin and Lakshmivarahan [6] derived necessary and sufficient conditions for absorption at the vertices of the unit simplex of proper dimension.

In section 2, the algorithm and the statement of problem are given. Necessary and sufficient conditions for strong absolute expediency and convergence of the

$(t)a_i, i \in I,$

of the individuals.

$a_i = a_i(m_{ii}-1),$

1 eigenvalue of  $Q$  is the growth (decay). If there are is the problem of maximizing

ality in average cost denume- problems, recurrency conditions 15 (1978), 356-373.

ation and stochastic control. ent Advances in Filtering February 1 to 6, 1982. ol and Information Sciences,

3 Markov Processes, Cambridge

imizing the characteristic Mathématiques pures et appl.

c programming recursions I, 526-547, 17 (1981), 310-328.



algorithm with probability one are established in section 3.  $\epsilon$ -optimality is derived in section 4.

## 2. LEARNING ALGORITHM and STATEMENT OF PROBLEM

There are  $M$  ( $2 \leq M < \infty$ ) coins. At time instant  $k$  ( $= 0, 1, 2, \dots$ ) the coin  $i$  ( $1 \leq i \leq M$ ) is chosen for tossing with probability  $P_i(k)$  where  $P(k) = (P_1(k), P_2(k), \dots, P_M(k))^T$ ,  $\sum_{i=1}^M P_i(k) = 1$ ,  $0 \leq P_i(k)$  and  $T$  denotes transpose. In tossing, the  $i$ th coin falls head (tail) with probability  $d_i$  ( $c_i = 1 - d_i$ ). It is assumed that (1)  $0 < d_i < 1$ ,  $i = 1, 2, \dots, M$ , (2) the  $d_i$ 's are all distinct, that is,  $d_i \neq d_j$  for all  $i$  and  $j$ . (3)  $d_i$ 's do not depend on  $k$  and (4)  $d_i$ 's are all unknown. The outcome of falling head is called success and falling tail is called failure. Let  $D = (d_1, d_2, \dots, d_M)^T$  and without loss of generality assume

$$d_1 > d_2 > d_3 > \dots > d_M \quad (1)$$

Any vector  $D$  that satisfies the above conditions is called an admissible  $D$ . The average probability of success at stage  $k$  denoted by  $n(k)$  is

$$n(k) = \sum_{i=1}^M P_i(k) d_i \quad (2)$$

Our basic problem is to make the average probability of success as close to its maximum (that is,  $d_1$ ) as possible. The primary interest in and the challenge of this problem arise basically from the fact that the success probabilities ( $d_i$ 's) of the coin are unknown.

As a first step towards the solution to this problem, in this paper, we propose to change  $p(k)$  in a learning algorithm. The key idea of the learning algorithm may be stated in words as follows: Increase (decrease) the probability of choosing the  $i$ th coin at time  $(k+1)$ , if it was chosen for tossing at time  $k$  and the toss resulted in success (failure). In particular, let

$$S_M = \{P | P = (P_1, P_2, \dots, P_M)^T, 0 \leq P_i, \sum_{i=1}^M P_i = 1\}$$

be the  $M$ -dimensional unit simplex and let

$$V_M = \{e_i | i = 1, 2, \dots, M\} \text{ where } e_i = (0, 0, \dots, 1, \dots, 0)^T$$

be the  $i$ th unit vector of dimension  $M$ . Further, let  $S_M^0 = \{P | P = (P_1, P_2, \dots, P_M)^T, 0 < P_i, \sum_{i=1}^M P_i = 1\}$ . Clearly,  $V_M$  corresponds to the corners of vertices of  $S_M$ .

Let  $f_s^i[.]$ ,  $g_s^i[.]$  be continuous functions such that

$$f_s^i : S_M \rightarrow [0, 1], g_s^i : S_M \rightarrow [0, 1], i, s = 1, 2, \dots, M$$

Formally the above algorithm may be stated as follows

$$\left. \begin{aligned} P_s(k+1) &= P_s(k) - \theta f_s^i[P(k)], s \neq i \\ P_i(k+1) &= P_i(k) + \theta f_i^i[P(k)] \end{aligned} \right\} \text{ if the toss of coin } i \text{ resulted in success} \quad (3)$$

and

$$\left. \begin{aligned} P_i(k+1) &= P_i(k) - \theta g_i^i[P(k)] \\ P_s(k+1) &= P_s(k) + \theta g_s^i[P(k)], s \neq i \end{aligned} \right\} \text{ if the toss of coin } i \text{ resulted in failure}$$

where  $0 < \theta \leq 1$  is called the step length parameter and the following condition (4):

$$\begin{aligned} &\text{either } f_s^i[.] = 0 \text{ for all } p \in S_M \\ &\text{or } f_s^i[p] \leq P_s, s \neq i \text{ and } f_i^i[p] = \sum_{s \neq i} f_s^i[p] \end{aligned} \quad (4)$$

and

$$\begin{aligned} &\text{either } g_s^i[p] = 0 \text{ for all } p \in S_M \\ &\text{or } g_s^i[p] \leq P_s \text{ and } g_i^i[p] = \sum_{s \neq i} g_s^i[p] \end{aligned}$$

for all  $i, s = 1, 2, \dots, M$ ,

imply that  $p(k+1) \in S_M$  if  $p(k)$  does. However in order to make the algorithm non trivial and interesting either  $f_s^i = 0$  or  $g_s^i = 0$  but not both.

Remark 1. If the coin  $i$  is chosen for tossing, it is clear from (3) that success (failure) increases (decreases) the probability of its choice. The increase and decrease are called reward and penalty and hence (3) is called reward-penalty algorithm. If  $f_s^i \neq 0$  but  $g_s^i = 0$  then (3) is called reward-inaction algorithm, if  $f_s^i = 0$  but  $g_s^i \neq 0$  then it is called inaction-penalty algorithm

Since the vector  $D$ , of success probabilities and the functions  $f_s^i[.]$  and  $g_s^i[.]$  (for all  $i$  and  $s$ ) are independent of  $k$ ,  $\{p(k)\}, k \geq 0$  is clearly a discrete time Markov process with stationary transition function over the state space  $S_M$ . To quantify the behaviour of the process  $\{p(k)\}$ , and hence of  $\{n(k)\}$  we introduce the following definitions. Let  $I = \{1, 2, \dots, M\}$ ,  $E = \{\text{success, failure}\}$ . The algorithm (3) defines a mapping  $T: S_M \times I \times E \rightarrow S_M$  where

$$p(k+1) = T[p(k), i(k), e(k)] \quad (5)$$

$i(k) \in I$  is the coin chosen for tossing and  $e(k) \in E$  is the outcome of the toss of that coin at time  $k$ .

Definition 1: A state  $p \in S_M$  is said to be absorbing state if and only if  $p = T[p, i, e]$  with probability one

(6)

Definition 2: A learning algorithm  $T$  is said to be absorbing if and only if there is at least one absorbing state.

For reasons that will become apparent in this paper our interests are in the class of learning algorithms for which  $V_M$  is the only set of absorbing states. We call such an algorithm absorbing barrier algorithms

Definition 3: A learning algorithm  $T$  is said to be

$$\text{a) optimal if } \lim_{k \rightarrow \infty} E[n(k)] = d_1 \quad (7)$$

$$\text{b) } \epsilon\text{-optimal if for all } p(0) = p \in S_M^0, \lim_{k \rightarrow \infty} |E[n(k)] - d_1| < \epsilon, \epsilon > 0 \quad (8)$$

$$\text{c) absolutely expedient if } E[n(k+1) | p(k) = p] \geq n(k) \text{ with probability one} \quad (9)$$

for all admissible  $D$  (satisfying the assumptions given at the beginning of this



section) and for all  $p \in S_M$ , with strict inequality in (9) holding good for all  $p \in S_M^0$ .

d) strongly absolutely expedient if

$$E[n(k+1) | p(k) = p] \geq n(k) \text{ with probability one} \quad (10)$$

for all admissible  $D$  and for all  $p \in S_M$  with strict inequality in (10) holding good for all  $p \in (S_M - V_M)$ .

**STATEMENT OF PROBLEM:** Our aim in this paper is to find conditions for  $\epsilon$ -optimality of the general class of learning algorithm given in (3).

**Remark 2:** Algorithm (3) is a generalization of many known algorithms in the sense that the functions in (3) used for updating the probabilities depend on the coin chosen as well. Such generalized algorithms are considered in [5] and [6]. In most of the earlier papers the functions that are used in updating are independent of the coin that is chosen for tossing, [1] [2] [3] [4] [7] [8] [9] [10] [12] [13].

### 3. Convergence with probability one:

As a first step towards  $\epsilon$ -optimality, in this section we derive conditions on the algorithm such that the Markov process  $\{p(k)\}$   $k \geq 0$  converges with probability one. To this end we begin by rewriting the consistency conditions (4) in a form more suitable for our analysis. Let

$$(C.1) \quad f_i^1[p] = \alpha[i, p] (1 - p_i), \quad f_s^1[p] = \alpha[i, s, p] p_s$$

and

$$\sum_{s \neq i} \alpha[i, s, p] p_s = \alpha[i, p] (1 - p_i)$$

and

$$(C.2) \quad g_i^1[p] = \gamma[i, p] p_i, \quad g_s^1[p] = \lambda[i, s, p] (1 - p_s)$$

and

$$\gamma[i, p] p_i = \sum_{s \neq i} \lambda[i, s, p] (1 - p_s)$$

where  $\alpha, \gamma: 1 \times S_M \rightarrow [0, 1]$  and  $\lambda: 1 \times 1 \times S_M \rightarrow [0, 1]$ .

The basic rule that governs the choice of functions in the above form is that if a term is subtracted from  $p_j$  then it is made proportional to  $p_j$  and if a term is added to  $p_j$  then it is made proportional to  $(1 - p_j)$  irrespective of which coin is chosen for the toss and whether the toss results in success or failure.

**Remark 3:** Our choice of the functions  $g_s^1[\cdot]$  is quite untraditional in the sense that in all most all the papers in Mathematical Psychology [18] [17] [13] [14]  $g_s^1[p]$  is made proportional to  $p_s$  for all  $i$  and  $s, s \neq i$ . Also in almost all the papers on Learning Automata [3] [4] [5] [7]  $g_i^1[p]$  is made proportional to  $(1 - p_i)$  and  $g_s^1[p]$  is made proportional to  $p_s$  for all  $s \neq i$ . Because of this there is a disparity in the behaviour of the algorithm (3) under success and failure. However, our present choice of functions  $g_s^1[\cdot]$  for all  $i$  and  $s$  given in (C.2) induce identical behaviour of the algorithm (3) under success and failure.

The following theorem is immediate

**Theorem 1:** The learning algorithm (3) is an absorbing barrier learning algorithm and only if (A.1) and (A.2) hold true.

(A.1) For all  $p \in S_M - V_M$  there exists  $1 \leq s \leq M$

such that  $\alpha[s, p] + \gamma[s, p] p_s > 0$

(A.2) For all  $1 \leq s \leq M, \gamma[s, e_s] = 0$

**Proof:** Refer [6]

A typical choice of functions in the algorithm (3) that satisfies the conditions (C.1) - (C.2) and (A.1) - (A.2) are given in the following example.

Let  $0 < C_1, C_2 < 1$ .

$$\alpha[i, s, p] = C_1 (1 - p_i) (1 - p_s), \quad \alpha[i, p] = C_1 \sum_{s \neq i} p_s (1 - p_s)$$

$$\lambda[i, s, p] = C_2 p_i p_s^2 (1 - p_i), \quad \gamma[i, p] = C_2 (1 - p_i) \sum_{s \neq i} p_s^2 (1 - p_s)$$

for all  $i, s = 1, 2, \dots, M$ .

**Remark 4:** The demonstration of the symmetry in the properties of the algorithm (3) under success and failure alluded to in Remark 3 is evidenced by Theorem 1.

Notice that (3) is an absorbing barrier learning algorithm if  $\alpha[i, p] \neq 0$  and  $\gamma[i, p] \neq 0$  or  $\alpha[i, p] = 0$  and  $\gamma[i, p] \neq 0$  or both  $\alpha[i, p]$  and  $\gamma[i, p] \neq 0$ . This is in sharp contrast with the properties of the currently available absolutely expedient learning algorithm [3] [4] [5] [7] wherein the reward-penalty and reward-inaction algorithms are absorbing barrier type but the inaction-penalty is not. In fact in all the inaction-penalty algorithms of the absolutely expedient type known so far [3], every state in  $S_M$  is an absorbing state. Our modified definition of strong absolute expediency is in fact motivated by the existence of the absorbing barrier algorithms of the reward-penalty, reward-inaction and inaction-penalty types.

**Remark 5:** For some  $1 \leq j \leq M$  if  $p_j(k) = 0$ , then it follows from (3) and (C.1) - (C.2) that  $p_j(k^*) = 0$  for all  $k \geq k$ . In other words, during learning process if  $p(k)$  reaches the boundary  $p_j = 0$  of the simplex  $S_M$ , then  $p(k)$  will continue to remain in that boundary.

Henceforth in this paper we will only be concerned with the absorbing barrier learning algorithms, that is, algorithm (3) under the conditions (A.1) - (A.2) of Theorem 1.

**Theorem 2:** Necessary and sufficient conditions for the absorbing barrier learning algorithm (3) to be strongly absolutely expedient are:

$$(S.1) \quad \sum_{j \neq i} p_j \alpha[i, j, p] = \sum_{j \neq i} p_j \alpha[j, i, p]$$

and

$$(S.2) \quad \sum_{j \neq i} p_j (1 - p_j) \lambda[i, j, p] = \sum_{j \neq i} p_j (1 - p_j) \lambda[j, i, p]$$

for all  $i = 1, 2, \dots, M$

**Sufficiency:** Define  $\Delta x(k) = x(k+1) - x(k)$  and let

$$\Delta n(k) = E[\Delta n(k) | p(k)] = \sum_{i=1}^M E[\alpha p_i(k) | p(k)] d_i \quad (11)$$

It can be seen by direct computation that

$$E [x_{p_i}(k) | p(k) = p] = p_i (1-p_i) d_i \alpha(i,p) - p_i^2 c_i \gamma(i,p) - \sum_{j \neq i} p_j p_i d_j \beta(j,i,p) + \sum_{j \neq i} p_j (1-p_i) c_j \lambda(j,i,p) \quad (12)$$

Substituting (12) in (11) and in view of (C.1) and (C.2) we obtain

$$\Delta \eta_1(k) = \Delta \eta_1(k) + \Delta \eta_2(k) \quad (13)$$

where

$$\Delta \eta_1(k) = \sum_{i=1}^M p_i d_i^2 \sum_{j \neq i} p_j \beta(j,i,p) - \sum_{i=1}^M p_i d_i \sum_{j \neq i} p_j d_j \beta(j,i,p) \quad (14)$$

and

$$\Delta \eta_2(k) = - \sum_{i=1}^M p_i d_i c_i \sum_{j \neq i} \lambda(j,i,p) (1-p_j) + \sum_{i \neq j} (1-p_i) d_i \sum_{j \neq i} p_j c_j \lambda(j,i,p) \quad (15)$$

After simplification it can be shown that [23]

$$\Delta \eta_1(k) = 1/2 \sum_{i=1}^M \sum_{j \neq i} p_i p_j (d_i - d_j)^2 (\alpha(i,j,p) + \lambda(j,i,p)) \geq 0 \text{ with equality holding only if } p \in V_m \quad (16)$$

Similarly it can be shown that [23]

$$\Delta \eta_2(k) = 1/2 \sum_{i=1}^M \sum_{j \neq i} p_i (1-p_j) (d_i - d_j)^2 (\alpha(i,j,p) + \lambda(j,i,p)) \geq 0 \text{ with equality holding good only if } p \in V_m \quad (17)$$

From (16) and (17) sufficiency follows.

Necessity:  $\Delta \eta_1(k)$  can be represented as a quadratic and linear term in the vector  $D$  as follows:

$$\Delta \eta_1(k) = D^T A D + D^T B \quad (18)$$

where  $A = [A_{ij}]$  and  $B = [B_1, B_2, \dots, B_m]^T$  with

$$A_{ii} = p_i (1-p_i) \alpha(i,p) + p_i^2 \gamma(i,p)$$

$$A_{ij} = -[p_i p_j \beta(j,i,p) + (1-p_i) p_j \lambda(j,i,p)]$$

and

$$B_i = p_i^2 \gamma(i,p) + (1-p_i) \sum_{j \neq i} p_j \lambda(j,i,p)$$

for all  $i, j = 1, 2, \dots, M$

From the definition of strong absolute expediency it follows that  $\Delta \eta_1(k)$  attains its minimum value zero either when all  $d_i$  are equal or when  $p \in V_m$ . Since every member of  $V_m$  is absorbing, on  $V_m$ , it is easily seen that  $\Delta \eta_1(k)$  attains its minimum value for all admissible  $D$  vectors. In the following we shall derive conditions for the minimum of  $\Delta \eta_1(k)$  when  $d_i = d$  for all  $i = 1, 2, \dots, M$ ,  $0 < d < 1$ . Necessary conditions for minimum are obtained by setting the derivative of  $\Delta \eta_1(k)$  (with respect to  $d_i$ ) at the point  $d_i = d$  for all  $i = 1, 2, \dots, M$  equal to zero, that is,

$$\frac{\partial \Delta \eta_1}{\partial d_i} \bigg|_{d_i=d} = 0 \text{ for all } i = 1, 2, \dots, M \quad (19)$$

From (18), the equation (19) take the form

$$(A + A^T) d \mathbf{1} + B = 0 \quad (20)$$

where  $\mathbf{1} = (1, 1, \dots, 1)^T$  is an  $M$  dimensional column vector of all ones.

Rewriting (20) we get

$$d K_i + L_i = 0$$

for all  $i = 1, 2, \dots, M$  and all  $0 < d < 1$

where

$$K_i = p_i (1-p_i) \alpha(i,p) + p_i^2 \gamma(i,p) - p_i \sum_{j \neq i} p_j \beta(j,i,p) - (1-p_i) \sum_{j \neq i} p_j \lambda(j,i,p)$$

and

$$L_i = p_i^2 \gamma(i,p) + (1-p_i) \sum_{j \neq i} p_j \lambda(j,i,p).$$

(21) is true for all  $0 < d < 1$  only if  $L_i = 0$  and  $K_i = 0$  for all  $i = 1, 2, \dots, M$ .

$L_i = 0$  leads to the condition

$$p_i^2 \gamma(i,p) = (1-p_i) \sum_{j \neq i} p_j \lambda(j,i,p). \quad (22)$$

Using (C.2), from (22) we obtain

$$p_i \sum_{j \neq i} (1-p_j) \lambda(j,i,p) = (1-p_i) \sum_{j \neq i} p_j \lambda(j,i,p). \quad (23)$$

which in fact is (S.2). Substituting (23) in  $K_i = 0$ , we obtain

$$p_i (1-p_i) \alpha(i,p) = p_i \sum_{j \neq i} p_j \beta(j,i,p).$$

Once again, using (C.1) we get  $p_i \sum_{j \neq i} p_j \beta(j,i,p) = p_i \sum_{j \neq i} p_j \lambda(j,i,p)$

which is the same as (S.1). Hence the theorem.

**Corollary 1:** The Markov process  $\{p(k)\}_{k \geq 0}$  as generated by the algorithm (3) under the conditions (C.1) - (C.2), (A.1) - (A.2) and (S.1) - (S.2) converges to  $V_m$  with probability one.

**Proof:** Since  $\{p(k)\}_{k \geq 0}$  is a Markov process, from theorem (2) we get

$$E [x_{p_i}(k) | P(r): 0 \leq r \leq k] = E [x_{p_i}(k) | P(k)] \geq 0$$

This in turn implies that  $x_{p_i}(k)_{k \geq 0}$  is a submartingale [22]. By martingale theorem  $\lim_{k \rightarrow \infty} x_{p_i}(k)$  exists and hence  $\lim_{k \rightarrow \infty} p(k) = p^*$  exists with probability one. As  $x_{p_i}(k) = 0$  only on  $V_m$ , it follows that  $p^* \in V_m$  with probability one.

**4. Optimality:** In the previous section it was established that  $p^* \in V_m$  with probability one. In this section we set out to quantify the distribution of  $p^*$ .

To this end, define  $r_i(p) = \text{prob}[P^* = e_i | P(0) = p]$  (24)

for  $i = 1, 2, \dots, M$ . Notice  $\sum_{i=1}^M r_i(p) = 1$  for all  $p \in S_M$

In view of the corollary 1 and (24) we obtain, for all  $p(0) = p \in S_M^0$

$$\lim_{k \rightarrow \infty} E [x_{p_i}(k)] = \sum_{i=1}^M r_i(p) d_i \quad (25)$$

To compute  $r_i(p)$  we need the following: Let  $C[S_M]$  be the class of all continuous functions from  $S_M$  to the real line. If  $f(\cdot) \in C[S_M]$  define the operator  $U$  as follows:

$$U f(p) = E [f(p(k+1)) | p(k) = p]$$

Clearly the operator  $U$  is linear and positive [12]

**Definition 4:** A function  $f: S_M \rightarrow \text{real line}$  is called super regular (regular, sub



regular) if

$$f(p) \geq (=, \leq) U f(p)$$

for all  $p \in S_M$ .

With these preliminaries, we now state two important propositions that lead an algorithm for quantifying  $\Gamma_1(p)$ .

**Proposition 1:**  $\Gamma_1(p)$  is the only continuous solution of the functional equation.

$$U \Gamma_1(p) = \Gamma_1(p) \quad (26)$$

satisfying the boundary condition

$$\Gamma_1(e_i) = 1 \text{ and } \Gamma_1(e_j) = 0 \text{ for all } i, j, i \neq j$$

Notice  $\Gamma_1(p)$  satisfying (26) by definition is a regular function. This functional equation is extremely difficult to solve. Hence in the following we establish upper and lower bounds on  $\Gamma_1(p)$ .

**Proposition 2:** Let  $f_1(p) \in C[S_M]$  be super (sub) regular function with  $f_1(e_i) = 1$  and  $f_1(e_j) = 0$  for all  $i, j, i \neq j$  then

$$f_1(p) \geq (\leq) \Gamma_1(p)$$

The proof of these propositions are rather involved and we refer the reader to Norman [12] for an elegant proof.

Thus, if we can find two functions  $h_1^{(1)}(p)$  and  $h_1^{(2)}(p)$  which are super and sub-regular functions respectively and satisfying the boundary conditions

$$h_1^{(1)}(e_i) = h_1^{(2)}(e_i) = 1, h_1^{(1)}(e_j) = 0 \text{ for all } j \neq i, i = 1, 2.$$

then from proposition 2 it follows that  $h_1^{(2)}(p) \leq \Gamma_1(p) \leq h_1^{(1)}(p)$

Consider a function

$$\Psi_1[x, p] = \exp[-x_1 p_1 / \theta]$$

where  $x_1 > 0$  is a parameter and  $\exp[-x] = e^{-x}$

Recognizing the fact that

$$\Psi_1[x, p] = \frac{1 - \exp[-x_1 p_1 / \theta]}{1 - \exp[-x_1 / \theta]}$$

is sub (super) regular whenever  $\Psi_1[x, p]$  is super (sub) regular, it follows from propositions 2 that

$$\Psi_1[Y_1, p] \leq \Gamma_1(p) \leq \Psi_1[Z_1, p]$$

where  $Y_1, Z_1$  are two constants such that  $\Psi_1[Y_1, p]$  and  $\Psi_1[Z_1, p]$  are

sub and super regular function respectively. Note that function  $\Psi_1[x, p]$  depends on just one component of  $p$  and thus leads to conservative results. However, it gives rise to expressions which are easily manageable.

The problem of getting bounds on  $\Gamma_1(p)$  now reduces to one of finding two positive constants  $Y_1$  and  $Z_1$  such that  $\Psi_1[Y_1, p]$  is sub and regular and  $\Psi_1[Z_1, p]$  is super regular. Further, from (25) and the inequality (1) it is clear that for  $c = \text{optimality}$  we need to concentrate only (on the lower bound) on  $\Gamma_1(p)$ .

It can be seen that (by dropping the subscript 1 from  $x$ , for convenience)

$$U \Psi_1[x, p] = \Psi_1[x, p] = -x F_1[x, p] + \Gamma_1(x, p)$$

where

$$F_1[x, p] = p_1 d_1 + [1, p] (1-p_1) \vee [-x, p] [1, p] (1-p_1)$$

$$-p_1 C_1 \vee [1, p] p_1 \vee [x, p] [1, p] p_1$$

$$- \sum_{j \neq 1} p_j d_j - p[j, 1, p] p_1 \vee [x, p] [1, 1, p] p_1$$

$$+ \sum_{j \neq 1} p_j C_j \wedge [j, 1, p] (1-p_1) \vee [-x, p] [j, 1, p] (1-p_1)$$

and  $V[x] = (\exp(x) - 1)/x$  if  $x \neq 0$  and 1 if  $x = 0$ .

Clearly,  $F_1[x, p] \geq 0$  implies  $\Psi_1[x, p]$  regular.

Since  $\alpha[1, p]$ ,  $\gamma[1, p]$ ,  $\delta[1, j, p]$  and  $\rho[1, j, p]$  are all bounded above by unity and as the  $V[\cdot]$  is strictly monotonically increasing, it follows

$$F_1[x, p] \geq 0 \text{ if } G[x, p] \geq \frac{V[-x(1-p_1)]}{V[x p_1]} \geq \frac{p_1^2 C_1 + [1, p] + \sum_{j \neq 1} p_j p_j d_j \rho[j, 1, p]}{\sum_{j \neq 1} p_j (1-p_1) C_j + [j, 1, p] + p_1 (1-p_1) d_1 + [1, p]}$$

It can be shown using the properties of the  $V[\cdot]$  function [4] [12] that

$$G[x, p] \geq \frac{1}{V[x]}$$

Further, using (C.1) - (C.2) and (S.1) - (S.2) it can be seen that

$$e^* \leq \frac{C_1 A + d_2 B}{C_M A + d_1 B} \geq \frac{p_1^2 C_1 + [1, p] + \sum_{j \neq 1} p_j p_j d_j \rho[j, 1, p]}{\sum_{j \neq 1} p_j (1-p_1) C_j + [j, 1, p] + p_1 (1-p_1) d_1 + [1, p]}$$

where  $A = \sum_{j \neq 1} p_j (1-p_1) C_j + [j, 1, p]$  and  $B = \sum_{j \neq 1} p_j p_j d_j \rho[j, 1, p]$

In view of the inequality (1) it follows that  $e^* \leq 1$ . From (27) and (28)

$$F_1[x, p] \geq 0 \text{ if } \frac{1}{V[x]} \geq e^* \quad (29)$$

Since  $e^* < 1$ ;  $V[x] = \frac{1}{e^*}$  has a unique solution  $x = y > 0$  so that  $\Psi[y, p]$  is subregular

**Remark 2:** For the reward-inaction algorithm (refer remark1)  $e^*$  reduces to  $\frac{d_2}{d_1}$  and

for the inaction - penalty algorithm  $e^*$  reduces to  $\frac{C_1}{C_M}$ . Thus in either of these special cases there exists a unique solution  $x=y > 0$  for the equation  $V[x] = \frac{1}{e^*}$  and

hence in both these special cases there exists a lower bound on  $\Gamma_1(p)$ . This is in sharp contrast with the existing results in the literature where in for the inaction-penalty (absolutely expedient) algorithms [4] no lower bound on  $\Gamma_1(p)$  has been established. One of the reasons for this anomaly is that none of the currently available inaction-penalty (absolutely expedient) algorithms [4] are of the absorbing barrier type.

Having established the lower bound on  $\Gamma_1(p)$  we now state our main result.

**Theorem 3:** For every  $c > 0$  and  $p(0) = p \in S_M^0$  there exists  $0 < \theta^* < 1$  such that for  $0 < \theta < \theta^*$ , the learning algorithm under the conditions (C.1) - (C.2), (A.1) - (A.2) and (S.1) - (S.2) is such that

$$\lim_{k \rightarrow \infty} |E[\eta(k)] - J_1| < \epsilon \quad (30)$$

**Proof:** From corollary 1 it follows that the limit in the left hand side of (30) exists

From (25)



That is,

$$\lim_{k \rightarrow \infty} E[n(k)] = \Gamma_1(p) d_1 + \sum_{j=1}^J \Gamma_j(p) d_j$$

$$\lim_{k \rightarrow \infty} E[n(k) - d_1] \leq (1 - \Gamma_1(p)) (d_2 - d_1)$$

Since  $\theta[y, p] \leq \Gamma_1(p)$ , Combining this with the above,

$$\lim_{k \rightarrow \infty} |E[n(k) - d_1]| \leq [1 - \theta[y, p]] (d_2 - d_1)$$

For any  $p \in S_M^0$ , we know that

$$\lim_{\theta \rightarrow 0} \theta[y, p] = 1$$

Combining these we see that for any  $\epsilon > 0$  there exists  $0 < \theta^* < 1$  such that for all  $0 < \theta < \theta^*$

$$\lim_{k \rightarrow \infty} |E[n(k) - d_1]| < \delta \quad |(d_2 - d_1)|$$

The theorem follows by choosing  $\delta = \epsilon / |(d_2 - d_1)|$ .

**5. Conclusions:** A new class of absorbing barrier absolutely expedient learning algorithms, whose behaviour under the action of success and failure are identical, is introduced. Conditions for the  $\epsilon$ -optimality of this class of learning algorithms are derived.

## 6. References:

- [1] K.S. Narendra and M.A.L. Thatachar. "Learning Automata - A survey." IEEE Transactions on Systems, Man, and Cybernetics. Vol. 4, pp. 323-334, 1974.
- [2] K.S. Narendra and S. Lakshmivarahan. "Learning Automata - A critique." Journal of Cybernetics and Information Sciences - Special Issue on Learning Automata, Vol. 1, pp. 53-66, 1978.
- [3] S. Lakshmivarahan and M.A.L. Thatachar. "Absolutely Expedient Learning Algorithms for Stochastic Automata." IEEE Transactions on Systems, Man and Cybernetics, Vol. 3, pp. 281-286, 1973.
- [4] S. Lakshmivarahan and M.A.L. Thatachar. "Bounds on the Probability of Convergence of Learning Automata." IEEE Transactions on Systems, Man and Cybernetics. Vol. 6, pp. 756-763, 1976.
- [5] H. Aso and M. Kimura. "Absolute Expediency of Learning Automata." Information Sciences, Vol. 17, pp. 91-112, 1979.
- [6] U. Herkenrath, D. Kalin and S. Lakshmivarahan. "On a General Class of Absorbing Barrier Learning Algorithms." School of EECS Technical report - 8001 University of Oklahoma, March, 1980. (Also in Information Sciences Vol. 24 pp. 255-263, 1981.
- [7] Y. Sawaragi and N. Baba. "Two  $\epsilon$ -optimal Non-Linear Reinforcement Schemes for Stochastic Automata." IEEE Transactions on Systems, Man and Cybernetics. Vol. 4, pp. 126-131, 1974.
- [8] K.S. Fu. "Stochastic Automata Models for Learning Systems." in Computers and Information Sciences II (Ed) by J.T. Tou, Academic Press, 1967.

- [9] G.J. McMurtry and K.S. Fu. "A Variable Structure Automaton Used in Multimodal Search Technique." IEEE Transactions on Automatic Control, Vol 11, pp. 379-387, 1966.
- [10] V.I. Varshavskii and I.P. Vorontsova. "On the Behaviour of Stochastic Automata with Variable Structure," Automation and Remote Control, Vol. 24, pp. 327-333, 1963.
- [11] M.L. Tsetlin. "Automaton Theory and Modelling of Biological Systems" Academic Press, 1973.
- [12] M.F. Norman. "On a Linear Model with Two absorbing Barriers." Journal of Mathematical Psychology. Vol. 5, pp. 225-241, 1968.
- [13] M.F. Norman. Markov Processes and Learning Models. Academic press, New York, 1972.
- [14] M. Iosifescu and R. Theodorescu. "Random Process and Learning." Springer Verlag, 1969.
- [15] S. Lakshmivarahan. " $\epsilon$ -optimal Learning Algorithms - Non Absorbing Barrier Type." Technical Report, EECS 7901, School of EECS, University of Oklahoma, February, 1979.
- [16] I.H. Witten. "The Apparent Conflict Between Estimation and Control - A Survey of Two Armad-Bandit Problem." Journal of Franklin Institute Vol. 301, pp. 161-189, 1976.
- [17] B. Chandrasekaran and D.W.C. Shen. "On Expediency and Convergence in Variable Structure Automata." IEEE Transactions on Systems Science and Cybernetics. Vol. 4, pp. 52-60, 1968.
- [18] R.R. Bush and F. Mosteller. Stochastic Models for Learning. John Wiley, New York, 1958.
- [19] R. Viswanathan and K.S. Narendra. "Expedient and Optimal Variable Structure Stochastic Automata." Dunham Lab. TR-37, 1970. Yale University.
- [20] I.J. Shapiro and K.S. Narendra. "Use of Stochastic Automata for Parameter Self Optimisation with Multimodal Performance Criteria." IEEE Transactions on Systems, Man and Cybernetics. Vol. 5, pp. 352-360, 1969.
- [21] S. Lakshmivarahan. Learning Algorithms: Theory and Applications. Springer Verlag, New York, 1981.
- [22] J.L. Dobb. Stochastic Processes. John Wiley, 1955.
- [23] M.R. Meybodi and S. Lakshmivarahan. " $\epsilon$ -Optimality of a general class of absorbing barrier learning algorithms." School of EECS Technical Report, August 1981. University of Oklahoma, Norman, Oklahoma, U.S.A.