

INCREMENTAL LEARNING ON RANDOM TRIALS

by

M. Frank Norman<sup>1</sup>

TECHNICAL REPORT NO. 62

December 9, 1963

PSYCHOLOGY SERIES

Reproduction in Whole or in Part is Permitted for  
any Purpose of the United States Government

INSTITUTE FOR MATHEMATICAL STUDIES IN THE SOCIAL SCIENCES

STANFORD UNIVERSITY

Stanford, California

0  
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99

## Abstract

When a subject is required to learn to make a certain response in the course of a sequence of consistently corrected trials, and when it is not clear which of success or failure is more efficacious for learning, the possibility suggests itself that acquisition might sometimes proceed incrementally but with a certain probability, constant over trials, that no learning occurs on a trial. In this paper many of the properties which one would expect acquisition of this type to possess if the incremental process were linear in the response probabilities are derived, and details of an application of this model to a paired-associate experiment are presented.

0  
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99

In a previous paper (Norman, 1963) an attempt was made to interpolate between the all-or-none model (see Bower, 1961) and the single-operator linear model (see Bush and Sternberg, 1959) by means of a "two-phase" model. In this model a random number of trials on which no change in response probability occurred preceded the beginning of the linear incremental learning process described by the single-operator model. In the present paper an interpolation of a slightly different kind is discussed; all applications of a linear operator to response probabilities are preceded by a random number of trials on which no learning occurs. This model will be referred to below as the random-trial incremental model. Many of its mathematical properties will be derived, and data from a paired-associate learning experiment will be presented to illustrate its usefulness.

A delightful feature of the single-operator and all-or-none models is their mathematical simplicity. It will emerge from the development below that this simplicity is shared in significant measure by the random-trial incremental model.

#### The Model

Let us assume that a subject's probability  $\underline{p}_{\underline{w}\underline{n}}$  of making the  $\underline{A}_2$  (error) response on trial  $\underline{n}$  satisfies the stochastic difference equation

$$(1) \quad \underline{p}_{\underline{w}\underline{n}+1} = \underline{a}^{\underline{y}\underline{n}} \underline{p}_{\underline{w}\underline{n}} \quad \underline{n} = 1, 2, \dots,$$

and the initial condition

$$\underline{p}_{w1} = \underline{p}$$

where  $\{y_{wn}\}$  is a sequence of identically and independently distributed random variables with

$$(2) \quad P(y_{w1} = 1) = \underline{c}$$

and

$$P(y_{w1} = 0) = 1 - \underline{c}.$$

It follows immediately that, for  $n \geq 1$ ,

$$(3) \quad \underline{p}_{wn} = \underline{a}^{Y_{wn}-1} \underline{p}$$

where

$$Y_{wk} = \sum_{i=1}^k Y_{wi} \quad \text{for } k \geq 1$$

and

$$Y_{w0} = 0.$$

For  $k \geq 1$ ,  $Y_{wk}$  clearly has the binomial distribution with parameters  $k$  and  $\underline{c}$ .

We must, of course, have  $0 \leq \underline{a}, \underline{c}, \underline{p} \leq 1$ . The reader will note that, when  $\underline{a} = 0$ , the random-trial incremental model reduces to the all-or-none model, whereas, if  $\underline{c} = 1$ , it reduces to the single-operator linear model. If the sequence  $\{y_{wn}\}$  is interpreted as the sequence of reinforcement indicator random variables for a non-contingent reinforcement schedule such that the probability of reinforcement

of  $A_1$  is  $c$  on every trial then the random-trial incremental model is just the identity-operator model (see Bush and Mosteller, 1955) used with such a schedule. Some of the mathematical results obtained below are new in this context and it is hoped that they will be useful in this connection even though this was not the motivation for the research discussed in this paper. While the mathematical work was being done I was thinking of the sequence  $\{y_{\underline{m}}\}$  as unobservable, and I was not able to identify the random process embodied in this sequence with any general psychological process or observable aspect of the experimental situation in the application of the model to be presented below.

#### The Basic Theorem

Most of the properties of the random-trial incremental model which I will discuss in this paper are readily obtained from the following fundamental result:

Let  $k$  be a positive integer and let  $i_1, i_2, \dots, i_k$  be a strictly increasing finite sequence of positive integers. Then

$$(4) \quad P\left(x_{\underline{w}1_1} = 1, x_{\underline{w}1_2} = 1, \dots, x_{\underline{w}1_k} = 1\right) = \underline{c}^k \prod_{j=1}^k b_{\underline{j}}^{i_{k-j+1} - i_{k-j}}$$

where

$$b_{\underline{j}} = \underline{a}^{\underline{j}} \underline{c} + (1 - \underline{c}) \quad \text{for } \underline{j} = 1, 2, \dots,$$

$i_0 = 1$ , and  $x_{\underline{w}m}$  is the error indicator random variable for trial  $\underline{m}$ .

The following proof proceeds by induction on  $\underline{k}$ .

1° If  $\underline{k} = 1$ , (4) reduces to

$$(5) \quad \underline{P}(\underline{x}_{\underline{w}_1} = 1) = \underline{p} \underline{b}_{\underline{1}}^{i-1}$$

That this is correct is seen as follows:

$$\begin{aligned} \underline{P}(\underline{x}_{\underline{w}_1} = 1) &= \sum_{\underline{n}=0}^{i-1} \underline{P}(\underline{x}_{\underline{w}_1} = 1 \mid \underline{Y}_{\underline{w}_1-1} = \underline{n}) \underline{P}(\underline{Y}_{\underline{w}_1-1} = \underline{n}) \\ &= \sum_{\underline{n}=0}^{i-1} \underline{a}^{\underline{n}} \underline{p} \binom{i-1}{\underline{n}} \underline{c}^{\underline{n}(1-\underline{c})} \underline{b}_{\underline{1}}^{i-1-\underline{n}} \\ &= \underline{p}(\underline{ac} + (1-\underline{c}))^{i-1} = \underline{p} \underline{b}_{\underline{1}}^{i-1} \end{aligned}$$

by the binomial theorem.

2° Suppose that (4) is correct for any collection of  $\underline{k}$  indices

$1 \leq \underline{m}_1 < \underline{m}_2 < \dots < \underline{m}_{\underline{k}}$  and all parameter values  $0 \leq \underline{a}, \underline{c}, \underline{p} \leq 1$ .

(In what follows, I will indicate the dependence of  $\underline{P}$  on the initial error probability explicitly by means of a subscript.) Let some set

$1 \leq \underline{i}_1 < \underline{i}_2 < \dots < \underline{i}_{\underline{k}+1}$  of  $\underline{k} + 1$  indices and some parame

values  $\underline{a}, \underline{c}$ , and  $\underline{p}$  be given. Then

$$\underline{P}_{\underline{p}}(\underline{x}_{\underline{w}_1} = 1, \dots, \underline{x}_{\underline{w}_1-\underline{k}} = 1, \underline{x}_{\underline{w}_1-\underline{k}+1} = 1)$$

$$= \sum_{\underline{m}=0}^{\underline{i}_1-1} \underline{P}_{\underline{p}}(\underline{x}_{\underline{w}_1} = 1, \dots, \underline{x}_{\underline{w}_1-\underline{k}+1} = 1 \mid \underline{Y}_{\underline{w}_1-1} = \underline{m}) \underline{P}(\underline{Y}_{\underline{w}_1-1} = \underline{m})$$

$$= \sum_{\underline{m}=0}^{\underline{i}_1-1} \underline{P}_{\underline{a}\underline{p}}^{\underline{m}}(\underline{x}_{\underline{w}_1} = 1, \underline{x}_{\underline{w}_1-\underline{i}_2+1} = 1, \dots, \underline{x}_{\underline{w}_1-\underline{i}_1+1} = 1) \underline{P}(\underline{Y}_{\underline{w}_1-1} = \underline{m})$$



for it is a property of the random-trial incremental model that after the error probability on trial  $i_1$  has been specified the sequence  $\frac{x_{i_1}}{w_{i_1}}, \frac{x_{i_1+1}}{w_{i_1+1}}, \frac{x_{i_1+2}}{w_{i_1+2}}$  is stochastically identical to the sequence  $\frac{x_1}{w_1}, \frac{x_2}{w_2}, \frac{x_3}{w_3}$  with the specified error probability on trial 1. Noting next that  $\frac{x_{i_1}}{w_{i_1}}$  is independent of the  $\frac{x_i}{w_i}$ 's with  $i > 1$  we obtain

$$\begin{aligned} P_p \left( \frac{x_{i_1}}{w_{i_1}} = 1, \dots, \frac{x_{i_1-k}}{w_{i_1-k}} = 1, \frac{x_{i_1-k+1}}{w_{i_1-k+1}} = 1 \right) &= \\ &= \sum_{m=0}^{i_1-1} \frac{p}{a-p} \left( \frac{x_{i_1}}{w_{i_1}} = 1 \right) \frac{p}{a-p} \left( \frac{x_{i_1-2}}{w_{i_1-2}} = 1, \dots, \frac{x_{i_1-k+1}}{w_{i_1-k+1}} = 1 \right) P \left( \frac{Y_{i_1-1}}{w_{i_1-1}} = m \right) = \\ &= \sum_{m=0}^{i_1-1} \frac{p}{a-p} \left( \frac{a-p}{p} \right)^k \left( \prod_{j=1}^k \frac{b_{i_1-k-j+1}^{i_1-k-j}}{b_{i_1-k-j+1}^{i_1-k-j}} \right) \binom{i_1-1}{m} \frac{c^m (1-c)^{i_1-1-m}}{a-p} \end{aligned}$$

by the induction hypothesis where  $\frac{i_h}{w_h} = \frac{i_{h+1}}{w_{h+1}} = 1$

$$\begin{aligned} &= p^{k+1} \left( \prod_{j=1}^k \frac{b_{i_1-k-j+1}^{i_1-k-j}}{b_{i_1-k-j+1}^{i_1-k-j}} \right) \sum_{m=0}^{i_1-1} \binom{i_1-1}{m} \frac{c^m (1-c)^{i_1-1-m}}{a-p} \\ &= p^{k+1} \prod_{j=1}^{k+1} \frac{b_{i_1-k-j+1}^{i_1-k-j}}{b_{i_1-k-j+1}^{i_1-k-j}} \end{aligned}$$

An application of the principle of induction completes the proof of (4).

#### Corollaries of the Basic Theorem

Let  $J$  be the number of errors before the first success.

Then  $J$  is nonnegative and for  $k \geq 1$ ,  $P(J \geq k) = P\left(\frac{x_1}{w_1} = 1, \dots, \frac{x_k}{w_k} = 1\right)$ .

Therefore,

$$(6) \quad P(\underline{J} \geq \underline{k}) = p^{\underline{k}} \prod_{j=1}^{\underline{k}-1} b_{\underline{j}}.$$

For  $\underline{k} \geq 1$  let  $u_{\underline{w}\underline{k}} = \sum_{n=1}^{\infty} \prod_{i=0}^{\underline{k}-1} x_{\underline{w}\underline{n}+i}$  be the total number of  $\underline{k}$  tuples of consecutive errors in an infinite sequence of trials. Then

$$\begin{aligned} E(u_{\underline{w}\underline{k}}) &= \sum_{n=1}^{\infty} P(x_{\underline{w}\underline{n}} = 1, \dots, x_{\underline{w}\underline{n}+\underline{k}-1} = 1) \\ &= \sum_{n=1}^{\infty} p^{\underline{k}} \left( \prod_{j=1}^{\underline{k}-1} b_{\underline{j}} \right) b_{\underline{k}}^{n-1}. \end{aligned}$$

Thus

$$(7) \quad E(u_{\underline{w}\underline{k}}) = \frac{p^{\underline{k}} \prod_{j=1}^{\underline{k}-1} b_{\underline{j}}}{1 - b_{\underline{k}}}$$

(where  $\prod_{j=1}^0 b_{\underline{j}}$  is defined to be 1;  $u_{\underline{w}\underline{1}}$ , the total number of error in an infinite sequence of trials will also be denoted by  $\frac{T}{\underline{w}}$  in what follows).

For  $\underline{k} \geq 1$ , let  $c_{\underline{w}\underline{k}} = \sum_{n=1}^{\infty} x_{\underline{w}\underline{n}} x_{\underline{w}\underline{n}+\underline{k}}$  be the autocovariance of errors of lag  $\underline{k}$  (clearly  $c_{\underline{w}\underline{1}} = u_{\underline{w}\underline{2}}$ ). Then

$$\begin{aligned} \mathbb{E}(c_{\frac{k}{n}}) &= \sum_{n=1}^{\infty} P(x_{\frac{n}{n}} = 1, x_{\frac{n+k}{n}} = 1) \\ &= \sum_{n=1}^{\infty} p^2 \frac{b_1^k}{b_1} \frac{b_2^{n-1}}{b_2} . \end{aligned}$$

Therefore

$$(8) \quad \mathbb{E}(c_{\frac{k}{n}}) = \frac{p^2 \frac{b_1^k}{b_1}}{1 - \frac{b_2}{b_1}} .$$

Clearly

$$T^2 = \sum_{i=1}^{\infty} x_{\frac{i}{i}}^2 + 2 \sum_{k=1}^{\infty} \sum_{n=1}^{\infty} x_{\frac{n}{n}} x_{\frac{n+k}{n}} = \mathbb{E}(T) + 2 \sum_{k=1}^{\infty} c_{\frac{k}{n}} .$$

Thus

$$\begin{aligned} \mathbb{E}(T^2) &= \mathbb{E}(T) + 2 \sum_{k=1}^{\infty} \mathbb{E}(c_{\frac{k}{n}}) \\ &= \mathbb{E}(T) + 2 \sum_{k=1}^{\infty} \frac{p^2 \frac{b_1^k}{b_1}}{1 - \frac{b_2}{b_1}} \end{aligned}$$

by (8)

$$= \mathbb{E}(T) + \frac{2p^2}{1 - \frac{b_2}{b_1}} \frac{b_1}{1 - \frac{b_2}{b_1}} = \mathbb{E}(T) + \left( \frac{p}{1 - \frac{b_2}{b_1}} \right) \left( \frac{2}{p} \right) \left( \frac{p^2 \frac{b_1}{b_1}}{1 - \frac{b_2}{b_1}} \right) .$$

We conclude that

$$(9) \quad \mathbb{E}(T^2) = \mathbb{E}(T) + \frac{2}{p} \mathbb{E}(T) \mathbb{E}(u_2) .$$

Similar computations yield the formulas

$$(10) \quad \mathbb{E}(T u_2) = \mathbb{E}(T) \mathbb{E}(u_2) \left[ \frac{2}{p} (1 - \frac{b_2}{b_1}) + \frac{b_2(1 - \frac{b_2}{b_1})}{1 - \frac{b_2}{b_3}} + \frac{b_2(1 - \frac{b_2}{b_1})}{1 - \frac{b_2}{b_3}} \right]$$

and

$$(11) \quad \underline{E}(u_{\underline{w}2}^2) = \underline{E}(u_{\underline{w}2}) + 2\underline{E}(u_{\underline{w}3}) + \frac{2}{1-b_2} \underline{E}(u_{\underline{w}4})$$

which we will use in evaluating our estimation procedure below.

For  $n \geq 1$  denote by  $\underline{T}(\underline{n})$  the total number of errors after the  $\underline{n}$  th trial in an infinite sequence of trials, i.e.,

$$\underline{T}(\underline{n}) = \sum_{k=\underline{n}+1}^{\infty} \underline{x}_{\underline{w}k} . \quad \text{Then}$$

$$\begin{aligned} \underline{E}(\underline{T}(\underline{n}) | \underline{x}_{\underline{w}n} = 1) &= \sum_{k=\underline{n}+1}^{\infty} \frac{P(\underline{x}_{\underline{w}k} = 1, \underline{x}_{\underline{w}n} = 1)}{P(\underline{x}_{\underline{w}n} = 1)} \\ &= \sum_{k=\underline{n}+1}^{\infty} \frac{p^2 b_1^{k-n} b_2^{n-1}}{p b_1^{n-1}} \\ &= p \left( \frac{b_2}{b_1} \right)^{n-1} \frac{b_1}{1-b_1} . \end{aligned}$$

Thus

$$(12) \quad \underline{E}(\underline{T}(\underline{n}) | \underline{x}_{\underline{w}n} = 1) = \underline{E}(\underline{T}) b_1 \left( \frac{b_2}{b_1} \right)^{n-1} .$$

From the derivations given above, it is clear how one would go about deriving from (4) the expectation of any statistic which can be defined as a sum of finite products of  $\underline{x}_{\underline{w}n}$  's. In this connection it should be noted that formula (4) does not require for its validity an infinite sequence of trials.

The Distributions of the Total Number of Errors and the Trial of the Last Error

I will now derive the probability generating function<sup>2</sup> of  $\underline{T}$ ,

i.e., the function  $\frac{g_T}{w}$  given by

$$\frac{g_T}{w}(s) = \sum_{k=0}^{\infty} P_{\frac{w}{w}}(T = k) s^k.$$

We can write  $\frac{T}{w}$  as

$$\frac{T}{w} = \sum_{m=0}^{\infty} \frac{T}{w^m}$$

where  $\frac{T}{w^m}$  is the number of errors made when the error probability is  $\frac{m}{a-p}$ . The  $\frac{T}{w^m}$  are obviously mutually independent. Thus, if  $\frac{g_T}{w^m}$  is the probability generating function of  $\frac{T}{w^m}$ , we have

$$(13) \quad \frac{g_T}{w}(s) = \prod_{m=0}^{\infty} \frac{g_T}{w^m}(s).$$

But we can write  $\frac{T}{w^m}$  as

$$\frac{T}{w^m} = \sum_{j=1}^{S_{\frac{w}{w^m}}} x_{\frac{w}{w^m}, j}$$

where  $x_{\frac{w}{w^m}, j}$  is the error indicator random variable for the  $j^{\text{th}}$  trial on which the error probability is  $\frac{m}{a-p}$  and  $S_{\frac{w}{w^m}}$  is the number of trials on which the error probability is  $\frac{m}{a-p}$ . Since the  $x_{\frac{w}{w^m}, j}$  for  $j = 1, 2, \dots, S_{\frac{w}{w^m}}$  are identically distributed, mutually independent, and independent of  $S_{\frac{w}{w^m}}$  it follows that

$$(14) \quad \frac{g_T}{w^m}(s) = \frac{g_{S_{\frac{w}{w^m}}}}{w^m} (g_{x_{\frac{w}{w^m}, 1}}(s)),$$

where  $\underline{g}_{\underline{S}}^{\underline{m}}$  and  $\underline{g}_{\underline{X}}^{\underline{m},1}$  are the probability generating functions of the subscript random variables. Since  $\underline{S}$  has the geometric distribution with parameter  $\underline{c}$  (thus  $\underline{g}_{\underline{S}}^{\underline{m}}(u) = \underline{c}u/1-(1-\underline{c})u$ ) and  $\underline{X}_{\underline{m},1}$  has the binomial distribution with parameters 1 and  $\underline{a}-\underline{p}$  (thus  $\underline{g}_{\underline{X}}^{\underline{m},1}(\underline{s}) = \underline{a}-\underline{p}\underline{s} + (1-\underline{a}-\underline{p})$ ), we obtain from (13) and (14)

$$(15) \quad \underline{g}_{\underline{T}}^{\underline{m}}(\underline{s}) = \prod_{\underline{m}=0}^{\infty} \frac{\underline{c}[\underline{a}-\underline{p}\underline{s} + (1-\underline{a}-\underline{p})]}{1 - (1-\underline{c})[\underline{a}-\underline{p}\underline{s} + (1-\underline{a}-\underline{p})]}.$$

In particular,

$$(16) \quad \underline{P}_{\underline{p}}^{\underline{m}}(\underline{T} = 0) = \underline{g}_{\underline{T}}^{\underline{m}}(0) = \prod_{\underline{m}=0}^{\infty} \frac{\underline{c}(1-\underline{a}-\underline{p})}{1 - (1-\underline{c})(1-\underline{a}-\underline{p})}.$$

The distribution of  $\underline{N}$ , the trial number of the last error in an infinite sequence of trials, is easily obtained from (16). For  $\underline{n} \geq 0$

$$\begin{aligned} \underline{P}_{\underline{p}}^{\underline{m}}(\underline{N} \leq \underline{n}) &= \underline{P}_{\underline{p}}^{\underline{m}}(\underline{x}_{\underline{m}+1} = 0, \underline{x}_{\underline{m}+2} = 0, \dots) \\ &= \sum_{\underline{k}=0}^{\underline{n}} \underline{P}_{\underline{p}}^{\underline{m}}(\underline{x}_{\underline{m}+1} = 0, \underline{x}_{\underline{m}+2} = 0, \dots | \underline{Y}_{\underline{m}} = \underline{k}) \underline{P}_{\underline{p}}^{\underline{m}}(\underline{Y}_{\underline{m}} = \underline{k}) \\ &= \sum_{\underline{k}=0}^{\underline{n}} \underline{P}_{\underline{a}-\underline{p}}^{\underline{k}}(\underline{x}_{\underline{m}+1} = 0, \underline{x}_{\underline{m}+2} = 0, \dots) \binom{\underline{n}}{\underline{k}} \underline{c}^{\underline{k}}(1-\underline{c})^{\underline{n}-\underline{k}}. \end{aligned}$$

So

$$(17) \quad \underline{P}_{\underline{p}}^{\underline{m}}(\underline{N} \leq \underline{n}) = \sum_{\underline{k}=0}^{\underline{n}} \left( \prod_{\underline{m}=0}^{\infty} \frac{\underline{c}(1-\underline{a}-\underline{p})}{1 - (1-\underline{c})(1-\underline{a}-\underline{p})} \right) \binom{\underline{n}}{\underline{k}} \underline{c}^{\underline{k}}(1-\underline{c})^{\underline{n}-\underline{k}}.$$

### Application to Paired-Associate Learning

In a recent experiment<sup>3</sup> conducted by Patrick Suppes and Madeleine Schlag-Rey, each of 40 college students learned a 12 item list of paired-associates (this part of the experiment will be referred to as "the first session" or S1 below) and then learned another such list 7 days later, on the average (during the "second session" or S2). The stimuli were CVC nonsense syllables and the appropriate response to each was a press on one of three keys available before the subject. (Thus the present experiment differs from conventional paired-associate learning experiments both in the number and nature of the response alternatives available to the subject.) The order of presentation of the 12 stimuli was randomized over successive presentations. The steps of the experimental routine were as follows: A nonsense syllable was displayed before a subject. As soon as he wished (the average latency was 1.2 seconds) the subject pressed a key. A light immediately flashed over one of the keys indicating which key was correct. Four seconds later the next stimulus appeared. Each session began with three practice items after which the list to be learned was presented 25 times.

In the analyses to follow all  $40 \times 12 = 480$  subject-items within each session are assumed to be stochastically independent and identical. Predictions of the random-trial incremental model for an infinite sequence of trials will be compared with data for the 25 trials of each experimental session. According to the model (consider Eq. 5 for  $i > 25$  and the parameter values used below) the error thus introduced is quite small.

Learning was noticeably faster in S2 than in S1 (the mean total numbers of errors for the two sessions, for instance, were 3.26 and 4.37, respectively) so the model was applied separately to the data from the two sessions. I will complete most of my discussion of the data from S2 before turning to that from S1.

The a priori estimate .6667 of  $\underline{p}$  was used in analyzing the data from S2. Defining the functions  $\underline{a}(\underline{x}, \underline{y})$  and  $\underline{c}(\underline{x}, \underline{y})$  by

$$(18) \quad \underline{a}(\underline{x}, \underline{y}) = \frac{\underline{p}(\underline{x}-\underline{p})}{\underline{y}} - 1$$

$$(19) \quad \underline{c}(\underline{x}, \underline{y}) = \frac{\underline{p}}{(2-\underline{p}(\underline{x}-\underline{p})/\underline{y})\underline{x}}$$

it can be shown by elementary computations that

$$(20) \quad \underline{a}(\underline{E}(\underline{T}), \underline{E}(\underline{u}_2)) \equiv \underline{a}$$

and

$$(21) \quad \underline{c}(\underline{E}(\underline{T}), \underline{E}(\underline{u}_2)) \equiv \underline{c}.$$

Therefore the statistics

$$(22) \quad \hat{\underline{a}}_{\overline{w}} = \underline{a}(\overline{\underline{T}}, \overline{\underline{u}_2})$$

and

$$(23) \quad \hat{\underline{c}}_{\overline{w}} = \underline{c}(\overline{\underline{T}}, \overline{\underline{u}_2})$$

(where  $\overline{\underline{T}}$  and  $\overline{\underline{u}_2}$  are the sample average total numbers of errors and pairs of errors) are moments estimators of  $\underline{a}$  and  $\underline{c}$ . For the data from S2, (22) and (23) give  $\hat{\underline{a}}_{\overline{w}} = .1655$  and  $\hat{\underline{c}}_{\overline{w}} = .2454$ . The small



value of  $\hat{\frac{a}{w}}$  is satisfying since one would not expect too much difference between two and three response paired-associate learning, and the success of the all-or-none model in predicting the course of two response paired-associate learning (see Bower, 1961) indicates that the value of  $\underline{a}$  appropriate for the random-trial incremental model in such learning is very small. Estimators like  $\hat{\frac{a}{w}}$  and  $\hat{\frac{c}{w}}$  are well understood mathematically and reasonably well behaved at least in the large sample case.<sup>4</sup> Specifically (see Wilks, 1962, Theorem 9.3.1a) p. 260), these estimators are asymptotically normally distributed with asymptotic means  $\underline{a}$  and  $\underline{c}$  and asymptotic variances

$$(24) \text{ asy. var. } \hat{\frac{a}{w}} = \frac{1}{n} \left( \frac{a^2}{\underline{x}} \text{var}\left(\frac{T}{w}\right) + \frac{a^2}{\underline{y}} \text{var}\left(\frac{u_2}{w^2}\right) + 2\frac{a}{\underline{x}\underline{y}} \text{cov}\left(\frac{T}{w}, \frac{u_2}{w^2}\right) \right)$$

and

$$(25) \text{ asy. var. } \hat{\frac{c}{w}} = \frac{1}{n} \left( \frac{c^2}{\underline{x}} \text{var}\left(\frac{T}{w}\right) + \frac{c^2}{\underline{y}} \text{var}\left(\frac{u_2}{w^2}\right) + 2\frac{c}{\underline{x}\underline{y}} \text{cov}\left(\frac{T}{w}, \frac{u_2}{w^2}\right) \right)$$

where the partial derivatives are evaluated at  $(\underline{E}\left(\frac{T}{w}\right), \underline{E}\left(\frac{u_2}{w^2}\right))$  and  $\underline{n}$  is the number of subject-items. The number of subject-items in the present experiment is sufficiently large to justify consideration of asymptotic means and variances. Replacing all moments in (24) and (25) by the corresponding predictions of the model using the parameter estimates  $\hat{\frac{a}{w}} = .1655$  and  $\hat{\frac{c}{w}} = .2454$  and replacing  $\underline{n}$  by 480, we obtain

$$\sqrt{\text{asy. var. } \hat{\frac{a}{w}}} = .0297 \text{ and } \sqrt{\text{asy. var. } \hat{\frac{c}{w}}} = .0145$$

as approximations to the asymptotic standard deviations of  $\hat{\frac{a}{w}}$  and  $\hat{\frac{c}{w}}$

for the data from S2. Certainly these values are larger than one might desire. On the other hand, they are not so large that I will feel too many qualms when I use them below.

An extensive set of comparisons of the model with the S2 data is presented in Table 1. The fit is generally excellent, though not uniformly so. The mild bimodality in the observed distribution of  $\frac{N}{w}$  (the same phenomenon appears in the data from the first session) is certainly not predicted by the random-trial incremental model (or any other model that I know of).<sup>5</sup> Also, all of the predicted  $\underline{E}(\frac{c}{w^n})$ ,  $n \geq 2$  are smaller than the corresponding observations (this was also observed in S1) and the effect is moderately large for  $n = 4, 5, \text{ and } 6$ .

The reader may have noticed earlier that the theoretical expressions in the model for  $\underline{P}(\frac{x}{w^n} = 1)$ ,  $\underline{E}(\frac{c}{w^n})$ , and  $\underline{E}(\frac{T^{(n)}}{w^n} | \frac{x}{w^n} = 1)$  are exponential functions of  $n$  (see Eqs. 5, 8, and 12). It is easily seen directly from (4) that the same is true of  $\underline{P}(\frac{x}{w^n} = 1, \frac{x}{w^{n+1}} = 1)$  and  $\underline{P}(\frac{x}{w^n} = 1, \frac{x}{w^{n+1}} = 1, \frac{x}{w^{n+2}} = 1)$ . It is therefore particularly interesting to plot the logarithms of these quantities, which are linear in  $n$ , and the logarithms of the corresponding sample averages as functions of  $n$ . Such plots appear in Fig. 1.

Two analyses were made of the data from the first session, one using the estimate .6667 of  $p$  as above and the other using the estimate .6937, the proportion of the subject-items which had an error on trial 1. While the qualitative conclusions to be reached below would have been the same under the two analyses, only the second will be discussed below since it yielded a somewhat better fit. The results

Table 1

Session 2 Data and Predictions of the  
Random-Trial Incremental Model

$$(\underline{p} = .6667, \underline{a} = .1655, \underline{c} = .2454)$$

$\underline{P}(x_{\underline{w}n} = 1)$	Data	Model
$\underline{n} = 1$	.6666	.6667
2	.48	.53
3	.43	.42
4	.37	.34
5	.25	.27
6	.22	.21
7	.18	.17
8	.14	.13
9	.11	.11
10	.09	.08
11	.08	.07
12	.05	.05
13	.03	.04
$\underline{P}(x_{\underline{w}n} = 1, x_{\underline{w}n+1} = 1)$		
$\underline{n} = 1$	.34	.35
2	.26	.27
3	.23	.20
4	.14	.16
5	.10	.12
$\underline{P}(x_{\underline{w}n} = 1, x_{\underline{w}n+1} = 1, x_{\underline{w}n+2} = 1)$		
$\underline{n} = 1$	.19	.18
2	.16	.14
3	.10	.10
4	.06	.08
$\underline{E}(T)$	3.2562	*
$\underline{\sigma}(T)$	2.90	2.67
$\underline{E}(u_{\underline{w}k})$		
$\underline{k} = 2$	1.4812	*
3	.72	.73
4	.35	.37
5	.18	.19
6	.10	.09

Table 1 (cont.)

$\frac{E(c_{\bar{w}k})}{\bar{w}k}$		Data	Model
	$\bar{k} = 2$	1.22	1.18
	3	.97	.94
	4	.82	.74
	5	.70	.59
	6	.54	.47
$P(\bar{J} = \bar{k})$	$\bar{k} = 0$	.33	.33
	1	.33	.31
	2	.15	.17
	3	.08	.09
	4	.05	.04
	5	.02	.02
	$\geq 6$	.03	.02
$\frac{E(T_{\bar{w}}^{(n)}   x_{\bar{w}1} = 1)}$	$\bar{n} = 1$	2.91	2.59
	2	2.86	2.48
	3	2.60	2.37
	4	2.09	2.27
	5	2.06	2.17
	6	2.37	2.08
	7	2.15	1.99
	8	1.78	1.91
	9	2.10	1.83
$P(\bar{N} = \bar{k})$	$\bar{k} = 0$	.13	.07
	1	.12	.13
	2	.09	.12
	3	.09	.10
	4	.12	.09
	5	.09	.08
	6	.06	.07
	$\geq 7$	.29	.34

\* Used to estimate parameters.

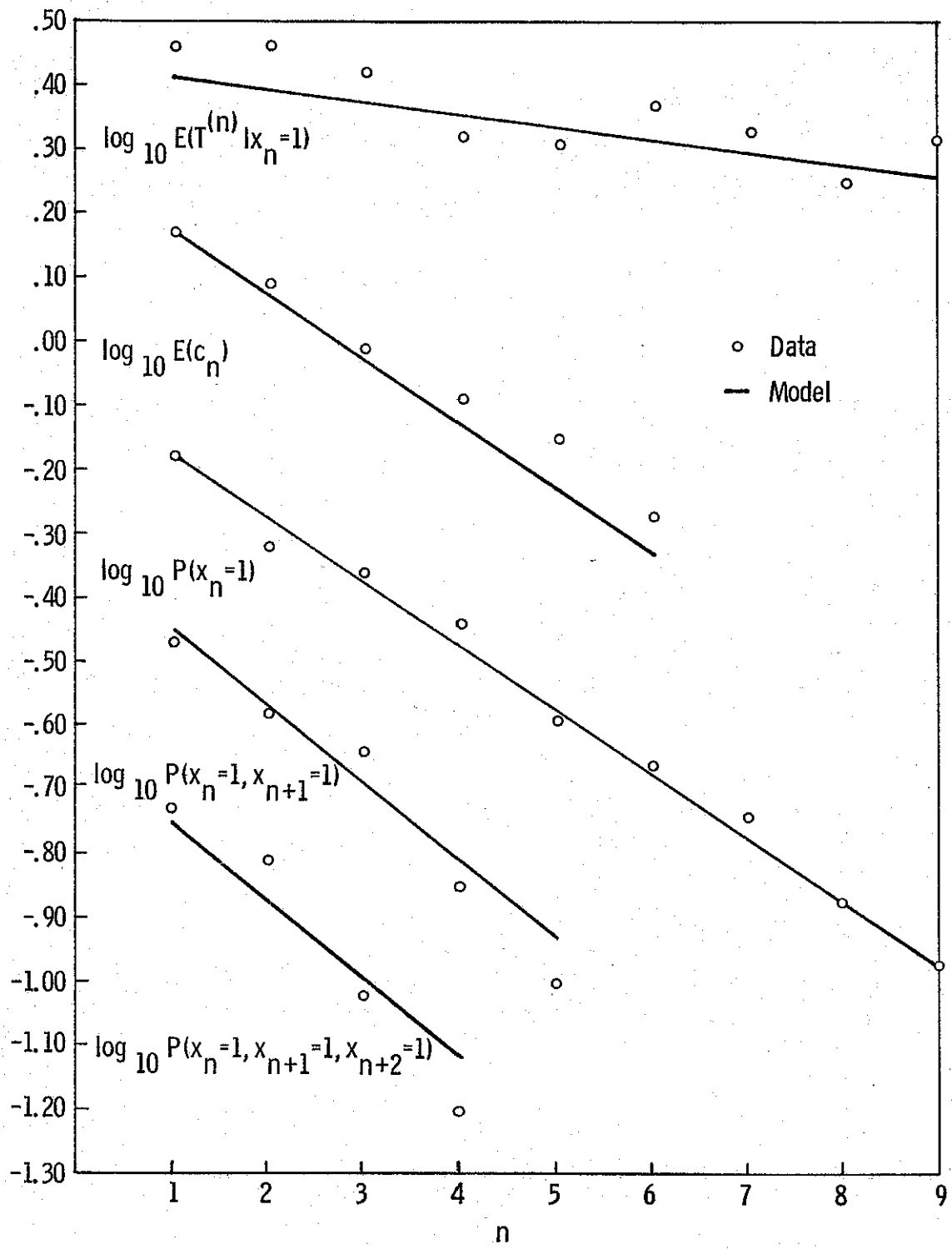
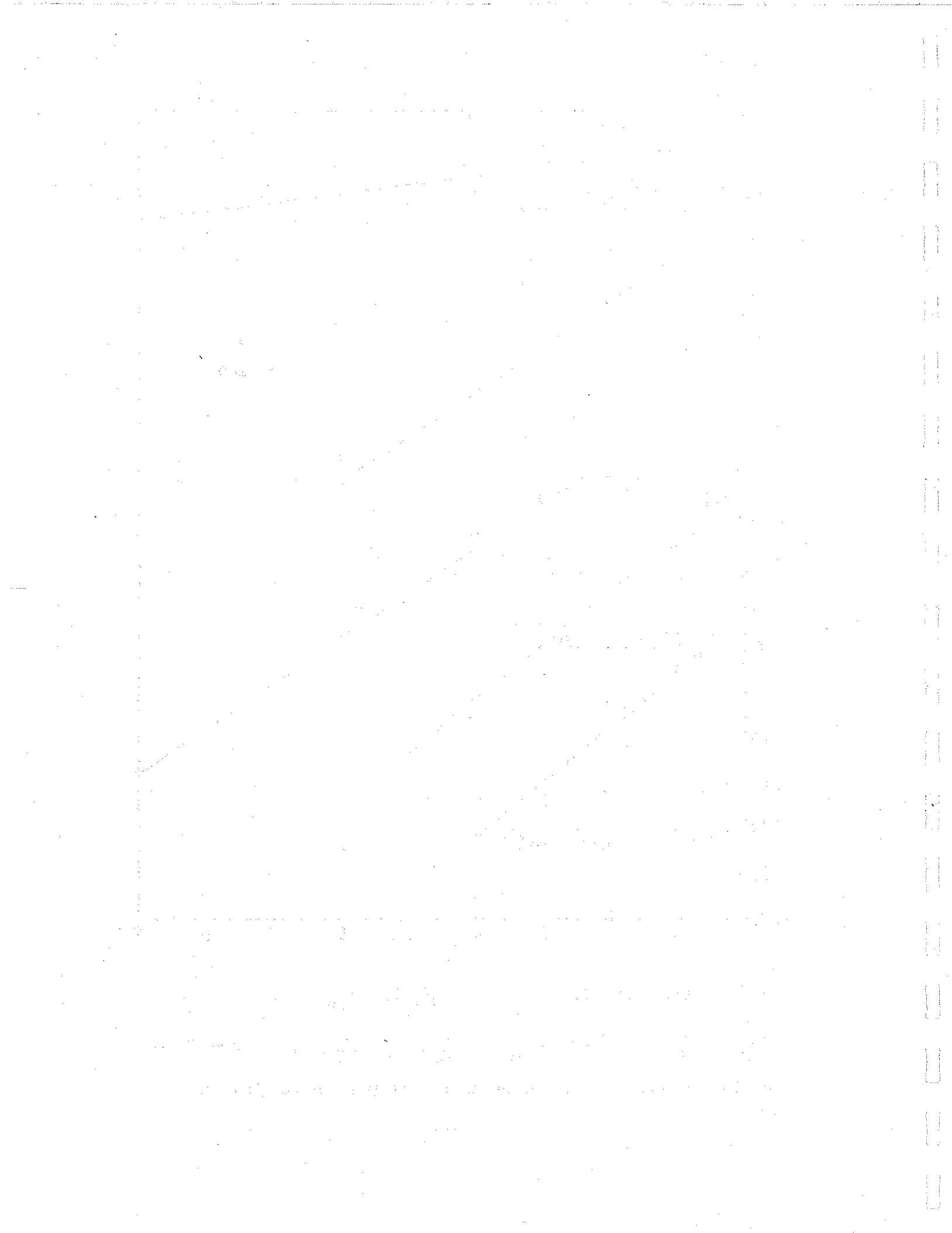


Fig. 1.  $\log_{10}$  of  $\underline{E}(T^{(n)} | \underline{x}_n = 1)$ ,  $\underline{E}(c_n)$ ,  $P(\underline{x}_n = 1)$ ,  $P(\underline{x}_n = 1, \underline{x}_{n+1} = 1)$  and  $P(\underline{x}_n = 1, \underline{x}_{n+1} = 1, \underline{x}_{n+2} = 1)$ : predictions of the random-trial incremental model and data for sessions 2.



are presented in Table 2.

Though the overall fit is not bad in the sense that the absolute values of the discrepancies between the model and the data are small in most cases, these discrepancies form a pronounced pattern. Examination of Table 2 shows that the model predicts too few errors, pairs of consecutive errors, and triples of consecutive errors early in learning. (Thus, as a consequence of the fact that the model predicts the total number of errors and pairs of consecutive errors throughout S1 exactly because of the way the parameters were estimated, it predicts too many errors and pairs of consecutive errors late in learning.) This generalization leads us to expect that the random-trial incremental model will overestimate the average number of errors following an error on trial  $n$  for  $n$  not too large. The denominator in the prediction will tend to be smaller than the comparable quantity in the data, and the numerator will tend to be larger than the comparable quantity in the data (see the first line in the derivation of Eq. 12). From Table 2 we see that this effect was obtained on trials 2-9.<sup>6</sup>

Such a pattern of errors would be expected if the random-trial incremental model were correct, but for the fact that its parameters were changing as learning proceeded. The average total number of errors in S2 was in fact a good bit smaller than the comparable statistic for S1. A natural inference is that the subjects were learning to learn gradually throughout the experiment. If such second order learning were negatively accelerated, a seemingly reasonable assumption, then its effects might have been negligible in S2. But under the further

Table 2

Session 1 Data and Predictions of the

Random-Trial Incremental Model

 $(p = .6937, a = .1215, c = .1805)$ 

	Data	Model
$\underline{P}(x_{\underline{w}n} = 1)$		
$\underline{n} = 1$	.6937	*
2	.56	.58
3	.50	.49
4	.44	.41
5	.41	.35
6	.30	.29
7	.29	.25
8	.21	.21
9	.19	.17
10	.14	.15
11	.14	.12
12	.11	.10
13	.08	.09
$\underline{P}(x_{\underline{w}n} = 1, x_{\underline{w}n+1} = 1)$		
$\underline{n} = 1$	.42	.40
2	.35	.33
3	.31	.27
4	.24	.22
5	.20	.18
$\underline{P}(x_{\underline{w}n} = 1, x_{\underline{w}n+1} = 1, x_{\underline{w}n+2} = 1)$		
$\underline{n} = 1$	.26	.23
2	.22	.19
3	.17	.16
4	.13	.13
$\underline{E}(T)$	4.3749	*
$\underline{\sigma}(T)$	3.40	3.64
$\underline{E}(u_{\underline{w}k})$		
$\underline{k} = 2$	2.2770	*
3	1.26	1.28
4	.71	.73
5	.40	.41
6	.24	.24



Table 2 (cont.)

		Data	Model
$\underline{E}(c_{\underline{w}\underline{k}})$	$\underline{k} = 2$	1.94	1.92
	3	1.64	1.61
	4	1.44	1.36
	5	1.19	1.14
	6	.97	.96
$\underline{P}(J = \underline{k})$	$\underline{k} = 0$	.31	.31
	1	.28	.29
	2	.16	.17
	3	.09	.10
	4	.08	.06
	5	.04	.03
	$\geq 6$	.05	.04
$\underline{E}(T_{\underline{w}}^{(n)}   x_{\underline{w}n} = 1)$	$\underline{n} = 1$	3.86	3.68
	2	3.61	3.60
	3	3.40	3.51
	4	3.17	3.43
	5	2.90	3.36
	6	3.03	3.28
	7	2.85	3.20
	8	2.62	3.13
	9	2.66	3.06
$\underline{P}(N = \underline{k})$	$\underline{k} = 0$	.08	.05
	1	.10	.10
	2	.07	.10
	3	.06	.09
	4	.09	.08
	5	.10	.07
	6	.06	.06
	$\geq 7$	.42	.45

\* Used to estimate parameters.

assumption that our parameter estimates for S1 were appropriate for approximately the middle of the session, and there is no reason to doubt the validity of this as a first order approximation, we would expect the model to predict too few errors early in learning and too many later in learning, just as was observed.

## References

- Bower, G. H. Application of a model to paired-associate learning. Psychometrika, 1961, 26, 255-280.
- Bush, R. R., & Mosteller, F. Stochastic models for learning. New York: Wiley, 1955.
- Bush, R. R., & Sternberg, S. H. A single operator model. In R. R. Bush & W. K. Estes (Eds.), Studies in mathematical learning theory. Stanford: Stanford Univer. Press, 1959. Pp. 204-214.
- Feller, W. An introduction to probability theory and its applications Vol. 1 (2nd ed.) New York: Wiley, 1957.
- Norman, M. F. A two-phase model and an application to verbal discrimination learning. In R. C. Atkinson (Ed.), Studies in mathematical psychology. Stanford: Stanford Univer. Press, 1963. Pp. 173-187.
- Wilks, S. S. Mathematical statistics. New York: Wiley, 1962.

## Footnotes

<sup>1</sup>I am grateful to Professors R. C. Atkinson, G. H. Bower, and P. Suppes for their encouragement during this research. The preparation of this report was supported by Air Force Contract AF 49(638)-1253.

<sup>2</sup>The reader may consult Chs. 11 and 12 of Feller (1957) for the properties of generating functions used in this section.

<sup>3</sup>The data needed for the analysis presented below were generously supplied by Professor Suppes.

<sup>4</sup>The comment about two similar estimators on p. 182 of Norman (1963) should be ignored.

<sup>5</sup>A referee has pointed out that this bimodality may be indicative of a violation of the assumption of homogeneity of subject-items.

<sup>6</sup> $P(\frac{x}{wn} = 1), \underline{n} = 1, 2, \dots, 14$  and  $\underline{E}(\frac{T^{(n)}}{wn} | \frac{x}{wn} = 1), \underline{n} = 1, 2, \dots, 8$  have been computed for the random-trial incremental model (parameter estimation as with the data from S2) for Bower's four response verbal discrimination learning experiment described in Norman (1963), and the same pattern is observed: the model predicts too few errors on the first few trials and too many on later trials while badly overestimating the average number of errors after an error on trial  $\underline{n}$  for  $\underline{n} = 3, 4, \dots, 8$ .

Erratum

The phrase "...and the numerator will tend to be larger than the comparable quantity in the data..." on p. 15 should be replaced by:

"...while the deficiencies of the model on the earlier and later trials will tend to cancel in the numerator..."