

MATHEMATICAL LEARNING THEORY

by

Richard C. Atkinson and Robert C. Calfee

TECHNICAL REPORT NO. 50

January 2, 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

Applied Mathematics and Statistics Laboratories

STANFORD UNIVERSITY

Stanford, California

U
N
I
T
E
D
N
A
T
I
O
N
A
L
A
R
M
Y
O
F
S
O
C
I
A
L
J
U
S
T
I
C
E

UNITED NATIONS

GENERAL ASSEMBLY

RESOLUTION

1981

1981

1981

1981

1981

1981

1981

1981

MATHEMATICAL LEARNING THEORY*

by

Richard C. Atkinson and Robert C. Calfee

Stanford University

Mathematical learning theory probably began in 1885 when Herman Ebbinghaus attempted to fit some data from an experiment on recall with a simple logarithmic function. However, the term has come to be associated closely with a number of recent developments in psychology, and it is these developments and their immediate historical antecedents that we will consider in this paper. In particular, we will discuss the role of mathematical models in contemporary learning theory, with special reference to the influence of such models on research and theory. A number of points that we wish to make clear will then be illustrated by a detailed consideration of a specific experiment. In this example, we will show how the design and analysis of experiments are related to a mathematical approach to learning.

Prior to 1950, the most significant attempt to formulate a mathematical theory of learning was that of Clark L. Hull. His theoretical system and variations of it (Hull, 1943; Spence, 1956; Logan, 1959) were based on the postulation of a set of unobservable intervening variables, psychological constructs such as habit and drive. These intervening variables were related to the observable dependent and independent variables by certain functions that were to be empirically determined. Within the Hullian framework behavior basically was a deterministic

*

The preparation of this document was supported by the National Institute of Health (Grant M-5184) and by the National Science Foundation (Grant 24264). The paper is a contribution to a forthcoming book edited by E. Nagel and B. B. Wolman entitled "Psychology and the Theory of Science" (Basic Books, Inc.).

process, though a probabilistic overlay was added. All response measures were functions of a single underlying factor, reaction potential, and it was usually assumed that, since all response measures should be correlated, the experimenter was free to choose the measure that he thought most appropriate.

In general, experiments designed to test theories in the pre-1950 period were of two types: the comparative experiment and the factorial experiment. (For a historical account of this period see Hilgard, 1956.) In the comparative experiment, a situation was arranged in which opposing predictions could be made by different theories. Few of the experiments proved to be as crucial as they were designed to be, since the protagonists were quite skillful at finding an interpretation of the theory that accounted for the results, and the theories themselves were quite resilient. The consequence of this experimentation was that it shortly became difficult to distinguish one theoretical system from another (Seward, 1956).

As an example of the factorial experiment, we may mention the efforts to determine whether drive and incentive combine additively or multiplicatively. If the latter condition holds, one would expect a significant interaction term when a factorial design is analyzed by the analysis of variance technique. (A factorial design is one in which several levels of each variable are represented in all possible combinations.) A non-significant interaction would be interpreted as evidence for the additive hypothesis. Clearly, the outcome depended on numerous conditions other than the assumption being tested, such as the choice of levels of the variables, the degree of experimental control, and probably the response measure chosen. In particular, to the extent that experimental control was poor, the additive hypothesis would be favored.

In the late 1940's and early 50's, there appeared a group of new developments that have come to be called mathematical learning theory. Let us at the outset state that we cannot hope to do justice to all the investigators whose work properly falls within this category. And we can only mention in passing that in addition to developing a number of new learning situations for their own purposes, experimenters have applied mathematical theories of one type or another to a wide variety of standard learning paradigms, such as classical conditioning, avoidance conditioning, discrimination learning, stimulus generalization, paired associate learning, memory processes and concept formation. (For recent reviews, see Bush, 1960; Estes, 1959, 1962; Restle, 1959.)

The movement has been characterized by a number of features. Behavior is seen as an essentially probabilistic phenomenon. The primary behavioral measure is taken to be the probability of occurrence of a member of some response class. Theories are stated in a way that has made mathematical development feasible. There has been a tendency to interpret behavioral phenomena, not by reference to underlying molecular processes, but by specification of the rules of operation (Estes, 1962). These rules are simple mathematical laws whose implications describe the overt response character of a behavioral system, much as Newton's laws describe the activity of the solar system.

In a sense, mathematical learning theory is a misnomer. One should not imagine that mathematical learning theory represents a position that is opposed to other learning theories. What is actually being expressed is an increased dependence upon the use of mathematics in the formulation of learning theory irrespective of whether the theory is

oriented toward stimulus-response notions, cognitive constructs, expectancies, or some other approach. As we shall see, issues that in the past have crucially differentiated opposing theoretical systems, when formulated in a precise mathematical fashion can live together quite comfortably within a single system. The use of mathematics has allowed the psychologist to analyze more adequately the content of his statements, and to determine whether a particular set of data are consistent with these statements. In particular, mathematical models have been of heuristic importance in the reformulation and extension of theory. A mathematical model consists of a set of axioms that are generated by the theory. By means of a calculus such as probability theory, predictions from the model are obtained by straightforward mathematical deduction. Even though the theory may be qualitatively stated, so that a given model is not the only one which might be interpretable from the theory, the model serves to make explicit the assumptions that are in fact being made. If additional or different assumptions become necessary, then the form of these, as well as their implications, become evident.

Two main lines of development in mathematical learning theory that appeared almost simultaneously are those associated with Bush and Mosteller (1951, 1955) and Estes (1950, 1959). Bush and Mosteller began with what Restle (1959) has called an abstract theory. For reasons of mathematical simplicity, they assumed that the probability of a given response on a trial could always be expressed as a linear function of the probability of the same response on the preceding trial. The form of the linear operator (i.e., the parameter values of the function) depended on the type of reinforcement event that intervened. Though the

theory was abstract, Bush and Mosteller showed that derivations and parameter estimation problems could be greatly simplified if certain restraints based on extra-theoretical considerations were imposed on the models. These considerations were of a sort that made sense psychologically. An example is the "equal alpha" condition, where for certain two-choice problems, (e.g., a T-maze) the symmetry of the situation permits the assumption that the learning rate parameters associated with the two responses are equal.

Estes' theoretical formulation which has come to be known as Stimulus Sampling Theory was of an entirely different form. The environment was represented by a large population of discrete, mutually exclusive conceptual entities which he called stimulus elements. Each element was conditioned to one and only one response class. The organism took a sample from the set of elements, and the probability of occurrence of any response class was simply the proportion of elements in the sample which were conditioned to that class. The reinforcement event acted upon the conditioning relations of the sample of elements in some specified fashion (e.g., all elements in the sample became conditioned to the response class which was designated as correct) and the sample was then returned to the population for resampling at a later time. In Estes' initial development of these notions, the function describing changes in response probability was a linear function similar to the Bush-Mosteller model.

Let us now mention a number of ways in which these systems represented advances over earlier formulations. A feature of psychological theories from Freud through Hull has been the postulation of multiple unobservable processes, that may interact in some complicated fashion

to either complement or oppose each other. At a qualitative level of analysis, by a suitable post facto weighting of such processes, one may account for virtually any experimental result. Mathematical models have allowed theorists to introduce and evaluate such notions in an unambiguous manner. One may assume more than one underlying process, and then determine the contribution of each process unequivocally.

Another advance brought about by mathematical developments in learning theory concerns changes in the organization and analysis of data. In this regard, perhaps the most important role of mathematical models has been to provide a framework within which the detailed trial-by-trial aspects of behavior can be scrutinized (Anderson, 1959). An experiment designed only to establish the existence of a gross relation between two variables, such as response speed and reward magnitude, ignores the many sequential properties of psychological phenomena. Examination of these properties is a significant step forward in that it provides a source of behavioral information that cannot be obtained from an analysis of average performance curves. Theories stated only in qualitative terms do not provide an adequate means for analyzing and interpreting such complex sequential phenomena.

In this connection we may note Estes' distinction (1959) between three levels of prediction from a mathematical model: extrapolation, overdetermination, and situational invariance. Extrapolation refers to the ability of a model to account for those statistics from which parameters are estimated. (This requirement is not as simple as it may seem. For example, no choice of parameters for a linear equation will give a satisfactory fit to the typical learning curve.) Overdetermination

refers to prediction, within the same body of data, of statistics that are independent of those yielding the parameter estimates. Finally, situational invariance is the degree to which parameter estimates made in one experimental situation can account for data collected in other experiments. The mixed record of successes and failures of mathematical models shows that these requirements are not trivial. For example, in numerous studies (Anderson, 1960; Suppes and Atkinson, 1960; Atkinson and Estes, 1963) the learning rate parameters that satisfy the mean learning curve requirements have proven inadequate in accounting for asymptotic sequential dependency statistics.

Among recent trends within the area of stimulus sampling theory, we may mention the introduction of models where the number of stimulus elements in the population is severely limited, leading in the limit to the one-element model. In many learning situations, it is reasonable to suppose that the subject does not sample randomly from a large population of different cues, but restricts his attention to a few homogeneous aspects of the environment. In particular, the subject may distinguish between stimulus events which consist of the same "elements", but that have different patterns. For example, in a paired-associate task, one can assume that each stimulus word is represented by a single pattern that is sampled with probability 1 when the stimulus word is displayed. In analyzing such tasks, there has been a shift from linear operator models to Markov models. In the latter models, the description of the organism on a particular trial is phrased in terms of the momentary state of each stimulus element; taken together, these descriptions constitute the state of the organism. It is usually assumed that the change

in response probabilities from trial n to trial $n+1$ of an experiment is dependent only on the state of the organism on trial n , and a transition matrix that specifies the change in states associated with each reinforcement event. This assumption, plus the restriction on the number of elements in the population, have served to reduce the number of states in the Markov process to a manageable number. The resulting models have proven mathematically tractable, and have given an excellent account of a wide array of data.

In addition to the specification on the stimulus side of the exact cues that are being sampled, there has been a relaxation of the original stimulus sampling assumption that each element is conditioned in an all-or-none fashion to some response class. For example, one may postulate neutral states, where if an element is sampled, the subject simply responds at random from among the available response alternatives (LaBerge, 1959). An element may be "strongly" conditioned to a response in which case, for example, at least two negative reinforcements must occur before the element changes conditioning to another response class, versus a "weak" state, where a single error may produce a change in conditioning, (Atkinson and Estes, 1963).

As an example of the effect of the use of models on experimental design, we may mention the Bower paired-associate experiments (1961, 1962). In these experiments an attempt was made to evaluate a one-element stimulus sampling model in which the learning is assumed to occur abruptly, in an all-or-none fashion; the ability of such a model to account for paired-associate data has been extremely good. It is important to note several features of Bower's experiments that are

relevant to the effects of a mathematical approach on research techniques. First, the experiments were designed explicitly to test a particular model. In the classical paired-associate task, both stimuli and responses frequently were verbal items such as nonsense syllables or familiar words. It seemed to Bower that at least two processes were taking place in the traditional situation--learning of the response set and the association of responses with the appropriate stimuli. The one-element model was designed to account only for the latter process. Hence response items were chosen that one could assume would already be part of the subject's response repertoire. Further in terms of the model it was desirable to treat each subject's protocol as though the several stimulus-response pairs were learned independently. Consequently, stimuli were chosen that, in other situations, showed minimum interference with each other. Thus the theory as interpreted in the model dictated the criteria for the selection of stimulus and response items that would be appropriate to test Bower's ideas about the associative phase of paired associate learning. The criticism sometimes made that such experiments are contrived and artificial fails to recognize the goal of laboratory research, which is to restrict the introduction of extraneous variables that are not relevant to the hypothesis being considered.

In the remaining part of this paper we shall try to give some concrete illustrations of the role of models in psychological research. In order to do this it will be necessary to describe a typical experimental problem, outline several alternative models, and then indicate some of the strategies and tactics involved in making a comparison

among the models. The task we select is a highly special case of paired-associate learning. The reason for selecting this experimental problem is that it illustrates many of the problems in psychological theorizing without introducing too much mathematical complexity.

The experiment involves a list of 18 different paired-associate items. The stimulus member of each pair is a single Greek letter and the response is the number 1 or 2. The subjects are told the response alternatives available to them, and each number occurs equally often as the to-be-learned response. Hence the probability of a correct response by guessing is $\frac{1}{2}$.

Two types of trials are defined. On a study trial the 18 letter-number pairs are exhibited singly in a random order. The subject is instructed simply to associate each letter with the appropriate number and is not required to make a response. On a test trial the letters alone are presented singly in a new random order and the subject attempts to give the correct number to each letter. The subject is required to respond to each letter on a test trial (even if he is uncertain and must guess), but he is not told whether his response is correct.

In the experiment we shall examine, two study trials were given followed by four test trials; the standard notation for this type of experiment is simply $R_1 R_2 T_1 T_2 T_3 T_4$ (Jones, 1962). If we represent a correct response by c and an error by e , then the response protocol for an individual stimulus item (i.e., a particular Greek letter) over the four test trials will consist of an ordered four-tuple of c 's and e 's. For example, the protocol $e_1 c_2 e_3 e_4$ would indicate a correct response on T_2 and incorrect responses on T_1, T_3 , and T_4 . The role of theory in this situation is to predict the types of sequences that will occur and their relative frequencies.

One feature of our experimental situation that has been established by several studies is that if we run enough test trials in sequence, then in time the subject will become consistent in his response to each stimulus. For some stimuli the stereotyped response is the correct one, for other stimuli it is incorrect.

The models that we shall examine are imbedded in the controversy regarding all-or-none learning versus incremental learning. Of late, there have been some particularly important studies dealing with this issue but we will not attempt to review them here. Rather, for illustrative purposes, we will take a naive approach and outline one model that might be viewed as characterizing the incremental position and another that typifies the all-or-none viewpoint.

The incremental model is in the spirit of Hullian theory and is very similar to the early work of Bush and Mosteller (1951, 1955). We assume that at the start of a trial there is a fixed number p associated with each stimulus item that specifies the probability that a correct response will be made to that item. The effect of a study trial is to increment that probability by a constant proportion θ of the total possible change. Specifically, if p is the probability before a study trial, then after a single study trial the new probability will be $p + \theta(1-p)$. That is, the new probability is the old one plus a constant θ of the possible increase. In mathematical terminology, we say that the effect of a study trial is to apply an operator Q to the operand p to yield a new quantity $Q(p)$; i.e., $Q(p) = p + \theta(1-p)$. As will be evident later, it will be more convenient to write this reinforcement operator in the following form:

$$Q(p) = (1-\theta)p + \theta \quad (1)$$

To obtain the new probability after two successive study trials we apply the operator Q twice, namely

$$\begin{aligned} Q^2(p) &= Q [Q(p)] = (1-\theta)^2 p + (1-\theta)\theta + \theta \\ &= 1 - (1-p)(1-\theta)^2 . \end{aligned}$$

By induction one can show that after n successive study trials

$$Q^n(p) = 1 - (1-p)(1-\theta)^n .$$

For our experimental situation the initial probability of guessing correctly is $\frac{1}{2}$ and hence we would set $p = \frac{1}{2}$. Thus, for this model the probability of a correct response on the first test trial following n successive study trials will be

$$\Pr(c_1) = 1 - \frac{1}{2}(1-\theta)^n . \quad (2)$$

The all-or-none learning process that we shall consider is one that has been actively investigated by Estes (1960, 1961), Bower (1961, 1962), Restle (1963), Suppes and Ginsberg (1963) and others. For this model we assume that each stimulus item is in one of two conditioning states: C or G. In state C the stimulus is conditioned to the correct response and on a test trial will elicit that response with probability 1. In state G the stimulus is not conditioned to any response, and in this state the probability of a correct response is $\frac{1}{2}$; i.e., a correct response will occur at the chance level. All items at the start of the experiment are in state G, but on each study trial there exists a probability θ that conditioning will occur. Thus,

the probability that a particular stimulus item is in state C after one study trial is θ , after two study trials $\theta + (1-\theta)\theta$, after three study trials $\theta + \theta(1-\theta) + \theta(1-\theta)^2$, etc. More generally the probability of being in state C after n successive study trials is

$$\Pr(\underline{C}_n) = 1 - (1-\theta)^n . \quad (3)$$

For this model the expected probability of a correct response on test trial T_1 after n successive study trials would be

$$\begin{aligned} \Pr(c_1) &= 1 - (1-\theta)^n + \frac{1}{2}(1-\theta)^n \\ &= 1 - \frac{1}{2}(1-\theta)^n . \end{aligned} \quad (4)$$

That is, the probability of being in state C plus $\frac{1}{2}$ times the probability of being in state G. The all-or-none character of this model is represented by the fact that for the underlying states the probability of a correct response can take on only two values; either $\frac{1}{2}$ if the subject is in state G, or 1 if the subject is in state C. Further, the transition from G to C occurs in an all-or-none fashion on a single trial.

To summarize to this point, for the incremental model two study trials generate a fixed number associated with each stimulus item that specifies the probability of a correct response on the first trial. We shall call this number ϕ , and it is given by Equation 2 when $n = 2$; i.e.,

$$\phi = \frac{1}{2}(1-\theta)^2 + \theta(1-\theta) + \theta . \quad (5)$$

For the all-or-none model, each stimulus item will be in either state C or state G. If the item is in state C a correct response occurs on a test trial; if the item is in state G a correct response occurs with probability $\frac{1}{2}$. The probability of being in state C after two study trials will be called x and is given by Equation 3; i.e.,

$$x = \theta + (1-\theta)\theta . \quad (6)$$

The next question is with regard to the events that occur on a test trial. As noted earlier, it is known that behavior eventually becomes stereotyped if sufficiently long series of test trials are run, and this observation suggests that systematic changes may be occurring over test trials. A plausible assumption that accounts for the changes is that in the absence of an experimenter-determined reinforcing event (i.e., the experimenter telling the subject which response was correct) the emitted response is the response reinforced. This last phrase characterizes much of the theoretical work of contiguity theorists such as Guthrie (1935); the idea being that the last response to take place in the presence of a stimulus will remain associated with it and will tend to reoccur when the stimulus is presented again.

The assumption that the emitted response is the one reinforced on a test trial delimits a class of qualitative theories that can be experimentally investigated. But this class of theories is large and difficult to characterize; also, too frequently new experimental findings can somehow be made to agree with almost any of the theoretical positions. Thus, much is to be gained by taking a qualitative assumption concerning test trial effects and examining the consequences of stating it

mathematically. To illustrate, assume that a qualitative theory predicts a difference between two experimental groups; an experiment is run to test for the difference and none is obtained. What conclusion can be drawn? Either that no difference exists or that if it exists it is too small to be detected by the experimental procedures utilized. In contrast, with a quantitative theory we know not only the direction of the predicted difference but also the exact magnitude. Consequently the equipment and experimental procedure can be designed so that they are sufficiently sensitive to detect the difference if it is present. Then if no difference is found there can be no alibi that the effect might be too small to be detected. Experiments that find no differences are ambiguous in evaluating qualitative predictions; for quantitative theories such results have an exact interpretation.

For the incremental model we shall assume the same reinforcing operator on test trials as on study trials. Specifically, if p is the probability of a correct response and that response occurs on a test trial, then the new probability will be $Q(p) = (1-\alpha)p + \alpha$ where α is the parameter describing learning under self reinforcement. If an incorrect response occurs (which has probability $q = 1-p$), then that response will be reinforced, which means that Q will be applied to q ; i.e., $Q(q) = (1-\alpha)q + \alpha$. By inspection of the last expression we see that reinforcing an incorrect response is equivalent to applying the operator $Q'(p) = (1-\alpha)p$. Stating our ideas exactly, if p is the probability of a correct response on a test trial, then p^* (the value at the end of the test trial) will be

$$p^* = \begin{cases} Q(p) = (1-\alpha)p + \alpha, & \text{if the correct response occurs.} \\ Q'(p) = (1-\alpha)p, & \text{if the incorrect response occurs.} \end{cases} \quad (7)$$

From these equations it can be shown that the probability of a correct response will approach 1 or 0 as the run of test trials becomes large. Thus, asymptotically some items will absorb on the correct response and others on an incorrect response. If ϕ denotes the probability of a correct response to a specific stimulus item at the start of the test sequence, then the probability that this item absorbs on the correct response will be ϕ .

Our assumption that on a test trial the emitted response is reinforced also has a natural interpretation in terms of the all-or-none model. As before, we assume that reinforcement of a response conditions the stimulus to that response with some probability, say β . If the stimulus is conditioned to the correct response on a test trial (i.e., in state C) then that response occurs and by reinforcing it we guarantee that it remains in state C. If the stimulus item is in state G, then with probability $\frac{1}{2}$ the correct response occurs, and by assumption this is reinforcing on a test trial; hence with probability β the item moves to state C. Now it is obvious that we must also allow for the occurrence of an incorrect response; if an item is in state G and an incorrect response occurs, then the item will become conditioned with probability β to the incorrect response. Therefore, in addition to states C and G which characterize study trials we also need a state E to denote conditioning to an incorrect response. These notions are embodied in the following transition matrix:

$$\begin{array}{c}
 \underline{\underline{C}} \\
 \underline{\underline{G}} \\
 \underline{\underline{E}}
 \end{array}
 \begin{bmatrix}
 \underline{\underline{C}} & \underline{\underline{G}} & \underline{\underline{E}} \\
 1 & 0 & 0 \\
 \frac{1}{2}\beta & 1-\beta & \frac{1}{2}\beta \\
 0 & 0 & 1
 \end{bmatrix}
 \quad (8)$$

The rows indicate the state at the start of a test trial and the columns the state at the end of the trial. Each entry denotes the probability of a transition from one state to another. In state $\underline{\underline{C}}$ (or $\underline{\underline{E}}$) no change can occur on a test trial. In state $\underline{\underline{G}}$ the item may become conditioned to the correct response if it occurs (with probability $\frac{1}{2}$) and conditioning is effective (with probability β); similarly in state $\underline{\underline{G}}$ the item may become conditioned to the incorrect response if it occurs and conditioning is effective. As in the case of the linear model each stimulus item eventually absorbs in either state $\underline{\underline{C}}$ or $\underline{\underline{E}}$. Thus, after a long run of test trials, a given stimulus item will eventually elicit either a correct or an incorrect response consistently.

These then are the two models we shall examine. From a qualitative viewpoint each represents the same psychological process; i.e., they both assume that a reinforcement tends to increase the likelihood of the reinforced response and they both assume the same subject-determined reinforcement schedules. However, the exact nature of the change that occurs following reinforcement is quite different for the two models. This is illustrated by the fact that the probability of a correct response in the incremental model may take on any value from 0 to 1, whereas for the all-or-none model it can take on only the values 0, $\frac{1}{2}$, or 1. This fact alone indicates that there are substantial

differences between the two models and it becomes important to determine which interpretation of the reinforcing event best approximates the actual learning process.

In order to compare these models we need to derive some predictions. Consider first the incremental model and the possible outcomes for a given stimulus item on trials T_1 and T_2 . What is the probability of two correct responses? Well, a correct response will occur on T_1 with probability ϕ (see Equation 5) and since the correct response occurred on T_1 , then the probability of a correct response on T_2 will be $(1-\alpha)\phi + \alpha$ (see Equation 7). The probability of two correct responses is simply the product of these two probabilities. Similarly, the probability of an incorrect followed by a correct is $(1-\phi)$ times $(1-\alpha)\phi$; i.e., the probability of an incorrect response on T_1 times the probability of a correct response on T_2 given that the preceding trial was incorrect. In this way we obtain the following expressions.

$$\begin{aligned}
 \Pr(c_1 c_2) &= \phi[(1-\alpha)\phi + \alpha] \\
 \Pr(c_1 e_2) &= \phi[1 - \{(1-\alpha)\phi + \alpha\}] \\
 \Pr(e_1 c_2) &= (1-\phi)[(1-\alpha)\phi] \\
 \Pr(e_1 e_2) &= (1-\phi)[1 - (1-\alpha)\phi]
 \end{aligned}
 \tag{9}$$

By similar methods expressions can be obtained for the probability of any sequence of responses over the four test trials. We display a few of these equations to indicate that even for a simple model of this sort the elaboration of the theory leads to predictions whose consequences are too complicated to be understood without the tools of mathematical analysis:

$$\begin{aligned}
\Pr(c_1 c_2 c_3 c_4) &= \phi[(1-\alpha)\phi+\alpha][[(1-\alpha)^2\phi+(1-\alpha)\alpha+\alpha][[(1-\alpha)^3\phi+(1-\alpha)^2\alpha+(1-\alpha)\alpha+\alpha] \\
\Pr(c_1 c_2 c_3 e_4) &= \phi[(1-\alpha)\phi+\alpha][[(1-\alpha)^2\phi+(1-\alpha)\alpha+\alpha][1-(1-\alpha)^3\phi-(1-\alpha)^2\alpha-(1-\alpha)\alpha-\alpha] \quad (10) \\
&\vdots \\
&\vdots \\
\Pr(e_1 e_2 e_3 e_4) &= (1-\phi)[1-(1-\alpha)\phi][1-(1-\alpha)^2\phi][1-(1-\alpha)^3\phi] \quad .
\end{aligned}$$

Predictions for the same quantities can be obtained for the all-or-none model. Consider first the possible outcomes for a given stimulus item on trials T_1 and T_2 . The probability of two correct responses in a row is $x + (1-x)\frac{1}{2}[\beta + (1-\beta)\frac{1}{2}]$. That is, with probability x the stimulus item is in state C before the first test trial and hence will generate correct responses on all subsequent trials; with probability $1-x$ the stimulus item starts in state G and a correct response occurs on T_1 with probability $\frac{1}{2}$, further a correct response can occur on T_2 if (1) conditioning was effective on T_1 (i.e., with probability β) or (2) if conditioning was not effective and by chance the subject again guessed correctly on the next trial. Proceeding in this way we obtain the following expressions:

$$\begin{aligned}
\Pr(c_1 c_2) &= x + \frac{1}{2}(1-x)[\beta + (1-\beta)\frac{1}{2}] \\
\Pr(c_1 e_2) &= \frac{1}{4}(1-x)(1-\beta) \\
\Pr(e_1 c_2) &= \frac{1}{4}(1-x)(1-\beta) \\
\Pr(e_1 e_2) &= \frac{1}{2}(1-x)[\beta + (1-\beta)\frac{1}{2}] \quad .
\end{aligned} \tag{11}$$

Using the same methods one can obtain expressions for the sequence of events over the four test trials. For example

$$\begin{aligned}
\Pr(c_1 c_2 c_3 c_4) &= x + \frac{1}{2}(1-x)\left[\beta + \frac{1}{2}\beta(1-\beta) + \frac{1}{4}(1-\beta)^2\left\{\beta + (1-\beta)\frac{1}{2}\right\}\right] \\
\Pr(c_1 c_2 c_3 e_4) &= \frac{1}{16}(1-x)(1-\beta)^3 \\
&\vdots \\
\Pr(e_1 e_2 e_3 e_4) &= \frac{1}{2}(1-x)\left[\beta + \frac{1}{2}\beta(1-\beta) + \frac{1}{4}(1-\beta)^2\left\{\beta + (1-\beta)\frac{1}{2}\right\}\right].
\end{aligned}
\tag{12}$$

Armed with these equations we now face the task of deciding which model provides the best account of our data. Table 1 presents observed frequencies for two identical experiments, one using 60 college students at Stanford University and the other 60 fourth grade children from the Oakland City schools*. As indicated earlier each subject was run on 18 paired-associates so that over a group of subjects we have information on $18 \times 60 = 1080$ test trial sequences. The table gives the frequency with which each sequence occurred in the two groups. In subsequent analyses we assume that the items in the list are of equal difficulty and that all subjects in a group learn at the same rate; i.e., we postulate that the values of the various parameters are the same for all subjects in the college group and for all subjects in the grade school group. Of course, this assumption is suspect, but for purposes of this paper it seems justified since it greatly simplifies subsequent analyses. For with this assumption the response sequence associated with any given stimulus item can be viewed as a sample of size one from a population of sequences all generated by the same underlying process. A discussion of the problems involved in treating individual data and group data is given in Suppes and Atkinson (1960).

* This study was conducted by Duncan Hansen of Stanford University.

TABLE 1

Observed frequencies for response sequences.

Sequence Code	T ₁	T ₂	T ₃	T ₄	College Group	Grade School Group
1	c	c	c	c	633	474
2	c	c	c	e	22	51
3	c	c	e	c	19	39
4	c	c	e	e	28	38
5	c	e	c	c	16	30
6	c	e	c	e	19	27
7	c	e	e	c	23	20
8	c	e	e	e	54	71
9	e	c	c	c	43	40
10	e	c	c	e	6	15
11	e	c	e	c	11	24
12	e	c	e	e	10	33
13	e	e	c	c	26	23
14	e	e	c	e	11	19
15	e	e	e	c	14	15
16	e	e	e	e	145	161

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
100

In order to make predictions for the data displayed in Table 1 we need estimates of the parameters ϕ and α for the incremental model and x and β for the all-or-none model. There are many ways of making these estimates, but for the present problem a simple method is to select the pair of parameter values that minimizes the χ^2 function*. To illustrate the method, let $p_i(\alpha, \phi)$ denote theoretical expressions for the four-response probabilities given in Equation 10, where the subscript i refers to code numbers assigned in Table 1. Further, let O_i ($i=1$ to 16) denote the observed frequencies for one of the groups in Table 1 and let $T = O_1 + O_2 + \dots + O_{16}$. Then we define the function

$$\chi^2(\alpha, \phi) = \sum_{i=1}^{16} \frac{[Tp_i(\alpha, \phi) - O_i]^2}{Tp_i(\alpha, \phi)} \quad (13)$$

and select our estimates of α and ϕ so that they jointly minimize the χ^2 function. Under the null hypothesis this minimum χ^2 has the usual limiting distribution with $16 - 3$ degrees of freedom. (If n parameters are estimated then there are $16 - n - 1$ degrees of freedom.)

Using this method we obtain parameter estimates for the incremental model with the following minimum χ^2 values:

	College	Grade School
α	.383	.331
ϕ	.705	.636
minimum χ^2	256.3	192.7

* For a review of some of these statistical methods as they apply to learning models see Chapter 2 in Suppes and Atkinson (1960).

Using the same minimum χ^2 method for the all-or-none model, we obtain the following parameter estimates:

	College	Grade School
β	.291	.191
x	.409	.297
minimum χ^2	40.4	67.0

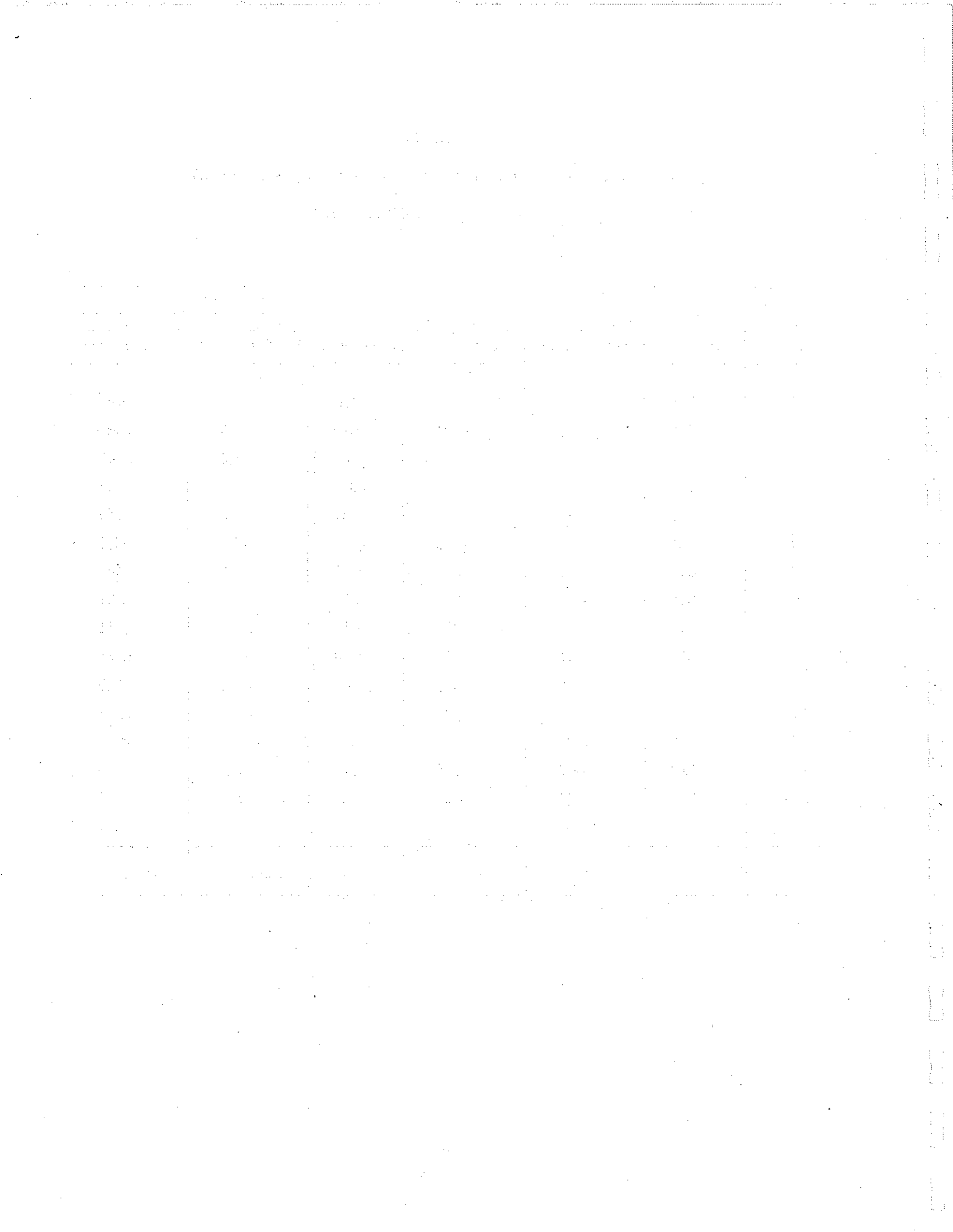
The above parameter estimates generate the predictions presented in Table 2. For both the College and the Grade School groups the all-or-none model provides the closest correspondence between predicted and observed proportions on the four-response data. The superiority of the all-or-none model is also reflected in the obtained minimum χ^2 's. For the College group the χ^2 for the incremental model is more than six times that for the all-or-none model; whereas, for the Grade School group the ratio is almost three to one in favor of the all-or-none model. Since all the χ^2 's are based on the same number of degrees of freedom, it seems reasonable to conclude that the all-or-none model provides the best account of these data. Further, independent of any comparative analysis of the models, the correspondence between the all-or-none predictions and these data is reasonably good in terms of the degree of accuracy that psychologists have come to expect in research of this sort. See Atkinson and Estes (1963) for a discussion of this point.

Finally, we should note that the parameter estimates for both the all-or-none model and the incremental model have the expected properties. For the all-or-none model the estimate of x is larger for the College

TABLE 2

Observed proportions and predictions for the Incremental
Model and the All-or-None Model.

Sequence Code	College			Grade School		
	Observed	Incremental Predictions	All-or-None Predictions	Observed	Incremental Predictions	All-or-None Predictions
1	.586	.476	.549	.439	.359	.426
2	.020	.035	.013	.047	.044	.023
3	.018	.060	.013	.036	.070	.023
4	.026	.004	.024	.035	.009	.034
5	.015	.045	.024	.028	.052	.034
6	.018	.020	.013	.025	.026	.023
7	.021	.020	.013	.019	.026	.023
8	.050	.044	.054	.066	.051	.061
9	.040	.066	.054	.037	.071	.061
10	.006	.018	.013	.014	.025	.023
11	.010	.018	.013	.022	.025	.023
12	.009	.027	.024	.031	.035	.034
13	.024	.025	.024	.021	.031	.034
14	.010	.020	.013	.018	.028	.023
15	.013	.020	.013	.014	.028	.023
16	.134	.102	.140	.149	.121	.129
χ^2		256.3	40.4		192.7	67.0



group than for the Grade School group, reflecting the fact that more learning occurred for the former group on the study trials. Also, the parameter β characterizing the conditioning rate on test trials is larger for the College group, indicating that the postulated self-reinforcing event had more effect for the College students. Similarly, for the incremental model the estimates of both ϕ and α are greater for the College group than for the Grade School group.

What conclusions can be drawn from our analysis of these two models? Should readers not familiar with psychology conclude that incremental models are clearly less satisfactory than all-or-none models and that subsequent theory construction should be along all-or-none lines of development? This might be one conclusion, but few psychologists inclined toward incremental theories of learning would be in agreement. They, of course, would want to see many more comparisons between the two models on a variety of experimental data. Further, they probably would argue that the incremental model presented here was too simple and that a more sophisticated interpretation of incremental theory would lead to much better results. On this last point, psychologists favoring an all-or-none position would agree, but they, in turn, would point out that similar improvements easily could be made for the all-or-none model. Further, they would emphasize that both models estimated the same number of parameters and are of similar mathematical difficulty, and that these facts are important when one evaluates the clear superiority of the all-or-none process.

It is not necessary to pursue such arguments to realize that a single comparison of this sort is unlikely to change anyone's theoretical disposition. Only as more evidence accumulates and more variations of each model are investigated will it become clear which approach is more parisimonious. Certainly since the time of Henri Poincaré we have known that no theory is correct in an absolute sense. Rather, some theories tend to be more useful than others in accomplishing the goals of the scientific enterprise if they (1) lead to natural and unambiguous interpretations of phenomena, (2) have a tractable logical structure, and (3) suggest new experimental dimensions. The success of the Relativity Theory was not due to the fact that the new concepts of space and time were in any sense more true than the old ones. For, any of the phenomena that could be explained by the new theory also could be explained on the basis of absolute dimensions of space and time. But such explanations became extremely cumbersome and artificial when contrasted to the explanations offered by Relativity Theory.

If we accept the notion that a single comparison of the sort offered in this paper cannot be regarded as crucial in selecting between models, then what is the next step? As indicated earlier, one obvious requirement is to extend the application of the two models to many other types of experiments. As more experimental comparisons are made a better understanding of the properties of each model will be obtained which will give us a measure of their relative power. However, in addition to extending the range of application it also is important to take each experiment and, by inspection of the discrepancies between theory and observation, obtain some clues for modifications that will lead to a

better model. This part of the scientific enterprise is extremely challenging. In effect the theorist, by scrutinizing unexpected perturbations in the data, attempts to come up with either a modification of the basic assumptions or a new interpretation of how the assumptions can be applied. We doubt if such revisions follow any clear or systematic pattern; more often than not the theorist tries many new schemes and then reports the one that seems most promising. From the viewpoint of understanding the scientific process, it is unfortunate that the trial and error stage between successive revisions of a model is not occasionally recorded. Looking at the end product tends to give the misleading impression that theory develops in a neat and orderly fashion.

In this next section we shall try to give the reader some idea of the way the psychologist may use a set of experimental results to suggest changes in the theory. We could present examples for either the incremental or the all-or-none model since both have natural extensions suggested by our data. However, it would be too lengthy to attempt both in this paper; consequently we will only examine possible revisions of the all-or-none theory. We select this model because the modifications to be considered are tractable and will not introduce any new mathematical techniques. The modifications illustrate two types of revisions that can occur in a theory: one modification represents a reinterpretation of how the axioms can be applied and the other, a basic change in the axioms.

Inspection of the correspondence between predictions for the all-or-none model and our data indicates two striking discrepancies. First, the theory predicts that the probability of a correct response on T_1

followed by an incorrect response on T_2 is the same as the probability of an incorrect response on T_1 followed by a correct response on T_2 ; i.e., $\Pr(c_1e_2) = \Pr(e_1c_2)$. However, the data clearly contradicts this prediction as indicated below:

	College	Grade School
$\Pr(c_1e_2)$.104	.137
$\Pr(e_1c_2)$.065	.104

These quantities are obtained directly from Table 2; i.e.,

$\Pr(c_1e_2) = p_5 + p_6 + p_7 + p_8$ and $\Pr(e_1c_2) = p_9 + p_{10} + p_{11} + p_{12}$. For both groups $\Pr(c_1e_2)$ is larger than $\Pr(e_1c_2)$.

Another discrepancy between theory and data is with regard to the probability of a correct response on test trial n . The theory predicts that this quantity, averaged over items, should be a constant; i.e., $\Pr(c_n) = x + (1-x)\frac{1}{2}$. The appropriate statistics are given below.

	College	Grade School
$\Pr(c_1)$.754	.694
$\Pr(c_2)$.715	.661
$\Pr(c_3)$.719	.629
$\Pr(c_4)$.727	.616

In both groups there is a tendency for the probability of a correct response to decrease over test trials. Parenthetically, the incremental model makes the same predictions; i.e., $\Pr(c_1e_2) = \Pr(e_2c_1)$ and $\Pr(c_n) = \text{constant}$.

Both of these observations suggest the need for a process in the model that will produce an increase in the error rate over test trials. One way of doing this is to assume that the conditioning parameter associated with a correct response is different from that associated with an incorrect response. Specifically, we can assume that when an item is in state \underline{G} on a test trial, then (1) if a correct response is made it becomes conditioned with probability μ and (2) if an incorrect response is made it becomes conditioned with probably δ . The transition matrix given in Equation 8 would now be rewritten as follows:

$$\begin{array}{c}
 \underline{C} \\
 \underline{G} \\
 \underline{E}
 \end{array}
 \begin{array}{c}
 \underline{C} \\
 \underline{G} \\
 \underline{E}
 \end{array}
 \begin{array}{c}
 \underline{E}
 \end{array}
 \left[\begin{array}{ccc}
 1 & 0 & 0 \\
 \frac{1}{2}\mu & 1 - \frac{1}{2}(\mu + \delta) & \frac{1}{2}\delta \\
 0 & 0 & 1
 \end{array} \right] \quad (14)$$

At a qualitative level of analysis this change provides for the discrepancies noted above. It can be shown that

$$\begin{aligned}
 \Pr(c_1 e_2) &= \frac{1}{4}(1-\mu)(1-x) \\
 \Pr(e_1 c_2) &= \frac{1}{4}(1-\delta)(1-x)
 \end{aligned} \quad (15)$$

which implies that $\Pr(c_1 e_2) > \Pr(e_1 c_2)$ when $\delta > \mu$. Similarly, it can be shown that the probability of a correct response on test trial n is

$$\Pr(c_n) = x + \frac{(1-x)}{2} \left[1 - \frac{1}{2}(\mu + \delta) \right]^{n-1} + \frac{(1-x)\mu}{\mu + \delta} \left[1 - \left\{ 1 - \frac{1}{2}(\mu + \delta) \right\}^{n-1} \right] \quad (16)$$

When $\delta > \mu$ this equation describes a function decreasing over test trials that approaches $x + \frac{(1-x)\mu}{(\mu+\delta)}$ in the limit.

What psychological rationale can be offered for postulating that $\delta > \mu$? There are several but we shall mention one that assumes differences in learning rates for various paired-associate items. Let us suppose that some stimulus-response connections are easier to learn than others; e.g., for some subjects the Greek letter ω may be easier to associate with response 1 than 2. If response 2 is assigned as the to-be-learned response, then this item should be learned more slowly and consequently would have a higher probability than other items of being in state G after two study trials. If the item still is in state G after the study trials then, since response 1 is the more compatible response, the item will be more likely to become conditioned to an error response over the series of test trials. In essence the assumption is that stimulus items not conditioned to a correct response after a series of study trials will tend to favor association with an error response on test trials.

The proposed modification that assumes differential conditioning parameters may be regarded as a revision of the basic axioms of the model. The next modification we offer simply requires a reinterpretation of how the model might be applied. Originally it was assumed that after a series of study trials each stimulus item was either in state C (with probability x) or state G (with probability $1-x$). A natural extension of these ideas is to permit the possibility that the item also may be in state E at the end of a study trial. That is, it seems reasonable to assume that at the start of T_1 there is a probability x'

that the item will be in state C, a probability y' that the item will be in state E, and a probability $1-x'-y'$ that it will still be in state G.

Introducing these two modifications yields a model that has four parameters to be estimated from the data (namely, μ, δ, x' , and y') as compared with two in the original all-or-none model. The question is whether the addition of two parameters yields a substantial improvement in the fit of the model.

The four parameters were estimated separately for each set of data by selecting that parameter vector (μ, δ, x', y') that minimized χ^2 . The obtained estimates were as follows.

	College	Grade School
μ	.207	.041
δ	.289	.270
x'	.478	.379
y'	0	0

The predictions for the revised model based on the above estimates are given in Table 3. The table also presents the associated minimum χ^2 's. It is evident, by comparing these predictions and χ^2 values with those given in Table 2, that the revised model does a better job. The improvement is not very large for the College data but is dramatic in the case of the Grade School data*. However, it is evident that the improvement in prediction is due entirely to but one of the two modifications. The change that permitted an item to be in state E at

* Statistical tests to evaluate these comparisons are available and are discussed in reference to learning models by Suppes and Atkinson (1960).

TABLE 3

Observed proportions and predictions for the revised
version of the All-or-None Model.

Sequence Code	College		Grade School	
	Observed	Predicted	Observed	Predicted
1	.586	.578	.439	.435
2	.020	.016	.047	.034
3	.018	.015	.036	.026
4	.026	.026	.035	.045
5	.015	.022	.028	.028
6	.018	.015	.025	.026
7	.021	.013	.019	.020
8	.050	.054	.066	.075
9	.040	.041	.037	.033
10	.006	.015	.014	.026
11	.010	.013	.022	.020
12	.009	.024	.031	.035
13	.024	.020	.021	.022
14	.010	.013	.018	.020
15	.013	.012	.014	.015
16	.134	.123	.149	.141
χ^2		30.2		21.7

the end of a study trial did not yield an improvement in prediction, because in both the College and the Grade School groups the estimated value of y' was 0. Thus, the revision that postulated differential conditioning parameters for correct and incorrect responses accounts entirely for the improvement in the fit of the modified all-or-none model.

Without pursuing this example further, we hope that the reader has a fairly clear picture of our view of the theoretical enterprise in psychology. However, there is one last point. If we scan the χ^2 's given in this paper and select the one associated with the best fit (i.e., the modified all-or-none model applied to the Grade School data) a value of 21.7 is obtained with 11 degrees of freedom. This value of χ^2 is significant at the .05 level, and therefore on the basis of statistical considerations we would reject the model. The sensible retort to this statement is the point we have tried to emphasize throughout the paper. We always assume that any model can be rejected on statistical grounds if enough observations are made. The goal is not to reject or accept a given model at some predetermined level of significance, but rather to make comparisons among models and ask how well a model performs relative to other models. Simply stated, a model will not be rejected on purely statistical grounds, but will be rejected only when there are other models that consistently do a better job of prediction.

In conclusion, we view the use of mathematical models as virtually synonymous with the construction of a quantitative theory of behavior. From a mathematical standpoint it is logically possible to have a theory of behavior that leads only to qualitative predictions. However,

it is difficult to find in the history of science, let alone in the history of psychology, theories of this sort that have had sustained empirical significance. From the systematic standpoint a theory or model based only on qualitative distinctions leads to a small number of testable predictions. Aristotle's physics and Lewin's topological field theory (1936) are good examples. The absence of precise systematization leads usually to pseudo-derivations from the theory. By pseudo-derivation we mean the derivation of prediction that requires many additional assumptions that are not part of the original theory. Further, as the set of phenomena that we study expand in complexity so also does the reasoning necessary for the design of experiments and the formulation of hypotheses. Ordinary logic becomes inadequate and the elaboration of the theory requires the powerful tool of mathematical analysis.

Finally, we remark that we have avoided a discussion of the general nature of models and theories. The words "model" and "mathematical model" are used in a variety of related senses by behavioral scientists and philosophers. Often the most reasonable interpretation is that the model is the set of mathematically formulated postulates that express in precise form the intuitive notions of the relevant psychological theory. Despite the dominance of this usage in the literature of the behavioral sciences, we prefer to use the more precise concept of model that has been adopted by mathematical logicians. Roughly speaking, a model is an abstract object that satisfies a theory. A theory given in axiomatic form can, if one chooses, be identified with its set of axioms. It is not to the point to enter into technical details here. What is important in our opinion is that models have a role to play whenever theory is

constructed. From our standpoint, if a theory has systematic content and is not simply a vague collection of heuristic ideas, then there exist models that satisfy the theory, and it is up to the experimenter to determine whether these models provide an adequate analysis of behavioral phenomena. We believe that the role of mathematical models in psychology is not really separate from the role of systematic theorizing.

REFERENCES

- Anderson, N. H. Effect of first-order conditional probability in a two-choice learning situation. J. exp. Psychol., 1960, 59, 71-93.
- Anderson, N. H. An analysis of sequential dependencies. In R. R. Bush and W. K. Estes (Eds.), Studies in mathematical learning theory. Stanford, California: Stanford University Press, 1959, 248-264.
- Atkinson, R. C. and Estes, W. K. Stimulus sampling theory. In R. R. Bush, E. Galanter, and R. D. Luce (Eds.), Handbook of mathematical psychology. Vol. 2. New York: Wiley, 1963, in press.
- Bower, G. H. A model for response and training variables in paired-associate learning. Psychol. Rev., 1962, 69, 34-53.
- Bower, G. H. Applications of a model to paired-associate learning. Psychometrika, 1961, 26, 255-280.
- Bush, R. R. A survey of mathematical learning theory. In R. Duncan Luce (Ed.), Developments in mathematical psychology. Glencoe, Illinois: The Free Press, 1960.
- Bush, R. R. and Mosteller, F. Stochastic models for learning. New York: Wiley, 1955.
- Bush, R. R. and Mosteller, F. A mathematical model for simple learning. Psychol. Rev., 1951, 58, 313-323.
- Estes, W. K. Learning theory. Annual Review of Psychology, 1962, 13, 107-144.
- Estes, W. K. New developments in statistical behavior theory: differential tests of axioms for associative learning. Psychometrika, 1961, 26, 73-84.
- Estes, W. K. Learning theory and the new mental chemistry. Psychol. Rev., 1960, 67, 207-223.
- Estes, W. K. The statistical approach to learning theory. In S. Koch (Ed.), Psychology: A study of a science. Vol. 2. New York: McGraw-Hill, 1959, 380-491.

- Estes, W. K. Toward a statistical theory of learning. Psychol. Rev., 1950, 57, 94-107.
- Guthrie, E. R. The psychology of learning. New York: Harper, 1935.
- Hilgard, E. R. Theories of Learning. (Rev. Edition). New York: Appleton-Century-Crofts, 1956.
- Hull, C. L. Principles of behavior: An introduction to behavior theory. New York: Appleton-Century-Crofts, 1943.
- Jones, J. E. All-or-none versus incremental learning. Psychol. Rev., 1962, 69, 156-160.
- LaBerge, D. A model with neutral elements. In R. R. Bush and W. K. Estes (Eds.), Studies in mathematical learning theory. Stanford, California: Stanford University Press, 1959, 53-64.
- Lewin, K. Principles of topological psychology. New York: McGraw-Hill, 1936.
- Logan, F. A. The Hull-Spence approach. In S. Koch (Ed.), Psychology: A study of a science. Vol. 2. New York: McGraw-Hill, 1959, 293-358.
- Luce, R. D. Individual choice behavior: A theoretical analysis. New York: Wiley, 1959.
- Restle, F. Sources of difficulty in learning paired associates. In Studies in mathematical psychology. Stanford, California: Stanford University Press, 1963, in press.
- Restle, F. A survey and classification of learning models. In R. R. Bush and W. K. Estes (Eds.), Studies in mathematical learning theory. Stanford, California: Stanford University Press, 1959, 415-428.
- Seward, J. P. Reinforcement and expectancy: two theories in search of a controversy. Psychol. Rev., 1956, 63, 105-113.
- Spence, K. W. Behavior theory and conditioning. New Haven, Conn.: Yale University Press, 1956.
- Suppes, P. and Atkinson, R. C. Markov learning models for multiperson interactions. Stanford, California: Stanford University Press, 1960.
- Suppes, P. and Ginsberg, R. A fundamental property of all-or-none models. Psychol. Rev., 1963, 70, in press.