

On the Complexity of Causal Models

Brian R. Gaines
Man-Machine Systems Laboratory
Department of Electrical Engineering Science
University of Essex, Colchester, England

It is argued that principle of causality is fundamental to human thinking, and it has been observed experimentally that this assumption leads to complex hypothesis formation by human subjects attempting to solve comparatively simple problems involving acausal randomly generated events. This correspondence provides an automata-theoretic explanation of this phenomenon by analyzing the performance of an optimal modeler observing the behaviour of a system and forming a minimal state model of it.

1. INTRODUCTION

The postulation of a principle of causality, “to every effect there is a cause,” has been a continuing central problem for philosophy (Popper, 1972). Its role as a source of contention in modern science (Jauch, 1973) is epitomized by Einstein’s remark that, “I can’t believe that God plays dice.” Many of the arguments about the application of the principle are very relevant to systems science and to problems of system identification and machine learning, on the one hand, and to epistemology and behavioural psychology, on the other. In current system science the theory of causal deterministic systems is most well developed and generally applied, while the theory of modeling with alternative structures, e.g., stochastic automata, indeterminate automata, products of asynchronous automata, etc., has not been developed to the same degree.

There is a significant parallel here in human psychology—some of Michotte’s (1963) experiments on perception demonstrate that causality can be a very low level percept, similar in status to shape and colour. The deep-felt assumption of causality, typified by Einstein’s remark, is apparent in cognitive tasks also. In our research on machine learning (Gaines and Andreae, 1966; Andreae and Cashin, 1969; Andreae, 1973) we used as one test environment a version of nim in which the learning machine, played against a partially random player of variable optimality. To demonstrate that the task was nontrivial for man we also built an automaton with four lamps (representing the number of matches left) and three keys (representing taking one to three matches) and a “reward” light indicating that the game had been won. In practice, human beings found the game virtually impossible to learn in these circumstances, whereas in the normal representation of nim it is trivial.

Much of the problem in learning seemed to stem from the partially random opponent in the automaton. Although this random element makes the game easier to win (the opponent makes fatal errors) it also makes the behaviour of the automaton partially indeterminate. A hypothesis of randomness or indeterminism was never introduced by human players against the automaton, and the longer they played the more confused they became.

To investigate this further we built a simple automaton with two lamps and two keys such that one lamp came on for a short period when a key was depressed. The problem was to press the key under whichever lamp was going to come on next. In fact when a key was depressed either lamp came on at random with equal probability. Again in tests with a wide range of subjects this random element was never recognized, and individuals were prepared to pit their wits against the automaton for many hours and most seemed to delight in doing so.

A number of records were made of verbal introspection as subjects played against the automaton. One mathematician, after some 30 minutes of play, came to the conclusion that the automaton had two modes of behaviour. In one the same light tended to come on as the key depressed while in the other the opposite light tended to come on—this is tautologous but was regarded as an empirical hypothesis. Most people tended to generate more elaborate models but felt that they were “not quite correct.”

One of the most interesting records is that of a behavioural scientist who played with the box for some 30 minutes and kept a tally of the occasions when he was correct less the occasions when he was wrong. After that period he was over 20 ahead on his count and announced that he had a good strategy but felt that he should drop it for a while and experiment with alternatives to seek a better one. His accumulated count then rapidly declined and when it went negative he announced that he was going back to his “good strategy.” The count continued to decline, however and he finally announced that the box contained a system that modeled his strategy and then structured itself to act in the opposite manner and thus outwit him!

Such a form of “frustration automaton” is possible and a construction of a universal form for deterministic players has been given (Gaines, 1971). However, the hypothesis of such a complex structure on the basis of interaction with a simple two-state stochastic automaton is a fascinating phenomenon in its own right. The remainder of this correspondence will demonstrate that such a elaborate models are necessarily generated if causality is postulated in modeling acausal phenomena.

2. EXPECTED SIZE OF CAUSAL MODEL

The following abstract situation captures all the relevant features of a modeler forming a causal model of some other system based on his observations. The observed behaviour is some sequence drawn from the union of an input alphabet X and an output alphabet Y , in which no two input symbols are adjacent (this constraint is applied to simplify the definition of the sequence “length”—it may easily be satisfied by introducing a null-output symbol). Only finite sequences of observations will be considered, and the length of such a sequence will be defined as the number of members of the output alphabet in it. For examples if $X = \{I,J\}$, $Y = \{A,B\}$, then a sequence of behaviour might be:-

AIBBAJBIAAJ

where the “length” would be seven. Note that no assumptions are made about the structure of the system generating the observed behaviour.

Now consider a modeler whose task is to explain the observed sequence of behaviour in terms of some model from a class of possible models. We are interested in the case in which the modeler assumes a priori that the observed behaviour is generated by a causal system and forms a deterministic model of it as a state-determined machine. For any finite sequence of observed behaviour there will be an unbounded number of finite-state machines whose input-output behaviour, starting from a specified state, is identical to that observed. Out of these possible models there will be some such that no other machine has a smaller number of states. An optimal causal modeler is defined as one who always chooses such a minimal-state machine as a model.

What we now wish to derive is some relationship between the length of the sequence of observations and the number of states in the model formed by an optimal causal modeler. The following simple properties may be stated;

- 1) the number of states is a monotonically non-decreasing function of the length of the observed sequence of behaviour;
- 2) the number of states in the model of a sequence of behaviour generated by an M-state deterministic automaton cannot exceed M;
- 3) the number of states cannot exceed the length of the sequence of observations;
- 4) there are observed behaviours of length N such that the model must have M states (For example, consider the sequence $A^{N-1}B$ or the infinite sequence $ABA^2B^2A^3B^3 \dots$).

The first two properties seem to express the intuitive expectations of a human causal modeler—after a sufficiently long sequence of observations his model should stabilize. The last two properties show that it is possible, however, for the size of model to grow at the highest possible rate, precisely as the number of observations. That is, in terms of the discussion of human behaviour, an “unlucky” observer may obtain data requiring a very complex model. However, if the system generating the behaviour is relatively simple, for example, finite state, even if acausal, we might still expect our optimal causal modeler to perform in a similar way to when modeling a causal system. Clearly a stochastic source will generate sequences requiring as many states in the model as the length of the sequence, but it is plausible to suppose that such “pathological” sequences might be generated only infrequently by a finite-state automaton.

Hence a more interesting characteristic of the observer’s behaviour would be, not the maximum complexity model he might generate, but instead the expected number of states in the model generated (the average complexity). This might have one of three possible behaviours, as a function of the length of the observation sequence N. The expected number of states in the model formed by an optimal causal modeler of the behaviour of a finite-state stochastic automaton might be the following:

- a) asymptotic to a finite number, i.e., closely similar to the corresponding situation when the behaviour modeled is generated by a finite-state deterministic automaton;
- b) grow without limit but slower than the number of observations N itself, for example, as $\log N$; one might hypothesize that at least the ratio:

$$R_N = \frac{\text{expected number of states}}{\text{number of observations}} \quad (1)$$

would tend to zero as the number of observations increased;

- c) grow without limit at a rate similar to the maximum possible N; this would imply that nearly all sequences generated were “pathological” requiring maximum-size models growing as fast as the number of observations.

The recent development of theories of probability based on computational complexity (Martin-Lof, 1966; Kolmogorov, 1968; Chaitin, 1969; Willis, 1970) shows that case c) in fact occurs, and it is probably a measure of the extent to which the result is counterintuitive that this basis for probability has been so late in developing. We have so long accepted that a random sequence

may have any structure, including possibly one that appears highly deterministic, that it comes as a shock when it is suggested that we can apply a test for randomness to an individual sequence rather than a distribution. The paradox is resolved, of course, because almost all randomly generated sequences may be shown to have high computational complexity (Willis, 1970).

The following additional property shows the significance of these results in the present context case c) occurs and the ratio of (1) can tend to unity:

- d) The expected number of states in the model formed by an optimal causal modeler observing a sequence of behaviour of length N may be at least $N - \log_2 N - 2$.

Derivation: Consider a Bernoulli sequence of outputs of the form $(A + B)^*$ generated by a two-state stochastic automaton. There are 2^N possible sequences of N observations and if A and B are equiprobable so are the sequences. Now enumerate the number of different S -state deterministic automata that are available as models of these sequences. We note that a given automaton, starting in a given state, can act as a model of only one of the sequences, and we are concerned to find an upper bound on the variety of models available for a given S . This can be derived by noting that the general form of a model will be a transient chain followed by a cycle. If there are S states to be filled with two symbols and the cycle can commence in any state, there are at most $S2^S$ automata (some of which may generate the same sequence making this an upper bound on the number of different models available).

Since automata of this form with less than S states are also included amongst those having S states, $S2^S$ is an upper bound on the number of different automata with S states or less. Thus to have enough models available we must have

$$S2^S > 2^N \tag{2}$$

Now consider the mean number of states in the minimal forms of this set of models. There are at most $R2^R$ models with R states, and hence a lower bound of the mean number of states in the set whose maximum number of states is S is given by

$$\bar{n}_s > \frac{1}{(S-1)2^s} \sum_{R=1}^{S-1} R(R2^R - (R-1)2^{R-1}) \tag{3}$$

$$= S-2 + \frac{2}{S-1} \sum_{R=1}^{S-1} \frac{2}{(S-1)2^{S-R}} \tag{4}$$

$$> S-2. \tag{5}$$

Also, by taking logarithms to base 2 of inequality (2) we have

$$S + \log_2 S > N \tag{6}$$

so that

$$S > N - \log_2 S \tag{7}$$

but we have $S < N$ so that

$$S > N - \log_2 N \tag{8}$$

hence from (5)

$$\bar{S} > N - \log_2 N - 2 \quad (9)$$

that is, the mean number of states in the ensemble of models necessary to account for all possible behaviours is greater than $N - \log_2 N - 2$, where N is the number of observations.

Note that the ratio of (1), of states to observations, is

$$R_N = 1 - \frac{(\log_2 N + 2)}{N} \quad (10)$$

which is asymptotic to unity as N increases.

3. COMPUTATIONAL RESULTS

The theoretical bound obtained could be sharpened, but it clearly indicates the closeness between the expected number of states and the maximum number. Sequences requiring complex models are the common ones in the example used in the derivation of Property 5. It is of interest to see how closely the derived lower bound follows the actual mean number of states in an optimal causal module Minimal-state model. Minimal-state models were computed for each of the set of 2^N binary sequences for N from 1 through 18. The case analyzed theoretically, of a Bernoulli sequence with generating probability $p = 0.5$, is intuitively a worst case since it has no “structure”. To examine the effect of a bias towards one symbol on the expected size of model the effect of varying the generating probability was also investigated.

Fig. 1 is a set of graphs of the expected number of states in an optimal causal model against the number of observations for binary Bernoulli sequences with generating probabilities ranging from 0.0 to 0.5. The maximal upper bound at which the number of states equals the length of the observed sequence is shown as a diagonal line. The dashed line is derived from (4) and (6) of Theorem S and represents the best bound derivable from the theoretical argument.

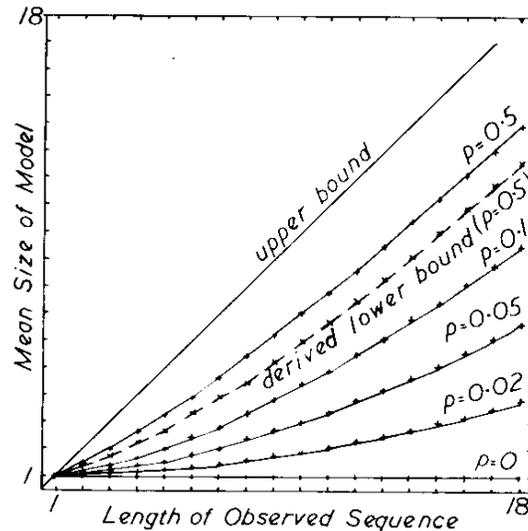


Fig. 1. Expected number of states in optimal causal model against number of observations for binary Bernoulli sequences of different generating probabilities.

When the generating probability p is zero the sequence is determinate and always represented by a one-state model. However, the introduction of the second symbol with a probability of one in

fifty ($p = 0.02$) gives rise to a steady growth in the size of a model with the number of observations. It is clear that the rate of growth of states for nonzero non-unity p must ultimately become at least $1/2$ since the only sequences of R A's and $N-R$ B's giving an average number of states less than $N/2$ are those with $R = 0$ and $R = N$. The probabilities of the occurrence of these sequences are p^N and $(1 - p)^N$, respectively, both of which decline to zero with increasing N unless p is zero or unity. Thus the slightest introduction of probabilistic acausality gives rise to models whose complexity is proportional to the length of observation.

4. CONCLUSIONS

The results of Sections 2 and 3 show that the assumption of causality when modeling acausal systems can lead to indefinitely complex models of comparatively simple systems. It is not just that perfect prediction has become impossible, but rather, that not accepting this leads to essentially meaningless "models" that cannot reflect any stochastic structure present. The example used in the derivation of Property 5 is precisely the stochastic automaton used in some of the experiments with human subjects described in Section 1. In Einstein's terms, if one assumes that God does not play dice, then when he does, one will obtain an over-complex view of the universe.

ACKNOWLEDGEMENT

I am grateful to Ian Witten of this Department for his critical comments on the first draft of this correspondence, and to Dr. John Andreae of the University of Christchurch, New Zealand, for stimulating this line of argument.

REFERENCES

- Andreae, J.H. (1973). Intelligent identification of signals from an unknown source. **Proceedings of National Electronics Conference**.
- Andreae, J.H. and Cashin, P.M. (1969). A learning machine with monologue. **International Journal of Man-Machine Studies** 1 1-20.
- Chaitin, G. (1969). On the length of programs for computing finite binary sequences. **Journal ACM** 16 145-159.
- Gaines, B.R. (1971). Memory minimization in control with stochastic automata. **Electronics Letters** 7(24) 710-711.
- Gaines, B.R. and Andreae, J.H. (1966). A learning machine in the context of the general control problem. **Proceedings of the 3rd Congress of the International Federation for Automatic Control**. London, Butterworths.
- Jauch, J.M. (1973). Determinism in classical and quantal physics. **Dialectica** 27 13-26.
- Kolmogorov, A.N. (1968). Logical basis for information theory and probability. **IEEE Transactions on Information Theory** IT-14 662-664.
- Martin-Lof, P. (1966). The definition of random sequences. **Information and Control** 9 602-619.
- Michotte, A.É. (1963). **The Perception of Causality**. New York, Basic Books.
- Popper, K.R. (1972). **Objective Knowledge: an Evolutionary Approach**. Oxford, Clarendon Press.
- Willis, D. (1970). Computational complexity and probability constructions. **Journal ACM** 17 241-259.