## THE ADEQUACY OF COMPLEXITY MODELS OF INDUCTION

R. Solomonoff, Zator Co., Camb., Mass., USA, 02138

We will discuss theorems showing that for a broad class of stochastic sources, Willis' model gives error rates close to minimum obtainable. Used as the basis of a betting system it again gives close to maximum yield. Cover's model, using Chaitin complexity, converges to the same optima but we are uncertain if the rate of convergence is as rapid.

Inductive inference can be defined to be the extrapolation of a binary sequence containing all the
data to be used in the induction. Almost all, if not
all, activities considered to be induction can be
put into this form. The Algorithmic complexity of a
sequence is related to the shortest programs for a
reference universal Turing machine producing that sequence as output.

Earliest application of algorithmic complexity was for inductive inference and it was shown that this complexity was relatively insensitive to choice of reference machine. Kolmogorov and Chaitin independently used somewhat different definitions of algorithmic complexity to define randomness of finite strings. Willis later developed the original induction system in more rigorous form and gave many theorems that greatly clarified its operation. Cover, apparently independently of Willis, used Chaitin's complexity for induction and employed a simple betting system to measure its efficiency.

A model equivalent to Willis' uses as reference a universal Turing machine having the "sequential property"- i.e. if input string a gives A as output, then any input of the form a b must have output of form A B. Willis assigns to string a, the probability

$$P_{w}(a) = \lim_{a \to \infty} N_{g}(a)/N_{g} \tag{1}$$

 $N_{\mathcal{Q}}(a)$  is the number of input strings of length  $\mathcal{Q}$  that result in output strings of the form a b, where b may be any finite (including the null) string.  $N_{\mathcal{Q}}$  is the number of input strings of length  $\mathcal{Q}$  that give outputs that are at least as long as a is.

Cover's model is based on Chaitin's complexity using a universal prefix machine as reference, assigning to a the probability

$$P_{c} = A \sum_{i} \tilde{z}^{c_{i}} \tag{2}$$

C; is the Chaitin complexity of the ith possible continuation of a, and A is a constant that assures that the total probability of all strings of a given length is unity.

While both of these definitions of probability are not effectively computable, there exist sequences of computable probability estimates that converge to equations (1) and (2).

Suppose we have a binary sequence that has been created by a stochastic generator. Through Bayes' theorem each of the two systems can be used to give the conditional probability of each bit in the sequence, given the antecedent subsequence.

First, we will use each system as the basis of a betting scheme at even odds. Suppose we start with a fortune of unity and bet a fraction  $p_i^l$  that the next bit will be 1 and a fraction  $p_i^l = 1-p_i^l$  that the next bit will be 0. At the end of n bets, our total fortune will be exactly  $2^n P_i$  (x(n)). Here x(n) is the binary sequence that has occurred, and  $P_i$  (x(n)) is the probability that our system has assigned to that sequence.

The log of Our fortune has a maximum expected value if

P<sub>i</sub> (x(n)) is the probability that would be assigned by the stochastic generator. If P<sub>i</sub> is a less perfect probability value the fortune will be somewhat less - depending on how good P<sub>i</sub> is. One goodness criterion of a system which we'll call b<sub>i</sub>, is the log of the ratio of its betting yield to the maximum possible.

$$b_i = log_2(P_i(x(n))/P(x(n))$$
 (3)

Another criterion considers the differences in conditional probability obtained by P and by P. Let  $\delta_n^i$  be the conditional probability by P, of the nth bit of x(n), given the previous n-l bits, i.e.

Then the expected value of 
$$\sum_{j=1}^{n} (\delta_{j}^{i} - \delta_{j}^{j})^{2}$$
 is a measure of the error in  $P_{i}$ . It can be shown that  $E(\sum_{j=1}^{n} (\delta_{j}^{i} - \delta_{j}^{i})^{2}) \equiv \sum_{n=1}^{n} (P(^{n} \times (n)) \sum_{j=1}^{n} (^{n} \delta_{j}^{i} - ^{n} \delta_{j}^{i})^{2}) \leqslant b_{i} \ln \sqrt{2}$  E is the expected value with respect to  $P$ .  $(n)$  is the kth sequence of length  $n$  (there are just  $2^{n}$  of them).  $(n)$  and  $(n)$  are the conditional probabilities of the jth bits of  $(n)$  for  $(n)$  and  $(n)$  respective-

ly. b<sub>i</sub> is the same as in equation (3).

If the underlying stochastic generator is describable in a finite number of bits, d, then for  $P_{i} = P_{w}$  (Willis' system),  $b_{w} = d$ .

If the generator is finitely describable, except for k differentiable parameters, (a differentiable parameter has an infinitely long description), then if the functional form is known and the parameters are known to within a certain error, b; will be of the form c log n, c being a constant. For such a generator, Willis' system will give b = c log n + d. c is the same constant as before, and d is the number of bits in the description of the functional form of the stochastic generator.

If the generator is ergodic, then Cover has shown that for his system,

$$\lim_{n\to\infty}\frac{b_{c}}{n}=0$$

In Willis' system  $\frac{b_{w}}{n}$  also approaches zero but at a known rate, probably as fast as is theoretically possible.

We don't yet have bounds on how rapidly  $\frac{b_c}{n}$  approaches zero - conceivably, it could do so very slowly. Preliminary analysis, however, suggests that the approach might be as rapid as that of Willis' method.

The systems described seem quite adequate for prediction and cast much light on some classic problems in inductive inference theory.

One possible approach to defining randomness of a finite sequence is that all future continuations of the sequence are about equally likely. The forgoing systems make it possible to put such a definition into an exact form and analyse its properties.

The problem of geometric probability is best understood by converting the data to digital form (using most any analog to digital conversion method) and analysing the data as a binary sequence. Various measure transformations have a simple interpretation in this light, and their significance (or lack of significance) in modifying probability values can be readily evaluated.

Goodman's paradoxes involving various linguistic transformation of inferential data, are also easy to analyse from this point of view.

The most important open problems in induction

theory are in the practical realization of systems such as have been described. We want to obtain the most accuracy in induction for a given amount of computing. This involves optimizing a function of algorithmic and computational complexities - a kind of problem that occurs in information retrieval and in some models of organic evolution.

## References

- 1. Solomonoff, R.J. "A Formal Theory of Inductive Inference." <u>Information and Control</u>, March 1964, pp. 1-22, June 1964, pp. 224-254.
- 2. Kolmogorov, A.N. "Logical Basis for Information Theory and Probability Theory." <u>IEEE Transactions on Information Theory IT-14</u>, pp. 662-664, 1968.

"Three Approaches to the Quantitative Definition of Information," <u>Information Transmission</u>, Vol. I, pp. 3-11, 1965.

- 3. Chaitin, G.J. "A Theory of Program Size Formal-ly Identical to Information Theory." To appear in Journal of the Association of Computing Machinery, July 1975 pp. 329-340.
- 4. Willis, D.G. "Computational Complexity and Probability Constructions," <u>Journal of the Assoc.</u>
  of Computing Machinery, April 1970, pp. 241-259.
- 5. Cover, T.M. "Universal Gambling Schemes and the Complexity Measures of Kolmogorov and Chaitin,"

  Report No. 12, Statistics Dept., Stanford University, Stanford, Calif., 1974. Submitted to Annals of Statistics.