# Centrum voor Wiskunde en Informatica

Centre for Mathematics and Computer Science

M. Li, P.M.B. Vitányi

Inductive reasoning and Kolmogorov complexity
(preliminary version)

1989

**CWI**

## Centrum voor Wiskunde en Informatica
Centre for Mathematics and Computer Science

M. Li, P.M.B. Vitányi

Inductive reasoning and Kolmogorov complexity
(preliminary version)

# Inductive Reasoning and Kolmogorov Complexity
## (Preliminary Version)

*Ming Li*

Department of Computer Science, York University
North York, Ontario M3J 1P3, Canada

*Paul M.B. Vitányi*

Centrum voor Wiskunde en Informatica, Kruislaan 413,
1098 SJ Amsterdam, The Netherlands
and
Universiteit van Amsterdam, Faculteit Wiskunde en Informatica

## ABSTRACT

Reasoning to obtain the "truth" about reality, from external data, is an important, controversial, and complicated issue in man's effort to understand nature. (Yet, today, we try to make machines do this.) There have been old useful principles, new exciting models, and intricate theories scattered in vastly different areas including philosophy of science, statistics, computer science, and psychology. We focus on inductive reasoning in correspondence with ideas of Solomonoff. While his proposals result in perfect procedures, they involve the noncomputable notion of Kolmogorov complexity. In this paper we develop the thesis that Solomonoff's method is fundamental in the sense that many other induction principles can be viewed as particular ways to obtain computable approximations to it. We demonstrate this explicitly in the cases of Gold's paradigm for inductive inference, Valiant's learning (by adding computational requirements), Rissanen's minimum description length (MDL) principle, Fisher's maximum likelihood principle, and Jaynes's maximum entropy principle. We present several new theorems and derivations to this effect. We also delimit what can be learned and what cannot be learned in terms of Kolmogorov complexity, and we describe an experiment in machine learning of handwritten characters.

*1980 Mathematics Subject Classification:* 68C25, 68F05, 68G05, 60A05, 62A10, 62F15, 62B10, 94A17.

*CR Categories:* F.2, F.4.

*Keywords and Phrases:* Solomonoff's Inductive Reasoning, Kolmogorov complexity, Algorithmic Information Theory, Bayesian Inference, Universal A Priori Probability, Solomonoff-Levin Distribution, Inductive Inference, Gold's Paradigm, Valiant learning, Rissanen's Minimum Description Length Principle, Fisher's Maximum Likelihood, Jaynes's Maximum Entropy Principle

*Note:* To appear in: *Proc. 4th Annual IEEE Structure in Complexity Theory Conference*, 1989.

# Inductive Reasoning and Kolmogorov Complexity

## (Preliminary Version)

*Ming Li*

Department of Computer Science, York University
North York, Ontario M3J 1P3, Canada

*Paul M.B. Vitányi*

Centrum voor Wiskunde en Informatica, Kruislaan 413,
1098 SJ Amsterdam, The Netherlands
and
Universiteit van Amsterdam, Faculteit Wiskunde en Informatica

*The eye of the understanding is like the eye of the sense; for as you may see great objects through small crannies or levels, so you may see great axioms of nature through small and contemptible instances.* [Francis Bacon, *Sylva Sylvarum* 337, 1627]

## 1. A Historical View of Inductive Reasoning

The Oxford English Dictionary gives as the meaning

of **induction:** *the process of inferring a general law or principle from the observations of particular instances.* This defines precisely what we would like to call *inductive inference.* On the other hand, we regard *inductive reasoning* as a more general concept than inductive inference, namely as a process of re-assigning a probability (or credibility) to a law or proposition from the observation of particular instances. In other words, in the way we use the notions, inductive inference draws conclusions that consist in *accepting or rejecting* a proposition, while inductive reasoning only changes the degree of our belief in a proposition. The former is a special case of the latter. In this paper we discuss inductive reasoning in correspondence with R. Solomonoff's ideas as expressed in e.g. [44]. However, Solomonoff's procedure is not effective, since it involves the non-computable Kolmogorov complexity of objects. We

shall show, however, that there is considerable structure in many different approaches proposed for induction, since they can be variously derived as computable approximations to Solomonoff's method.

The history of inductive inference, which is as old as empirical science itself, dates back at least to the Greek philosopher of science Epicurus (342? - 270? B.C). To reason by induction is nothing but to learn from experience. As the sun rises day by day, our belief in that the sun will rise tomorrow increases, and we eventually infer the *truth* that the sun *will* rise every morning. As human history evolves, man tries to understand and explain the events that happen around him: this takes the form of different induction methods to formulate scientific theories from positive and negative, fortunate and unfortunate, lucky and unlucky, happy and miserable experiences. Two metaphysical principles stand out and prevail today: the principle of Epicurus' multiple explanations (or indifference) and Occam's principle of simplest explanation (Occam's razor).

The **Principle of Multiple Explanations:** *If more than one theory is consistent with the data, keep them all.*

The source of the following material is [4]. Epicurus, in his *Letter to Pythocles*, explains that: There are cases, especially of events in the heavens such as the risings and settings of heavenly bodies and eclipses, where it is sufficient for our happiness that several explanations be discovered. In these cases, the events "have multiple causes of coming into being and a multiple predication of what exists, in agreement with the perceptions." Epicurus maintains that, if several explanations are in agreement with the (heavenly) phenomena, then we must keep all of them for two reasons. Firstly, the degree of precision achieved by multiple explanations is sufficient for human happiness. Secondly, it would be unscientific to prefer one explanation to another when both are equally in agreement with the phenomena. This, he claims, would be to "abandon physical inquiry and resort to myth." His follower Lucretius (95 - 55 B.C.) illustrates the inevitablity of the use of the multiple explanation principle by the following example:

There are also some things for which it is not enough to state a single cause, but several, of which one, however, is the case. Just as if you were to see the lifeless corpse of a man lying far away, it would be fitting to state all the causes of death in order that the single cause of this death may be stated. For you would not be able to establish conclusively that he died by the sword or of cold or of illness or perhaps by poison, but we know that there is something of this kind that happened to him. [Lucretius 6. 703-11]

Based on the same intuition, in the calculus of probabilities it has been customary to postulate the "principle of indifference" or the "principle of insufficient reason". When there is no other evidence, because of the absolute lack of knowledge concerning the conditions under which a die falls, we have no reason to assume that a certain face has higher probability of turning up. Hence we assume that each side of the die has the probability 1/6. The **principle of indifference** considers events to be equally probable if we have not the slightest knowledge of the conditions under which each of them is going to occur. For the case of a die, this actually coincides with the so-called "maximum entropy principle", we will discuss later, which states that we should choose probabilities $p_i$ for face $i$ to be the outcome of a trial, $i = 1, 2, ..., 6$, such that $-\sum_{i=1}^{6} p_i \ln p_i$ is maximised under the only constraint $\sum_{i=1}^{6} p_i = 1$. We obtain precisely $p_i = 1/6$ for $i = 1, 2, ..., 6$.

The second and more sophisticated principle is the celebrated Occam's razor principle commonly attributed to William of Ockham (1290? - 1349?). This enters the scene about 1500 years after Epicurus. In sharp contrast to the principle of multiple explanations, it states:

**Occam's Razor Principle:** *Entities should not be multiplied beyond necessity.*

This is generally interpreted as: Among the several theories that are all consistent with the observed phenomena, one should pick the simplest theory. (According to Bertrand Russell, the actual phrase used by Ockham was: "It is vain to do with more what can be done with fewer.") Surely Occam's razor principle is easily understood from an 'utilitarian' point of view: if both theories explain the same set of facts, why not use the simpler theory?! However things become more intricate when we want to know whether a simpler theory is really better than the more complicated one. This also raises another question which has been a bone of contention in Philosophy ever since the razor's inception. For what is the proper measure of simplicity? Is $x^{100} + 1$ more complicated than $ax^{17} + bx^2 + cx + d$? E.g., the distinguished contemporary philosopher K. Popper pronounced that the razor is without sense or use on such grounds. However, it is interesting to notice that the principle can be given objective contents, and has recently been very successfully applied in many different forms in computational learning theory.

To explain this, let us consider an oversimplified example of inferring a finite automaton with one-letter input using Occam's razor principle.
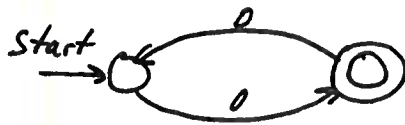
Accepted inputs: 0, 000, 00000, 000000000;

Rejected inputs: $\epsilon$, 00, 000000;

For these data, there exist many consistent finite automata. The most trivial one would be:

The smallest automaton is:



Intuitively, the first machine does not tell us anything. We therefore do not expect that machine to predict unknown data. On the other hand the second machine makes an *inference* that the language accepted consists of strings of an odd number of 0's, and this can be used for predicting future data. However, a too simplistic application of Occam's razor principle does not have a very good explanatory value, as the following story illustrates.

Once upon a time, there was a little girl named Emma. Emma had never eaten a banana, nor had she been on a train. One day she went for a journey from New York to Pittsburgh by train. To relieve Emma's anxiety, her mother gave her a large bag of bananas. At Emma's first bite of a banana, the train plunged into a tunnel. At the second bite, the train broke into daylight again. At the third bite, Lo! into a tunnel; the fourth bite, La! into daylight again. And so on all the way to Pittsburgh and to the bottom of her bag of bananas. Our bright little Emma (applying Occam's razor principle?) told her grandpa: "Every odd bite of a banana makes you blind; every even bite puts things right again." [After N.R. Hanson, "Perception and Discovery", Freeman, Cooper & Co, 1969, p.359.]

Let us consider how the idea of 'simplicity' affects a scientist's thinking. We refer to a beautiful study of simplicity by Kemeny [25]. Initially, there were no new facts that failed to be explained by the Special Theory of relativity. The incentive to invent the General Theory of Relativity, by Albert Einstein, was his conviction that the Special Theory was not the *simplest* theory that can explain all the observed facts. Reducing the number of independent variables obviously simplifies a theory. By the requirement of general covariance Einstein succeeded in replacing the previous independent 'gravitational mass' and 'inertial mass' by a single concept.

In spite of the apparent universal acceptance of Occam's razor, consciously or unconsciously, the concept of simplicity remains highly controversial. Generally speaking, it has remained a crude non-precise idea. Things are subtler than they appear. Is the following formulation precise?

**Occam's Razor Rule:** *Select a hypothesis which is as well in agreement with the observed value as possible; if there is any choice left, choose the simplest possible hypothesis.*

**Example.** Consider the problem of fitting $n$ points by a polynomial. The above rule tells us to choose the polynomial of lowest degree passing through all the $n$ points. But due to measurement precision and possible noise, whatever degree polynomial the points originated from, we will end up a polynomial of degree $n-1$ which fits the data precisely. But this polynomial most likely does not help us to predict future points.

**Example.** Consider another example given by Kemeny: Let there be unknown number of white balls and black balls in a sealed urn. Through an opening you randomly pick one ball at a time, note its color and replace it, and shake the urn thoroughly. After $n$ draws you must decide what fraction of the balls in the urn is white. The possible hypotheses state that some rational fraction $r$ of balls in the urn is white, where $0 \leqslant r \leqslant 1$. By the above rule, if in $n$ draws, $m$ white balls are selected, then we should formulate the hypothesis $r = m/n$. Let there be 1/3 white and 2/3 black balls. However the probability of getting the true hypothesis $m = n/3$ is zero if $n$ is not divisible by 3, and it tends to zero, even under the assumption that $n$ is divisible by 3. On the other hand we know that to obtain a hypothesis $1/3 - \epsilon \leqslant r \leqslant 1/3 + \epsilon$, for any $\epsilon$, has probability tending to 1 exponentially fast, by the so-called Chernoff formula. (For Chernoff's formula see e.g. [2]. ) Even when the process converges, $n$ may be too large for practical use.

**Kemeny's Rule.** *Select the simplest hypothesis compatible with the observed values.*

Here 'compatible' is defined as follows. The hypothesis $H_i$ is *compatible* with data $D$ if, assuming the truth of $H_i$, there was at most one percent chance of getting a deviation as great as $m(H_i, D)$ for some measure function $m$. This is related to Valiant's learning theory to be discussed later.

But how does one define simplicity? Is 1/4 simpler than 1/10? Is 1/3 simpler than 2/3? Saying that an urn contains 1/3rd part white balls comes down to the same thing as saying that it contains a 2/3rd part black balls. Kemeny warned: "do not use more precision in your theories than is necessary." But what is necessary and what is not? All these issues are very subjective. Does a simple theory generate a hypothesis which is good for predicting future outcomes? How do we achieve fast convergence? How does one trade between 'simplicity' and 'truth' ('compatibility')? Kemeny actually asked for "a criterion combining an optimum of simplicity and compatibility" [crediting Nelson Goodman for this suggestion].

## 1.1. Combining Epicurus, Ockham, and Bayes

The study of inductive reasoning predates artificial intelligence or computer science by more than 2000 years. There is tremendous amount of literature in many different fields under diverging terminologies. Our goal is to extract a common core of simple ideas underlying all these approaches, in the spirit of Occam's Razor principle. We will start with Bayesian inference theory.

To apply Bayesian type reasoning one has to assign *a priori* probabilities (prior probability) to each possible hypothesis. Since the posthumous publication in 1763 of Th. Bayes's (??- 1761) famous memoir 'An essay towards solving a problem in the doctrine of chances' by his friend Richard Price, [5], there has been continuous bitter debate on the controversial prior probability in the Bayesian formula.

The invention of Kolmogorov complexity, by its first inventor R. Solomonoff, was as an auxilliary notion to resolve this particular problem. Namely, using Kolmogorov complexity he found a single 'universal' prior distribution which can be substituted for any particular actually valid distribution (as long as it is computable) in Bayes's formula, and obtain approximately as good results as if the actually valid distribution had been used. It sounds like magic, but Solomonoff's approach does give a more or less satisfactory solution to this unlikely objective.

The elegant idea of a universal prior is a combination of Occam's razor and modern computability theory. However, the universal prior is uncomputable, since it involves Kolmogorov complexity. In this paper we develop the thesis that many theories, models, and principles for inductive reasoning that were formulated both before and after Solomonoff's inception, can be rigorously derived as particular computable approximations to it.

In this paper we first describe the basics of Bayesian theory and how to apply Kolmogorov complexity to obtain the Universal prior probability distribution. We then derive the Gold paradigm and its principles. Combination with ideas from computational complexity theory leads to Valiant's model of deductive learning. We derive a form of Rissanen's Minimum Description Length (MDL) principle. From the MDL principle Rissanen derives the Fisher's Maximum Likelihood principle and Jaynes Maximum Entropy principle. This paper contains a review of all these theories and principles. We were not satisfied with the fact that some experts say the connections as claimed are obvious, while some other experts deny those connections exist. Hence, in this paper we explicitly establish them. That is, we forge some new links, connections, explanations, and provide explicit derivations together with the appropriate related theorems. We also present a theorem about what can be inferred correctly, and what cannot be inferred correctly, in terms of Kolmogorov complexity. We describe an experiment we have performed in machine learning of recognition of handwritten characters using the MDL principle. A more extensive treatment of this material will be given in our forth-coming textbook [34].

## 2. The Universal Prior Distribution

### 2.1. Bayesian Inference

In the following discussion of probability we assume the usual so-called Kolmogorov Axioms, see e.g. [11]. For our purpose we need the following. We have a *hypothesis space*, $\mathrm{H} = \{H_1, H_2, ...\}$, which consists of a countable set of hypotheses, which are mutually exclusive (in the sense that at most one of them is right), and exhaustive (in the sense that at least one of them is right). With each hypothesis $H$ we associate a *probability* $P(H)$ such that $\sum P(H) = 1$. The student is supplied with some data $D$, providing information about which hypothesis is correct. We provide the *conditional probability* $P(D \mid H)$ that the data $D$ arises given that $H$ is the correct hypothesis. We assume that $P(D \mid H)$ is computable from $D$ and $H$. From the definition of conditional probability, i.e., $P(A \mid B) = P(A \cap B)/P(B)$, it is easy to derive **Bayes's formula** (rewrite $P(A \cap B)$ in two different ways):

$$P(H_i \mid D) = \frac{P(D \mid H_i) P(H_i)}{P(D)} \qquad (1)$$

where

$$P(D) = \sum_i P(D \mid H_i) P(H_i).$$

We interpret the different variables in the formula as follows. The $H_i$'s represent the possible alternative hypotheses concerning the phenomenon we wish to discover. The term $D$ represents the empirically or otherwise known data concerning this phenomenon. The term $P(D)$, the probability of data $D$, can be considered as a normalizing factor so that $\sum_i P(H_i \mid D) = 1$. The term $P(H_i)$ is called the *prior* probability or *initial* probability of hypothesis $H_i$, i.e., the probability that $H_i$ is true before we have seen any evidence. The term $P(H_i \mid D)$ is called the *final, a posteriori*, or *inferred* probability, which reflects the probability of $H_i$ modified from the prior probability $P(H_i)$ after seeing the data $D$. The term $P(D \mid H_i)$, the conditional probability of seeing $D$ when $H_i$ is true, is assumed to be computable from $D$ and $H_i$. In many learning situations, data can only be consistent with an hypothesis $H_i$ in the sense of being forced by it such that $P(D \mid H_i) = 1$. If the data is inconsistent with hypothesis $H_i$ then $P(D \mid H_i) = 0$. In such a situation, the data either is determined by a hypothesis, or disqualifies it. (We assume there is no noise that distorts the data.) (For example, the hypothesis is datum $x \in L$ or $x \notin L$.)

The most interesting term is the prior probability $P(H_i)$. In the context of machine learning, $P(H_i)$ is often considered as the learner's *initial degree of belief* in hypothesis $H_i$. In essence Bayes's rule is a *mapping* from *a priori* probability $P(H_i)$ to *a posteriori* probability $P(H_i \mid D)$, where the

mapping is determined by data $D$. In general, the problem is not so much that in the limit the inferred probability would not 'condense' on the 'true' hypothesis, but that the inferred probability gives as much information as possible about the possible hypotheses from only a limited number of data, cf. example below. In fact, the continuous bitter debate between the Bayesian and non-Bayesian opinions centered on the *prior* probability. The controversy is caused by the fact that Bayesian theory does not say how to initially derive the prior probabilities for the hypotheses. Rather, Bayes's rule only says how they are to be *updated*. However, in each actual case the prior probabilities may be unknown, uncomputable, or conceivably do not exist. (What is the prior probability of use of words in written English? There are many different sources of different social backgrounds living in different ages.) This problem is solved if we can find a *single* probability distribution to use as the prior distribution in each different case, with approximately the same result as if we had used the real distribution. Surprisingly, this turns out to be possible up to some mild restrictions.

**Historical Note.** Properly speaking, formula (1) is not due to Bayes, but it is due to P.S. Laplace (1749 - 1827) who stated the formula and attached Bayes's name to it [29]. Actually, Bayes in his original paper [5] assumed the uniform distribution for the a priori probability, hence he has essentially derived $P(H_i | D) = P(D | H_i) / \sum_i P(D | H_i)$. Although this formula can be derived from (1) by simply assuming that all $P(H_i)$ are the same, Bayes did not state his result in the general form as in (1), nor did he derived his result through a formula similar to (1). Despite the fact that Bayes's rule is just a rewriting of the definition of conditional probability and nothing more, it is its interpretation and applications that are most profound and caused much bitter controversy during the past two centuries.

**Example 1.** We use an example of von Mises [37]. Let an urn contain many dice with different attributes. (For convenience, the attribute space is discrete.) A die with attribute $p$ has probability $p$ showing "6" in a random throw. If we randomly draw a die from the urn, how can we determine its attribute? Let $H_p$ be the event of drawing a die with attribute $p$ from an urn. We draw a die from the urn and throw $n$ times independently. Let $D_{m,n}$ be the experimental data such that $m$ successes (6's) were observed out of $n$ throws. So

$$P(H_p | D_{m,n}) = \frac{P(D_{m,n} | H_p) P(H_p)}{P(D_{m,n})},$$

where $P(D_{m,n}) = \sum_p P(D_{m,n} | H_p) P(H_p)$. According to Chernoff's formula (see [2] ), for $\alpha < 1$, we have:

$$P(m - np > \alpha np | H_p) < e^{-\frac{1}{2}\alpha^2 np},$$

$$P(np - m > \alpha np | H_p) < e^{-\frac{1}{3}\alpha^2 np}.$$

Hence, if $p_0$ is the attribute of the die we have drawn, then $P(D_{m,n} | H_p)$ goes to 0 exponentially fast when the number of experiments increases, for all terms with $p \neq p_0$; And $P(H_{p_0} | D_{m,n})$ goes to 1. Hence in polynomial time with respect to the minimum distance between the $p$'s we can derive $p_0$ with high probability.

**Example 2.** We explain a simple version of Solomonoff's theory of inductive inference. Consider theory formation in science as the process of obtaining a compact description of the past observations together with predictions of future ones. The investigator observes increasingly larger initial segments of an infinite binary sequence as the outcome of an infinite sequence of experiments on some aspect $X$ of nature. To describe the underlying regularity of this sequence, the investigator tries to formulate a theory that governs $X$, on the basis of the outcome of past experiments. Candidate theories (hypotheses) are identified with computer programs that compute binary sequences starting with the observed initial segment.

First assume the existence of prior probability distribution $P$ over $\Omega = \{0,1\}^*$. Define the measure $\mu$ over $\Omega$ by $\mu(x) = \sum \{P(xy): y \in \Omega\}$. (Actually, because of the subadditive property below, $\mu$ is not a proper measure. One may call $\mu$ a 'semimeasure'.) Thus, $\mu(x)$ is the probability of a sequence starting with $x$. Given a previously observed data string $S$, the inference problem is to predict the next symbol in the output sequence, i.e., extrapolating the sequence $S$. In terms of the variables in Formula (1), $H_i$ is the hypothesis that the sequence under consideration starts with initial segment $S a$. The data $D$ consist in the assertian that the sequence in fact starts with initial segment $S$. Thus, for $P(H_i)$ and $P(D)$ in Formula (1) we substitute $\mu(S a)$ and $\mu(S)$, respectively, and obtain, $a = 0$ or $a = 1$,

$$P(Sa | S) = \frac{P(S | Sa)\mu(Sa)}{\mu(S)}.$$

Here
$\mu(S) = P(S | S0)\mu(S0) + P(S | S1)\mu(S1) + P(S)$.
When the string $S$ is generated by a deterministic process, we must have $P(S | Sa) = 1$ for any $a$, hence,

$$P(Sa | S) = \frac{\mu(Sa)}{\mu(S)}. \tag{2}$$

In terms of inductive inference or machine learning, the final probability $P(Sa | S)$ is the probability of the next symbol being $a$, given the initial sequence $S$. The goal of inductive inference in general is to be able to infer an underlying effective process (in the most general case, a Turing machine, according to the Church-Turing thesis) that generated $S$, and hence to be able to predict (extrapolate) the next symbol. Obviously we now only need the prior probability to evaluate $P(Sa | S)$.

In order to solve the problem for unknown

prior probability, Solomonoff proposed what he called a **universal prior distribution**. We now carefully define the universal prior distribution and prove several fundamental theorems due to Solomonoff and Levin, and afterwards continue this example. The definitions and theorems are so fundamental that our approach totally rests upon them. However since these results have only been published either by Solomonoff in a veiled fashion, or by Levin in the Russian literature or in a complicated form [50], or are unpublished by Gàcs [17], they are almost unknown except to a few people doing research in this area. First we need the basic definitions of Kolmogorov complexity.

## 2.2. Kolmogorov Complexity

Inductive reasoning was the midwife that stood at the cradle of Kolmogorov complexity. Nowadays, Kolmogorov complexity has been applied in many areas of computer science and mathematics (see [32] for a general survey), and few realize that Kolmogorov complexity was at first invented for the purpose of inductive inference. In this essay, we go back to this origin.

We are interested in defining the complexity of a concrete individual finite string of zeros and ones. Unless otherwise specified, all strings will be binary and of finite length. All logarithms in this paper are base 2, unless it is explicitly noted they are not. If $x$ is a string, then $l(x)$ denotes the *length* (number of zeros and ones) of $x$. We identify throughout the $x$th finite binary string with the natural number $x$, according to the correspondence:

$$(\epsilon, 0), (0, 1), (1, 2), (00, 3), (01, 4), (10, 5),...$$

Intuitively, we want to call a string simple if it can be described in a few words, like "the string of a million ones"; A string is considered complex if it cannot be so easily described, like a "random" string which does not follow any rule and hence we do not know how to describe apart from giving it literally. A description of a string may depend on two things, the decoding method (the machine which interprets the description) and outside information available (input to the machine). We are interested in descriptions which are effective, and restrict the decoders to Turing machines. Without loss of generality, our Turing machines use binary input strings which we call *programs*. More formally, fixing a Turing machine $T$, we would like to say that $p$ is a description of $x$ if, on input $p$, $T$ outputs $x$. It is also convenient to allow $T$ to have some extra information $y$ to help to generate $x$. We write $T(p,y)=x$ to mean that Turing machine $T$ with input $p$ and $y$ terminates with output $x$.

**Definition 1.** The descriptional complexity $C_T$ of $x$, *relative* to Turing machine $T$ and binary string $y$, is defined by

$$C_T(x\mid y) = \min\{l(p): p \in \{0,1\}^*, \ T(p,y)=x\},$$

or $\infty$ if such $p$ do not exist.

The complexity measure defined above is useful and makes sense only if the complexity of a string does not depend on the choice of $T$. Therefore the following simple theorem is vital. This Invariance Theorem is given by Solomonoff [44], Kolmogorov [26], and Chaitin [8].

**Theorem 1.** *There exists a universal Turing machine $U$, such that, for any other Turing machine $T$, there is a constant $c_T$ such that for all strings $x,y$, $C_U(x\mid y) \leq C_T(x\mid y)+c_T$.*

*Proof.* Fix some standard enumeration of Turing machines $T_1,T_2,\cdots$. Let $U$ be the Universal Turing machine such that when starting on input $0^n 1p$, $p \in \{0,1\}^*$, $U$ simulates the $n$th Turing machine $T_n$ on input $p$. For convenience in the proof, we choose $U$ such that if $T_n$ halts, then $U$ first erases everything apart from the halting contents of $T_n$'s tape, and also halts. By construction, for each $p \in \{0,1\}^*$, $T_n$ started on $p$ eventually halts iff $U$ started on $0^n 1p$ eventually halts. Choosing $c_T=n+1$ for $T_n$, finishes the proof. □

Clearly, the Universal Turing machine $U$ that satisfies the Invariance Theorem is *optimal* in the sense that $C_U$ minorizes each $C_T$ up to a fixed additive constant (depending on $U$ and $T$). Moreover, for each pair of Universal Turing machines $U$ and $U'$, satisfying the Invariance Theorem, the complexities coincide up to an additive constant (depending only on $U$ and $U'$), for all strings $x, y$:

$$|C_U(x\mid y) - C_{U'}(x\mid y)| \leq c_{U,U'}.$$

Therefore, we set the canonical *conditional Kolmogorov* complexity $C(x\mid y)$ of $x$ under condition of $y$ equal to $C_U(x\mid y)$, for some fixed optimal $U$. We call $U$ the *reference* Turing machine. Hence the Kolmogorov complexity of a string does not depend on the choice of encoding method and is well-defined. Define the *unconditional Kolmogorov* complexity of $x$ as $C(x) = C(x\mid \epsilon)$, where $\epsilon$ denotes the empty string ($l(\epsilon) = 0$).

**Definition 2.** For technical reasons, we need in the following the so-called *prefix* complexity or *self-delimiting* complexity, rather than $C(x)$ from Definition 1. The two complexities are approximately the same, namely, for each $x$ they coincide to within an additive term of $2\log C(x)$. We use a special model of Turing machine, viz. one with three tapes: a one-way input tape, a one-way output tape, and a two-way work tape. Initially, the input tape contains an indefinitely long sequence of bits. If the machine halts, then the initial segment on the input tape it has read up till that time is considered the *input* or *program*, and the contents of the output tape is the *output*. Clearly, the set of programs of each such machine is a prefix-code. (Recall that if $p$ and $q$ are two code words of a *prefix-code*, then $p$ is not a proper prefix of $q$.) Such a mchine is called a *prefix machine*. We can give an effective enumeration of all prefix machines in the standard way. Then

the *self-delimiting* descriptional complexity of $x \in \{0,1\}^*$, with respect to prefix machine $T$, and binary string $y$, is similarly defined as

$$K_T(x \mid y) = \min\{l(p): p \in \{0,1\}^*, \ T(p,y)=x\},$$

or $\infty$ if such $p$ do not exist. One can prove an Invariance Theorem for self-delimiting complexity, and define the conditional and unconditional *self-delimiting Kolmogorov complexity*, by fixing some reference optimal prefix machine, in exactly the same way as before, so we do not repeat the construction.

*Remark.* The self-delimiting Kolmogorov complexity of string $x$, is the length of the shortest self-delimiting program that outputs $x$. Mathematically, self-delimiting Kolmogorov complexity has nicer properties. In this exposition, we will use $K(x)$ to denote the self-delimiting Kolmogorov complexity of $x$. $C(x)$ and $K(x)$ differ by at most a $2\log K(x)$ additive term. In some applications this does not make any difference. But in some other applications, for example inductive inference, this is vital.

**Definition 3.** A binary string $x$ is *incompressible* if $K(x) \geqslant l(x)$.

*Remark.* Since Martin-L:*of [36] has shown that incompressible strings pass all effective statistical tests for randomness, we will also call incompressible strings *random* strings. A simple counting argument shows that most strings are random. The theory of computability shows that the function $K(x)$ is noncomputable, but can be approximated to any degree of accuracy by a computable function. However, at no point in this approximation process can we know the error. Cf. also the surveys [33, 50].

## 2.3. Semicomputable Functions and Measures

We consider recursive functions with values consisting of pairs of natural numbers. If $<p, q>$ is such a value then we interpret this value as the rational number $p/q$, and say that the recursive function is *rational valued*.

**Definition.** A real function $f$ is *semicomputable from below* iff there exists a recursive function $g(x, k)$ with rational values (or, equivalently, a computable real function $g(x, k)$), nondecreasing in $k$, with $f(x) = \lim_{k \to \infty} g(x, k)$. A function $f$ is semicomputable from above, if $-f$ is semicomputable from below.

(An equivalent definition: $f$ is a function that is semicomputable from below if the set $\{(x, r): r \leqslant f(x), r \text{ is rational}\}$ is recursively enumerable.)

A real function $f$ is computable iff there is a recursive function $g(x, k)$ with rational values, and $|f(x) - g(x, k)| < 1/k$.

Obviously, all recursive functions are computable, and all computable functions are semicomputable. However, not all semicomputable functions

are computable, and not all computable functions are recursive. Nontrivial examples of functions that are semicomputable from above but not computable are $C(x)$, $C(x \mid y)$, $K(x)$, and $K(x \mid y)$.

The following analysis is a simplified version over the discrete space $N$ (or the set of finite binary strings), of Zvonkin and Levin [50]. We follow to some extent [17]. Functions $\mu: N \to [0, 1]$ that satisfy the usual properties of probability distributions except that

$$\sum_x \mu(x) \leqslant 1.$$

we shall call *measures*. We say that a measure $\mu$ (multiplicatively) *dominates* a measure $\mu'$ if there exists a constant $c$ such that, for all $x$ in $N$, we have $\mu'(x) \leqslant c\,\mu(x)$. It is known from the calculus that *no* measure $\mu$ dominates all measures: for each measure $\mu$ there is a measure $\mu'$ such that $\lim \mu'(x)/\mu(x) = \infty$. However, if we restrict ourselves to the class of semicomputable measures, then it turns out that this class contains an "absorbing" element, a measure that dominates all measures in the class. We call the measure that dominates all other measures in a given class a *universal measure* for that class. This important observation that such a measure exists was first made by Levin [50].

**Theorem 2.** *The class of measures that are semicomputable from below contains a universal measure.*

*Proof.* First we consider the standard enumeration of all partial recursive functions $\phi_1, \phi_2, \ldots$. We change each $\phi$ into a partial recursive function $\psi$ with the same range as $\phi$ but with, for each $x$, the value of $\psi(<x, k>)$ is defined only if $\psi(<x, 1>), \psi(<x, 2>), \ldots, \psi(<x, k-1>)$ are defined. (Assign values to arguments in enumeration order.) We use each $\psi$ to define a semicomputable function $s$ by approximations $s^k(x)$, $k = 1, 2, \ldots$, from below:

$$s(x) = \sup \{s^k(x): s^k(x) = p/q,$$

$$\psi(<x, k>) = <p, q>, k = 1, 2, \ldots\}.$$

The resulting $s$-enumeration contains all semicomputable functions. Next we use each semicomputable function $s$ to compute a measure $\mu$ from below. Initially, set $\mu(x) = 0$ for all $x$. If $s(1)$ is undefined then $\mu$ will not change anymore and it is trivially a measure. Otherwise, for $k = 1, 2, \ldots$, if $s^k(1) + s^k(2) + \ldots + s^k(k) \leqslant 1$ then set $\mu(i) := s^k(i)$ for $i = 1, 2, \ldots, k$, else the computation of $\mu$ is finished.

There are three mutually exclusive ways the computation of $\mu$ can go, exhausting all possibilities. Firstly, $s$ is already a measure and $\mu := s$. Secondly, for some $x$ and $k$ with $x \leqslant k$ the value $s^k(x)$ is undefined. Then the values of $\mu$ do not change anymore from $\mu(i) = s^{k-1}(i)$ for $i = 1, 2, \ldots, k-1$, and $\mu(i) = 0$ for $i \geqslant k$, even though the computation of $\mu$ goes on forever. Thirdly, there is a first $k$ such

that $s^k(1) + s^k(2) + ... + s^k(k) > 1$, that is, the new approximation of $\mu$ violates the condition of measure. Then the approximation of $\mu$ is finished as in the second case. But in this case the algorithm terminates, and $\mu$ is even computable.

Thus, the above procedure yields an effective enumeration $\mu_1, \mu_2, ...$ of all semicomputable measures. Define the function $\mu_0$ as:

$$\mu_0(x) = \sum_n 2^{-n} \mu_n(x).$$

It follows that $\mu_0$ is a measure since

$$\sum_x \mu_0(x) = \sum_n 2^{-n} \sum_x \mu_n(x) \leqslant \sum_n 2^{-n} = 1.$$

The function $\mu_0$ is also semicomputable from below, since $\mu_n(x)$ is semicomputable from below in $n$ and $x$. (Use the universal partial recursive function $\phi_0$ and the construction above.) Finally, $\mu_0$ dominates each $\mu_n$ since $\mu_0(x) > 2^{-n} \mu_n(x)$. Therefore, $\mu_0$ is a universal semicomputable measure. $\square$

Obviously, there are countably infinite universal semicomputable measures. We now fix a *reference* universal semicomputable measure $\mu_0(x)$, and denote it by $\mathbf{m}(x)$. It will turn out that function $\mathbf{m}(x)$ adequately captures Solomonoff's envisioned universal *a priori* probability. It represents Levin's approach to the great (but technically unsatisfactory) original idea of Solomonoff.

If a semicomputable measure is also a probability distribution then it is computable. Namely, if we compute an approximation $\mu^k$ of the function $\mu$ from below for which $\sum_x \mu^k(x) > 1 - \epsilon$, then we have $|\mu(x) - \mu^k(x)| < \epsilon$ for all $x$.

Any positive function $w(x)$ such that $\sum_x w(x) \leqslant 1$ must converge to zero. Hence $\mathbf{m}(x)$ converges to zero as well. However, it converges to zero slower than any positive computable function that converges to zero. That is, $\mathbf{m}(x)$ is not computable, and therefore it is not a proper probability distribution: $\sum_x \mathbf{m}(x) < 1$. There is no analogous result to Theorem 2 for computable measures: amongst all computable measures there is no universal one. This fact is one of the reasons for introducing the notion of semicomputable measures.

## 2.4. The Solomonoff-Levin Distribution

The original incentive to develop a theory of algorithmic information content of individual objects was Ray Solomonoff's invention of a *universal a priori* probability that can be used instead of the actual (but unknown) *a priori* probability in applying Bayes's Rule. His original suggestion was to set the *a priori* probability $P(x)$ of a finite binary string $x$ to $\sum 2^{-l(p)}$, the sum taken over all programs $p$ with $U(p) = x$ where $U$ is the reference Turing machine of Theorem 1 for the $C$-complexity. However, using plain Turing machines this is improper, since not only does $\sum_x P(x)$ diverge, but for some $x$ even $P(x)$ itself diverges. To counteract this defect,

Solomonoff in 1960 and 1964 used normalizing terms, but the overall result was unconvincing. Levin [50] succeeded in 1970 to find a proper mathematical expression of the *a priori* probability, of which we present the simpler version over the discrete domain $N$. This was elaborated by Levin in 1973 and 1974 [30, 31], and Levin and Gács in 1974 [16] and independently by Chaitin in 1975 [9].

**Definition.** The *Solomonoff-Levin distribution* (actually a measure) on the positive integers is defined by $P_U(x) = \sum 2^{-l(p)}$, where the sum is taken over all programs $p$ for which the reference prefix-machine $U$ of Theorem 1 outputs $x$. This is a measure because of the following.

**Kraft's Inequality.** If $l_1, l_2, ...$ is a sequence of positive integers such that $\sum_n 2^{-l_n} \leqslant 1$ then there is a prefix-code $c : N \to \{0, 1\}^*$ (i.e., if $n \neq m$ are positive integers, then $c(n)$ is not a prefix of $c(m)$), with $l(c(n)) = l_n$. Conversely, if $c : N \to \{0, 1\}^*$ is a prefix-code, such that the sequence $l_1, l_2, ...$ with $l_n = l(c(n))$, $n = 1, 2, ...$ satisfies the inequality above. See e.g. [12].

Hence, by the Kraft Inequality, for the prefix-code formed by the programs $p$ of $U$ we have $\sum_p 2^{-l(p)} \leqslant 1$. Therefore, the combined probability $\sum_x P_U(x)$, summed over all $x$'s, sums up to less than one, no matter how we choose reference $U$, because for some program $q$ there is no output at all.

Another way to conceive of $P_U(x)$ is as follows. We think of the input to the reference prefix machine $U$ as being provided by indefinite long sequences of fair coin flips. Thus, the probability of generating a program $p$ for $U$ is $P(p) = 2^{-l(p)}$ where $P$ is the standard 'coin-flip' uniform measure. (Namely, presented with any infinitely long sequence starting with $p$, the machine $U$, being a prefix-machine, will read exactly $p$ and no further.) Due to the halting problem, for some $q$ the reference $U$ does not halt. Therefore, the *halting probability* $\Omega$ satisfies

$$\Omega = \sum_x \mathbf{m}(x) < 1.$$

Now we are ready to state the remarkable and powerful fact that Levin's universal semicomputable measure $\mathbf{m}(x)$, the Solomonoff-Levin universal *a priori* probability $P_U(x)$, and the simpler expression $2^{-K(x)}$, all coincide up to an independent fixed multiplicative constant. It is a consequence of currently universally accepted views in mathematical logic (Church's Thesis), that the widest possible effective notion of simplicity of description of an object $x$ is quantified by $K(x)$.

The Solomonoff-Levin distribution can be interpreted as an recursively invariant notion that is the formal representation of "Occam's Razor": the statement that one object is simpler than the other is equivalent to saying that the former object has higher probability than the latter.

**Theorem 3.** *There is a constant c such that for each x, up to additive constant c, we have* $-\log \mathbf{m}(x) \stackrel{+}{=} -\log P_U(x) = K(x)$.

*Proof.* Since $2^{-K(x)}$ represents the contribution to $P_U(x)$ by a shortest program for $x$, $2^{-K(x)} \leqslant P_U(x)$ for all $x$. Since $P_U(x)$ is semicomputable from below by enumerating all programs for $x$, we have by the universality of $\mathbf{m}(x)$ that there is a fixed constant $c$ such that for all $x$ we have $P_U(x) \leqslant c\, \mathbf{m}(x)$.

It remains to show that $\mathbf{m}(x) = O(2^{-K(x)})$. This is equivalent to showing that for some constant $c$ we have $-\log \mathbf{m}(x) \geqslant K(x) + c$. It suffices to exhibit a prefix code such that for some other fixed constant $c'$, for each $x$ there is a code word $p$ such that $l(p) \leqslant -\log \mathbf{m}(x) + c'$, together with a prefix-machine $T$ such that $T(p) = x$. Then, $K_T(x) \leqslant l(p)$ and hence by the Invariance Theorem 1 also $K(x) \leqslant l(p)$ up to a fixed additive constant. First we recall a construction for the Shannon-Fano code.

*Claim.* If $\mu$ is a measure on the integers, $\sum_x \mu(x) \leqslant 1$, then there is a binary prefix-code $r: N \to \{0, 1\}^*$ such that the code words $r(1), r(2),\ldots$ are in lexicographical order, such that $l(r(x)) \leqslant -\log \mu(x) + 2$. This is the Shannon-Fano code.

*Proof.* Let $[0, 1)$ be the half open unit real interval. The half open interval $[0.x, 0.x + 2^{-l(x)})$ corresponding to the set (cylinder) of reals $\Gamma_x = \{0.y : y = xz\}$ ($x$ finite and $y$ and $z$ infinite binary strings) is called a *binary interval*. We cut off disjoint, consecutive, adjacent (not necessarily binary) intervals $I_n$ of length $\mu(n)$ from the left end of $[0, 1)$, $n = 1, 2,\ldots$. Let $i_n$ be the length of the longest binary interval contained in $I_n$. Set $r(n)$ equal to the binary word corresponding to the first such interval. It is easy to see that $I_n$ is covered by at most four binary intervals of length $i_n$, from which the claim follows. $\square$

Since $\mathbf{m}(x)$ is semicomputable from below, there is a partial recursive function $\phi(t, x)$ such that $\phi(t, x) \leqslant \mathbf{m}(x)$ for all $t$, and $\lim_{t \to \infty} \phi(t, x) = \mathbf{m}(x)$. Let $\psi(t, x) = 2^{-k}$, with $k$ is a positive integer, be the greatest partial recursive lower bound of this form on $\phi(t, x)$. We can assume that $\psi$ enumerates its range without repetition. Then,

$$\sum_{x, t} \psi(t, x) = \sum_x \sum_t \psi(t, x) \leqslant \sum_x 2\,\mathbf{m}(x) \leqslant 2.$$

(The series $\sum_t \psi(t, x)$ can only converge to precisely $2\,\mathbf{m}(x)$ in case there is a positive integer $k$ such that $\mathbf{m}(x) = 2^{-k}$.)

Similar to before, we chop off consecutive, adjacent, disjoint half open intervals $I_{t,x}$ of length $\psi(t, x)/2$, in order of computation of $\psi(t, x)$, starting from the left side of $[0, 1)$. This is possible by the last displayed equation. It is easy to see that we can construct a prefix-machine $T$ as follows. If $\Gamma_p$ is the largest binary interval of $I_{t,x}$, then $T(p) = x$.

Otherwise, $T(p)$ is undefined (e.g., $T$ doesn't halt).

By construction of $\psi$, for each $x$ there is a $\psi(t, x) > \mathbf{m}(x)/2$. By the construction in the Claim, each interval $I_{t,x}$ has length $\psi(t, x)/2$. Each $I$-interval contains a binary interval $\Gamma_p$ of length at least one quarter of that of $I$. Therefore, there is a $p$ with $T(p) = x$ such that $2^{-l(p)} \geqslant \mathbf{m}(x)/16$. This implies $K_T(x) \leqslant -\log \mathbf{m}(x) + 4$. The proof of the theorem is finished. $\square$

Theorem 3 demonstrates a particularly important instance of the two conceptually different, but equivalent, definitions of the semicomputable measures. We analyse this equivalence in some detail. Let $P_1, P_2,\ldots$ be the effective enumeration of all semicomputable probability distributions constructed in Theorem 2. Let $T_1, T_2,\ldots$ be the standard enumeration of prefix-machines. For each prefix-machine $T$, define

$$Q_T(x) = \sum_{T(p) = x} 2^{-l(p)}.$$

In other words, $Q_T(x)$ is the probability that $T$ computes output $x$ if its input $p$ is generated by successive tosses of a fair coin. In other words, the inputs $p$ are uniformly distributed with the probability of $p$ occurring equal $2^{-l(p)}$. It is easy to see that each $Q_T$ satisfies

$$\sum_x Q_T(x) \leqslant 1.$$

Equality holds iff $T$ halts for all inputs (proper programs). Let $Q_1, Q_2,\ldots$ (where we do not require equality to hold) be the probability distributions associated with $T_1, T_2,\ldots$.

*Claim.* There are recursive functions $\sigma, \pi$ such that $Q_n = P_{\sigma(n)}$ and $P_n = \Theta(Q_{\pi(n)})$, for $n = 1, 2,\ldots$.

*Proof.* Omitted. $\square$

**Remark.** The Coding Theorem tells us that there is a constant $c > 0$ such that $-\log P_U(x) - K(x) \leqslant c$. We recall from the definition of the Solomonoff-Levin distribution that $-\log P_U(x) = -\log \sum_{U(p) = x} 2^{-l(p)}$, and $K(x) = \min\{l(p): U(p) = x\}$. A priori an outcome $x$ may have high probability because it has many long descriptions. But these relations show that in that case it must have a short description too. In other words, the *a priori* probability of $x$ is governed by the shortest program for $x$.

**Remark.** Let $P$ be any probability distribution (not necessarily computable). The $P$-expected value of $\mathbf{m}(x)/P(x)$ is

$$\sum_x P(x) \frac{\mathbf{m}(x)}{P(x)} < 1.$$

We find by Chebychev's first Inequality[1] that

$$\sum \{P(x): \mathbf{m}(x) \leqslant k\,P(x)\} \geqslant 1 - \frac{1}{k}. \qquad (3)$$

Since $\mathbf{m}(x)$ dominates all semicomputable

---

[1] Recall that *Chebychev's First Inequality* says the follow-

measures multiplicatively, for all $x$ we have

$$P(x) \leqslant c_P \mathbf{m}(x), \qquad (4)$$

for a fixed positive constant $c_P$ independent of $x$ (but depending on the index of $P$ in the effective enumeration $\mu_1, \mu_2, ...$ of semicomputable measures).

Inequalities (3) and (4) have the following consequences:[2]

(i) If $x$ is a random sample from a simple computable distribution $P(x)$, then $\mathbf{m}(x)$ is a good estimate of $P(x)$.

(ii) If we know or believe that $x$ is random with respect to $P$, and we know $P(x)$, then we can use $P(x)$ as an estimate of $\mathbf{m}(x)$.

**Example 2, Continued.** We continue Solomonoff's approach to inductive inference. The problem was that the proper a priori probabilities $\mu$ in formula (2) are not known. Define

$$\mathbf{M}(x) = \sum \{\mathbf{m}(xy) : y \in \{0, 1\}^*\},$$

i.e., $\mathbf{M}(x)$ is the *a priori probability* that the output of the reference prefix-machine $U$ starts with $x$. (This is slightly different from the more sophisticated approach of Levin in, e.g., [50], but it suffices for our purposes.) Now instead of using formula (2), we estimate the conditional probability $P(xy \mid x)$ the next segment after $x$ is $y$ by the expression

$$\frac{\mathbf{M}(xy)}{\mathbf{M}(x)}. \qquad (5)$$

Now let $\mu$ in Formula (2) be an arbitrary computable measure. This case includes all computable sequences, as well as many Bernouilli sequences. If the length of $y$ is fixed, and the length of $x$ grows to infinity, then we have

$$\frac{\mathbf{M}(xy)/\mathbf{M}(x)}{\mu(xy)/\mu(x)} \rightarrow 1,$$

with $\mu$-probability one. In other words, the conditional a priori probability is almost always

asymptotically equal to the conditional probability. One can show that convergence is very fast, cf. [45] which shows that if we use (5) instead of the real value (2), then our inference is almost as good. One can also show that:

$$-\log \mathbf{M}(x) = K(x) + O(\log K(x)). \qquad (6)$$

We now come to the punch line: Bayes's rule using the universal prior distribution yields Occam's razor principle. Namely, if several programs could generate $S0$ then the shortest one is used (for the prior probability), and further if $S0$ has a shorter program than $S1$ then $S0$ is preferred (*i.e.* predict 0 with higher probability than predicting 1 after seeing $S$). Bayes's rule via the universal prior distribution also gives Epicurus's multiple choice principle in case $S0$ and $S1$ have roughly equal length shortest programs which 'explain' $S0$ and $S1$, respectively.

The universal prior distribution $\mathbf{m}(x)$ is not computable because the Kolmogorov complexity is not computable. However, we can compute approximations to $K(x)$ and $\mathbf{m}(x)$. It turns out that using Solomonoff's inference principles with such computable approximations yields many other known inference models or principles. In the next few sections, we derive or establish connections with various well-known machine learning models and inductive inference paradigms or principles. Thus we provide an alternative view of these models and principles from the lofty perspective of Kolmogorov complexity.

### 3. Gold's Inductive Inference Paradigm

There are many different ways of formulating concrete inductive inference problems in the real world. We will try to simplify matters as much as possible short of losing significance.

(i) **The class of rules** we consider can be various classes of languages or functions, where we restrict ourselves to classes of recursive sets, context-free languages, regular sets and sets of finite automata, and sets of Boolean formulae. We treat a language $L$ as a *function* $f$ using its characteristic function, *i.e.*, $f(x) = \chi_L(x) = 1$ if $x \in L$, and 0 otherwise.

(ii) **The hypothesis space** or **rule space** denoted by $\mathbf{R}$ specifies syntactically how each rule in (i) should be represented. We fix a standard enumeration of the representations for a class of rules, $\mathbf{R} = \{R_1, R_2, ...\}$, and assume that each rule $f$ has at least one description in the corresponding hypothesis space. For example, the hypothesis space can be a set of standard Turing machine codes of halting Turing machines, a set of standard encodings of context-free grammars, or a set of standard encodings of Finite Automata. In any case, it is assumed that the hypothesis space is effectively enumerable (so it cannot be the set of all halting Turing machine codes). For convenience, this enumeration of hypotheses $R_1, R_2, ...$ consists of codes for

---

ing. Let $P$ be any probability distribution, $f$ any nonnegative function with expected value $E_P(f) = \sum_x P(x)f(x) < \infty$. For $\lambda \geqslant 0$ we have $\sum \{P(x) : f(x) > \lambda\} < \lambda^{-1} E_P(f)$. Here we use it with $k E_P(f)$ substituted for $\lambda$.

[2] We shortly remark, without further explanation, that in both cases the degree of approximation depends on the index of $P$, and the randomness of $x$ with respect to $P$, as measured by the randomness deficiency $\delta_0(x \mid P) = \log(\mathbf{m}(x)/P(x))$. If $\delta_0(x \mid P) = O(1)$ then $x$ is random, otherwise $x$ is not random. For example, for the Uniform Distribution $\delta_0(x \mid P) = n - K(x \mid n) + O(1)$, where $n = l(x)$. Such a (universal Martin-Löf) test is needed, since otherwise we cannot distinguish, for instance, between randomness and nonrandomness of samples from the uniform distribution. (Clearly, the word C o n s t a n t i n o p l e is not a random 14-letter word. The probability of seeing it somewhere written is decidedly greater than $128^{-14}$, say, for a randomly selected fourteen letter ASCII word.)

algorithms to compute recursive functions $f_1, f_2,...$ (languages are represented by their characteristic functions).

(iii) **The presentation of examples** is vital to the inference process. We choose the simplest, and yet most general, form of data presentation. For a function $f$ to be inferred, there is a fixed infinite sequence of *examples* $(s_1, f(s_1)), (s_2, f(s_2)),...$. When $f = \chi_L$, we have $\chi_L(s) = 1$ iff $s \in L$ in which case we say $s$ is a positive example of $L$. We say $s$ is a negative example of $L$ iff $s \notin L$. For instance, $0,1,1,2,3,5,8,...$ (the fibonacci sequence) can be represented by a sequence of pairs $(0,0), (1,1), (2,1), (3,2), (4,3), (5,5), (6,8),...$. Then, $(2,1)$ is a positive instance of the corresponding characteristic function $f$, $f(2,1) = 1$, and $(3,1)$ is a negative instance, $f(3,1) = 0$.

A rule (or function) $f$ is said to be *consistent* with the initial segment of examples

$$S = (s_1, a_1), \ldots, (s_n, a_n), \qquad (7)$$

if $f(s_i) = a_i$, $i = 1,..,n$. We require that *all* strings will eventually appear as first component in a pair in $S$. The last assumption is strong, but is essential to the Gold paradigm.

The prediction problem is now: given $Ss$, where $S$ is the initial segment (7), predict $f(s)$. In other words, our task is to learn $f$. For example, in case of inference of language $L$, given $Ss$, we must predict $\chi_L(s)$ (whether or not $s \in L$). Let us formulate a procedure to infer a rule according to Solomonoff.

By (ii), there is an effective enumeration $f_1, f_2,...$ of partial recursive functions corresponding to the enumeration of hypotheses. The a priori probability of $f_k$ is $\mathbf{m}(f_k) = \mathbf{m}(k)$. (Actually, $\mathbf{m}(f_k) = c\,\mathbf{m}(k)$, for some constant $c$ depending on the effective enumeration involved, but not depending on $n$. To assume that $c = 1$ makes no difference in the following discusssion.) We are given an infinite sequence of examples representing the rule or function $f$ to be learned. Let $f(s) = a$. (E.g., if we want to infer a language $L$ then $f = \chi_L$, and $a = 1$ for $s \in L$, and $a = 0$ for $s \notin L$.) According to Bayes's rule (1), for $k = 1, 2,...$, the inferred probability of $f_k$ after the sequence of examples (7) is given by:

$$P(f_k = f \mid f(s_i) = a_i, i = 1,..,n) \qquad (8)$$

$$= \frac{P(f(s_i) = a_i, i = 1,..,n \mid f_k = f)\,\mathbf{m}(k)}{\sum \{\mathbf{m}(j): f_j(s_i) = a_i, i = 1,..,n\}}.$$

In the numerator of the right-hand term, the first factor is zero or one depending on whether $f_k$ is consistent with $S$, and the second factor is the a priori probability of $f_k$. The denominator is a normalizing term giving the combined a priori probability of all rules consistent with $S$. With increasing $n$, the denominator term is monotonically nonincreasing. Since all examples eventually appear in $S$, the

denominator converges to a limit, say $d \leqslant 1$. For each $k$, the inferred probability $f_k$ is monotonically nondecreasing with increasing $n$, until $f_k$ is inconsistent with a new example, in which case it falls to zero and stays there henceforth. In the limit, only the $f_k$'s that are consistent with the sequence of presented examples have positive inferred probability $\mathbf{m}(k)/d$. By Theorem 3, since $\mathbf{m}(k) \leqslant c\,2^{-K(k)}$, with equality for most $k$, for some positive constant $c$, the highest inferred probability is carried by the rule $f_k$ with least Kolmogorov complexity among the remaining ones. Similar statements hold after each initial segment of $n$ examples, $n = 1,2,...$.

Reasoning inductively, we transform the *a priori* probability according to Formula (8), inferring a new *posterior* probability by the evidence of each initial segment of examples. At each step, we can select the rule with the highest inferred probability, and in the limit we have selected the proper rule. At each step we predict the rule with the highest inferred probability. Reformulating, if we want to infer a language $L$ using this procedure, then:

(a) The Bayesian *a posteriori* probability for the correct answer converges to $c\,2^{-l(p)}/d$, where $p$ is the shortest program which the reference machine uses to simulate $M_0$, where $M_0$ is the smallest TM that accepts $L$. This correct answer will have the highest probability in the limit. That is, inferred probability distribution over the underlying machines converges to a highest probability for $M_0$ *in the limit*. In other words, after $n$ steps for some $n$, all the machines smaller than $M_0$ violate some data pair in $S$, and $M_0$ is the choice forever after step $n$.

(b) It is interesting to notice that the *a posteriori* probability decreases monotonically until it converges to $c\,2^{-l(p)}/d$ for $p$ the program with which $U$ simulates $M_0$. Smaller machines are chosen first and then canceled because they violate some data.

(c) If we want to infer $f(s)$, rather than $f$, given the sequence of examples $S$, then using formulas (2) and (5), the inferred probability that $f(s) = a$ is (denoting a string $s_1 s_2 \cdots s_n$ as $s_{1:n}$):

$$P(f(s) = a \mid f(s_i) = a_i, i = 1,..,n) \qquad (9)$$

$$= \frac{\sum \{\mathbf{m}(j): f_j(s_i) = a_i, i = 1,..,n, f_j(s) = a\}}{\sum \{\mathbf{m}(j): f_j(s_i) = a_i, i = 1,..,n\}}$$

$$= \frac{M(S(s,a) \mid s_{1:n}s)}{M(S \mid s_{1:n})}.$$

(Here we use that the probability of an initial segment of examples $S$, conditional to the sequence of queries $s_{1:n}$, denoted by $M(S \mid s_{1:n})$, is equal (as always up to a constant) to the summed probability of all functions $f$ that are consistent with $S$.)

The well-known Gold paradigm of inductive inference [14, 15], can be viewed simply as a computable approximation to Equation (8). The

fundamental idea of the Gold paradigm is the idea called *identification in the limit* and a universal method of implementing the identification in the limit is called "identification by enumeration". These are contained in facts (a) and (b), as a computable analogue of Solomonoff's approach. We now investigate the correspondence between these two basic ideas in some detail.

**Identification in the Limit** views inductive inference as an infinite process. Formally, let $M$ be an inductive inference method in order to derive some unknown rule $R$. If $M$ receives a larger and larger set of examples (bigger and bigger initial segment $S$), a larger and larger sequence of $M$'s conjectures is generated, say, $f_1, f_2, f_3, \cdots$. If there is some integer $m$ such that $f_m$ is a correct description of $R$ and for all $n > m$

$$f_m = f_n,$$

then $M$ *identified* $R$ in the limit. Two facts deserve mentioning: $M$ cannot determine whether it has converged and therefore stop with a correct hypothesis. $M$ may be viewed as learning more and more information about the unknown rule $R$ and monotonically increasing its approximation to $R$ until the correct identification. Gold gave the best explanation to his definition:

I wish to construct a precise model for the intuitive notion "able to speak a language" in order to be able to investigate theoretically how it can be achieved artificially. Since we cannot write down the rules of English which we require one to know before we say he can "speak English", an artificial intelligence which is designed to speak English will have to learn its rules from implicit information. That is, its information will consist of examples of the use of English and/or of an informant who can state whether a given usage satisfies certain rules of English, but cannot state these rules explicitly.

$\cdots$ A person does not know when he is speaking a language correctly; there is always the possibility that he will find that his grammar contains an error. But we can guarantee that a child will eventually learn a natural language, even if it will not know when it is correct.

**Identification by enumeration** is a method to implement identification in the limit. It refers to the following guessing rule: Enumerate the class of rules in rule space. At step $t$, guess the unknown rule to be the first rule of the enumeration which agrees with data received so far. Formally speaking, in our setting, if we have received an initial segment $S$, then, given $s$, predict as the next example $(s, f(s))$ for $f$ is the first rule that is consistent with $S$. Now if this can be done *effectively*, identification in the limit will be achieved. In fact, let $G$ and $G'$ be two guessing methods. $G$ will be said to be *uniformly faster* than $G'$ if the following two conditions hold: (1) Given any $R$ from the rule space, $G$ will identify $R$ correctly at least as soon as $G'$, expressed in the number of examples needed, for all sequences of examples; and (2) for some $R$, some sequence of examples, $G$ will identify $R$ sooner than $G'$. It is easy to prove that the identification-by-enumeration method will identify a hypothesis in the limit if this

hypothesis can be identified in the limit at all. Further if $G_0$ is an identification-by-enumeration guessing rule, then there is no guessing rule uniformly faster than $G_0$. Indeed,

**Theorem 4.** *Identification-by-enumeration is a computable approximation to inductive inference (Solomonoff's inference) associated with Formula (8). Neither method is uniformly faster than the other.*

*Proof.* An effective enumeration for the identification-by-enumeration method, can be viewed as a computable approximation to Solomonoff's procedure according to formula (8) as follows. Let the effective enumeration of the rule space be: $R_1, R_2, R_3 \cdots$. Convert this to an effective prefix-free description of each rule $R_i$ in the rule space. For instance, if $x = x_1,...,x_n$ is a binary string, then $\bar{x} = x_1 0 x_2 0...0 x_n 1$ is a *prefix-code* for the $x$'s. Similarly, $x' = l(x)x$ is a prefix-code. Note that $l(x') = l(x) + 2 \log l(x)$. We encode each rule $R_i$ (a binary string) as $p\, i'$, where $p$ is a (self-delimiting) program that enumerates the rule space. The resulting code for the $R_i$'s is an effective prefix-code. Denoting the length of the description of $R_i$ by $|R_i|$, we have:

(a)  if $i < j$, then $|R_i| \leq |R_j|$; and

(b)  $\sum_i 2^{-|R_i|} \leq 1$ (by Kraft's inequality).

Assign *a priori* probability $P(R_i) = 2^{-|R_i|}$ to rule $R_i$, $i = 1, 2, ....$ (This is possible because of (b).) Using Formula (8) with $P(R_i)$ instead of $m(i)$ yields a computable approximation to Solomonoff's inductive inference procedure. Formula (8) chooses the shortest encoded consistent rule which coincides with the first consistent rule in the effective enumeration. This shows that identification by enumeration can be formulated as an computable approximation to Solomonoff's procedure. It remains to show that neither method is uniformly faster than the other.

Let $G_1, G_2,...$ be an effective enumeration of the hypotheses space by a Gold procedure, and let $H_1, H_2,...$ be the (noneffective) enumeration of the hypotheses space by decreasing *a priori* probability according to Solomonoff. In other words, $K(H_1) \leq K(H_2) \leq \cdots$. In both cases we deal with identification-by-enumeration, so it is known that there is no guessing rule uniformly faster than either of them. $\square$

**Remark.** What about non-uniform speed comparison? In case the particular rule $f$ to be inferred is sufficiently *simple* (has low Kolmogorov complexity) then Solomonoff's procedure can be much faster than Gold's enumeration. Let $f$ be the function we want to infer, and let $f = f_m$, with $m$ minimal, in Gold's enumeration order. Let also $f = f_n$, for $n$ with $K(n)$ minimal. To infer the correct $f$, in Gold's approach we must eliminate all $f_k$ with $k < m$. But in Solomonoff's approach, we only need to eliminate all $f_k$ with $K(k) < K(n)$. Now necessarily there are

many $f$'s that are 'simple' in the sense that $K(n) \ll l(m)$, for which e.g. Solomonoff's procedure works much (sometimes noncomputably) faster than Gold's method.

The following theorem sets limits on the number of examples needed to infer a particular function $f$.

**Theorem 5.** *Let $f_1, f_2, \ldots$ be an effective enumeration of the rule space. Suppose we want to infer $f = f_i$, with $i$ minimal, from a set of $n$ examples $S$ as in (7). Let $c$ be an appropriate large enough constant.*
*(a) If $K(i) > K(f(s_1)...f(s_n) \mid s_1...s_n) - c$, then it is impossible by any effective deterministic procedure to infer $f$ correctly.*
*(b) If we can infer $f$ correctly by computable approximation to Solomonoff's method (8) using only $S$, and $c$ extra bits of information, then $K(i \mid S) \leq c$.*
*(c) If $K(i \mid S) \leq c$ then we can compute $f_i$ from $S$ and $c$ bits extra information.*

*Proof.* (a) Otherwise we would be able to compute $i$ from a program of length significantly shorter than $K(i)$: contradiction. Items (b) and (c) are obvious. □

There is an enormous amount of research in the area under the title of Gold paradigm. We refer the readers to the articles [3, 7] and the book [39]. We present three examples of reasoning by means of Gold's paradigm in order to give a flavor of this research direction.

**Example.** [Gold, 1967] We can learn a function in the set of primitive recursive functions.

*Proof.* Effectively enumerate the set of all primitive recursive functions by $\psi_1, \psi_2, \cdots$. On any initial input segment $(x_1, y_1) \cdots (x_k, y_k)$, our inference machine just prints the least $i$ such that $\psi_i$ is consistent with the input, i.e., $\psi_i(x_k) = y_k$ for $k = 1, \cdots, n$.

**Example.** [Gold, 1967] We cannot learn in general a function in the set of all total recursive functions.

*Proof.* By diagonalization. Suppose $M$ can identify all recursive functions. But then one can define a recursive function $f$ so that the guesses of $M$ will be wrong on $f$ infinitely often. We construct $f$ by simply simulating $M$. Let $f(0) = 0$, Suppose the value of $f(0), f(1), \cdots, f(n-1)$ have been constructed. On input $n$, simulate $M$ on initial input $f(0), f(1), \cdots, f(n-1)$. Then define $f(n)$ equal 1 plus the guess of $M$ (modulo 1). So $M$ never guesses $f$ correctly. □

**Example.** One of the first studied problem was extrapolating a sequence. A machine $M$ *extrapolates* a sequence $f(1), f(2), \cdots$ as follows. It makes an initial guess $f'(0)$. Then it inputs the real $f(0)$. At step $i$, based on previous inputs $f(1), f(2), \cdots, f(i-1)$, it guesses $f'(i)$. If there is a $i_0$ such that for all $i > i_0$ $f'(i) = f(i)$, then we say $M$

*extrapolates $f$.* Bringing everything in our setting, the initial segment before step $i$ is a sequence of pairs $(1, f(1))(2, f(2))...(i-1, f(i-1))$, and $M$ extrapolates with the pair $(i, f'(i))$. It is not surprising that the class of functions computable by a Turing machine running in time $t(n)$, for any computable function $t$, can be extrapolated (by identification by enumeration).

## 4. The Valiant Model of Deductive Learning

In practice, as we have mentioned, the Solomonoff paradigm, or the Gold paradigm, do not offer feasible solutions in two senses. Firstly, both require all strings to appear eventually and the learning program never knows when to stop. Secondly, in many cases even within a single learning step, just finding the first rule that is consistent with the data in the 'identification by enumeration' method may require exponential time in terms of the size of data. The goal of this section is to do machine learning *fast*.

Obviously, this task is impossible in general if we are asked to *precisely* infer a law of nature, even if we assume that this law is a computable function. In fact, the infinite behavior for the whole process and the exponential behavior for each inference step are inherent. However, in practice, learning usually appears to be fast. A child recognizes an apple after seeing several apples. A kitty learns to avoid skunks from only a few bad hunting adventures. After several bad experiences, a newcomer to the city of Boston finds that his car will get clamped after receiving five unpaid parking tickets.

These examples all bear some common characteristics: (1) The learning process does not depend on the number of apples in the world, the number of skunks in the jungle, or the number of tickets in the Boston police department. (2) The learning process is fast. (3) The results usually converge with remarkable precision. As Valiant noticed,

Human learning often shows remarkable properties of convergence. In large populations there is a high degree of agreement on the meaning of thousands of words as they relate to everyday situations. This suggests that highly reliable program acquisition may be feasible even in the absence of explicit programming. Hence it is reasonable to insist that our study of learning be disciplined by insistence on fast convergence.

Use the same setting as in the last section. Again we consider only the case where the function $f$ to be learned is deterministic. Let there be an arbitrary, fixed, distribution $D$, according to which the examples $(s, f(s))$ are drawn. That is, $D(s)$ is the probability of drawing example $(s, f(s))$. Assume that we want to learn a rule $R$ from an effectively enumerated rule space $\mathbf{R} = \{R_1, R_2, ...\}$ with corresponding functions $\{f_1, f_2, ...\}$. Assume that, given some sequence $S$ of $n$ examples, as in Formula (7), we can compute $R(S)$, defined as a minimum length rule $R$ that is consistent with $S$. Using formula (8) we can show:

**Theorem 6.** *Let $f$ be the function to be learned,*

*and let* $l(R(S)) = m$. *After drawing a sequence of $n$ examples, say $S$ as in Formula (7), from distribution $D$, and $n$ is a polynomial $n(m, 1/\epsilon)$, then with probability $\geq 1 - \epsilon$ we have $\sum D(s) < \epsilon$, the sum taken over all $s$ such that $f_k(s) \neq f(s)$ and $k$ is the index that maximizes $P(f_k = f \mid S)$. (The proof of this theorem is similar to that of Occam's Razor theorem below.)*

However, computing the minimum representation is hard in many cases. For example, given positive and negative examples, computing the smallest consistent DFA is NP-complete [1, 15, 40]. Computing the minimum $k$-DNF (DNF formula with each term containing at most $k$ literals) from a set of examples is also NP-complete. Do we really need to compute the minimum representation? What about just approximation? Adding this feature to the approach above brings us to the **Valiant deductive learning model**. (A similar model was also studied by Vapnik and Chervonenkis [48] ).

In 1983, Valiant [46] introduced a learning model which requires that the computational process by which the machine deduces the desired programs requires a feasible (i.e., polynomial) number of steps. One version of the Valiant model is given here: Consider an effectively enumerated class of rules **R** with associated functions, as before. In this context we call a rule a *concept*. For each concept $f$ (associated with $R$) in **R**, if $v$ is a positive example to $f$ we write $f(v) = 1$, else $v$ is a negative example to $f$ and denoted by $f(v) = 0$. The learning algorithm has available two buttons labeled POS and NEG. If POS (NEG) is pushed, a positive (negative) example is generated according to some fixed but unknown probability distribution $D^+$ ($D^-$) according to nature. We assume nothing about the distributions $D^+$ and $D^-$ except that for each concept $f \in \mathbf{R}$ we have $\sum_{f(v)=1} D^+(v) = 1$ and $\sum_{f(v)=0} D^-(v) = 1$.

**Definition V1.** **R** is (Valiant) learnable from examples iff there exists a polynomial $n$ and a (possibly randomized) learning algorithm $A$ such that, for each $R \in \mathbf{R}$ and $\epsilon > 0$, algorithm $A$ halts in $n(l(R), 1/\epsilon)$ time and number of examples, and outputs a rule $R' \in \mathbf{R}$ that with probability at least $1 - \epsilon$ has the following properties: the function $g$ associated with $R'$ satisfies $\sum_{g(v)=0} D^+(v) < \epsilon$ and $\sum_{g(v)=1} D^-(v) < \epsilon$.

Examples for class **R** are: Set of finite automata; Set of DNF Boolean formulae; Set of $k$-DNF formulae;

By Theorem 6, a class of rules **R** is (Valiant) learnable if we can compute fast the shortest rule in rule space that is consistent with the given examples. As said before, this may be infeasible (in the sense of NP-completeness). However, it turns out that it suffices to identify a rule of which the length is sufficiently near to that of the shortest one. In other words, instead of the minimum consistent hypothesis we can use a consistent hypothesis that is sufficiently

short to let the analogue of Theorem 6 hold. Informally, we can learn in the sense of Valiant if we can find a consistent rule that is sufficiently short in terms of the minimal length consistent rule and the number of examples. The precise form of this statement is due to Blumer, Ehrenfeucht, Haussler, Warmuth, [6], who called it 'Occam's Razor Theorem'.

**Theorem 7.** *Let **R** be an effectively enumerated class of rules, and let $S$ as in Formula (7) be a sequence of $n$ independent examples of a rule (rather the associated function) $f$ to be inferred. Let $R(S)$ be a minimum length rule consistent with $S$, and denote $s = l(R(S))$. The class **R** is (Valiant) learnable, if there is an algorithm $A$ which produces a rule $R'$ in **R** consistent with $S$ and with $l(R') \leq s^c n^\alpha$, where $c \geq 1$ and $0 \leq \alpha < 1$, and the running time of $A$ is polynomial in $n$. The sample size $n = n(s, 1/\epsilon)$ is required to be polynomial in the two arguments.*

*Proof.* Fix an error tolerance $\epsilon$. We first prove that the probability that a consistent rule has error $\leq \epsilon$ is large. Precisely,

*Claim.* Given any rule $f$ in a class of $r$ hypotheses. Let $g$ be any hypothesis in this class, for which the probability that it disagrees with $f$ on a single example $s$ exceeds $\epsilon$ (notation: $Pr(g(s) \neq f(s) \mid g, f) \geq \epsilon$). The probability that $g$ is consistent with a sample of $f$ of $n$ examples is less than $(1 - \epsilon)^n r$.

*Proof of Claim.* Let $E_g$ be the event that hypothesis $g$ agrees with all $n$ examples of $f$. Also let $Pr(g(s) \neq f(s) \mid g, f) \geq \epsilon$, i.e. $g$ and $f$ disagree on one example with probability at least $\epsilon$. By independence of the examples:

$$Pr(E_g) \leq (1 - \epsilon)^n,$$

and

$$P\left( \bigcup_{Pr(g(s) \neq f(s) \mid g, f) > \epsilon} E_g \right) \leq (1 - \epsilon)^n r. \quad \Box$$

The rule $R'$ under consideration is given by a binary string of length at most $s^c n^\alpha$. The number $r$ of such rules satisfies:

$$\log r \leq s^c n^\alpha \leq -\frac{n}{2} \log(1 - \epsilon).$$

By the claim, the probability of producing a hypothesis with error larger than $\epsilon$ is less than

$$(1 - \epsilon)^n r \leq (1 - \epsilon)^{\frac{n}{2}} \leq \epsilon,$$

for large $n$ polynomial in $s$ and $1/\epsilon$. $\quad \Box$

Put in words, this theorem shows that given a set of positive and negative data, any consistent concept of size 'reasonably' shorter than the size of data is an 'approximately' correct concept. If one can find a shorter representation of data, then one learns. *The shorter the conjecture is, the more efficiently it explains the current data in the learning examples.* This provide a rigorous support of our use of Occam's razor principle.

This theorem turns out to be a very useful tool for identifying polynomial time learnable classes in the Valiant model. There are many problems for which it is actually NP-complete to derive their minimum representation, but they can be approximated so that above theorem can be applied. An important special case of this theorem occurs when $\alpha = 0$, in which case the theorem says that if one can compress the data to length independent of $n$, then one learns a rule to within polynomial length.

**Example** [Valiant, 1983]. It is known that it is NP-complete to find the minimum $k$-DNF formula consistent with the given positive and negative data. We now show that the set of $k$-DNF of $r$ variables $v_1, \cdots, v_n$ is (*polynomial time*) learnable. The learning algorithm for a $k$-DNF formula (associated function) $f$ is as follows:

$$g := \sum_{m \in M_k} m \text{ , } M_k \text{ is the set of monomials of}$$
size $k$;

**repeat**

$(s, f(s)) := $ a negative example for $f$ (i.e., $f(s) = 0$);

for each term $c_i$ in $g$ delete $c_i$ if $s$ implies $c_i$;

**until** satisfied $((2n)^{k+2}$ loops suffice)

In the beginning, $g$ contains all conjunctive terms of $k$ literals. There are at most $(2n)^{k+1}$ of them, where $k$ is considered to be a constant here. We claim that by the end of the computation, with high probability $g$ will satisfy Definition VI. Fix any $\epsilon > 0$. Initially, all positive examples satisfy $g$, *i.e.* $\sum_{g(v)=1} D^+(v) = 1$, but also all negative examples satisfy $g$. After each step, some clauses in $g$ which imply some negative examples are necessarily deleted. It remains to show that $\sum_{g(v)=0} D^+(v) = 0$, that is, no error on the positive side is introduced. The size of $g$ is at most $(2n)^{k+1}$. By Occam's Razor Theorem, after a polynomial number of, say $(2n)^{k+2}$, steps and examples, with probability greater than $1 - \epsilon$, we have $Pr(g(s) \neq f(s) \mid g, f) \leqslant \epsilon$.
□

For related results, see [6, 18, 23, 24, 35, 38, 41, 43, 46, 47].

## 5. Rissanen's Minimum Description Length Principle

Solomonoff's ideas about inductive reasoning have explicitly served as guiding principle in Rissanen's development of Minimum description length (MDL) principle. Let us derive Rissanen's MDL principle from Solomonoff's induction principle. For simplicity, we deal with only non-adaptive models. A non-adaptive model is a model $P(D \mid \theta)$ where the parameter vector $\theta = \theta(D)$ is estimated from $n$ observed data points denoted by $D$.

Scientists formulate their theories in two steps: firstly a scientist must, based on scientific observations or given data, formulate alternative hypotheses, and secondly he selects one definite hypothesis. This is the subject of inference in statistics. Statisticians have developed many different principles to do this, like Occam's razor principle, the Maximum Likelihood principle, various ways of using Bayesian formula with different prior distributions. No single principle turned out to be satisfactory in all situations. Philosophically speaking, Solomonoff's approach presents an ideal way of solving induction problems using Bayes's rule with the universal prior distribution. However, due to the non-computability of the universal prior function, such a theory cannot be directly used in practice. Some approximation is needed in the real world applications. Further, from theory to inductive inference and statistical practice, there is still a big distance, for example, concrete formulae are needed.

Gold's principle was a particularly simple approximation to Solomonoff's induction - the sophisticated notion of probability distribution is replaced by linear enumeration. In Valiant's learning theory we have added computability requirements and the notion of approximate solution to this general theory. Now we will closely follow Solomonoff's ideas, but substitute a 'good' computable approximation to $m(x)$. This results in Rissanen's **Minimum Description Length Principle** (MDL principle). He not only gives the principle, more importantly he also gives the detailed formulas on how to use this principle. This makes it possible to use the MDL principle in real problems. The principle can be intuitively stated as follows:

**Minimum Description Length Principle.** *The best theory to explain a set of data is the one which minimizes the sum of*

(1) *the length, in bits, of the description of the theory;*

(2) *the length, in bits, of data when encoded with the help of the theory.*

We now develop this MDL principle from Bayes's rule, Formula (1), using the Universal distribution $m(x)$. Recall Bayes's formula:

$$P(H \mid D) = \frac{P(D \mid H) P(H)}{P(D)}.$$

Here $H$ is an hypothesis, here a probability distribution, which we assume to be computable or anyway semicomputable, nd $D$ the observed data. We must choose the hypothesis $H$ such that $P(H \mid D)$ is maximized. First we take the negative logarithm on both sides of the formula:

$$-\log P(H \mid D) = -\log P(D \mid H) - \log P(H) + \log P(D).$$

Since $P(D)$ can be considered as a normalizing factor, we ignore it in the following discussion. Since we are only concerned with maximizing the term $P(H \mid D)$ or, equivalently, *minimizing* the term

$-\log P(H \mid D)$, this is equivalent to minimizing

$$- \log P(D \mid H) - \log P(H).$$

Now to get the minimum description length principle, we only need to explain the above two terms in the sum properly. The term $- \log P(H)$ is straightforward. Following Solomonoff, we set $P(H) = \mathbf{m}(H) = 2^{-K(H) \pm O(1)}$ where $K(H)$ is the self-delimiting Kolmogorov complexity of $H$. That is, $- \log P(H)$ is the *length* of a minimum *prefix code*, or program, of the hypothesis $H$.

We now consider the term $- \log P(D \mid H)$. Since $P$ is computable, using (the conditional version of) Theorem 3 and its corollaries, Formulas 3 and 4, we know that the universal semimeasure $\mathbf{m}(x)$ has the following properties.

(a) There is a constant $c$, such that $\mathbf{m}(D \mid H) \geqslant cP(D \mid H)$.

(b) The $P$-probability that $\mathbf{m}(D \mid H) \leqslant kP(D \mid H)$ is at least $1 - 1/k$.

By a conditional version of Theorem 3, $\mathbf{m}(D \mid H) = 2^{-K(D \mid H) \pm O(1)}$. Hence again $2^{-K(D \mid H)}$ is a reasonable approximation of $P(D \mid H)$, and minimizing $-\log P(D \mid H)$ can be considered as minimizing $K(D \mid H)$, *i.e.*, finding an $H$ such that the description length, or the Kolmogorov complexity, of $D$ given $H$ is minimized. The term $-\log P(D \mid H)$ can also be thought as the *ideal code length* for describing data $D$, given hypothesis $H$. Such prefix code length can be achieved by the Shannon-Fano code. The term $-\log P(D \mid H)$, also known as the *self-information*, in information theory, and the negative log likelihood in statistics, can now be regarded as the number of bits it takes to redescribe or encode $D$ with an ideal code relative to $H$.

In the original Solomonoff approach, $H$ in general is a Turing machine. In practice we must avoid such an overly general approach in order to keep things computable. In different applications, the hypothesis $H$ can mean many different things. For example, if we infer decision trees, $H$ is a decision tree; In case of learning finite automata, $H$ can be a finite automaton; In case we are interested in learning Boolean formulae, then $H$ may be a Boolean formula; If we are fitting a polynomial to a set of points, then $H$ may be a polynomial of some degree; In general statistical applications, one assumes that $H$ is some model $H(\theta)$ with a set of parameters $\theta = \{\theta_1, \cdots, \theta_k\}$, where the number $k$ may vary and influence the (descriptional) complexity of $H(\theta)$. In such case, from

$$- \log P(D \mid \theta) - \log P(\theta),$$

using Taylor expansion at the point of optimal $\hat{\theta}$ (for best maximum likelihood estimator), and taking only dominant terms, Rissanen has derived a formula for the minimum description length as

$$\min_{\theta, k} \{ -\log P(D \mid \theta) + \frac{1}{2} k \log n \},$$

where $k$ is the number of parameters in $\theta = \{\theta_1, \cdots, \theta_k\}$, and $n$ is number of observations (or data points) contained in $D$. At the optimal $k$ and $\theta$, the term $1/2\, k \log n$ is called the optimum model cost.

Since $K(H)$ is not computable and hard to approximate, Rissanen suggested the following approach. First convert (or encode) $H$ to a positive integer in $N = \{1, 2, \cdots\}$. Then we try to assign prior distribution to each integer in $N$. Jeffreys [22] suggested to assign probability $1/n$ to integer $n$. But this results an improper distribution since the series $\sum 1/n$ diverges. We modify Jeffreys distribution. It is possible, by iterating the idea of encoding $n$ (viewed as the corresponding $n$th binary string) as $n' = l(n)n$, to obtain a prefix-code such that $L(n)$ denotes the length of the code for $n$, with $L(n)$ defined by

$$l^*(n) = \log n + \log\log n + \cdots,$$

all positive terms, and

$$L(n) = l^*(n) + \log c,$$

where $c = 2.865064 \cdots$. Viz, it can be proved [42] that:

$$\sum_{n=1}^{\infty} 2^{-l^*(n)} = c.$$

Therefore, the existence of a prefix-code as claimed follows from Kraft's Inequality.

Assign prior probability $P(n) = 2^{-L(n)}$ to each integer $n$. We obtain the following desired properties: (a) $\sum_{n=1}^{\infty} 2^{-L(n)} = 1$; and (b) integers $n$ are coded by a prefix code. Hence, descriptions of two integers, $n_1$ and $n_2$, can be just concatenated to produce the code for the pair $(n_1, n_2)$, and so on. The decoding process is trivial.

Using the MDL principle, Wax and Rissanen (according to Wax) and Quinlan and Rivest [41] have developed procedures to infer decision trees. Other work by Wax [49] and by Gao and Li [13] applied MDL principle to recognition problems.

**Example.** We sketch an initial experiment we [13] have performed in on-line handwritten character learning using the MDL principle. Inputting Chinese characters into computers is a difficult task. There are at least 5,000 characters in daily use, all of different shapes. Many methods have been invented for key-board input. Some have been successful in the limited sense that they can be used by highly trained typists only. Handwriting input is an alternative choice. Many such systems have been built with various recognition rates.

We [13] have implemented such a system that learns handwritten characters from examples under the guidance of the MDL principle. We now sketch a simple experiment we have performed. An input character is drawn on a digitizer board with 200/inch resolution in both horizontal and vertical

directions. The system learns a character from examples. The basic algorithm involves low level preprocessing, scaling, forming a prototype of a character (for learning), elastic matching (for recognizing), and so on. At the stage of forming a prototype of a character, we have to decide on the *feature extraction intervals*. Then we code a character into a prototype so that future inputs are classified according to their (elastic Hamming) distance to the prototypes.

Handwritten characters are usually quite arbitrary and prone to lots of noise. If the feature extraction interval is very small, then the algorithm will be very sensitive to errors and slight changes in the recognition phase, causing low recognition rate. If the feature extraction interval is very large, then it becomes less likely that we extract the essential features of a character and hence we get a low recognition rate again. We must compromise. The compromise is on the basis of minimum description length of prototypes.

We proceeded as follows to establish an optimal feature selection interval. A set of 186 characters drawings by one subject, exactly 3 examples for each of the 62 alphanumerical characters, were recorded. The character drawings were stored in a standardized integer coordinate system ranged from 0 to 30 in both x and y directions. These character drawings were then input to the system to establish a knowledge base, which formed the collection of prototypes with normalized real coordinates, based on some selected feature extraction interval. After the construction of knowledge base was finished, the system was tested by having it classify the same set of character drawings. If a character is misclassified, it is encoded using extra bits (i.e., the term $P(D \mid H)$). The *error code length* is the sum of the total number of points for all the incorrectly classified character drawings. The *model code* length is the total number of points in all the prototypes in the machine's knowledge base multiplied by 2. The factor of 2 comes from the fact that the prototype coordinates are stored as real numbers. This takes twice as much memory (in C) as the character drawing coordinates which are in integer form. The prototype coordinates are real instead of integer numbers, to facilitate the elastic matching process to give small resolution for comparisons of classification.

Thus, both the model code length and the error code length are directly related to the feature extraction interval. The smaller this interval, the more complex the prototypes, but the smaller the error code length. The effect is reversed if the feature extraction interval goes toward larger values. Since the *total code length* is the sum of the two code lengths, there should be a value of the feature extraction interval gives a minimum for the total code length. This feature extraction interval is considered to be the 'best' one in the spirit of the MDL

principle. The corresponding model, or knowledge base, is considered to be optimal in the sense that it contains enough of the essence of the raw data but eliminates most redundancy of the noise component from the raw data. This optimal feature extraction interval can be found empirically by carrying out the above described build-and-test procedure repeatedly. That is, build the knowledge base, and then test it based on the same set of characters for which it was built. Repeat this for a number of different extraction intervals.

In fact, this actual optimization process is implemented on the system and is available whenever the user wants to call it. For our particular set of characters, the results
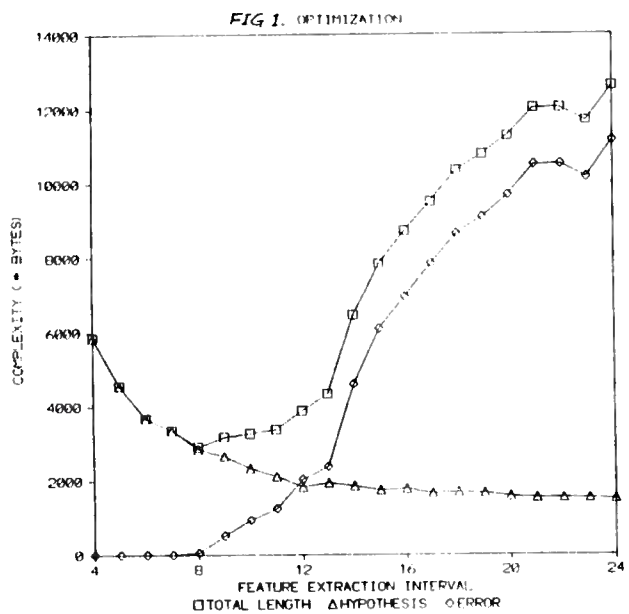


**Figure 1.** Optimization

of this optimization are given in Figure 1, which depicts three quantities: the model code length, the error code length, and the total code length versus feature extraction interval (SAMPLING INTERVAL in the Figure). For larger feature extraction intervals, the model code length is small but most of the character drawings are misclassified, giving a very large total code length. On the other hand, when the feature extraction interval is at the small end of the scale, all the training characters get correctly classified, and the error code length is zero. However the model code length reaches its largest value, resulting in a larger total code length again. The minimum code length occurred at extraction interval of 8, which gives 98.2 percent correct classification. Figure 2 illustrates the fraction of correctly classified character drawings for the training data.

Whether the resulting 'optimal' model really performs better than the models in the same class, the knowledge bases established using different
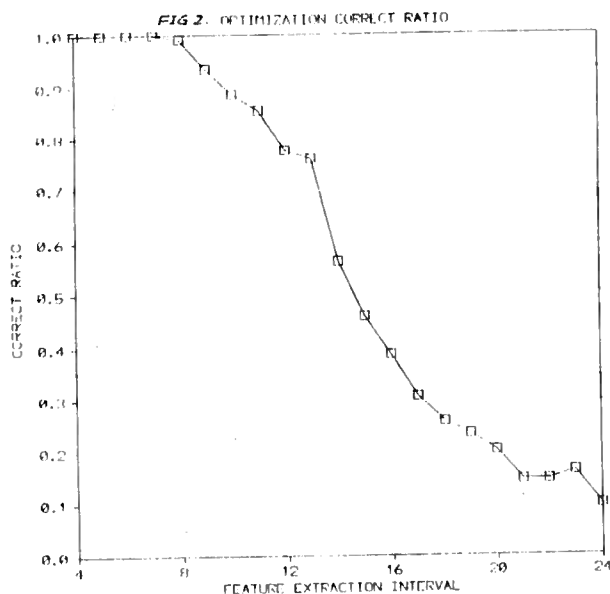
**Figure 2.** Optimization correct ratio

feature extraction intervals, is subject to testing it on new character drawings. For this purpose, the set of 62 handwritten characters were drawn again by the same person who provided the initial data to build the knowledge base. Thus the new data can be considered to be from the same source as the previous data set. The new data were classified by the system using the knowledge base built from the former data set of 186 character drawings, based on different feature extraction intervals. The testing result is
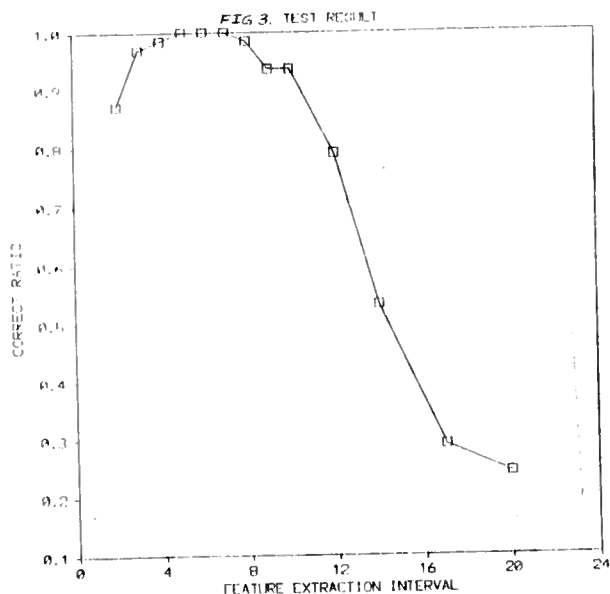


**Figure 3.** Test result

plotted in Figure 3 in terms of the fraction of correct classification (CORRECT RATIO) versus

feature extraction interval (SAMPLING INTERVAL). It is interesting to see that 100% correct classification occurred at feature extraction intervals 5, 6 and 7. These values of feature extraction intervals are close to the optimized value 8. At the low end of feature extraction interval scale the correct classification drops, indicating disturbance caused by too much redundancy in the model. The recommended working feature extraction interval is thus either 7 or 8 for this particular type of character drawings. For more information on this research, see [13] (preprint available from the first author).

## 6. Fischer's Maximum Likelihood Principle

Rissanen [42] has argued that Fisher's maximum likelihood principle is a special case of the MDL principle. By our treatment of MDL it is therefore a more restricted computable approximation to Solomonoff's induction. Notice that minimizing the first term $-\log P(D \mid H)$ is precisely the Maximum Likelihood (ML) principle. ML is sound in statistics, subject to the implicit assumption that each hypothesis $H$ consists of a probability distribution $\theta = (\theta_1, \cdots, \theta_k)$ with the same number $k$ of parameters, each parameter $\theta_i$ with the same fixed precision. In other words, in the probability distribution $P(D \mid H = \theta)$ the number $k$ of parameters of $\theta$, and the precision of each of them, is the same for each $H$. (We use $H$ and $\theta$ interchangeably, statisticians like to use $\theta$) Hence, one assumes that the descriptions of *all hypotheses (models $\theta$)* are of equal length; that is, the complexity of the models is considered to be fixed. This is, obviously, an *objective* assumption. In contrast, the MDL principle minimizes the sum of $-\log P(D \mid H)$ and $-\log P(H)$ Intuitively, if one increases the description length of the hypothesis $H$, it may fit the data better and therefore decrease the description of data given $H$. In the extreme case, if one encodes all the data information into the model $H$ precisely, $P(H)$ is minimized and $-\log P(H)$ is maximized. In that case, no code is needed to describe the data; that is, $P(D \mid H)$ is maximized (equals 1) and $-\log P(D \mid H)$ is minimized (equals 0).

On the other hand, if one decreases the description length of $H$, then this may be penalized by the increasing description length of the data, given $H$. In the extreme case say, $H$ is a trivial hypothesis that contains nothing, then one needs 0 bits to describe $H$. But then, one gains no insight of data and has to "plainly" describe the data without help from any hypothesis.

Hence one may consider the MDL principle as a more general principle than the ML principle in the sense that it considers the trade-off between the complexity of the model $H$ and the power of the model to describe the data $D$, whereas the ML principle does not take the hypothesis complexity into account.

Yet the rationale behind the ML principle was

to be objective by avoiding the 'subjective' assumption of the prior probability. The ML principle is equivalent with selecting the probabilistic model $P(D \mid \theta)$ which permits the shortest ideal code length for the observed sequence, provided that the model used in the encoding, i.e., the parameter $\theta$ is given, too. Thus, the ML principle is just a special case of the MDL principle under the assumption that hypotheses are equally likely and the number of parameters in $\theta$ are fixed and small (so they do not make $P(D \mid \theta) = 1$). The shortcoming of the ML principle is that it cannot handle the situation where we do not know the number (and precision) of the parameters. For example, in the fitting polynomial example, the ML principle does not work well when the degree of the polynomial is not fixed. On the other hand the MDL principle works naturally for this example.

## 7. Jaynes's Maximum Entropy Principle

Rissanen [42] and M. Feder [10] have shown that Jaynes's Maximum Entropy (ME) principle [19-21] can also be considered as a special case of the MDL principle. This is interesting since it is known in statistics that there are a number of important applications where the ML principle fails but where the maximum entropy formalism has been successful, and vice versa. In order to apply Bayes's theorem, we need to decide what the prior probability $p_i = P(H_i)$ is subject to condition

$$\sum_i p_i = 1,$$

and certain other constraints provided by empirical data or considerations of symmetry, probabilistic laws, and so on. Usually these constraints are not sufficient to determine the $p_i$'s uniquely. Jaynes proposed to use the estimated values $\hat{p}_i$ which satisfy said constraints and maximize the entropy function

$$H = -\sum_i p_i \ln p_i$$

subject to the constraints. This is called the *maximum entropy (ME)* principle.

We now demonstrate the rationale behind the ME principle, its use, and its connection with the MDL principle following discussions in [10, 20, 42]. Consider a random experiment with $k$ possible outcomes in each trial, thus $k^n$ possible outcomes in $n$ trails. Let $n_i$ be the number of times the $i$th value appears in an outcome $D$ of $n$ trials. Let frequency $f_i = n_i / n$, $i = 1, 2, ..., k$. The *entropy* of outcome $D$ is:

$$H(f_1, \ldots, f_k) = -\sum_{i=1}^{k} f_i \ln f_i. \tag{10}$$

Let there be $m < k$ linearly independent constraints of the form

$$\sum_{i=1}^{k} a_{ji} f_i = d_j, \quad 1 \leq j \leq m, \text{ and} \tag{11}$$

$$\sum_{i=1}^{k} f_i = 1 \tag{12}$$

where the set $\mathbf{D} = \{d_1, ..., d_m\}$ is related to the observed data, measuring as it were $m$ 'physical quantities' subject to the matrix $A = \{a_{ji}\}$.

**Example.** Consider a loaded die, $k = 6$. If we do not have any information about the die, then using the Epicurus's multiple explanation principle, we may assume that $p_i = 1/6$ for $i = 1, ..., 6$. This actually coincides with the ME principle, since $H(p_1, ..., p_6) = \sum_{i=1}^{6} p_i \ln p_i$ subject to (12) achieves maximum value $\ln 6 = 1.7917595$ for $p_i = 1/6$ for all $i$. Now suppose some experiments on the die have been performed, and it is observed that the die is biased and the average throw gives 4.5. That is,

$$\sum_{i=1}^{6} i p_i = 4.5.$$

In terms of Equation (11), we have $m = 1$, $\mathbf{D} = \{4.5\}$, and $a_{j1} = (1,2,3,4,5,6)$. Maximizing the expression in Equation (10), subject to constraints (11) and (12) gives estimates:

$$\hat{p}_i = e^{-\lambda i} (\sum e^{-\lambda i})^{-1}, \quad 1 = 1, ..., 6,$$

where $\lambda = -0.37105$. Hence $(\hat{p}_1, ..., \hat{p}_6) = (0.0543, 0.0788, 0.1142, 0.1654, 0.2398, 0.3475)$. The maximized entropy $H(\hat{p}_1, \cdots, \hat{p}_6)$ equals 1.61358. How dependable is the ME principle? Jaynes has proven an 'entropy concentration theorem' which, for example, implies that in an experiment of $N = 1000$ trails, 99.99% of all outcomes satisfying the constraints of Equations (11) and (12) have entropy

$$1.602 \leqslant H(\frac{n_1}{n}, \ldots, \frac{n_6}{n}) \leqslant 1.614.$$

Now we turn to the MDL principle to deal with the same problem. The following argument can be derived from probabilistic assumptions. But Kolmogorov [27, 28] advocated a purely combinatorial approach, such as we give below, which does not need any such assumptions. Let $\theta = (p_1, \cdots, p_k)$ be the actual prior distribution of a random variable. We perform a sequence of $n$ independent trials. Kolmogorov observed that the real substance of Formula (10) is that we need approximately $n H(\theta)$ bits to record the sequence of $n$ outcomes. Namely, it suffices to state that each outcome appeared $n_1, ..., n_k$ times, respectively, and afterwards give the index of which one of the

$$C(n_1, \cdots, n_k) = \frac{n!}{n_1! \cdots n_k!}$$

possible sequences $D$ of $n$ outcomes actually took place. For this no more than

$$k \log n + \log C(n_1, \cdots, n_k) + O(\log \log n)$$

bits are needed. The first term corresponds to $-\log P(\theta)$, the second term corresponds to

$-\log P(D \mid \theta)$, and the third term represents the cost of encoding separators between the individual items. Using Stirling's approximation for the factorial function, we find that for large $n$ this is approximately

$$n\left(-\sum_{i=1}^{k} \frac{n_i}{n} \log \frac{n_i}{n}\right) = n H\left(\frac{n_1}{n}, \cdots, \frac{n_k}{n}\right).$$

Since $k$ and $n$ are fixed, the least upper bound on the minimum description length for an arbitrary sequence of $n$ outcomes under certain given constraints $D$ is found by maximizing the term $\log C(n_1,...,n_k)$ subject to said constraints. This is equivalent to maximizing the entropy function (10) under constraints $D$. (Such constraints may be derived, for example, from the laws of large numbers: in case of independent experiments with a probability distribution $\theta$, we have $n_i/n \sim p_i$, and we have a certain rate of convergence with certain probability.)

## 8. Acknowledgement

## References

1. Angluin, D., "On the complexity of minimum inference of regular sets," *Information and Control* **39**, pp. 337-350 (1978).

2. Angluin, D. and L. Valiant, "Fast probabilistic algorithms for hamiltonian circuits and matchings," *J. Computer and System Sciences* **18** (1979).

3. Angluin, D. and C.H. Smith, "Inductive inference: theory and methods," *Computing Surveys* **16**, pp. 239-269 (1983).

4. Asmis, E., *Epicurus Scientific Method*, Cornell University Press, Ithaka, N.Y. (1984).

5. Bayes, Th., "An essay towards solving a problem in the doctrine of chances," *Philosophical Transactions Royal Society*, London **53**, pp. 376-398, (*Ibid.*, 54(1764), pp. 298-310, R. Price, Ed.) (1763).

6. Blumer, A., A. Ehrenfeucht, D. Haussler, and M. Warmuth, "Classifying Learnable Geometric Concepts With the Vapnik-Chervonenkis Dimension," pp. 273-282 in *Proceedings 18th ACM Symp. on Theory of Computing* (1986).

7. Case, J. and C. Smith, "Comparison of identification criteria for machine inductive inference," *Theoret. Comp. Sci.* **25**, pp. 193-220 (1983).

8. Chaitin, G.J., "On the length of programs for computing finite binary sequences: statistical considerations," *J. Assoc. Comp. Mach.* **16**, pp. 145-159 (1969).

9. Chaitin, G.J., "A theory of program size formally identical to information theory," *J. Assoc. Comp. Mach.* **22**, pp. 329-340 (1975).

10. Feder, M., "Maximal entropy as special case of the minimum length description criterion," *IEEE Transactions on Information Theory* **IT-32**, pp. 847-849 (1986).

11. Feller, W., *An Introduction to Probability Theory and Its Applications*, Wiley, New York (1968). (Third edition)

12. Gallager, R.G., *Information Theory and Reliable Communication*, Wiley, New York (1968).

13. Gao, Q. and M. Li, "An application of minimum description length principle to online recognition of handprinted alphanumerals," in *11th International Joint Conference on Artificial Intelligence*, Detroit, MI (1989).

14. Gold, E.M., "Language identification in the limit," *Information and Control* **10**, pp. 447-474 (1967).

15. Gold, E.M., "Complexity of automaton identification from given data," *Information and Control* **37**, pp. 302-320 (1978).

16. Gács, P., "On the symmetry of algorithmic information," *Soviet Math. Dokl.* **15**, p. 1477, (Correction, *Ibid.*, 15(1974)) (1974).

17. Gács, P., "Lecture notes on descriptional complexity and randomness," Manuscript, Boston University, Boston, Mass. (October 1987). (Unpublished)

18. Haussler, D., N. Littlestone, and M Warmuth, "Expected mistake bounds for on-line learning algorithms," Manuscript (1988).

19. Jaynes, E.T., "Prior probabilities," *IEEE Transactions on Systems Science and Cybernetics* **SSC-4**, pp. 227-241 (1968).

20. Jaynes, E.T., "On the rationale of maximum entropy methods," *Proceedings of the IEEE* **70**, pp. 939-952 (1982).

21. Jaynes, E.T., "Prior information in inference," *J. American Statist. Assoc.* (1982).

22. Jeffreys, Z., *Theory of Probability*, Oxford at the Clarendon Press (1961). (Third Edition)

23. Kearns, M., M. Li, L. Pitt, and L. Valiant, "On the learnability of Boolean formulae," pp. 285-295 in *Proc. 19th ACM Symp. on Theory of Computing* (1987).

24. Kearns, M., M. Li, L. Pitt, and L. Valiant, "Recent results on Boolean concept learning," pp. 337-352 in *Proc. 4th Workshop on Machine Learning*, ed. T. Mitchell, Morgan Kaufmann (1987).

25. Kemeny, J.G., "The use of simplicity in induction," *Philos. Rev.* **62**, pp. 391-408 (1953).

26. Kolmogorov, A.N., "Three approaches to the quantitative definition of information," *Problems in Information Transmission* **1**(1), pp. 1-7 (1965).

27. Kolmogorov, A.N., "On the logical foundations of information theory and probability theory," *Problems in Information Transmission* **5**, pp. 1-4 (1969).

28. Kolmogorov, A.N., "Combinatorial foundations of information theory and the calculus of probabilities," *Russian Mathematical Surveys* **38**, pp. 29-40 (1983).

29. Laplace, P.S., *A Philosophical Essay on Probabilities*, Dover, New York (1819). (Date is of original publication.)

30. Levin, L.A., "On the notion of a random sequence," *Soviet Math. Dokl.* **14**, pp. 1413-1416 (1973).

31. Levin, L.A., "Laws of information conservation (non-growth) and aspects of the foundation of probability theory," *Problems in Information Transmission* **10**, pp. 206-210 (1974).

32. Li, M. and P.M.B. Vitányi, "Kolmogorov Complexity and Its Applications," in *Handbook of Theoretical Computer Science*, ed. J. van Leeuwen, North-Holland. (To appear)

33. Li, M. and P.M.B. Vitányi, "Two decades of applied Kolmogorov complexity: In memoriam A.N. Kolmogorov 1903 - 1987," pp. 80-101 in *Proc. 3rd IEEE Conference on Structure in Complexity Theory* (1988).

34. Li, M. and P.M.B. Vitányi, *An Introduction to Kolmogorov Complexity and Its Applications*, Addison-Wesley (To appear).

35. Littlestone, N., "Learning quickly when irrelevant attributes abound: A new linear threshold algorithm," pp. 68-77 in *Proc. 28th IEEE Symposium on Foundations of Computer Science* (1987).

36. Martin-Lof, P., "The definition of random sequences," *Information and Control* 9, pp. 602-619 (1966).

37. Mises, R. von, *Probability, Statistics and Truth*, MacMillan, New York (1939). (Reprint: Dover, 1981)

38. Natarajan, B.K., "Learning functions from examples," Manuscript, Computer Science Dept., Carnegie-Mellon University (1988(?)).

39. Osherson, D., M. Stob, and S. Weinstein, *Systems That Learn*, MIT Press (1986).

40. Pitt, L. and M. Warmuth, "The minimum consistent DFA problem cannot be approximated within any polynomial," in *Proc. 21th ACM Symposium on Theory of Computing* (1989).

41. Quinlan, J. and R. Rivest, "Inferring decision trees using the minimum description length principle," *Information and Computation* 80, pp. 227-248 (1989).

42. Rissanen, J., "A universal prior for integers and estimation by minimum description length," *Annals of Statistics* 11, pp. 416-431 (1982).

43. Rivest, R., "Learning Decision-Lists," Unpublished manuscript, MIT Lab. for Computer Science (December, 1986).

44. Solomonoff, R.J., "A formal theory of inductive inference, Part 1 and Part 2," *Information and Control* 7, pp. 1-22, 224-254 (1964).

45. Solomonoff, R.J., "Complexity-based induction systems: comparisons and convergence theorems," *IEEE Transactions on Information Theory* IT-24, pp. 422-432 (1978).

46. Valiant, L. G., "A Theory of the Learnable," *Comm. ACM* 27, pp. 1134-1142 (1984).

47! Valiant, L. G., "Deductive Learning," *Phil. Trans. Royal Soc. Lond.* A 312, pp. 441-446 (1984).

48. Vapnik, V.N. and A.Ya. Chervonenkis, "On the uniform convergence of relative frequencies of events to their probabilities," *Theory of Probability and Its Applications* 16:2, pp. 264-280 (1971).

49. Wax, M., "Detection of signals by information theoretic criteria," *IEEE Trans. on Accoustics, Speech, and Signal Processing* ASSP-33:2, pp. 387-392 (1985).

50. Zvonkin, A.K. and L.A. Levin, "The complexity of finite objects and the development of the concepts of information and randomness by means of the Theory of Algorithms," *Russ. Math. Surv.* 25, pp. 83-124 (1970).