

# Zeta Distribution and Transfer Learning Problem

Eray Özkural

Celestial Intellect Cybernetics  
celestialintellet.com

**Abstract.** We explore the relations between the zeta distribution and algorithmic information theory via a new model of the transfer learning problem. The program distribution is approximated by a zeta distribution with parameter near 1. We model the training sequence as a stochastic process. We analyze the upper temporal bound for learning a training sequence and its entropy rates, assuming an oracle for the transfer learning problem. We argue from empirical evidence that power-law models are suitable for natural processes. Four sequence models are proposed. Random typing model is like no-free lunch where transfer learning does not work. Zeta process independently samples programs from the zeta distribution. A model of common sub-programs inspired by genetics uses a database of sub-programs. An evolutionary zeta process samples mutations from Zeta distribution. The analysis of stochastic processes inspired by evolution suggest that AI may be feasible in nature, countering no-free lunch sort of arguments.

## 1 Introduction

Although power-law distributions have been analyzed in depth in physical sciences, little has been said about their relevance to Artificial Intelligence (AI). We introduce the zeta distribution as an analytic device in algorithmic information theory and propose using it to approximate the distribution of programs. We have been inspired by the empirical evidence in complex systems, especially biology and genetics, that show an abundance of power-law distributions in nature. It is well possible that the famous universal distribution in AI theory is closely related to power-law distributions in complex systems.

The transfer learning problem also merits our attention, as a general model of it has not been presented in machine learning literature. We develop a basic formalization of the problem using stochastic processes and introduce temporal bounds for learning a training sequence of induction problems, and transfer learning. The entropy rate of a stochastic process emerges as a critical quantity in these bounds. We show how to apply the bounds by analyzing the entropy rates of simple training sequence models that generate programs. Two models are close to what critics of AI have imagined, and easily result in unsolvable problems, while two models inspired by evolution suggest that there may be stochastic processes in nature on which AGI algorithms may be quite effective.

## 2 Approximating the Distribution of Programs

Solomonoff's universal distribution depends on the probability distribution of programs. A natural model is to consider programs, the bits of which are generated by a fair coin. Solomonoff defined the probability of a program  $\pi \in \{0, 1\}^+$

as:

$$P(\pi) = 2^{-|\pi|} \quad (1)$$

where  $|\pi|$  is the program length in bits. The total probability of all programs thus defined unfortunately diverges if all bit-strings  $\pi \in \{0, 1\}^*$  are considered valid programs. For constructing probability distributions, a convergent sum is required. Extended Kraft inequality shows that the total probability is less than 1 for a prefix-free set of infinite programs [2]. Let  $M$  be a reference machine which runs programs with a prefix-free encoding like LISP. The algorithmic probability that a bit-string  $x \in \{0, 1\}^*$  is generated by a random program of  $M$  is:

$$P_M(x) = \sum_{M(\pi)=x^*} P(\pi) \quad (2)$$

which conforms to Kolmogorov's axioms [9].  $P_M$  is also called the universal prior for it may be used as the prior in Bayesian inference, as any data can be encoded as a bit-string.

## 2.1 Zeta Distribution of Programs

We propose the zeta distribution for approximating the distribution of programs of  $M$ . The distribution of (1) is already an approximation, even after normalization, since it contains many programs that are semantically incorrect, and those that do not generate any strings. A realistic program distribution requires us to specify a detailed probability model of programs, which is not covered by the general model, however, the general model, which is approximate, still gives excellent bounds on the limits of Solomonoff's universal induction method. Therefore, other general approximations may also be considered.

Additionally, the zeta function is universal, which encourages us to relate algorithmic information theory to zeta distribution [12].

Let us consider a program bit-string  $\pi = b_1 b_2 b_3 \dots b_k$ . Let  $\phi : \{0, 1\}^+ \rightarrow \mathbb{Z}$  define the arithmetization of programs represented as bit-strings, where the first bit is the most significant bit.

$$\phi(\pi) = \sum_{i=1}^{i \leq |\pi|} b_i \cdot 2^{|\pi|-i} \quad (3)$$

Thus arithmetized, we now show a simple, but interesting inequality about the distribution of programs:

$$P(\pi) = 2^{-\lceil \log_2(\phi(\pi)+1) \rceil} \quad (4)$$

$$(2a)^{-1} \leq 2^{-\lceil \log_2 a \rceil} \leq a^{-1}, \text{ for } a \geq 4 \quad (5)$$

$$(2(\phi(\pi) + 1))^{-1} \leq P(\pi) \leq (\phi(\pi) + 1)^{-1}, \text{ for } \phi(\pi) \geq 3 \quad (6)$$

which shows an approximation that is closer than a factor of 2. Program codes  $\phi(\pi) < 3$  are discarded.

Zipf's law  $f_n \propto n^{-1}$  manifests itself as the Zipf distribution of ranked discrete objects  $\{o_1, o_2, \dots, o_n\}$  in order of increasing rank  $i$

$$P(Z_s^{(n)} = o_i) \triangleq \frac{1}{i^s Z} \quad (7)$$

where  $Z_s^{(n)}$  is a random variable,  $Z$  is the normalization constant and  $s \geq 1$  (we used the notation  $Z_s^{(n)}$  simply to avoid confusion with exponentiation,  $Z_s$  is a standard notation for the zeta random variable). Zeta distribution is the countably infinite version of Zipf distribution with parameter  $s > 1$

$$P(Z_s = k) = \frac{1}{k^s \cdot \zeta(s)} \quad (8)$$

where  $Z_s$  is a random variable with co-domain  $\mathbb{Z}^+$  and the zeta function is defined as

$$\zeta(s) = \sum_{n=1}^{\infty} \frac{1}{n^s} \quad (9)$$

Note that Zeta distribution is a discrete variant of Pareto distribution. It is much involved to work with a prefix-free set, therefore we will suggest an alternative device to approximate  $P(\pi)$ .

**Theorem 1.** *A program distribution may be approximated by the Zipf distribution with  $s = 1$ , or by the zeta distribution with a real  $s$  close to 1 from above.*

*Proof.* (a) Zeta distribution is undefined for  $s = 1$ . However, if we use the Zipf distribution instead, and model programs up to a fixed program-length, we can approximate the program distribution from above using  $(\phi(\pi) + 1)^{-1}$  and from below using  $(2\phi(\pi) + 2)^{-1}$  due to the sandwich property (6).

(b) We can approximate the program distribution from below using  $(2\phi(\pi) + 2)^{-1}$ . Since

$$\forall \epsilon > 0, (2\phi(\pi) + 2)^{-(1+\epsilon)} \leq (2\phi(\pi) + 2)^{-1} < P(\pi),$$

we can also approximate it with the Zeta distribution (8) for  $s$  close to 1.

In either case, the need for a prefix-free set of programs is obviated. Of the simplified distribution, we investigate if the approximations are usable.

**Theorem 2.** *The program distribution  $P(\pi)$  asymptotically obeys a power law with exponent  $-1$  as program size grows.*

*Proof.* The probability of arithmetized program  $\pi$  is sandwiched between  $(\phi(\pi) + 1)^{-1}$  and  $(2\phi(\pi) + 2)^{-1}$ , therefore as  $|\pi|$  grows, Zipf's law grows closer to  $P(\pi)$ .

$$\lim_{|\pi| \rightarrow \infty} (\phi(\pi) + 1)^{-1} - (2\phi(\pi) + 2)^{-1} = 0 \quad (10)$$

$$\lim_{|\pi| \rightarrow \infty} 2^{-|\pi|} - (2\phi(\pi) + 2)^{-1} = \lim_{|\pi| \rightarrow \infty} (\phi(\pi) + 1)^{-1} - 2^{-|\pi|} = 0 \quad (11)$$

Combining Theorem 1 and Theorem 2, we propose using a Zeta distribution with a parameter close to 1. Obviously, lower and upper bounds vary only by a factor of 2 within each other, therefore the error in the approximation of program distribution is at most by 1 bit (this property will be analyzed in detail in an extended version of the present paper). Substituting into (2), we propose an approximation

**Definition 1.**

$$P_M(x) \cong \sum_{M(\pi)=x^*} \frac{1}{(\phi(\pi) + 1)^{1+\epsilon} \cdot \zeta(1+\epsilon)} \quad (12)$$

where  $\zeta(1+\epsilon) \geq 2$  ( $\zeta(1.7) \cong 2$ ). Definition 1 may be useful for machine learning theorists wherever they must represent a priori program probabilities, as it allows them to employ number theory. See Elias Gamma Code [3] for an alternative integer code.

### 3 Training Sequence as a Stochastic Process

Although Solomonoff has theoretically described how the transfer learning problem might be solved in [10], a detailed theoretical model of transfer learning for the universal induction setting is missing in the literature. Here, we attempt to fill this gap. In his treatise of incremental learning, Solomonoff approached the transfer learning problem by describing an update problem which improves the guiding conditional probability distribution (GCPD) of the system as an inductive inference problem of the type that the system usually solves. Solomonoff's modular approach started with a number of problem solving methods and invented new such methods as the system progressed. The initial methods, however, are not fully specified, and we leave it as an open problem in this paper. Instead, we attempt at describing the space of training sequences using the zeta distribution, showing an interesting similarity to our world, whereas most problems in a sequence may be solved, but rarely they are not solvable at all. For instance, a mathematician may solve most problems, but stall at a conjecture that requires the invention of a new, non-trivial axiom indefinitely.

In usual Solomonoff induction (with no transfer learning component), a computable stochastic source  $\mu$  is assumed. The stochastic source may generate sequences, sets, functions, or other structures that we please, the general law of which may be induced via Solomonoff's method. We extend Solomonoff's induction model to a training sequence of induction problems, by considering a stochastic process  $\mathcal{M}$  of  $n$  random variables.

$$\mathcal{M} = \{\mu_1, \mu_2, \mu_3, \dots, \mu_n\} \quad (13)$$

The transfer learning problem thus is constituted from solving  $n$  induction problems in sequence which are generated from the stochastic process  $\mathcal{M}$ . It does not matter which type of induction problem these problems are, as long as they are generated via  $\mathcal{M}$ .

#### 3.1 Entropy Rate of a Training Sequence

A critical measurement of a stochastic process is its entropy rate, which is defined as the following for  $\mathcal{M}$ :

$$H(\mathcal{M}) = \lim_{n \rightarrow \infty} \frac{H(\mu_1, \mu_2, \mu_3, \dots, \mu_n)}{n} \quad (14)$$

and the conditional entropy rate,

$$H'(\mathcal{M}) = \lim_{n \rightarrow \infty} \frac{H(\mu_n | \mu_1, \mu_2, \mu_3, \dots, \mu_{n-1})}{n} \quad (15)$$

which gives the entropy given past observations. Observe that there is a well-known relation between average Kolmogorov complexity and the entropy of an i.i.d. stochastic process (Equation 5 in [1]):

$$\lim_{n \rightarrow \infty} \frac{K_M(X_1, X_2, X_3, \dots, X_n)}{n} = H(X) + O(1) \quad (16)$$

where  $X$  is a stochastic process and  $X_i$  its random variables. We assume that the relation extends to conditional entropy without proof due to lack of space.

### 3.2 Training Time

Let  $\pi_i^*$  be the minimal program for exactly simulating  $\mu_i$  on  $M$ . The most general expression for  $\pi_i^*$  is given in the following

$$\pi_i^* = \arg \min_{\pi_j} (|\pi_j| \mid \forall x, y \in \{0, 1\}^* : M(\pi_j, x, y) = P(\mu_i = x|y)) \quad (17)$$

where the pdf of stochastic source  $\mu_i$  is simulated by a program  $\pi_j$ . The conditional parameter  $y$  is optional. Let us note the following identity

$$K_M(\mu_i) = |\pi_i^*| \quad (18)$$

since arguments  $x, y$  are extraneous input to the pdf specified by  $\pi_i^*$ . Let  $t(\mu_i)$  denote the time taken to solve  $\mu_i$ , and  $t(\pi)$  denote the time taken by program  $\pi$  on  $M$ . Assume that  $t(\mu_i) < \infty$ . We know that the running time of extended Levin Search is bias-optimal [10], and

$$\frac{t(\pi_i^*)}{P(\pi_i^*)} \leq t(\mu_i) \leq \frac{2t(\pi_i^*)}{P(\pi_i^*)} \quad (19)$$

for a computable stochastic source  $\mu_i$  ( $K_M(\mu_i) < \infty$ ). The lower bound in (19) has been named conceptual jump size by Solomonoff, because it refers to the solution of individual induction problems within a training sequence, quantifying how much conceptual innovation is required for a new problem [10]. We cannot exactly predict  $t(\mu_i)$  due to the incomputability of algorithmic probability. Extended Levin Search will keep running indefinitely. It is up to the user to stop execution, which is usually bounded only by the amount of computational resources available to the user. We should also mention that Levin himself does not think that any realistic problems can be solved by Levin search or created on a computer [8]. In the present paper, we run counter to Levin's position, by arguing that Levin search can work in an evolutionary setting, assuming an  $O(1)$  oracle for the transfer learning problem.

We substitute the relation between  $K_M(x)$  and  $P_M(x)$  in the upper bound for  $t(\mu_i)$ ,

$$K_M(\pi_i^*) = -\log_2 P(\pi_i^*) \quad (20)$$

obtaining the following fact due to (18) and (20):

**Lemma 1.**  $t(\mu_i) \leq 2t(\pi_i^*)2^{K_M(\mu_i)}$

The inequality translates to the time for the training sequence  $\mathcal{M}$  as

**Theorem 3.**

$$t(\mathcal{M}) \leq \sum_{i=1}^n t(\pi_i^*)2^{K_M(\mu_i)+1} \quad (21)$$

which is a simple sum of Lemma 1.

The conditional entropy rate is useful when the stochastic process has interdependence. Let us define conditional Kolmogorov complexity for the training sequence  $\mathcal{M}$ ,

$$K'(\mathcal{M}_{<k}) \triangleq K(\mu_k | \mu_1, \mu_2, \mu_3, \dots, \mu_{k-1}) \quad (22)$$

where  $\mathcal{M}_{<k} \triangleq \{\mu_i | i \leq k\}$ . We define likewise for the stochastic process probabilities.

$$P'(\mathcal{M}_{<k}) \triangleq P(\mu_k | \mu_1, \mu_2, \mu_3, \dots, \mu_{k-1}) \quad (23)$$

$K'(\mathcal{M}_{<k})$  captures new algorithmic information content for the  $k^{\text{th}}$  variable of the stochastic process given the entire history.

As  $n$  grows, the transfer learning oracle has to add  $H'(\mathcal{M})$  bits of information to its memory on the average in the stochastic process  $\mathcal{M}$  as Kolmogorov-Shannon entropy relation (16) holds in the limit for conditional entropy, as well. Since the upper temporal bound grows exponentially, (22) only relates loosely to the solution time  $t(\mu_i)$  of a particular problem. We instead define the conditional expected training time upper bound with respect to  $\mathcal{M}$ :

$$\mathbb{E}'[t(\mathcal{M}_{<k})] \triangleq \mathbb{E}_{\mathcal{M}}[t(\mu_k) | \mu_1, \dots, \mu_{k-1}] \leq \sum_{\forall \mu_k \in \{0,1\}^*} 2t(\pi_k^*) 2^{K'(\mathcal{M}_{<k})} P'(\mathcal{M}_{<k}) \quad (24)$$

### 3.3 Random Typing Model

Let us start by considering the well-known model of random typing. If each  $\mu_i$  is regarded as a random  $m$ -bit program out of  $2^m$  such programs, the programs are independent, and the entropy rate is  $m$  bits exactly (under usual i.i.d. assumptions, e.g., we are using fair coin tosses, and we construct programs using a binary alphabet). Assume  $2^m \gg n$ .

In the random typing model, all  $\mu_i$  are algorithmically independent, therefore there is no saving that can be achieved by transfer learning. The time it takes for any problem is therefore:

$$t(\mu_i) \leq t(\pi_i^*) 2^{m+1} \quad (25)$$

for any of the  $2^m$  programs. Since  $m$  can be arbitrarily large, this model is compatible with Levin's conjecture that AI is impossible. Note that this simplistic model is reminiscent of various no-free lunch theorems that were heralded as mathematical proof that general-purpose machine learning was impossible. However, this scenario is highly unrealistic. It is extremely difficult to find problems that are completely independent, as this would require us to be using true random number generators to generate any problem. In other words, we are only showing this "model" to demonstrate how far removed from reality no-free lunch theorems are. In a physical world, this model would correspond to the claim that quantum randomness saturates every observation we may make. However, we already know this claim to be false, since our observations do not consist of noise. On the contrary, there is a lot of dependable regularity in the environment we inhabit, which is sometimes termed "common sense" in AI literature.

### 3.4 Power-law in Nature

A more realistic model, however, uses the zeta distribution for programs instead of uniform distribution. We propose this indeed to be the case since zeta distribution is empirically observed in a multitude of domains, and has good theoretical justification for the abundance of power-law in nature. Theorem 2 gives some

weak and indirect justification as to why we might observe fractions of the zeta distribution of programs in a computable universe. However, there are more direct and appealing reasons why we must expect to see the zeta distribution in highly evolved complex systems. First, it is a direct consequence of the power-law ansatz, and scale-invariance [1] or preferential attachment in evolutionary systems [13]. Second, it follows from an application of maximum entropy principle where the mean of logarithms of observations is fixed [11]. Third, biologists have observed the zeta distribution directly in genetic evolution, thus strengthening the case that our  $\pi_i^*$ 's are likely to conform to zeta distributions. For instance, gene family sizes versus their frequencies follow a power-law distribution [5] and the gene expression in various species follows Zipf's law [4]. Universal regularities in evolution have been observed, for instance in the power-law relation between the number of gene families and gene family size, and number of genes in a category versus number of genes in genome, and power-law like distribution of network node degree [6]. Therefore, there is not only a highly theoretical heuristic argument that we are following, but there exist multiple theoretical and empirical justifications for expecting to observe the zeta distribution of programs in nature. The material evolution of the environment in a habitat, is not altogether different from biological evolution. Except in the case of rare natural catastrophes, the material environment changes only gradually in accord with the dynamic flow of natural law (surprise is small), and is dependent mostly on the actions of organisms in a complex habitat, which may be considered to be programs from an information-theoretic point of view. In that sense, the entire ecology of the habitat in question may be considered to be an evolutionary system, with program frequencies similar to the case of genes in a single organism. In the following, we introduce novel models of training sequences inspired by these empirical justifications.

### 3.5 Identical Zeta Random Variables

Let  $\mathcal{M}$  be i.i.d. generated from zeta distribution according to Theorem 2. Then,

$$H'(\mathcal{M}) = H(\mu_1) = H(Z_s) \quad (26)$$

indicating that the constant entropy rate depends only on the entropy of the zeta distribution. We thus analyze the running time. Let  $t_{max} = \max \{t(\mu_i)\}$ .

$$\mathbb{E}'[t(\mathcal{M}_{<k})] \leq \frac{2t_{max}}{\zeta(s)} \sum_{k=1}^{\infty} 2^{\lceil \log_2 k \rceil} k^{-s} \leq \frac{4t_{max}}{\zeta(s)} \sum_{k=1}^{\infty} \frac{k}{k^s} \quad (27)$$

For the first 1 trillion programs,  $t_{max} \sum_{k=1}^{10^{12}} 4k/k^{1.001} \zeta(1.001) \approx 3.89 \times 10^9 t_{max}$  for  $s = 1.001$ , which is a feasible factor for a realistic program search limit.

Note that AI theorists interpret i.i.d. assumptions as the main reason why no free-lunch theorems are unrealistic [7]. Our i.i.d. zeta process here may be interpreted as an elaboration of that particular objection to no free-lunch theorems. Therefore, we follow the heuristic argument that the right description of the environment which we observe must be something else than the random typing model since agents succeed in transfer learning. The constant zeta process leans towards feasibility, but it does not yet model transfer learning in complex environments.

### 3.6 Zipf Distribution of Sub-programs

Based upon the observations of genetic evolution above and the fact that the whole ecology is an evolutionary system, we may consider a process of programs that has the following property. Each  $\pi_i^*$  that corresponds to  $\mu_i$  is constructed from a number of sub-programs (concatenated). The joint distribution of sub-programs is  $Z_s^{(n)}$ . This is a model of gene frequencies observed in chromosomes, where each chromosome corresponds to a program, and each gene corresponds to a sub-program. Such a distribution would more closely model a realistic distribution of programs by constraining possible programs, as in the real-world the process that generates programs is not ergodic. The total entropy of the process therefore depends on the sub-programs that may be assumed to be random, and program coding. Let each sub-program be a  $k$ -bit random program for the sake of simplicity. The sub-programs that correspond to instructions are specified in a database of  $2^k$  bits. Instructions are not equiprobable, however, as in the random typing model. Let each program have  $m$  instructions drawn from the set of  $2^k$  instructions:

$$A = \{a_1, a_2, a_3, \dots, a_{2^k}\}. \quad (28)$$

Then, we can model each optimal program  $\pi_i^*$  as

$$\pi_i^* = \pi_{i,1}^* \pi_{i,2}^* \pi_{i,3}^* \dots \pi_{i,m}^* \quad (29)$$

which make up a matrix of instructions  $P^* = \pi_{i,j}^*$  where  $\pi_{i,j}^*$  is drawn from the set  $A$  of instructions. The total entropy is due to the database of sub-programs, and the entropy of the global distribution of sub-programs  $Z_s^{(n)}$  which determines the entropy of  $P^*$ . The total entropy is then approximately,

$$H(\mu_1, \mu_2, \dots, \mu_n) \approx \log_2 k + k \cdot 2^k + \log_2 n + \log_2 m + H(Z_s^{(2^k)}) \quad (30)$$

where we show the significant terms for  $k, n, m$ , parameters.

**Lemma 2.** *For the Zipf distribution of sub-programs,*

$$H'(\mathcal{M}) \approx \lim_{n \rightarrow \infty} \frac{1}{n} \left( k \cdot 2^k + \frac{s}{H_{2^k, s}} \sum_{l=1}^{2^k} \frac{\ln(l)}{l^s} + \ln(H_{2^k, s}) + \log_2 k + \log_2 n + \log_2 m \right) \quad (31)$$

due to (30).

which is to say that, the entropy rate, and thus running time, critically depends on the choice of  $k$  and  $n$ .

### 3.7 An Evolutionary Zeta Process

Another process of programs may be determined by mimicking evolution, by considering random mutations of programs in a training sequence. Let us set

$$\pi_1^* = \wedge \quad (32)$$

$$\pi_i^* = \begin{cases} M(Z_s, \pi_{i-1}^*), & \text{if } Z_s \text{ is a valid transformation} \\ \pi_{i-1}^*, & \text{otherwise} \end{cases} \quad (33)$$

which would apply a random transformation sampled from  $Z_s$  in sequence to an initially null program. Such mutations are unlikely to be too complex. The resulting process has small conditional entropy rate, which is wholly dependent on  $Z_s$ .

$$\lim_{n \rightarrow \infty} H'(\mathcal{M}) = H(Z_s) = \log(\zeta(s)) - \frac{s\zeta'(s)}{\zeta(s)} \quad (34)$$

**Lemma 3.**

$$H(Z_{1.1}) = 13.8 \quad H(Z_{1.05}) = 24.5 \quad (35)$$

$$H(Z_{1.01}) = 106.1 \quad H(Z_{1.001}) = 1008.4 \quad (36)$$

The lemma suggests that if an evolutionary process evolves slowly enough, then an AI can easily learn everything there is to learn about it provided that the time complexity of random variables is not too large. We can also employ  $Z_s^{(k)}$  instead of  $Z_s$  in (33). For a universal induction approximation,  $Z_{1.001}$  may be difficult to handle, however, for efficient model-based learning algorithms such as gradient descent methods, digesting new information on the order of a thousand bits is not a big challenge given sufficiently many samples for a problem  $\mu_i$  in the sequence.

## 4 Concluding Remarks

We have shown novel relations between Zipf's law and program distribution by means of the arithmetization of programs. We have shown that zeta distribution may be used for approximating program distributions. We have proposed using the conditional entropy rate as an informative quantity for transfer learning. We have extended Solomonoff's induction model to a training sequence of problems as a stochastic process. We have proposed that the entropy rate of a stochastic process is informative. We have defined conditional Kolmogorov complexity and probability for the sequence, and have used these quantities to define a conditional expected upper bound of training time assuming an  $O(1)$  transfer learning oracle. We introduced sequence models to show that there is a wide range of possible stochastic processes that may be used to argue for the possibility of general purpose AI. The random typing model is a sensible elaboration of no-free lunch theorem kind of arguments, and demonstrate how artificial and unlikely they are since everything is interconnected in nature and pure randomness is very hard to come by, which we therefore rule out as a plausible model of transfer learning. We have shown several empirical justifications for using a power-law model of natural processes. Independent Zeta process tends to be feasible, but does not explain transfer learning. The models that were inspired by natural evolution allow general purpose learning to be feasible. In particular, the model of common sub-programs which is inspired by empirical evidence in genetics supports a view of evolution of natural processes that allows incremental learning to be effective. The evolutionary Zeta process applies random mutations, which can be slow enough for a machine learning algorithm to digest all the new information.

A more detailed analysis of the transfer learning problem will be presented in an extended journal paper. Open problems include analyzing the complexity of the optimal update algorithm, time complexity analysis for the evolutionary processes, and accounting for the time complexity of individual programs.

## Acknowledgements

The paper was substantially improved owing to the extensive and helpful comments of anonymous AGI 2014 and AGI 2018 reviewers.

## References

1. Corominas-Murtra, B., Solé, R.V.: Universality of zipfs law. *Phys. Rev. E* 82, 011102 (Jul 2010), <http://link.aps.org/doi/10.1103/PhysRevE.82.011102>
2. Cover, T.M., Thomas, J.A.: *Elements of Information Theory*. Wiley-Interscience, New York, NY, USA (1991)
3. Elias, P.: Universal codeword sets and representations of the integers. *IEEE Transactions on Information Theory* (1975)
4. Furusawa, C., Kaneko, K.: Zipfs law in gene expression. *Phys. Rev. Lett.* 90, 088102 (Feb 2003), <http://link.aps.org/doi/10.1103/PhysRevLett.90.088102>
5. Huynen, M.A., van Nimwegen, E.: The frequency distribution of gene family sizes in complete genomes. *Molecular Biology and Evolution* 15(5), 583–589 (1998), <http://mbe.oxfordjournals.org/content/15/5/583.abstract>
6. Koonin, E.V.: Are there laws of genome evolution? *PLoS Comput Biol* 7(8), e1002173 (08 2011), <http://dx.doi.org/10.1371/journal.pcbi.1002173>
7. Lattimore, T., Hutter, M.: No free lunch versus occams razor in supervised learning. In: Dowe, D. (ed.) *Algorithmic Probability and Friends. Bayesian Prediction and Artificial Intelligence*, Lecture Notes in Computer Science, vol. 7070, pp. 223–235. Springer Berlin Heidelberg (2013), [http://dx.doi.org/10.1007/978-3-642-44958-1\\_17](http://dx.doi.org/10.1007/978-3-642-44958-1_17)
8. Levin, L.A.: Forbidden Information. eprint arXiv:cs/0203029 (Mar 2002)
9. Levin, L.A.: Some theorems on the algorithmic approach to probability theory and information theory. *CoRR* abs/1009.5894 (2010)
10. Solomonoff, R.J.: A system for incremental learning based on algorithmic probability. In: *Proceedings of the Sixth Israeli Conference on Artificial Intelligence*. pp. 515–527. Tel Aviv, Israel (December 1989)
11. Visser, M.: Zipf’s law, power laws and maximum entropy. *New Journal of Physics* 15(4), 043021 (2013), <http://stacks.iop.org/1367-2630/15/i=4/a=043021>
12. Voronin, S.M.: Theorem on the “universality” of the riemann zeta-function. *Izv. Akad. Nauk SSSR Ser. Mat.* 39(3), 475–486 (1975)
13. Yule, G.U.: A mathematical theory of evolution, based on the conclusions of dr. j. c. willis, f.r.s. *Philosophical Transactions of the Royal Society of London. Series B, Containing Papers of a Biological Character* 213(402-410), 21–87 (1925), <http://rstb.royalsocietypublishing.org/content/213/402-410/21.short>