

A Safe Approximation for Kolmogorov Complexity

Peter Bloem¹, Francisco Mota², Steven de Rooij¹,
Luís Antunes², and Pieter Adriaans¹

¹ System and Network Engineering Group
University of Amsterdam, Amsterdam, The Netherlands
`uva@peterbloem.nl`, `steven.de.rooij@gmail.com`, `p.w.adriaans@uva.nl`

² CRACS & INESC-Porto LA and Institute for Telecommunications
University of Porto, Porto, Portugal
`fmota@fmota.eu`, `lfa@dcc.fc.up.pt`

Abstract. Kolmogorov complexity (K) is an uncomputable function. It can be approximated from above but not to arbitrary given precision and it cannot be approximated from below. By restricting the source of the data to a specific model class, we can construct a computable function $\bar{\kappa}$ to approximate K in a probabilistic sense: the probability that the error is greater than k decays exponentially with k . We apply the same method to the normalized information distance (NID) and discuss conditions that affect the safety of the approximation.

The Kolmogorov complexity of an object is its shortest description, considering all computable descriptions. It has been described as “the accepted absolute measure of information content of an individual object” [1], and its investigation has spawned a slew of derived functions and analytical tools. Most of these tend to separate neatly into one of two categories: the platonic and the practical.

On the platonic side, we find such tools as the normalized information distance [2], algorithmic statistics [1] and sophistication [3, 4]. These subjects all deal with uncomputable “ideal” functions: they optimize over all computable functions, but they cannot be computed themselves.

To construct practical applications (ie. runnable computer programs), the most common approach is to take one of these platonic, uncomputable functions, derived from Kolmogorov complexity (K), and to approximate it by swapping K out for a computable compressor like GZIP [5]. This approach has proved effective in the case of normalized information distance (NID) [2] and its approximation, the normalized compression distance (NCD) [6]. Unfortunately, the switch to a general-purpose compressor leaves an analytical gap. We know that the compressor serves as an upper bound to K —up to a constant—but we do not know the difference between the two, and how this error affects the error of derived functions like the NCD. This can cause serious contradictions. For instance, the normalized information distance has been shown to be non-approximable [7], yet the NCD has proved its merit empirically [6]. Why this

should be the case, and when this approach may fail has, to our knowledge, not yet been investigated.

We aim to provide the first tools to bridge this gap. We will define a computable function which can be said to approximate Kolmogorov complexity, with some practical limit to the error. To this end, we introduce two concepts:

- We generalize resource-bounded Kolmogorov complexity (K^t) to *model-bounded Kolmogorov complexity*, which minimizes an object’s description length over any given enumerable subset of Turing machines (a *model class*). We explicitly assume that the source of the data is contained in the model class.
- We introduce a probabilistic notion of approximation. A function approximates another *safely*, under a given distribution, if the probability of them differing by more than k bits, decays at least exponentially in k .¹

While the resource-bounded Kolmogorov complexity is computable in a technical sense, it is never computed practically. The generalization to model bounded Kolmogorov complexity creates a connection to *minimum description length* (MDL) [8, 9, 10], which does produce algorithms and methods that are used in a practical manner. Kolmogorov complexity has long been seen as a kind of platonic ideal which MDL approximates. Our results show that MDL is not just an upper bound to K , it also approximates it in a probabilistic sense.

Interestingly, the model-bounded Kolmogorov complexity itself—the smallest description using a single element from the model class—is not a safe approximation. We can, however, construct a computable, safe approximation by taking into account all descriptions the model class provides for the data.

The main result of this paper is a computable function $\bar{\kappa}$ which, under a model assumption, safely approximates K (Theorem 3). We also investigate whether a $\bar{\kappa}$ -based approximation of NID is safe, for different properties of the model class from which the data originated (Theorems 5, 6 and 7).

1 Turing Machines and Probability

Turing Machines

Let $\mathbb{B} = \{0, 1\}^*$. We assume that our data is encoded as a finite binary string. Specifically, the natural numbers can be associated to binary strings, for instance by the bijection: $(0, \epsilon)$, $(1, 0)$, $(2, 1)$, $(3, 00)$, $(4, 01)$, etc, where ϵ is the empty string. To simplify notation, we will sometimes conflate natural numbers and binary strings, implicitly using this ordering.

We fix a canonical prefix-free coding, denoted by \bar{x} , such that $|\bar{x}| \leq |x| + 2\log|x|$. See [11, Example 1.11.13] for an example. Among other things, this gives us a canonical pairing function to encode two strings x and y into one: \bar{xy} .

¹ This consideration is subject to all the normal drawbacks of asymptotic approaches. For this reason, we have foregone the use of big-O notation as much as possible, in order to make the constants and their meaning explicit.

We use the Turing machine model from [11, Example 3.1.1]. The following properties are important: the machine has a read-only, right-moving input tape, an auxiliary tape which is read-only and two-way, two read-write two-way work-tapes and a read-write two-way output tape.² All tapes are one-way infinite. If a tape head moves off the tape or the machine reads beyond the length of the input, it enters an infinite loop. For the function computed by TM i on input p with auxiliary input y , we write $T_i(p \mid y)$ and $T_i(p) = T_i(p \mid \epsilon)$. The most important consequence of this construction is that the programs for which a machine with a given auxiliary input y halts, form a prefix-free set [11, Example 3.1.1]. This allows us to interpret the machine as a probability distribution (as described in the next subsection).

We fix an effective ordering $\{T_i\}$. We call the set of all Turing machines \mathcal{C} . There exists a universal Turing machine, which we will call U , that has the property that $U(\bar{i}p \mid y) = T_i(p \mid y)$ [11, Theorem 3.1.1].

Probability

We want to formalize the idea of a probability distribution that is *computable*: it can be simulated or computed by a computational process. For this purpose, we will interpret a given Turing machine T_q as a probability distribution p_q : each time the machine reads from the input tape, we provide it with a random bit. The Turing machine will either halt, read a finite number of bits without halting, or read an unbounded number of bits. $p_q(x)$ is the probability that this process halts and produces x : $p_q(x) = \sum_{p: T_q(p)=x} 2^{-|p|}$. We say that T_q *samples* p_q . Note that if p is a semimeasure, $1 - \sum_x p(x)$ corresponds to the probability that this sampling process will not halt.

We model the probability of x conditional on y by a Turing machine with y on its auxiliary tape: $p_q(x \mid y) = \sum_{p: T_q(p|y)=x} 2^{-|p|}$.

The *lower semicomputable semimeasures* [11, Chapter 4] are an alternative formalization. We show that it is equivalent to ours:

Lemma 1. *** The set of probability distributions sampled by Turing machines in \mathcal{C} is equivalent to the set of lower semicomputable semimeasures.*

The distribution corresponding to the universal Turing machine U is called m : $m(x) = \sum_{p: U(p)=x} 2^{-|p|}$. This is known as a universal distribution. K and m dominate each other, ie. $\exists c \forall x : |K(x) - \log m(x)| < c$ [11, Theorem 4.3.3].

2 Model-Bounded Kolmogorov Complexity

In this section we present a generalization of the notion of resource-bounded Kolmogorov complexity. We first review the unbounded version:

² Multiple work tapes are only required for proofs involving resource bounds.

** Proof in the appendix.

Definition 1. Let $k(x \mid y) = \arg \min_{p: U(p|y)=x} |p|$. The prefix-free, conditional Kolmogorov complexity is

$$K(x \mid y) = |k(x \mid y)|$$

with $K(x) = K(x \mid \epsilon)$.

Due to the halting problem, K is not computable. By limiting the set of Turing machines under consideration, we can create a computable approximation.

Definition 2. A model class $C \subseteq \mathcal{C}$ is a computably enumerable set of Turing machines. Its members are called models. A universal model for C is a Turing machine U^C such that $U^C(\bar{i}p \mid y) = T_i(p \mid y)$ where i is an index over the elements of C .

Definition 3. For a given C and U^C we have $K^C(x) = \min \{|p| : U^C(p) = x\}$, called the model-bounded Kolmogorov complexity.

K^C , unlike K , depends heavily on the choice of enumeration of C . A notation like K_{U^C} or $K^{i,C}$ would express this dependence better, but for the sake of clarity we will use K^C .

We define a model-bounded variant of m as $m^C(x) = \sum_{p: U^C(p)=x} 2^{-|p|}$, which dominates all distributions in C :

Lemma 2. For any $T_q \in C$, $m^C(x) \geq c_q p_q(x)$ for some c_q independent of x .

Proof.

$$m^C(x) = \sum_{i, p: U^C(\bar{i}p)=x} 2^{-|\bar{i}p|} \geq \sum_{p: U^C(\bar{q}p)=x} 2^{-|\bar{q}|} 2^{-|p|} = 2^{-|\bar{q}|} p_q(x). \quad \square$$

Unlike K and $-\log m$, K^C and $-\log m^C$ do not dominate one another. We can only show that $-\log m^C$ bounds K^C from below ($\sum_{U^C(p)=x} 2^{-|p|} > 2^{-|k^C(x)|}$). In fact, as shown in Theorem 1, $-\log m^C$ and K^C can differ by arbitrary amounts.

Example 1 (Resource-Bounded Kolmogorov Complexity [11, Ch. 7]).

Let $t(n)$ be some time-constructible function.³ Let T_i^t be the modification of $T_i \in \mathcal{C}$ such that at any point in the computation, it halts immediately if more than k cells have been written to on the output tape and the number of steps that have passed is less than $t(k)$. In this case, whatever is on the output tape is taken as the output of the computation. If this situation does not occur, T_i runs as normal. Let $U^t(\bar{i}p) = T_i^t(p)$. We call this model class C^t . We abbreviate K^{C^t} as K^t .

Since there is no known means of simulating U^t within $t(n)$, we do not know whether $U^t \in C^t$. It can be run in $ct(n) \log t(n)$ [11, 12], so we do know that $U^t \in C^{ct \log t}$.

Other model classes include Deterministic Finite Automata, Markov Chains, or the exponential family (suitably discretized). These have all been thoroughly investigated in coding contexts in the field of Minimum Description Length [10].

³ I.e. $t: \mathbb{N} \rightarrow \mathbb{N}$ and t can be computed in $O(t(n))$ [13].

3 Safe Approximation

When a code-length function like K turns out to be uncomputable, we may try to find a lower and upper bound, or to find a function which dominates it. Unfortunately, neither of these will help us. Such functions invariably turn out to be uncomputable themselves [11, Section 2.3].

To bridge the gap between uncomputable and computable functions, we require a softer notion of approximation; one which states that errors of any size may occur, but that the larger errors are so unlikely, that they can be safely ignored:

Definition 4. *Let f and f_a be two functions. We take f_a to be an approximation of f . We call the approximation b -safe (from above) for a distribution (or adversary) p if for all k and some $c > 0$:*

$$p(f_a(x) - f(x) \geq k) \leq cb^{-k}.$$

Since we focus on code-length functions, usually omit “from above”. A safe function is b -safe for some $b > 1$. An approximation is safe for a model class C if it is safe for all p_q with $T_q \in C$.

While the definition requires this property to hold for all k , it actually suffices to show that it holds for k above a constant, as we can freely scale c :

Lemma 3. *If $\exists_c \forall_{k:k > k_0} : p(f_a(x) - f(x) \geq k) \leq cb^{-k}$, then f_a is b -safe for f against p .*

Proof. First, we name the k below k_0 for which the ratio between the bound and the probability is the greatest: $k_m = \arg \max_{k \in [0, k_0]} [p(f_a(x) - f(x) \geq k) / cb^{-k}]$. We also define $b_m = cb^{-k_m}$ and $p_m = p(f_a(x) - f(x) \geq k_m)$. At k_m , we have $p(f_a(x) - f(x) \geq k_m) = p_m = \frac{p_m}{b_m} cb^{-k_m}$. In other words, the bound $c'b^{-k}$ with $c' = \frac{p_m}{b_m}c$ bounds p at k_m , the point where it diverges the most from the old bound. Therefore, it must bound it at all other $k > 0$ as well. \square

Safe approximation, domination and lowerbounding form a hierarchy:

Lemma 4. *Let f_a and f be code-length functions. If f_a is a lower bound on f , it also dominates f . If f_a dominates f , it is also a safe approximation.*

Proof. Domination means that for all x : $f_a(x) - f(x) < c$, if f_a is a lower bound, $c = 0$. If f_a dominates f we have $\forall p, k > c : p(f_a(x) - f(x) \geq k) = 0$. \square

Finally, we show that safe approximation is transitive, so we can chain together proofs of safe approximation; if we have several functions with each safe for the next, we know that the first is also safe for the last.

Lemma 5. *The property of safety is transitive over the space of functions from \mathbb{B} to \mathbb{B} for a fixed adversary.*

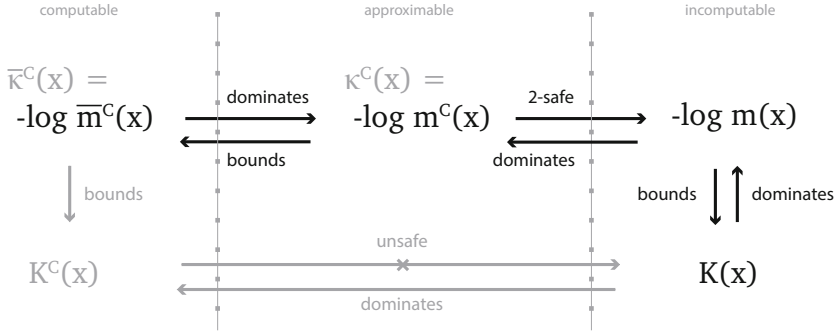


Fig. 1. An overview of how various code-length functions relate to each other in terms of approximation safety. These relations hold under the assumption that the data is generated by a distribution in C and that C is sufficient and complete.

Proof. Let f , g and h be functions such that

$$p(f(x) - g(x) \geq k) \leq c_1 b_1^{-k} \text{ and } p(g(x) - h(x) \geq k) \leq c_2 b_2^{-k}.$$

We need to show that $p(f(x) - h(x) \geq k)$ decays exponentially with k . We start with

$$p(f(x) - g(x) \geq k \vee g(x) - h(x) \geq k) \leq c_1 b_1^{-k} + c_2 b_2^{-k}.$$

$\{x : f(x) - h(x) \geq 2k\}$ is a subset of $\{x : f(x) - g(x) \geq k \vee g(x) - h(x) \geq k\}$, so that the probability of the first set is less than that of the second:

$$p(f(x) - h(x) \geq 2k) \leq c_1 b_1^{-k} + c_2 b_2^{-k}.$$

Which gives us

$$p(f(x) - h(x) \geq 2k) \leq c b^{-k} \quad \text{with } b = \min(b_1, b_2) \text{ and } c = \max(c_1, c_2),$$

$$p(f(x) - h(x) \geq k') \leq c b'^{-k'} \quad \text{with } b' = \sqrt{b}. \quad \square$$

4 A Safe, Computable Approximation of K

Assuming that our data is produced from a model in C , can we construct a computable function which is safe for K ? An obvious first choice is K^C . For it to be computable, we would normally ensure that all programs for all models in C halt. Since the halting programs form a prefix-free set, this is impossible. There is however a property for prefix-free functions that is analogous. We call this *sufficiency*:

Definition 5. A sufficient model T is a model for which every infinite binary string contains a halting program as a prefix. A sufficient model class contains only sufficient models.

We can therefore enumerate all inputs for U^C from short to long in series to find $k^C(x)$, so long as C is sufficient. For each input, U^C either halts or attempts to read beyond the length of the input.

In certain cases, we also require that C can represent all $x \in \mathbb{B}$ (ie. $m^C(x)$ is never 0). We call this property *completeness*:

Definition 6. *A model class C is called complete if for any x , there is at least one p such that $U^C(p) = x$.*

We can now say, for instance, that K^C is computable for sufficient C . Unfortunately, K^C turns out to be unsafe:

Theorem 1. *There exist model classes C so that $K^C(x)$ is an unsafe approximation for $K(x)$ against some p_q with $T_q \in C$.*

Proof. We first show that K^C is unsafe for $-\log m^C$.

Let C contain a single Turing machine T_q which outputs x for any input of the form $\bar{x}p$ with $|p| = x$ and computes indefinitely for all other inputs.

T_q samples from $p_q(x) = 2^{-|\bar{x}|}$, but it distributes each x 's probability mass uniformly over many programs much longer than $|\bar{x}|$.

This gives us $K^C(x) = |\bar{x}| + |p| = |\bar{x}| + x$ and $-\log m^C(x) = |\bar{x}|$, so that $K^C(x) + \log m^C(x) = x$. We get

$$m^C(K^C(x) + \log m^C(x) \geq k) = m^C(x \geq k) = \sum_{x: x \geq k} 2^{-|\bar{x}|} \geq \sum_{x: x \geq k} 2^{-2 \log x} \geq k^{-2}$$

so that K^C is unsafe for $-\log m^C$.

It remains to show that this implies that K^C is unsafe for K . In Theorem 2, we prove that $-\log m^C$ is safe for K . Assuming that K^C is safe for K (which dominates $-\log m^C$) implies K^C is safe for $-\log m^C$, which gives us a contradiction. \square

Note that the use of a model class with a single model is for convenience only. The main requirement for K^C to be unsafe is that the prefix tree of U^C 's programs distributes the probability mass for x over many programs of similar length. The greater the difference between K^C and $-\log m^C$, the greater the likelihood that K^C is unsafe.

Our next candidate for a safe approximation of K is $-\log m^C$. This time, we fare better. We first require the following lemma, called the *no-hypercompression theorem* in [10, p103]:

Lemma 6. *Let p_q be a probability distribution. The corresponding code-length function, $-\log p_q$, is a 2-safe approximation for any other code-length function against p_q . For any p_r and $k > 0$: $p_q(-\log p_q(x) + \log p_r(x) \geq k) \leq 2^{-k}$.*

Theorem 2. *$-\log m^C(x)$ is a 2-safe approximation of $K(x)$ against any adversary from C .*

Proof. Let p_q be some adversary in C . We have

$$\begin{aligned} p_q(-\log m^C(x) - K(x) \geq k) \\ &\leq cm^C(-\log m^C(x) - K(x) \geq k) && \text{by Lemma 2,} \\ &\leq c2^{-k} && \text{by Lemma 6.} \quad \square \end{aligned}$$

While we have shown m^C to be safe for K , it may not be computable, even if C is sufficient (since it is an infinite sum). We can, however, define an approximation, which, for sufficient C , is computable and dominates m^C .

Definition 7. Let the model class D be the union of C and some arbitrary sufficient and complete distribution from \mathcal{C} .

Let $\overline{m}_c^C(x)$ be the function computed by the following algorithm: Dovetail the computation of all programs on $U^D(x)$ in cycles, so that in cycle n , the first n programs are simulated for one further step. After each such step we consider the probability mass s of all programs that have stopped (where each program p contributes $2^{-|p|}$), and the probability mass s_x of all programs that have stopped and produced x . We halt the dovetailing and output s_x if $s_x > 0$ and the following stop condition is met:

$$\frac{1-s}{s_x} \leq 2^c - 1.$$

Note that if C is sufficient so is D , so that s goes to 1 and s_x never decreases. Since all programs halt, the stop condition must be reached. The addition of a complete model is required to ensure that s_x does not remain 0 indefinitely.

Lemma 7. If C is sufficient, $\overline{m}_c^C(x)$ dominates m^C with a constant multiplicative factor 2^{-c} (ie. their code-lengths differ by at most c bits).

Proof. We will first show that \overline{m}_c^C dominates m^D . Note that when the computation of \overline{m}_c^C halts, we have $\overline{m}_c^C(x) = s_x$ and $m^D(x) \leq s_x + (1-s)$. This gives us:

$$\frac{m^D(x)}{\overline{m}_c^C(x)} \leq 1 + \frac{1-s}{s_x} \leq 2^c.$$

Since $C \subseteq D$, m^D dominates m^C (see Lemma 9 in the appendix) and thus, \overline{m}_c^C dominates m^C . \square

The parameter c in \overline{m}_c^C allows us to tune the algorithm to trade off running time for a smaller constant of domination. We will usually omit it when it is not relevant to the context.

Putting all this together, we have achieved our aim:

Theorem 3. For a sufficient model class C , $-\log \overline{m}^C$ is a safe, computable approximation of $K(x)$ against any adversary from C

Proof. We have shown that, under these conditions, $-\log m^C$ safely approximates $-\log m$ which dominates K , and that $-\log \overline{m}^C$ dominates $-\log m^C$. Since domination implies safe approximation (Lemma 4), and safe approximation is transitive (Lemma 5), we have proved the theorem. \square

Figure 1 summarizes this chain of reasoning and other relations between the various code-length functions mentioned.

The negative logarithm of m^C will be our go-to approximation of K , so we will abbreviate it with κ :

Definition 8. $\kappa^C(x) = -\log m^C(x)$ and $\bar{\kappa}^C(x) = -\log \bar{m}^C(x)$.

Finally, if we violate our model assumption we lose the property of safety. For adversaries outside C , we cannot be sure that κ^C is safe:

Theorem 4. *There exist adversaries p_q with $T_q \notin C$ for which neither κ^C nor $\bar{\kappa}^C$ is a safe approximation of K .*

Proof. Consider the following algorithm for sampling from a computable distribution (which we will call p_q):

- Sample $n \in \mathbb{N}$ from some distribution $s(n)$ which decays polynomially.
- Loop over all x of length n return the first x such that $\kappa^C(x) \geq n$.

Note that at least one such x must exist by a counting argument: if all x of length n have $-\log \bar{m}^C(x) < n$ we have a code that assigns 2^n different strings to $2^n - 1$ different codes.

For each x sampled from q , we know that $\bar{\kappa}(x) \geq |x|$ and $K(x) \leq -\log p_q(x) + c_q$. Thus:

$$\begin{aligned} p_q(\bar{\kappa}^C(x) - K(x) \geq k) &\geq p_q(|x| + \log p_q(x) - c_q \geq k) \\ &= p_q(|x| + \log s(|x|) - c_q \geq k) = \sum_{n: n + \log s(n) - c_q \geq k} s(n). \end{aligned}$$

Let n_0 be the smallest n for which $2n > n + \log s(n) - c_q$. For all $k > 2n_0$ we have

$$\sum_{n: n + \log s(n) - c_q \geq k} s(n) \geq \sum_{n: 2n \geq k} s(n) \geq s\left(\frac{1}{2}k\right). \quad \square$$

For C^t (as in Example 1), we can sample the p_q constructed in the proof in $O(2^n \cdot t(n))$. Thus, we know that κ^t is safe for K against adversaries from C^t , and we know that it is unsafe against C^{2^t} .

5 Approximating Normalized Information Distance

Definition 9 ([2, 6]). *The normalized information distance between two strings x and y is*

$$NID(x, y) = \frac{\max[K(x | y), K(y | x)]}{\max[K(x), K(y)]}.$$

The information distance (ID) is the numerator of this function. The NID is neither lower nor upper semicomputable [7]. Here, we investigate whether we can safely approximate either function using κ . We define ID^C and NID^C as the ID and NID functions with K replaced by $\bar{\kappa}^C$. We first show that, even if the adversary only combines functions and distributions in C , ID^C may be an unsafe approximation.

Definition 10. ⁴A function f is a (b -safe) model-bounded one-way function for C if it is injective, and for some $b > 1$, some $c > 0$, all $q \in C$ and all k :

$$p_q(\kappa^C(x) - \kappa^C(x \mid f(x)) \geq k) \leq cb^{-k}.$$

Theorem 5. ^{**} Under the following assumptions:

- C contains a model T_0 , with $p_0(x) = 2^{-|x|}s(|x|)$, with s a distribution on \mathbb{N} which decays polynomially or slower,
- there exists a model-bounded one-way function f for C ,
- C is normal, ie. for some c and all x : $\kappa^C(x) < |\bar{x}| + c$

ID^C is an unsafe approximation for ID against an adversary T_q which samples x from p_0 and returns $\bar{x}f(x)$.

If x and y are sampled from C independently, we can prove safety:

Theorem 6. ^{**} Let T_q be a Turing machine which samples x from p_a , y from p_b and returns $\bar{x}y$. If $T_a, T_b \in C$, $ID^C(x, y)$ is a safe approximation for $ID(x, y)$ against any such T_q .

The proof relies on two facts:

- $\bar{\kappa}^C(x \mid y)$ is safe for $K(x \mid y)$ if x and y are generated this way.
- Maximization is a *safety preserving operation*: if we have two functions f and g with safe approximations f_a and g_a , $\max(f_a(x), g_a(x))$ is a safe approximation of $\max(f(x), g(x))$.

For *normalized* information distance, which is dimensionless, the error k in bits as we have used it so far does not mean much. Instead, we use f/f_a as a measure of approximation error, and we introduce an additional parameter ϵ :

Theorem 7. ^{**} We can approximate NID with NID^C with the following bound:

$$p_q\left(\frac{NID(x, y)}{NID^C(x, y)} \notin \left(1 - \frac{k}{c}, 1 + \frac{k}{c}\right)\right) \leq c'b^{-k} + 2\epsilon$$

with

$$p_q(ID^C(x, y) \geq c) \leq \epsilon \text{ and } p_q(\max[\kappa^C(x), \kappa^C(y)] \geq c) \leq \epsilon$$

for some $b > 1$ and $c' > 0$, assuming that p_q samples x and y independently from models in C .

⁴ This is similar to the Kolmogorov one-way function [14, Definition 11].

6 Discussion

We have provided a function $\bar{\kappa}^C(x)$ for a given model class C , which is computable if C is sufficient. Under the assumption that x is produced by a model from C , $\bar{\kappa}^C(x)$ approximates $K(x)$ in a probabilistic sense. We have also shown that $K^C(x)$ is not safe. Finally, we have given some insight into the conditions on C and the adversary, which can affect the safety of NCD as an approximation to NID.

Since, as shown in Example 1, resource-bounded Kolmogorov complexity is a variant of model-bounded Kolmogorov complexity, our results apply to K^t as well: K^t is not necessarily a safe approximation of K , even if the data can be sampled in t and κ^t is safe if the data can be sampled in t . Whether K^t is safe ultimately depends on whether a single shortest program dominates among the sum of all programs, as it does in the unbounded case.

For complex model classes, κ^C may still be impractical to compute. In such cases, we may be able to continue the chain of safe approximation proofs. Ideally, we would show that a model which is only locally optimal, found by an iterative method like gradient descent, is still a safe approximation of K . Such proofs would truly close the circuit between the ideal world of Kolmogorov complexity and modern statistical practice.

Acknowledgement. We would like to thank the reviewers for their insightful comments. This publication was supported by the Dutch national program COMMIT, the Netherlands eScience center, the ERDF (European Regional Development Fund) through the COMPETE Programme (Operational Programme for Competitiveness) and by National Funds through the FCT (Fundação para a Ciência e a Tecnologia, the Portuguese Foundation for Science and Technology) within project *FCOMP-01-0124-FEDER-037281*.

References

- [1] Gács, P., Tromp, J., Vitányi, P.M.B.: Algorithmic statistics. *IEEE Transactions on Information Theory* 47(6), 2443–2463 (2001)
- [2] Li, M., Chen, X., Li, X., Ma, B., Vitányi, P.M.B.: The similarity metric. *IEEE Transactions on Information Theory* 50(12), 3250–3264 (2004)
- [3] Vitányi, P.M.B.: Meaningful information. *IEEE Transactions on Information Theory* 52(10), 4617–4626 (2006)
- [4] Adriaans, P.: Facticity as the amount of self-descriptive information in a data set. *arXiv preprint arXiv:1203.2245* (2012)
- [5] Gailly, J., Adler, M.: The GZIP compressor (1991)
- [6] Cilibrasi, R., Vitányi, P.M.B.: Clustering by compression. *IEEE Transactions on Information Theory* 51(4), 1523–1545 (2005)
- [7] Terwijn, S.A., Torenvliet, L., Vitányi, P.M.B.: Nonapproximability of the normalized information distance. *J. Comput. Syst. Sci.* 77(4), 738–742 (2011)
- [8] Rissanen, J.: Modeling by shortest data description. *Automatica* 14(5), 465–471 (1978)

- [9] Rissanen, J.: Universal coding, information, prediction, and estimation. IEEE Transactions on Information Theory 30(4), 629–636 (1984)
- [10] Grünwald, P.D.: The Minimum Description Length Principle. Adaptive computation and machine learning series. The MIT Press (2007)
- [11] Li, M., Vitányi, P.M.B.: An introduction to Kolmogorov complexity and its applications, 2nd edn. Graduate Texts in Computer Science. Springer (1997)
- [12] Hennie, F.C., Stearns, R.E.: Two-tape simulation of multitape Turing machines. J. ACM 13(4), 533–546 (1966)
- [13] Antunes, L.F.C., Matos, A., Souto, A., Vitányi, P.M.B.: Depth as randomness deficiency. Theory Comput. Syst. 45(4), 724–739 (2009)
- [14] Antunes, L.F.C., Matos, A., Pinto, A., Souto, A., Teixeira, A.: One-way functions using algorithmic and classical information theories. Theory Comput. Syst. 52(1), 162–178 (2013)

A Appendix

A.1 Turing Machines and lsc. Probability Semimeasures (Lemma 1)

Definition 11. A function $f : \mathbb{B} \rightarrow \mathbb{R}$ is lower semicomputable (lsc.) iff there exists a total, computable two-argument function $f' : \mathbb{B} \times \mathbb{N} \rightarrow \mathbb{Q}$ such that: $\lim_{i \rightarrow \infty} f'(x, i) = f(x)$ and for all i , $f'(x, i+1) \geq f'(x, i)$.

Lemma 8. If f is an lsc. probability semimeasure, then there exists a function $f^*(x, i)$ with the same properties of the function f' from Definition 11, and the additional property that all values returned by f^* have finite binary expansions.

Proof. Let x_j represent $x \in \mathbb{D}$ truncated at the first j bits of its binary expansion and x^j the remainder. Let $f^*(x, i) = f'(x, i)_i$. Since $f'(x, i) - f^*(x, i)_i$ is a value with $i+1$ as the highest non-zero bit in its binary expansion, $\lim_{i \rightarrow \infty} f^*(x, i) = \lim_{i \rightarrow \infty} f'(x, i) = f(x)$.

It remains to show that f^* is nondecreasing in i . Let $x \geq y$. We will show that $x_j \geq y_j$, and thus $x_{j+1} \geq y_j$. If $x = y$ the result follows trivially. Otherwise, we have $x_j = x - x^j > y - x^j = y_j + y^j - x^j \geq y_j - 2^{-j}$. Substituting $x = f'(x, i+1)$ and $y = f'(x, i)$ tells us that $f^*(x, i+1) \geq f^*(x, i)$ \square

Theorem 8. Any TM, T_q , samples from an lsc. probability semimeasure.

Proof. We will define a program computing a function $p'_q(x, i)$ to approximate $p_q(x)$: Dovetail the computation of T_q on all inputs $x \in \mathbb{B}$ for i cycles.

Clearly this function is nondecreasing. To show that it goes to $p(x)$ with i , we first note that for a given i_0 there is a j such that, $2^{-j-1} < p_q(x) - p_q(x, i_0) \leq 2^{-j}$. Let $\{p_i\}$ be an ordering of the programs producing x , by increasing length, that have not yet stopped at dovetailing cycle i_0 . There is an m such that $\sum_{i=1}^m 2^{-|p_i|} \geq 2^{-j-1}$, since $\sum_{i=1}^{\infty} 2^{-|p_i|} > 2^{-j-i}$. Let i_1 be the dovetailing cycle for which the last program below p_{m+1} halts. This gives us $p_q(x) - p_q(x, i_1) \leq 2^{-j-1}$. Thus, by induction, we can choose i to make $p(x) - p'(x, i)$ arbitrarily small. \square

Theorem 9. *Any lsc. probability semimeasure can be sampled by a TM.*

Proof. Let $p(x)$ be an lsc. probability semimeasure and $p^*(x, i)$ as in Lemma 8. We assume—without loss of generality—that $p^*(x, 0) = 0$. Consider the following algorithm:

```

initialize  $s \leftarrow \epsilon, r \leftarrow \epsilon$ 
for  $c = 1, 2, \dots$ :
  for  $x \in \{b \in \mathbb{B} : |b| \leq c\}$ 
     $d \leftarrow p^*(x, c - i + 1) - p^*(x, c - i)$ 
     $s \leftarrow s + d$ 
    add a random bit to  $r$  until it is as long as  $s$ 
  if  $r < s$  then return  $x$ 

```

The reader may verify that this program dovetails computation of $p^*(x, i)$ for increasing i for all x ; the variable s contains the summed probability mass that has been encountered so far. Whenever s is incremented, mentally associate the interval $(s, s + d]$ with outcome x . Since $p^*(x, i)$ goes to $p(x)$ as i increases, the summed length of the intervals associated with x goes to $p(x)$ and s itself goes to $\overline{s} = \sum_x p(x)$. We can therefore sample from p by picking a number r that is uniformly random on $[0, 1]$ and returning the outcome associated with the interval containing r . Since s must have finite length (due to the construction of p^*), we only need to know r up to finite precision to be able to determine which interval it falls in; this allows us to generate r on the fly. The algorithm halts unless r falls in the interval $[\overline{s}, 1]$, which corresponds exactly to the deficiency of p : if p is a semimeasure, we expect the non-halting probability of a TM sampling it to correspond to $1 - \sum_x p(x)$. \square

Theorems 8 and 9 combined prove that the class of distributions sampled by Turing machines equals the lower semicomputable semimeasures (Lemma 1).

A.2 Domination of Model Class Supersets

Lemma 9. *Let C and D be model classes. If $C \subseteq D$, then m^D dominates m^C :*

$$\frac{m^D(x)}{m^C(x)} \geq \alpha$$

for some constant α independent of x .

Proof. We can partition the models of D into those belonging to C and the rest, which we'll call \overline{C} . For any given enumeration of D , we get $m^D(x) = \alpha m^C(x) + (1 - \alpha)m^{\overline{C}}(x)$. This gives us:

$$\frac{m^D(x)}{m^C(x)} = \alpha + (1 - \alpha) \frac{m^{\overline{C}}(x)}{m^C(x)} \geq \alpha.$$

\square

A.3 Unsafe Approximation of ID (Theorem 5)

Proof.

$$\begin{aligned}
 p_q(\text{ID}^C(x, y) - \text{ID}(x, y) \geq k) &= \\
 p_0(\max[\bar{\kappa}^C(x \mid f(x)), \bar{\kappa}^C(f(x) \mid x)] - \max[K(x \mid f(x)), K(f(x) \mid x)] \geq k) & . \\
 p_q(|x| - \text{ID}^C(x, y) \geq 2k) &\leq p_0(|x| - \bar{\kappa}^C(x \mid f(x)) \geq 2k) \\
 &\leq p_0(|x| - \kappa^C(x) \geq k \vee \kappa^C(x) - \bar{\kappa}^C(x \mid f(x)) \geq k) \\
 &\leq p_0(|x| - \kappa^C(x) \geq k) + p_0(\kappa^C(x) - \kappa^C(x \mid f(x)) \geq k) \leq 2^{-k} + cb^{-k} .
 \end{aligned}$$

K can invert $f(x)$, so

$$\text{ID}(x, y) = \max[K(x \mid f(x)), K(f(x) \mid x)] = \max[|f^*|, |f_{\text{inv}}^*|] < c_f$$

where f^* and f_{inv}^* are the shortest program to compute f on U and the shortest program to compute the inverse of f on U respectively.

$$\begin{aligned}
 p_q(\text{ID}^C(x, y) - \text{ID}(x, y) \geq k) + p_q(|x| - \text{ID}^C(x, y) \geq k) \\
 \geq p_q(\text{ID}^C(x, y) - \text{ID}(x, y) \geq k \vee |x| - \text{ID}^C(x, y) \geq k) \\
 \geq p_q(|x| - \text{ID}(x, y) \geq k) \geq p_0(|x| - c_f \geq k) = \sum_{i \geq k - c_f} s(i) .
 \end{aligned}$$

Which gives us:

$$\begin{aligned}
 p_q(\text{ID}^C(x, y) - \text{ID}(x, y) \geq k) \\
 \geq -p_q(|x| - \text{ID}^C \geq k) + \sum_{i \geq k - |f|} s(i) \geq -cb^{-k} + \sum_{i \geq k - |f|} s(i) \\
 \geq s(k - |f|) - cb^{-k} \geq c's(k) \quad \text{for the right } c' . \quad \square
 \end{aligned}$$

Corollary 1 *Under the assumptions of Theorem 5, $\bar{\kappa}^C(x \mid y)$ is an unsafe approximation for $K(x \mid y)$ against q .*

Proof. Assuming $\bar{\kappa}^C$ is safe, then since max is safety-preserving (Lemma 11), ID^C should be safe for ID. Since it isn't, $\bar{\kappa}^C$ cannot be safe. \square

A.4 Safe Approximation of ID (Theorem 6)

Lemma 10. *If q samples x and y independently from models in C , then $\kappa^C(x \mid y)$ is a 2-safe approximation of $-\log m(x \mid y)$ against q .*

Proof. Let q sample x from p_r and y from p_s .

$$\begin{aligned}
 p_q(-\log m^C(x \mid y) + \log m(x \mid y) \geq k) &= p_q(m(x \mid y)/m^C(x \mid y) \geq 2^k) \\
 &\leq 2^{-k} E[m(x \mid y)/m^C(x \mid y)] = 2^{-k} \sum_{x, y} p_s(y) m(x \mid y) \frac{p_r(x)}{m^C(x \mid y)} \\
 &\leq c2^{-k} \sum_{x, y} p_s(y) m(x \mid y) \frac{m^C(x \mid y)}{m^C(x \mid y)} \leq c2^{-k} \sum_{x, y} p_s(y) m(x \mid y) \leq c2^{-k} . \quad \square
 \end{aligned}$$

Since m and K mutually dominate, $-\log m^C$ is 2-safe for $K(x \mid y)$, as is $\bar{K}(x \mid y)$.

Lemma 11. *If f_a is safe for f against q , and g_a is safe for g against q , then $\max(f_a, g_a)$ is safe for $\max(f, g)$ against q .⁵*

Proof. We first partition \mathbb{B} into sets A_k and B_k :

$A_k = \{x : f_a(x) - f(x) \geq k \vee g_a - g(x) \geq k\}$ Since both f_a and g_a are safe, we know that $p_q(A_k)$ will be bounded above by the sum of two inverse exponentials in k , which from a given k_0 is itself bounded by an exponential in k .

$B_k = \{x : f_a(x) - f(x) < k \wedge g_a - g(x) < k\}$ We want to show that B contains no strings with error over k . If, for a given x the left and right max functions in $\max(f_a, g_a) - \max(f, g)$ select the outcome from matching functions, and the error is below k by definition. Assume then, that a different function is selected on each side. Without loss of generality, we can say that $\max(f_a, g_a) = f_a$ and $\max(f, g) = g$. This gives us: $\max(f_a, g_a) - \max(f, g) = f_a - g \leq f_a - f \leq k$.

We now have $p(B_k) = 0$ and $p(A_k) \leq cb^{-k}$, from which the theorem follows. \square

Corollary 2 ID^C is a safe approximation of ID against sources that sample x and y independently from models in C .

A.5 Safe Approximation of NID (Theorem 7)

Lemma 12. *Let f and g be two functions, with f_a and g_a their safe approximations against adversary p_q . Let $h(x) = f(x)/g(x)$ and $h_a(x) = f_a(x)/g_a(x)$. Let $c > 1$ and $0 < \epsilon \ll 1$ be constants such that $p_q(f_a(x) \geq c) \leq \epsilon$ and $p_q(g_a(x) \geq c) \leq \epsilon$. We can show that for some $b > 1$ and $c > 0$*

$$p_q \left(\left| \frac{h(x)}{h_a(x)} - 1 \right| \geq \frac{k}{c} \right) \leq cb^{-k} + 2\epsilon.$$

Proof. We will first prove the bound from above, using f_a 's safety, and then the bound from below using g_a 's safety.

$$\begin{aligned} p_q \left(\frac{h}{h_a} \leq 1 - \frac{k}{c} \right) &\leq p_q \left(\frac{h}{h_a} \leq 1 - \frac{k}{c} \ \& \ c < f_a \right) + \epsilon \leq p_q \left(\frac{h}{h_a} \leq 1 - \frac{k}{f_a} \right) + \epsilon \\ &= p_q \left(\frac{f}{f_a} \frac{g_a}{g} \leq 1 - \frac{k}{f_a} \right) + \epsilon \leq p_q \left(\frac{f}{f_a} \leq 1 - \frac{k}{f_a} \right) + \epsilon \\ &= p_q \left(\frac{f + k}{f_a} \leq 1 \right) + \epsilon = p_q(f_a - f \geq k) + \epsilon \leq c_f b_f^{-k} + \epsilon. \end{aligned}$$

The other bound we prove similarly. Combining the two, we get

$$p_q(h/h_a \notin (k/c - 1, k/c + 1)) \leq c_f b_f^{-k} + c_g b_g^{-k} + 2\epsilon \leq c' b'^{-k} + 2\epsilon. \quad \square$$

Theorem 7 follows as a corollary.

⁵ We will call such operations *safety preserving*.