# Kolmogorov Complexity and Distance Sets: Two Notions of Set Complexity

Elana Elman

April 26, 2024

### Abstract

I introduce Kolmogorov complexity and the Erdős distinct distance problem, describe an intuitive connection between these topics, and explore whether they are actually related. I construct sets in $\mathbb{R}^2$ with arbitrary Kolmogorov complexity and distance set size in order to show that the Kolmogorov complexity and the size of the distance set for nite sets of xed size are independent quantities. I consider what alternative descriptions of set complexity might agree more closely with my geometric intuition.

## 1    Introduction

I am interested in capturing the \structuredness" of a nite set of points. I study two notions of complexity, Kolmogorov complexity and the size of the distance set. Kolmogorov complexity is the length of the shortest program that outputs a given nite string; in the context of sets, I consider

De nition 2.3 (Big-O). Given two functions $f$ and $g$: $\mathbb{R} \to \mathbb{R}$, $f$ is said to be $\in O(g)$ or $= O(g)$ if $\exists C \in \mathbb{R}$ and $x_0 \in \mathbb{R}$ such that
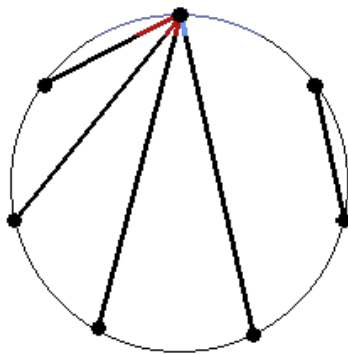
$$|f(x)| \le Cg(x) \,\& x > x_0:$$

Erd os' original paper [1] shows that $M(n)$ is between $\Omega(\sqrt[p]{n})$ and $O(\frac{n}{\sqrt{\log n}})$. The lower bound was gradually improved until 2015, when Guth and Katz proved $M(n)$ is at least $\Omega(\frac{n}{\log n})$.

### 2.0.1 Upper Bounds

**Example.** Let $P$ be $n$ random points in the unit square. It is overwhelmingly likely that all distances are unique and $|(P)| = \binom{n}{2} \in O(n^2)$.

**Example.** Let $P$ be $n$ points arranged evenly on a circle. Pick any of these points to call $a$: it is the same distance from each of its neighbors, and the same distance from the second point on its right as from the second point on its left, etc: all distances from $a$ to other points in $P$ come in pairs, except for the single distance from $a$ to the point directly opposite it if $n$ is even. By symmetry, any distance between non-$a$ points is the same as some distance involving $a$. So $|(P)| = \lfloor \frac{n}{2} \rfloor \in O(n)$.

- Even the green triangle contains some repeated distances. Any triangle where all three side lengths are integers$\leq \sqrt{n}$ is present within the green triangle, and any integer distance$\leq \sqrt{n}$ also occurs along the top edge of the grid, so the hypotenuse of such a triangle repeats a distance found along the top. The purple lines marked on Figure 2 are two different lines of length 5: one sits on the top of the grid and the other is the hypotenuse of a 3-4-5 triangle.

  Sets of integers $x, y, h$ with $x^2 + y^2 = h^2$ are called Pythagorean triples. Any Pythagorean triple is of the form

$$x = k(2ab)$$
$$y = k(a^2 - b^2)$$
$$h = k(a^2 + b^2)$$

  where $k$ is any positive integer and $a > b$ [2]. Restricting to $k = 1$, h is a repeated distance in the grid when $2ab \leq \sqrt{n}$ and $a^2 - b^2 \leq \sqrt{n}$. The number of integer grid points in this region is $O(\sqrt{n} \log n)$[1], so the number of distinct distances in the grid is at most $O(\sqrt{n} \log n)$.
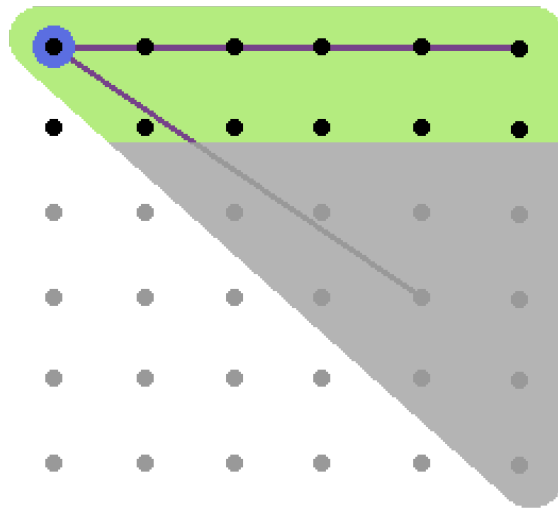


Figure 2: A 6 × 6 square grid. All distances outside the green triangle are also present in the green triangle. The green triangle contains some repeating distances: two lines of length 5 are drawn in purple.

### 2.0.2 Lower Bounds

Finding any lower bound on $M$ is less obvious. Here is Erdős's $\sqrt{n}$ bound:

*Proof.* Let $P$ be any set of $n$ points in the plane. Take any $z \in P$. For each other point $p$ in $P$, draw a circle through $p$ centered at $z$. Call the number of distinct circles drawn $m$.

- If $m < \sqrt{n}$, some circle must have at least

$$\frac{n}{m} \geq \sqrt{n}$$

points, so some hemisphere has at least $p$

| Time | Tape | Current State | Next State | Write | Move |
|------|------|---------------|------------|-------|------|
| 0 | :::B 101B::: | SEEK | SEEK | 1 | R |
| 1 | :::B 101B::: | SEEK | SEEK | O | R |
| 2 | :::B 101B::: | SEEK | SEEK | 1 | R |
| 3 | :::B 101B::: | SEEK | CARRY | B | L |
| 4 | :::B 101B::: | CARRY | CARRY | O | L |
| 5 | :::B 100B::: | CARRY | DONE | 1 | L |
| 6 | :::B 110B::: | DONE | | | |

Table 1: An example run of the increment machine on the input 101

Turing machines for any particular problem aren't unique. In fact, it's possible to translate any Turing machine into a machine with only two states by increasing the number of symbols, or into a machine with two symbols by increasing the number of states. However, the product of the number of states in symbols stays within a constant factor after either of these translations, so we can comfortably use this state symbol complexity to describe any Turing machine's complexity.

### 3.1.2 Programming Languages

Compared to Turing machines, real computers have additional features such as random-access memory and limitations such as finite memory. Turing machines are also difficult to design and understand. Programming languages allow simpler, human-readable descriptions of algorithms. Almost all programming language are Turing-complete, meaning they are capable of simulating Turing machines. A language is called Turing-equivalent if a Turing machine is capable of simulating running its programs, and nobody has ever found a Turing-complete language which is not Turing-equivalent. It seems extremely likely that no programming language can be more powerful than Turing machines.

For a Turing-equivalent language $L$, let $T_L$ be a Turing machine that simulates a program written in $L$, and let $L_T$ be a program in $L$ that simulates execution of a Turing machine. Then language $L_2$ can simulate a program written in language $L_1$ by simulating the execution of $T_{L_1}$. Since these programs and Turing machines are finite, there exist finite translations between Turing-equivalent languages. This shows that finite compilers { programs that translate Turing machines are finite, there exist finitely that

Definition 3.5 (Universal Functions). A function $f$ is universal for a set of functions $F$ if $f$ minorizes all functions in $F$.

If $f$ and $g$ are both universal for $F$, then their difference on any $x$ is bounded by a constant, and any non-universal function $h$ in $F$ is no more than a constant less than $f$. This means we can consider only optimal descriptions relative to universal functions without worrying that other functions provide shorter descriptions.

Definition 3.6 (Computable Functions). A partial function $f : \mathbb{N} \to \mathbb{N}$ is called computable if there is some Turing machine which terminates with output $f(n)$ on each input $n$ for which $f$ is defined.

Theorem 1 (A universal computable function exists). *Proof.* Let U be a Turing machine which simulates other Turing machines. U must take two pieces of information, a description of the machine to simulate and the input string to simulate running. U expects this input in the format $1^{l(g)} 0 g p$ this is $l(g)$ 1s, then one 0, then the literal description $g$ of the desired Turing machine, then the literal input string $p$. U may then use the prefix $1^{l(g)} 0$ to separate $g$ from $p$ and run its simulation. Let $f_U$ be the function executed by U. $\square$

Definition 3.7 (Kolmogorov Complexity). The Kolmogorov complexity of a finite string $s$ is defined as $C(s) = C_{f_U}(s)$. We require that the program doesn't take arguments, or equivalently that any input to the program counts toward its length.

Theorem 2. Most strings are incompressible.

*Proof.* The set of $n$-length binary strings has $2^n$ elements, and the set of shorter strings has $2^n$ elements, so there are at most $2^n - 1$ different outputs of shorter programs. Even if we optimistically assume that each of these outputs has length at least $n$, $2^n - 2^{n-1} = 2^{n-1}$ of the $2^n$ strings with length $n$ cannot be more concisely represented. $\square$

Theorem 3. Kolmogorov complexity is uncomputable.

*Proof.* Suppose a program K with length $a$ takes as input a finite string $s$ and returns its Kolmogorov complexity. Note that some shortest program always exists because it is at longest $\log s$. Write a new function C which decides if a string $s$ is compressible:

```
from math import log
def C(S):
    if K(s) >= log(s, 2):
        return false
    return true
```

This program is only a constant number of bits longer than $C$: $B$ is $a + c_C + c_B$ bits long.

Choose $min = 2^{a + c_C + c_B}$. Now $B(min)$ returns a number $b$ with $K(b) > 2^{a + c_C + c_B}$. On the other hand, we just saw that $b$ was generated by the program $B(min)$, which, including the length of its argument, was at most $a + c_C + c_B + \log 2^{a + c_C + c_B} = 2(a + c_C + c_B)$. This is a description of $b$ which is shorter than than $K(b)$, the shortest possible description. This is a contradiction, so the $K$ function cannot actually exist.

$\square$

The above definitions and facts are from Li and Vitanyi's book on Kolmogorov complexity [3], except for the proof of uncomputability, for which I referenced Peter Miltersen's course notes [5]. The below definition is my own, for purposes of comparison with the Erdős distance problem:

**Definition 3.8** (Kolmogorov Complexity of Sets) Let $A$ be a finite set of integers. Index by $a_i$ for $i$ from 0 through $|A|$, where $a_1$ is the smallest value in $A$, $a_2$ is the next smallest, etc. Define a string representation of $A$ to be the string $S = "a_1, a_2, \ldots, a_{|A|}"$ with each $a_i$ replaced by its literal value. Define $C(A)$ to be $C(S)$.

Let $B$ be a finite set of pairs of integers. Index these pairs by $(a_{i1}, a_{i2})$ (for $i$ from 1 through $|B|$, where the $a_i$s in dictionary order. Define a string representation of $B$ to be the string $S = "(a_{11}, a_{12}) \n (a_{21}, a_{22}) \n \ldots (a_{n.c[(n)]TJ/F27\ 673iw312}$

Of course a more time-efficient sort is possible, but this Python program lets us more easily picture a corresponding Turing machine.

Appending programs comes at a length cost of log of the shorter program length, which is at most a constant for this program.

Running this sorting program after the shorter program selected by $C_b$ gives us a program for $S$ with the order $a$ that has length $C_b + c$ for some constant $c$ independent of $S$. Since $C_a$ is defined to be the length of the shortest program producing $S$ with order $a_i$, $C_a$ can't be larger than $C_b + c$. So we can see that printing any order of the elements of $S$ is a problem in the same class of Kolmogorov complexity as $C_a$.

$\square$

### 3.1.4 Kolmogorov Complexity versus Algorithm Complexity

There are many ways other ways of measuring how "hard" a problem is.

1. The usual metric of interest is time complexity: as the size of the input increases, roughly how many CPU cycles does the program take to run? The problem of factoring large primes is hard in this sense. While time complexity is defined for specific programs, it is also commonly used to describe the best known solution to a problem.

2. Another important metric is space complexity: as the size of the input increases, roughly how much active memory does the program require? Working with adjacency list representations of large matrices is hard in this sense.

3. The complexity of a particular program is sometimes described by how many branching points (conditional jumps) it has. High complexity in this sense indicates that a program is hard to maintain and debug.

4. Mathematicians and programmers are frequently interested in the difficulty of coming up with any solution at all to a problem, informally measuring complexity by years left unsolved.

Kolmogorov complexity is independent of all of the above metrics (note that Kolmogorov complexity does not count memory used during computation, so it is not the same as space complexity).

Kolmogorov complexity is not frequently used in the field of computer science, possibly because it is inconvenient. While any program gives an upper bound for the Kolmogorov complexity of the problem it solves, finding the true value is generally impossible. In addition, computers thankfully have enough memory these days that program size isn't much of a limitation. Finally, it seems to me that programmers simply refuse to write substantial programs which grow linearly with their inputs.

While Kolmogorov complexity doesn't have much influence on concrete programs, it has value as an abstract measure of problem complexity.

**Definition 3.9** (Computable Numbers) Computable numbers are real numbers which a Turing machine can approximate to any desired precision. All rational numbers and some irrationals, such as $\pi$, are computable. However, the computable numbers are countable because the set of Turing machines is countable, so most real numbers are uncomputable.

Because I am comparing Kolmogorov complexity to distance set size, approximations of real numbers are not precise enough for my purposes. I want to consider only numbers which are precisely finitely representable. These include the natural numbers and numbers $\frac{k}{2} \in \mathbb{N}$ because their

These two representations of numbers have analogs in actual computers, which store numbers in either of the forms $k \cdot 2^l : k, l \in \mathbb{Z}$ or $\frac{k}{m}$ for a fixed large natural $m$ to provide both a large range of values alongside a good density of small numbers.

## 4 An Intuitive Relationship?

**Proof.** If $n = 5$, the program is 55 characters long, which translates to 55 bytes or 440 bits. The program length increases by one character (or one byte or eight bits) for each extra digit in n, so the length of the whole program is at least the number of digits in n. The rest of the program runs correctly without modification for any n, so the remaining 54 bits of the program are constant. So, the length of this program for arbitrary n is $54 + \log_{10} n$. The Kolmogorov complexity of printing S is at worst this length: $C(S) \leq 54 + \log_{10} n \in O(\log n)$. $\square$

## 5.2   No upper bound on Kolmogorov complexity from distance set size

Let $s$ be an incompressible string of $n$ bits. Divide $s$ up into $\sqrt{n}$ segments: $s_1$ is the first $\sqrt{n}$ bits, $s_2$ is the next $\sqrt{n}$ bits, etc.

For each string $s_i$, define a set $S_i \subseteq \mathbb{N} \setminus [1; \sqrt{n}]$ such that $k$ is in $S_i$ exactly if the $k$th bit in $s_i$ is 1. Then take $S = \{f \times S_i : f \neq i, g\}$. S is an incompressible random subset of the integers $(\mathbb{N} \cap [1; \sqrt{n}])^2$, since if $S$ was compressible we could reverse this construction and compress $s$. So any program for $S$ requires at least $n$ bits. This is much larger than the $\log n$ bits required to encode the complete $\sqrt{n} \times \sqrt{n}$ integer grid.

On the other hand, $S$ is a subset of the $\sqrt{n} \times \sqrt{n}$ integer grid, Erdős' original bound says that $S$ has at most $O(\frac{n}{\sqrt{\log n}})$ distances.

# 6   Sets of arbitrary Kolmogorov complexity and distance set size exist

In this section I will construct sets with minimal distance sets and arbitrary complexity, and sets with small complexity but maximal distance sets.

## 6.1   Sets of arbitrary complexity exist with $|j| \in O(n)$

There is a program which prints $(i, 0)$ for all but $m$ points, and either $(i, 0)$ or $(i, 0)$ at random for the remaining $m$ points. This program has a minimum program length $O(\log n)$ bits.

For $n$ points, suppose we want an arrangement which has complexity $O(m)$, with $0 \leq m \leq n$. Place all n points along the x-axis. Assign the first $m$ points each the y-coordinate 0 or 1 at random, and assign all the remaining points the y-coordinate 0. Since the first $m$ points have random y-coordinates which take 1 bit each to describe, any program returning this set must include at least $m$ bits.

Here is a Python program which prints this set:

```python
rand = [b1, b2, ..., bm]
for b in rand:
    print(f"({b},0)")
for i in range(n-m):
    print("(0,0)")
```

For each of the random coordinates $b_i$, we decide arbitrarily whether $b_i$ in the program is literally 1 or 0.

Since most strings of any given length are incompressible, we can choose rand incompressible, which means there is no more efficient way to remember which points are shifted. Then we can't avoid spending at least $m$ bits.

Surprisingly, this set has very few distances. All points lie on an integer grid between $(0, 0)$ and $(n; 1)$, so its distance set is a subset of this grid section's distance set. The possible distances are those between points with $y = 1$, those between points with $y = 0$, and those between a point with $y = 1$ and a point with $y = 0$:

its derivative with respect to $n$ is

From the other direction, Kolmogorov complexity is closely related to Hausdorff dimension, and a few researchers have developed a constructive Hausdorff dimension which might serve as a more useful description of complexity [4].

Either an energy-based metric or constructive Hausdorff dimension seem likely to work for infinite sequences and real numbers.

## 10    Acknowledgements