

# Two Notions of Set Complexity

Elana Elman

April 16, 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, and I show that the Kolmogorov complexity and the size of the distance set for finite sets of fixed 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 idea of how structured a finite set of points is. 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 finite string; in the context of sets, I consider the shortest program which lists the coordinates of all the set's points. The distance set size, taken from Erdős' distinct distances problem, counts the number of unique distances between points in a set. Each of these intuitively seems to measure how regular a set is. However, I find that they are unrelated concepts: knowing a set's Kolmogorov complexity gives no information about its distance set size, and vice versa.

## 2 The Erdős Distinct Distance Problem

**Definition 2.1** ( $\Delta(P)$ ). For a set  $S \in \mathbb{R}^2$ , the distance set  $\Delta(P)$  is the set of unique distances between points in  $P$ .

Note that  $|\Delta(P)|$  is the number of distinct distances between points in  $P$ .

Erdős' distinct distance problem asks the smallest possible number of distinct distances in any arrangement of  $n$  points. For example, three points may be arranged into a regular triangle, producing one unique distance, but four points can at best be arranged into a square, producing two unique distances (the square's side length and its diagonal length).

Let  $\mathbf{P}_n$  be the collection of all sets of  $n$  different points in  $\mathbb{R}^2$ .

**Definition 2.2** ( $M(n)$ ). For a positive integer  $n$ , let  $M(n)$  be the smallest possible number of distinct distances in any arrangement of  $n$  different points on the plane:

$$M(n) = \min_{P \in \mathbf{P}_n} |\Delta(P)|$$

The precise value of  $M(n)$  is only known for very small  $n$ , and work on the distinct distance problem seeks to bound  $M$  asymptotically:

**Definition 2.3** (Big-O). Given two functions  $f$  and  $g: \mathbb{R} \mapsto \mathbb{R}$ ,  $f$  is said to be  $\in O(g)$  or  $= O(g)$  if  $\exists C \in \mathbb{R}$  and  $x_0 \in \mathbb{R}$  such that

$$|f(x)| \leq Cg(x) \forall x > x_0.$$

Erdős' original paper [1] shows that  $M(n)$  is between  $O(\sqrt{n})$  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  $O(\frac{n}{\log n})$ .

### 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  $|\Delta(P)| = \binom{n}{2} \in O(n^2)$ .

**Example.** Let  $P$  be  $n$  points arranged evenly on a circle. Pick any of these points to call  $a$ .  $a$  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  $|\Delta(P)| = \lfloor \frac{n}{2} \rfloor \in O(n)$ .

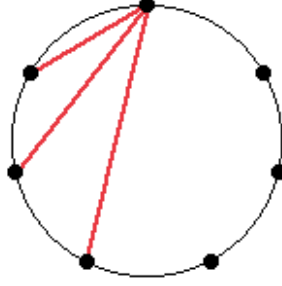


Figure 1: 7 points on a circle. One instance of each of this set's distances is drawn in red.

**Example** (Erdős' upper bound). For the sake of simplicity, assume  $n$  is a square integer. Let  $P$  be  $n$  points arranged on a  $\sqrt{n} \times \sqrt{n}$  integer grid, with each coordinate running from 0 through  $\sqrt{n} - 1$ . Then

$$|\Delta(P)| = \left| \left\{ \sqrt{(x_{p_i} - x_{p_j})^2 + (y_{p_i} - y_{p_j})^2} \right\} \right| = \left| \left\{ (x_{p_i} - x_{p_j})^2 + (y_{p_i} - y_{p_j})^2 \right\} \right|$$

$$\leq \left| \left\{ a^2 + b^2 \mid a, b \in \mathbb{Z}_+ \cap [0, \sqrt{n} - 1] \right\} \right| \leq \left| \left\{ a^2 + b^2 \mid a, b \in \mathbb{Z}_+, a^2 + b^2 \leq 2n \right\} \right| \leq \frac{bn}{\sqrt{\log n}} \in O\left(\frac{n}{\sqrt{\log n}}\right)$$

The last inequality is a limit by Landau.

### 2.0.2 Lower Bounds

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

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

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

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

points, so some hemisphere has at least  $\frac{\sqrt{n}}{2}$  points. The distance from the leftmost point on this hemisphere to each other point on the hemisphere is different, so there are at least  $\frac{\sqrt{n}}{2} \in O(\sqrt{n})$  distances in  $P$ .

- Otherwise  $m \geq \sqrt{n}$ , and each circle represents a distance from  $a$  to another point of  $p$ , so there are at least  $\frac{\sqrt{n}}{2} \in O(\sqrt{n})$  distances in  $P$ .

□

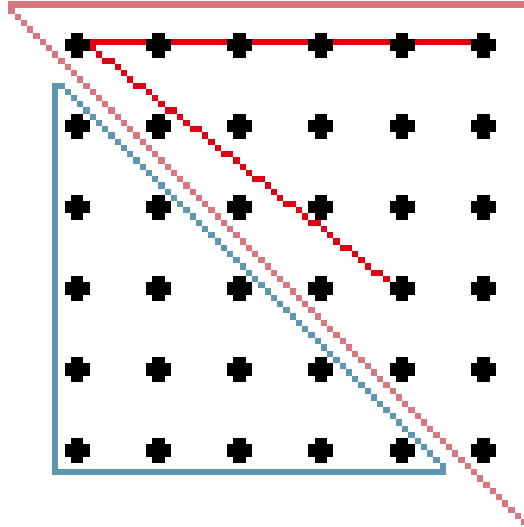


Figure 2: A  $6 \times 6$  square grid. All distances in the bottom triangle, marked in blue, are also present in the top triangle, marked in red. The red triangle contains some repeating distances: two lines of length 5 are drawn in bright red.

### 3 Kolmogorov Complexity

Kolmogorov complexity is an idea from information theory designed to measure randomness. Randomness is usually a property ascribed to generating processes, not individual objects, but Kolmogorov complexity allows us to argue that the string “000000000000” is less random than “128683487345”. The Kolmogorov complexity of a string is the length of the shortest program that outputs it. A precise definition requires defining a “program” and leaves us with classes of complexity instead of exact bit lengths.

#### 3.1 Formalizing Computation

##### 3.1.1 Turing Machines

**Definition 3.1** (Turing machines). A Turing machine is a finite set of states with a finite collection of transitions between states which accesses an infinite tape through a pointer. On transitions, the machine may overwrite the currently accessible tape cell with a new symbol and/or move the pointer right or left. The Turing machine must also specify a set of acceptable states to end in and an input string composed of finitely many non-blank characters from a finite alphabet which is written on the tape.

Any Turing machine may be finitely described by a tuple  $\langle \Sigma, b, S, s_0, \delta, F \rangle$ , where  $\Sigma$  is the alphabet of the input string,  $b \in \Sigma$  is the blank character,  $S$  is the set of states,  $s_0$  is the starting state,  $\delta$  is the transition function  $(S - F) \times \Sigma \rightarrow S \times \Sigma \times \{L, R, N\}$ , and  $F$  is any subset of  $S$ .  $\delta$  may be a partial function, in which case a global state with no successor halts.

**Example** (Increment). Let  $\Sigma = \{0, 1, B\}$ ,  $b = B$ ,  $S = \{\text{SEEK}, \text{CARRY}, \text{DONE}\}$ ,  $s_0 = \text{SEEK}$ ,  $F = \{\text{DONE}\}$ , and  $\delta$  is the function

Intuitively, this Turing machine moves to the right side of the input (**SEEK**), then flips bits from right to left until it’s able to put a 1 in a previously-0 cell (**CARRY**). We can now (inefficiently) add  $n$  to positive integers by writing  $n$  1s to the left of our input and erasing one per run of this increment routine. Positive integer addition can be (extremely inefficiently) built into integer multiplication, and

(state, symbol)	next state	write	move
(SEEK, 0)	SEEK	0	R
(SEEK, 1)	SEEK	1	R
(SEEK, B)	CARRY	B	L
(CARRY, 0)	DONE	1	N
(CARRY, 1)	CARRY	0	L
(CARRY, B)	DONE	1	N

so on; this capacity for arithmetic plus the recursive structure provided by state transitions allows for arbitrary computation. Notably, there exists a Turing machine which takes as input descriptions of another Turing machines and simulates their execution, returning the result the specified Turing machine would return.

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  $\times$  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  $L_1$  can simulate a program written in language  $L_2$  by simulating the execution of  $T_{L_2}$ . Since these programs and Turing machines are finite, there exist finite translations between Turing-equivalent languages. This shows that finite compilers – programs that translate between languages – exist. We may run a program in a different language by appending a compiler to the program. This adds a constant length to the program which depends on the two languages in use.

For the rest of this paper I will write programs in Python, but they could be written in any other language with only a constant penalty to length.

### 3.1.3 Defining Kolmogorov Complexity

**Definition 3.2.** The length  $l(p)$  of a natural number  $p$  is the number of digits in its binary representation.

**Definition 3.3** (Kolmogorov complexity with respect to a specifying function). Let  $S$  be a countable set and enumerate it with natural numbers  $n(x)$ . The Kolmogorov complexity of  $x \in S$  with respect to the partial function  $f$  on natural numbers is

$$C_f(x) = \min_{p \in \mathbb{N}} (\{l(p) \mid f(p) = n(x)\} \cup \{\infty\}).$$

In terms of computing,  $S$  is a set of outputs,  $f$  is a programming language,  $p$  are programs, and  $C_f$  reports the length of the shortest program that returns  $x$  when interpreted in language  $f$ .  $C_f$  returns infinity if no program in  $f$  returns  $x$ .

**Definition 3.4** (Minorizing functions). A function  $f$  minorizes a function  $g$  if there is a constant  $c$  such that  $C_f(x) \leq C_g(x) + c \forall x \in S$ .

**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  $f : \mathbb{N} \rightarrow \mathbb{N}$  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.** Kolmogorov complexity is uncomputable.

*Proof.* Suppose a Turing machine  $T$  takes as input a Turing machine  $K$  [Proof in progress, but this is roughly the halting problem...] □

**Theorem 2** (A universal computable function  $f_U$  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 on.  $U$  expects this input in the format  $11 \dots 10gp$ : 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 of 1s to separate  $g$  from  $p$  and run its simulation. Let  $f_U$  be the function executed by  $U$ . □

**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.

The above definitions and facts are from Li and Vitányi's book on Kolmogorov complexity [2]. 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  $A$  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  $S = "a_1, a_2, \dots, a_{|A|}"$  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  $S = "(a_{11}, a_{12}) \backslash \backslash (a_{21}, a_{22}) \backslash \backslash \dots (a_{n1}, a_{n2})"$  where each  $(a_{i1}, a_{i2})$  is replaced by its literal values. Define  $C(B)$  to be  $C(S)$ .

**Theorem 3.** The Kolmogorov complexity of a set  $S$  does not depend on the order  $a_i$ .

*Proof.* Suppose that any set  $S \in \mathbb{R}^2$  with  $|S| = n$  has an order  $b_i$  such that  $C(S)$  with respect to  $b_i$  (call this value  $C_b$ ) is less than  $C(S)$  with respect to  $a_i$  (call this value  $C_a$ ).

If a list of numbers is already in the computer's memory at `list1`, sorting it is possible with a constant-length addition to the program:

```
list2 = []
while list1: #this means "while list1 is not empty:"
    minimum = list1[0]
    for k in list1:
        if (k[0] < minimum[0] or
            (k[0] = minimum[0] and k[1] < minimum[1])):
            minimum = k
    list1.remove(k)
```

```
list2.append(k)
```

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  $a_i$  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  $C_a$ . □

To extend the above definitions and facts to  $\mathbb{Q}^2$ , store rationals as pairs of integers delimited by a new character in  $\Sigma$  and modify the comparison in Theorem 3's program to be the dictionary comparison for 4-tuples of integers.

### 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 in terms of 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  $e$ , 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  $\frac{k}{2^l}$  for  $k, l \in \mathbb{N}$  because their

binary representations are finite. Any other rational number may be represented as fractions by delimiting pairs of integers with a new tape symbol, since we can perform arithmetic operations perfectly well on fractions.

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.

### 3.1.5 Kolmogorov Complexity versus Randomness

**Theorem 4.** Most strings are incompressible.

*Proof.* The set of  $n$ -length binary strings has  $2^n$  elements, and the set of shorter strings has  $2^{n-1}$  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  $n$ , at least  $2^n - 2^{n-1} = 2^{n-1}$  of the  $2^n$  strings with length  $n$  cannot be more concisely represented.  $\square$

## 4 Kolmogorov Complexity and Distance Set Size are Independent

### 4.1 No upper bound on distance set size from Kolmogorov Complexity

Given an integer  $n$ , consider the following Python program:

```
for i in range(0, n):
    x = 2**i
    print(f"({x}, 0)\n")
```

After substituting in the actual integer for  $n$ , this program prints a list of  $n$  points  $(2^m, 0)$  for  $m \in \mathbb{N}, m < n$ .

The set  $S = \{(2^m, 0) | m \in \mathbb{N}, m < n\}$  has no duplicate distances between its points:

*Proof.* Of course a set of one point has no duplicate distances.

Suppose the set  $S^{m-1} = \{(2^k, 0) | k \in \mathbb{N}, k < m - 1\}$  has no duplicate distances.  $(2^m, 0)$  is  $2^{m-1}$  away from  $(2^{m-1}, 0)$  and further from each other point in  $S^{m-1}$ . Since the maximum distance within  $S^{m-1}$  is  $2^{m-1} - 1$ , all distances between  $(2^m, 0)$  and points in  $S^{m-1}$  are greater than all distances within  $S^{m-1}$ . Also,  $(2^m, 0)$  is a distinct distance  $2^m - 2^k, k < m$  from each point in  $S^{m-1}$ . The distances in  $S^m = \{(2^k, 0) | k \in \mathbb{N}, k < m\}$  are the distances in  $S^{m-1}$  and the distances between  $(2^m, 0)$  and points in  $S^{m-1}$ , and since there are no repeated distances within or between these groups,  $S^m$  has no duplicate distances.  $\square$

So  $S$  has the maximum possible number of distances in any set of  $n$  points in  $\mathbb{R}^2$ , which is  $\binom{n}{2} \in O(n^2)$ .

On the other hand, the program above has the minimum possible length for any program containing the number  $n$ , so the problem of listing the points in  $S$  has very low Kolmogorov complexity:

*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:  $C(S) \leq 54 + \log_{10} n \in O(\log n)$ .  $\square$

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

Let  $S_1$  be the  $\sqrt{n} \times \sqrt{n}$  integer grid, and let  $S_2 \subset S_1$  be an incompressible subset of  $S_1$ . Then  $|\Delta(S_2)| < |\Delta(S_1)| \in O(\frac{n}{\sqrt{\log n}})$ , but a description of  $S_2$  requires at least  $n$  bits.

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

### 5.1 Sets of arbitrary complexity exist with $|\Delta| \in O(n)$

There is a program which prints  $(0,0)$  for all but  $m$  points, and either  $(1,0)$  or  $(0,0)$  at random for the remaining  $m$  points. This program has a minimum program length of  $O(\log m)$  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  $\lfloor m \rfloor$  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  $n$  bits.

Here is a Python program which prints this set:

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

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

What stops us from picking a sequence like  $(0,0,0,\dots,0)$ , for which our program's list of bits is not the most efficient way to describe this set? Few sequences of  $n$  bits are describable in much fewer than  $n$  bits: since there are  $2^n$  sequences of  $n$  bits and only  $2^{n-1}$  sequences of less than  $n$  bits, if more than half of the  $n$ -length sequences were perfectly compressible we'd run out of shorter-length strings to compress them into. So most  $n$ -length sequences are incompressible, which means it's possible to pick coordinates such that we can't avoid spending a bit per  $b_i$ .

So some sequence of  $m$  arbitrary coordinates results in a set which cannot be described in less than  $m$  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  $x = 1$ , those between points with  $x = 0$ , and those between a point with  $x = 1$  and a point with  $x = 0$ :

- Between points which both have  $x = 0$ , only integer distances  $\leq n$  are possible. The same is true between points which each have  $x = 1$ . Each of these cases produce at most  $n$  distinct distances.
- Any distance between a point with  $x = 0$  and a point with  $x = 1$  is the same as the distance from  $(0,0)$  to some point with  $x = 1$ , so this case also produces at most  $n$  distinct distances.

So  $|\Delta|$  here is at most  $3n \in O(n)$ .

#### 5.1.1 Sets of arbitrary complexity exist with $|\Delta| \in O(\frac{n}{\sqrt{\log n}})$

Place  $n$  points on an  $\sqrt{n} \times \sqrt{n}$  grid. Then decide arbitrarily whether to shift each of the first  $m$  points right by  $\sqrt{n}$ . The total number of distances in this arrangement is no more than the number of distances in a  $2\sqrt{n} \times 2\sqrt{n}$  grid which has  $4n$  points and  $O(\frac{4n}{\sqrt{\log 4n}}) \leq O(\frac{4n}{\sqrt{\log n}}) = O(\frac{n}{\sqrt{\log n}})$  distances.

A program for this set is



```

import math
rand = [t1, t2, ..., tm]
s = math.sqrt(n)

index = 0
for i in range(1, s):
    for j in range(1, s):
        if index < m:
            if t[index]:
                print(f"({i+s},{j})")
            else:
                print(f"({i},{j})")
            index += 1
        else:
            print(f"({i},{j})")

```

Where each `ti` is a boolean representing whether its corresponding point is shifted.

As in the linear case, most bit strings of length  $m$  are incompressible, so some sequence of `tis` gives a set with a complexity of at least  $m$  bits. We can also see from this program that every sequence of `tis` gives a set with a complexity of at most  $O(m)$ , since the whole program is constant with respect to  $m$  except for the length of `rand`.

## 5.2 Sets with arbitrary $|\Delta|$ exist with complexity $O(\log n)$

For any  $n \in \mathbb{N}$ , suppose we want an arrangement of  $n$  points with  $O(n^k)$  distances for  $k \in [1, 2]$ .

Let  $M \subset \mathbb{R}^2$  be  $m$  of the  $n$  points evenly-spaced around a circle of radius  $\frac{1}{4}$  centered at the origin so that there are  $\frac{m}{2}$  distinct distances in  $M$ . Let  $K \subset \mathbb{R}^2$  be the remaining  $n - m$  points placed at  $(2^i, 0)$  for  $i \in \mathbb{N}, i < m - n$ , so that there are  $\binom{n-m}{2}$  distinct distances in  $K$ .

How many distinct distances exist between points in  $M$  and points in  $K$ ? There are at most  $m \cdot (n - m)$  distances between these sets if there are no repeats.

Consider the distance  $d$  between some  $a \in M$  and some  $b \in K$ . Points in  $M$  are all on a circle, the set of points in  $\mathbb{R}^2$  which are  $d$  away from  $b$  form a circle, and there are at most two points of intersection between two distinct circles, so at most one other point in  $M$  is  $d$  away from  $b$ . (We know the circles are distinct because the largest x-coordinate among points in  $M$  is  $\frac{1}{4}$  and the smallest x-coordinate among points in  $K$  is 1, so the circles are centered at different points.) If more than half the distances between  $b$  and points in  $M$  are duplicates, then there are fewer than  $\frac{n}{2}$  distinct distances and some group of more than two points are the same distance from  $b$ , so at most half of the distances between  $b$  and points in  $M$  are distinct.

For two different points  $b$  and  $b'$  in  $K$ , if  $d = |b - a|$  for some  $a \in M$  then  $d \neq |b' - a'|$  for any  $a' \in M$ :

WLOG assume  $b < b'$ .

$$|a| \leq \frac{1}{4}, |a'| \leq \frac{1}{4}$$

$$b_x = |b| \geq 1$$

$$|a - b| = \sqrt{(a_x - b_x)^2 + a_y^2}$$

$$a_x \in [-\frac{1}{4}, \frac{1}{4}] \text{ and } b_x \geq 1, \text{ so } |a_x| < |b_x| = b_x, \text{ so } |a_x - b_x| \leq b_x + \frac{1}{4} \text{ and } \geq b_x - \frac{1}{4}:$$

$$|a - b| \in [\sqrt{(b_x - \frac{1}{4})^2 + a_y^2}, \sqrt{(b_x + \frac{1}{4})^2 + a_y^2}]$$

and  $|a_y| \in [0, 1/4]$ , so

$$|a - b| \in [\sqrt{(b_x - \frac{1}{4})^2}, \sqrt{(b_x + \frac{1}{4})^2 + (\frac{1}{4})^2}] = [b_x - \frac{1}{4}, \sqrt{(b_x + \frac{1}{4})^2 + \frac{1}{16}}].$$

By the same argument,  $|a' - b'| \in [b'_x - \frac{1}{4}, \sqrt{(b'_x + \frac{1}{4})^2 + \frac{1}{16}}]$ . We know that  $b_x + 1 \leq b'_x$ , so  $b'_x - \frac{1}{4} > b_x + \frac{3}{4}$ . Now if  $b_x + \frac{3}{4} \geq \sqrt{(b_x + \frac{1}{4})^2 + \frac{1}{16}}$ , it follows that  $b'_x - \frac{1}{4} > \sqrt{(b_x + \frac{1}{4})^2 + \frac{1}{16}}$ :

$$\begin{aligned} b_x + \frac{3}{4} &\geq \sqrt{(b_x + \frac{1}{4})^2 + \frac{1}{16}} \iff \\ (b_x + \frac{3}{4})^2 &\geq (b_x + \frac{1}{4})^2 + \frac{1}{16} \iff \\ b_x^2 + \frac{3}{2}b_x + \frac{9}{16} &\geq b_x^2 + \frac{b_x}{2} + \frac{1}{16} + \frac{1}{16} \iff \\ b_x &\geq -\frac{7}{16} \end{aligned}$$

Which is true because  $b_x = |b| \geq 0$ . Now  $|a-b| \leq \sqrt{(b_x + \frac{1}{4})^2 + \frac{1}{16}} < b'_x - \frac{1}{4} \leq |a'-b'|$  so  $|a-b| \neq |a'-b'|$ .

This shows that no two points in  $K$  repeat distances with any points in  $M$ . So there are no duplicate distances between points in  $M$  and points in  $K$ , and the number of distinct distances between points in  $M$  and points in  $K$  is  $m \cdot (n - m)$ .

Now the total number of distinct distances  $|\Delta|$  in  $M \cup K$  is

$$\begin{aligned} &\frac{m}{2} + \binom{n-m}{2} + m \cdot (n-m) \\ &= \frac{m}{2} + \frac{(n-m)(n-m-1)}{2} + mn - m^2 \\ &= \frac{m}{2} + \frac{n^2 - 2mn - n + m^2 + m}{2} + mn - m^2 \\ &= \frac{n^2}{2} - \frac{n}{2} + m - \frac{m^2}{2} \end{aligned}$$

And we want  $|\Delta| = n^k$ :

$$\begin{aligned} \frac{n^2}{2} - \frac{n}{2} + m - \frac{m^2}{2} &= n^k \\ \frac{m^2}{2} - m - (\frac{n^2}{2} - \frac{n}{2} - n^k) &= 0 \\ m &= 1 \pm \sqrt{1 + n^2 - n - 2n^k} \end{aligned}$$

Which exists for large enough  $n$  since  $O(n^k) \in O(n^2)$ .

Since  $0 \leq m$ , we can't actually have  $m = 1 - \sqrt{1 + n^2 - n - 2n^k}$ , since it will become negative if  $n$  is large.

Also, since  $m \leq n$ , we need  $n \geq 1 + \sqrt{1 + n^2 - n - 2n^k}$ :

$$\begin{aligned} (n-1)^2 &\geq n^2 - 2n^k - n + 1 \\ n^2 - 2n + 1 &\geq n^2 - 2n^k - n + 1 \\ 2n^k - n &\geq 0 \end{aligned}$$

Which is true for all  $k > 1$ .

So,  $|\Delta| = n^k$  for  $1 \leq k < 2$  when  $m = 1 + \sqrt{n^2 - 2n^k - n + 1}$ .

Unfortunately  $m$  might not be an integer. Note that the derivative of  $|\Delta|$  with respect to  $m$  is  $1 - m$ , so  $|\Delta|$  is decreasing for  $m \in [1, n]$ . Any  $m$  is within  $\frac{1}{2}$  of some integer  $l$ , and  $|\Delta|_{m-\frac{1}{2}} \geq |\Delta|_{\lfloor m \rfloor} \geq |\Delta|_m = n^k \geq |\Delta|_{\lceil m \rceil} \geq |\Delta|_{m+\frac{1}{2}}$  for  $m > \frac{3}{2}$ , so there is some integer  $l$  with  $|\Delta|_l - n^k \in [|\Delta|_{m-\frac{1}{2}} - n^k, |\Delta|_{m+\frac{1}{2}} - n^k] = [\frac{n^2-n}{2} + m - \frac{1}{2} - \frac{m^2-m+\frac{1}{4}}{2} - (\frac{n^2-n}{2} + m - \frac{m^2}{2}), \frac{n^2-n}{2} + m + \frac{1}{2} - \frac{m^2+m+\frac{1}{4}}{2} - (\frac{n^2-n}{2} + m - \frac{m^2}{2})] = [-\frac{1}{2} - \frac{-m+\frac{1}{4}}{2}, \frac{1}{2} - \frac{m+\frac{1}{4}}{2}]$  where both bounds are  $O(n)$ . So  $|\Delta|_l \in O(n^k) \pm O(n) \in O(n^k)$ .

## 6 Capturing Set Complexity

Describing the Kolmogorov complexity of an arbitrary set of real numbers is impossible because of most real numbers' uncomputability. Beyond that, relating distance set size to Kolmogorov complexity seems to fail because Kolmogorov complexity is too sensitive to randomness. When choosing random points on a finite integer grid, distance set size says that losing a couple of points can't increase complexity regardless of which points are chosen, while Kolmogorov complexity suggests that (by order of magnitude) "forgetting" a point on a regular grid is as difficult as adding in an arbitrary rational point.

Kolmogorov also doesn't capture symmetry, which feels like a relevant part of complexity: the set of points  $2^i$  for  $0 \leq i < n$  has the same Kolmogorov complexity as the set of points  $i$  for  $0 \leq i < n$ .

When discussing set complexity, I would like to have the property that subsets are no more complicated than their supersets. [...]

## References

- [1] P. Erdos. "On Sets of Distances of n Points". In: *The American Mathematical Monthly* 53.5 (1946), pp. 248–250. ISSN: 00029890, 19300972. URL: <http://www.jstor.org/stable/2305092> (visited on 04/15/2024).
- [2] Ming Li and Paul M.B. Vitnyi. *An Introduction to Kolmogorov Complexity and Its Applications*. 3rd ed. Springer Publishing Company, Incorporated, 2008. ISBN: 0387339981.