

Reconstructing Sets From Interpoint Distances

Paul Lemke

Steven S. Skiena

Warren D. Smith

Abstract

Which point sets realize a given distance multiset? Interesting cases include the “turnpike problem” where the points lie on a line, the “beltway problem” where the points lie on a loop, and multidimensional versions. We are interested both in the algorithmic problem of determining such point sets for a given collection of distances and the combinatorial problem of finding bounds on the maximum number of different solutions. These problems have applications in genetics and crystallography.

We give an extensive survey and bibliography in an effort to connect the independent efforts of previous researchers, and present many new results. In particular, we give improved combinatorial bounds for the turnpike and beltway problems. We present a pseudo-polynomial time algorithm as well as a practical $O(2^n n \log n)$ -time algorithm that find all solutions to the turnpike problem, and show that certain other variants of the problem are NP-hard. We conclude with a list of open problems.

1 Introduction

A set of n points in some space defines a set of distances between all pairs of points. In this paper we consider the inverse problem of constructing all point sets which realize a given distance multiset. The complexity of an algorithm to generate all such point sets depends upon the number of solutions, and so we are also interested in bounds on the maximum number of distinct solutions in a given space, as a function of n .

The problem dates back to the origins of X-ray crystallography in the 1930’s [patt35] [picc39] [patt44]. More recently it has arisen in restriction site mapping of DNA, and was independently posed by M.I. Shamos [sham77] as a computational geometry problem. We encourage the reader to consult the recent thesis by Dakic [daki00] for the most recent results, including efforts to apply semi-definite programming to the problem. Pandurangan and Ramesh [pand01] have recent work on a variant of our problem which assumes additional information.

Spaces of particular interest include restricting the points to a line or a circular loop. The analogy of these points as exits on a road lead us to call these cases the “turnpike” and “beltway” problems, respectively. A turnpike problem instance consists of a multiset of $\binom{n}{2}$ distances; a beltway instance consists of a list of $(n-1)n$ distances.

It should be made clear that the correspondences between the distances and point pairs are *not* known and the entire difficulty of the reconstruction problem is to deduce such labeling information. If the labels are known, then given the $\binom{n}{2}$ labeled distances among n points in d -space, a suitable set of coordinates may be determined in $\mathbf{O}(n^2d)$ time. Let the n th point lie at $\vec{0}$ and let the coordinates of the $n-1$ nonzero points be given by the columns of a $d \times (n-1)$ upper triangular matrix A . Define B by $B = A^T A$. Then $B_{ij} = \frac{1}{2}(q_{0i} + q_{0j} - q_{ij})$ for $1 \leq i, j, n$, and $B_{ii} = q_{0i}$ where q_{ij} is the squared distance between points i and j . We may solve for A in terms of B consecutively column by column in time $\mathbf{O}(dn)$ per column, for a total runtime $\mathbf{O}(dn^2)$. If $d = n$, this algorithm is called the “Cholesky factorization” [gol83]. Alternately, we may find the “eigendecomposition” of $B = Q^T \Lambda Q$ where Λ is the diagonal matrix of $n-1$ real eigenvalues of B (in decreasing order; only the first d can be nonzero) and the columns of Q are the eigenvectors of B [gol83]. Q has orthonormal rows and columns. Then $\Lambda^{1/2}Q$ has d nonzero rows. Its $n-1$ columns give coordinates for our $n-1$ points. This approach has numerical advantages in situations in which our distances are contaminated by noise or roundoff error, because, e.g. the best approximation, in the Frobenius norm, of a symmetric matrix by a rank- d positive definite symmetric matrix is precisely its eigendecomposition with all eigenvalues besides the d largest ones, artificially zeroed ([gol83]; this is the symmetric case of the SVD approximation theorem).

1.1 Notation

\mathbf{Z} is the set of integers and \mathbf{R} the set of real numbers, while \mathbf{S}^d denotes the unit d -sphere $\{\vec{x} \in \mathbf{R}^{d+1}: |\vec{x}| = 1\}$. We will use $(\mathbf{R}/\mathbf{Z})^d$ to denote a *flat d -torus*, i.e. the d -cube $[0, 1]^d$ with opposite faces equivalenced. In X-ray crystallography, one must reconstruct a point set from its vector-differences modulo some d -paralleliped unit cell; by taking an affine transformation this paralleliped may be transformed to $(\mathbf{R}/\mathbf{Z})^d$. Our asymptotic notation $\mathbf{O}()$, $o()$, $\theta()$, \sim follows [knu76].

An algorithm is said to run in *pseudo-polynomial time* if it runs in time polynomial in the size of its input, when this input consists of integers written as unary numbers [gare78]. Similarly, a problem is *strongly NP-complete* if it is still NP-complete even if the input is required to consist of unary integers. For convenience, we adopt the “real RAM” [prep85] as our model of computation in this paper, although we have taken care not to abuse its excessive power. All of our lower bounds, NP-completeness proofs, and

undecidability proofs have been obtained by the explicit use of weaker models of computation polynomially equivalent to a Turing machine.

Two noncongruent n -point sets are *homometric* if the multisets of $\binom{n}{2}$ distances they determine are the same. A set of points in \mathbf{R}^d are in *general position* if their coordinates do not satisfy any nontrivial algebraic equation with integer coefficients.

The function $H_d(n)$ denotes the maximum possible number of mutually noncongruent and homometric n -point sets which can exist in \mathbf{R}^d . $S_d(n)$ denotes the maximum possible number of such sets for which all points must lie on the sphere \mathbf{S}^d . $H_d^*(n)$ and $S_d^*(n)$ are defined in the same way, except that we allow 0-distances or, equivalently, overlapping points. Thus $H_d(n) \leq H_d^*(n)$.

1.2 Summary of results

Our results on the worst-case computational complexity of n -point reconstruction problems may be summarized as follows. All d -dimensional vector difference turnpike reconstructions can be found in $\mathbf{O}(2^n n \log n)$ time and all d -torus vector difference beltway reconstructions in $\mathbf{O}(dn^n \log n)$ time. These include the one-dimensional turnpike and beltway problems respectively as special cases. The turnpike problem is also soluble in pseudo-polynomial time by a different algorithm. The *decision* problem of whether a multiset of $\binom{n}{2}$ distances is realized by n points in \mathbf{R}^d is NP-complete. Also shown NP-complete is the question of whether there exist n points in \mathbf{R} whose distances are contained in $\binom{n}{2}$ given distance intervals.

Our bounds on $H_d(n)$ are summarized below. The lower bounds hold for an infinite number of values of n , while the upper bounds are valid for all n .

$$\frac{1}{2}n^{0.8107144} \leq H_1(n) \leq \frac{1}{2}n^{1.2324827}, \quad H_1^*(n) \leq \frac{1}{2}n^{2.4649654}$$

$$n^{(d/2-1-\epsilon)n} \leq H_d(n) \leq H_d^*(n) \leq n^{(2d-1)n}, \quad d \geq 2, \epsilon > 0$$

Further, $H_1(n)$ is always a power of 2. The maximal order of $\log H_d(n)$ is known to within a constant factor as $n \rightarrow \infty$ except when $d = 2$. We also have similar results on the number of solutions for the case when the points lie on a d -dimensional sphere or torus. If the requirement that the points be distinct is relaxed, then we can show

$$\exp \left\{ \theta(n \exp[-0.7044985\sqrt{\ln n}]) \right\} \leq H_2^*(n)$$

Concerning the maximum possible number $S_1(n)$ of n -point beltway reconstructions, the upper bound $S_1(n) \leq H_2(n)$ can be shown by embedding a circle in the plane. We also have

$$\exp(2\frac{\ln n}{\ln \ln n} + o(1)) \leq S_1(n) \leq S_1^*(n) \leq \frac{1}{2}n^{n-2};$$

this lower bound is far stronger asymptotically. But it is valid only for a certain infinite set of n , while the upper bound is valid for all n .

A preliminary version of this paper appears in [skie90]. This paper is organized as follows. Section 1.3 discusses applications of reconstruction problems to DNA sequencing and x-ray crystallography. Section 2 discusses bounds on $H_d(n)$ and related functions. Section 3 presents reconstruction algorithms and their analysis. Section 4 presents computational complexity results including NP-completeness proofs. We conclude by presenting a list of open problems in Section 5.

1.3 Application in DNA sequencing

The genetic code of an organism can be thought of as a string on an alphabet $\{A, C, G, T\}$ (adenine, cytosine, guanine, or thymine). *Sequencing* is the process of determining this pattern for a given strand of DNA. *Restriction enzymes* are chemicals which recognize and cut DNA molecules at particular patterns. For example, the enzyme *Eco RI* cuts at *GAATTC*. Over 100 restriction enzymes are known and each may cut a given strand of DNA many times. The lengths of DNA fragments may be measured by means of differing forced diffusion speeds in electrophoresis; length measurement accuracies approaching $\approx 0.1\%$ are feasible.

Restriction site mapping is the process of determining where all the restriction sites lie by measuring the lengths of DNA fragments. By partially digesting DNA with some restriction enzymes, fragments of all possible lengths are produced. Additional information may be obtained by using different combinations of enzymes, for intervals containing sites for more than one enzyme, although we not consider this possibility. Then determining where the restriction sites for those enzymes lie is a turnpike or beltway problem, depending on whether the original DNA was linear or circular. In this paper, we give a backtracking algorithm for the turnpike problem that in practice is capable of handling almost all distance sets of up to several hundred points.

Additional work extending the backtracking algorithm discussed in this paper to biological data is reported in [skie94] [zhan94]. Related combinatorial bounds on the probed partial digestion problem are reported in [newb93]. For additional information on restriction site mapping, see [ally88] [bell88] [dixk88] [rhee88] [stef78] [tuff88].

1.4 Applications in X-ray crystallography

In Fraunhofer monochromatic X-ray diffraction [hose62] [patt44] [patt35] there is a simple correspondence, up to certain angle-dependent factors, between the far-field amplitude pattern of X-rays scattered from the sample, and the modulus

$$\left| \int \int \int_{\text{sample}} e^{i\vec{k} \cdot \vec{x}} \rho(\vec{x}) d^3\vec{x} \right|$$

of the Fourier transform of the X-ray scattering density $\rho(\vec{x})$ within the sample. Our object is to determine the function $\rho(\vec{x})$ from this scattering data. This problem would be solvable directly by an inverse Fourier transform, except that only the *modulus* and not the (complex) phase of the Fourier transform is known. By the Fourier convolution identity, inverse Fourier transforming the square of this modulus gives us the autocorrelation function or *Patterson function*

$$Q(\vec{y}) = \left| \int \int_{\text{sample}} \rho(\vec{x}) \rho(\vec{y} - \vec{x}) d^3 \vec{x} \right|^2$$

which constitutes exactly the information that is recoverable from the scattering data.

If the atoms are modeled as identical Dirac-delta masses, then Q gives us exactly the vector-difference multiset for the n atoms in the sample. Reconstructing the coordinates of the atoms from this multiset is a three-dimensional vector-difference problem, while considering each coordinate separately gives a turnpike problem. Generally one is dealing with a crystal, in which the atoms lie in an essentially infinite periodic pattern. X-ray scattering then gives the vector differences among the atoms in a single (parallelepiped) unit cell, *modulo* fundamental translations. To reconstruct each cell, a beltway problem must be separately solved in each coordinate, while the entire problem is a vector difference problem in a 3-torus.

Other instances of distance reconstruction problems occur in astronomy [fink83], pattern recognition [mcla61], and the psychology of vision [gilb74].

2 Bounds on the Number of Homometric Sets

In this section we analyze the behavior of the functions $H_d(n)$, $S_d(n)$, $H_d^*(n)$, $S_d^*(n)$. We remark that the maximum number of n -point sets in \mathbf{R}^d (respectively a d -torus) having the same *vector-difference* multiset, is exactly $H_1(n)$ (resp. $S_1(n)$) by projection onto a general rational line, so that we need not consider this problem separately.

2.1 Homometric Sets for the Turnpike Problem

We will prove power law upper and lower bounds on $H_1(n)$.

Lemma 2.1 *If $n \leq 5$, then $H_1(n) = 1$, but if $n \geq 6$, then $H_1(n) \geq 2$.*

Proof: For $n \geq 6$, the following construction by Lemke and Werman [lemk88], based on a 6-point pair found by Hosemann and Bagchi [hose62] gives a homometric pair. Observe that for all $n \geq 6$, the two n -point sets X and X' given by $X = n+1, n+3 \cup S \cup T$ and $X' = 2, n+2 \cup S \cup T$ where S is the set of integers i with $5 \leq i \leq n-2$, and $T = \{0, 1, n, n+5\}$, are homometric and noncongruent.

Table 1. Some examples of homometric sets.

n	the sets	source
6	$\{0,1,2,6,8,11\}, \{0,1,6,7,9,11\}$	[hose54], Lemma 2.1
6	$\{0,1,4,10,12,17\}, \{0,1,8,11,13,17\}$	[bloo77]
7	$\{0,1,5,7,8,10,12\}, \{0,1,2,5,7,9,12\}$	[ross82]
8	$\{0,1,5,6,8,9,11,13\}, \{0,1,2,5,6,8,10,13\}$	Lemma 2.1
9	$\{0,1,2,3,4,6,7,8,11\}, \{0,1,4,5,6,7,8,9,11\}$	Lemma 2.2 and $\{0,1,4\}, \{0,2,7\}$
13	$\{0,2,3,5,7,9,10,13,16,17,18,22,28\},$ $\{0,2,7,9,13,14,15,17,18,19,22,25,28\},$ $\{0,2,6,11,12,13,15,17,18,20,21,25,28\},$ $\{0,2,3,5,6,9,11,13,15,16,20,21,28\}$	computer search
14	$\{0,1,5,10,11,12,13,15,17,19,20,22,23,26\},$ $\{0,1,3,4,7,9,10,11,12,14,16,21,22,26\}$	computer search
15	$\{0,1,3,4,5,8,10,11,13,14,15,20,22,26\},$ $\{0,1,7,9,11,12,13,16,18,19,21,22,23,26\}$ $\{0,1,4,5,8,9,10,11,12,13,19,23,25,26,28\},$ $\{0,1,2,3,10,13,15,16,17,19,20,21,24,25,28\}$ $\{0,1,2,3,4,5,6,10,13,14,18,19,21,25,28\},$ $\{0,1,7,8,10,11,15,19,20,22,23,24,25,26,28\}$	computer search

We now show that $H_1(n) = 1$ if $n = 2, 3, 4$, or 5 . The cases $n = 2$ and $n = 3$ are trivial. For the case $n = 4$, place the two furthest-apart points $x_1 = 0, x_4$; then the second-largest distance determines the third point (up to a mirror reflection). Finally the fourth point is then completely determined by the three remaining distances.

For the case $n = 5$, place the two furthest apart points $x_1 = 0$ and $x_5 = 1$ (say). Of the remaining 9 distances, there are the three pairs of distances each of which sum to 1; call this set of 6 distances S and the remaining 3 distances T . The set $T = \{a, b, c\}$ must have the property that $a + b = c$. This partitioning into sets S and T is necessarily unique because since $a + b = c$, if $a + c = 1$ then b uniquely determines T , while if $a + c \neq 1$ (and $b + c \neq 1$), then S is determined uniquely. By the case $n = 3$, T determines the three remaining points x_2, x_3, x_4 uniquely up to a translation and a reflection. The set S then determines the entire 5-point set up to a reflection about the midpoint $\frac{1}{2}$. Note that this idea will not suffice to prove uniqueness for $n \geq 6$, because there are then enough degrees of freedom to make the selection of the $(n - 2)$ -pair set S ambiguous. ■

Some examples of homometric sets appear in Table 1. Homometric sets of larger multiplicity can be constructed with the following observation:

Lemma 2.2 *If $b \geq a \geq 3$, then $H_1(ab) \geq 2 H_1(a) H_1(b)$.*

Proof: To each n -point set $a_i, 1 \leq i \leq n$, we associate the generating function

$$P(x) = \sum_i x^{a_i} \quad (1)$$

There is a correspondence between the distance multiset $|a_i - a_j|, 1 \leq i < j \leq n$ and the distance generating function $P(x)P(1/x)$. Any two sets with the

same distance generating function must be homometric. Given a set of $H_1(a)$ homometric a -point sets whose generator polynomials are $P_i(x)$, $1 \leq i \leq H_1(a)$, and a set of $H_1(b)$ homometric b -point sets whose generator polynomials are $Q_j(x)$, $1 \leq j \leq H_1(b)$, we can construct $2H_1(a)H_1(b)$ mutually homometric ab -point sets as follows. Construct the sets whose generating functions are $P_i(x)Q_j(x)$ and $P_i(x)Q_j(1/x)$. If the a -point sets have been appropriately pre-scaled, e.g. by dilation with some constant incommensurable with the b -point sets, then no point overlaps can occur. ■

In the proof above, we have implicitly assumed that the original sets from $P_i(x)$ and $Q_j(x)$ were not reflection symmetric; this assumption is justified because:

Lemma 2.3 *Lemma 1 Any point set in \mathbf{R} that is invariant under a reflection, cannot be homometric with any incongruent set.*

Proof: Follows immediately from either of the algorithms of Section 3 for finding all point sets with a given distance set (see second paragraph of Section 3.3). ■

Theorem 2.4 *For an infinite number of values of n , $\frac{1}{2}n^\alpha \leq H_1(n)$, where*

$$\alpha = \ln(8)/\ln(13) \approx 0.8107144.$$

Proof: Table 1 proves that $H_1(13) \geq 4$. By iterative application of Lemma 2.2, $H_1(13^k) \geq 2^{3k-1}$ for all $k \geq 1$. If $n = 13^k$, this may be rewritten $2H_1(n) \geq n^{\ln 8 / \ln 13} n^{0.8107144}$, giving the result. More generally, if $H_1(a) \geq r$, then whenever $n = a^k$,

$$H_1(n) \geq \frac{1}{2}n^{\log_a(2r)}.$$

■

This result provides incentive to determine the least n such that $H_1(n) \geq 2^r$, which is open for $r \geq 2$. In particular, demonstrating that $H_1(n) = 8$ for some $n \leq 30$ would improve the result of Theorem 2.4.

Before presenting our asymptotic upper bound on $H_1(n)$, let us first present a few results concerning the structure of homometric sets in \mathbf{R}^1 .

Consider the set of all equations of the form $d_1 + d_2 = d_3$ satisfied by the distances d_i . Two distance sets are called *equivalent* if they satisfy exactly same system of equations of this form (up to some relabeling), because any solution of the turnpike problem for one distance set, is immediately converted into a solution for an equivalent set by substituting corresponding distances. In theory, it is possible to characterize all possible inequivalent types of distance sets for n -point sets, and hence determine any desired value

of $H_1(n)$, in $\mathbf{O}(n^{6n})$ time. This is by investigating every possible set of $\leq n$ linearly independent equations (in addition to the usual set of triangle equalities) and solving the resulting linear systems for the point sets. Another result along these lines is

Lemma 2.5 *Given any n -point turnpike problem, there is an equivalent turnpike problem whose distances are all integers bounded by $6^{(n-2)/2}$. This reduction may be carried out in polynomial time.*

Proof: Consider the set of all equations of the form $d_1 + d_2 = d_3$ satisfied by the distances d_i , where d_i is the i th largest distance. Since its coefficients are integers, this system must have a rational solution – indeed an integer solution by scaling – satisfying no other linearly independent equations of this form. Simply find such a solution in polynomial(n) operations.

This linear system has rank $\leq n - 2$ because all the distances are determined as differences among the coordinates of the n points and we may fix the two furthest-apart points at 0 and 1 without loss of generality. We may use the coordinates of the middle $n - 2$ points as (at least) a basis. When all equations are rewritten in these variables, every equation has norm (sum of the squares of its coefficients) ≤ 6 . Hence Hadamard's inequality and Cramer's rule, applied to a spanning set of $\leq n - 2$ equations, shows that the numerator and (common) denominator of each (rational) coordinate is $\leq 6^{(n-2)/2}$. The bound follows upon removing the denominators. ■

A theorem given by Rosenblatt and Seymour [ross82] nicely characterizes homometric pairs using generating functions.

Theorem 2.6 *Two point sets with generator polynomials $P(x)$ and $Q(x)$ are homometric if and only if there exist generating functions $A(x)$ and $B(x)$ and numbers μ and ν such that $P(x) = x^\mu A(x)B(x)$ and $Q(x) = x^\nu A(x)B(1/x)$.*

Theorem 2.7 *The number of finite subsets of \mathbf{R}^1 that have a given distance multiset, is always a power of 2. Consequently, $H_1(n)$ is always a power of 2.*

The proof will require a few intermediate lemmas.

Lemma 2.8 *If $F(x)$ is a polynomial with integer coefficients which divides a polynomial whose coefficients are 0's and 1's only, then $F(x)$ has first and last coefficients that are either both +1 or both -1.*

Proof: Clearly the first and last coefficients are ± 1 . If their signs differed, then $F(x)$ would have a positive real root, and 0-1 polynomials cannot. ■

Lemma 2.9 *If $F(x)$ and $G(x)$ are monic polynomials with integer coefficients, and if $F(x)F(\frac{1}{x}) = G(x)G(\frac{1}{x})$, then if either $F(x)$ or $G(x)$ is 0-1, then the other is also.*

Proof: Let $F(x) = \sum_k f_k x^k$, $G(x) = \sum_k g_k x^k$. We have $F(1) = \pm G(1)$; assume the $+$ sign for the moment. It follows that $\sum_k g_k = \sum_k f_k$. Equating the constant terms in $F(x)F(\frac{1}{x}) = G(x)G(\frac{1}{x})$ yields $\sum_k (g_k)^2 = \sum_k (f_k)^2$. Hence

$$\sum_k g_k - (g_k)^2 = \sum_k f_k - (f_k)^2. \quad (2)$$

If all f_k are 0 or 1, then both sides of this equation are zero; but if any g_k were not 0 or 1, then the left hand side of (2) would be negative – a contradiction. If, on the other hand, $F(1) = -G(1)$, then by the same argument $-G(x)$ would be 0-1, contradicting the assumption that it is monic. ■

A corollary is that if $P(x) = F(x)G(x)$ is 0-1, then so is $Q(x) = F(x)G(\frac{1}{x})$, because $P(x)P(\frac{1}{x}) = Q(x)Q(\frac{1}{x})$. An alternate proof follows from Filastta's discussion of the factorization of 0-1 polynomials into reciprocal and nonreciprocal parts [fila99].

A polynomial $P(x)$ of degree k is said to be *reciprocal* if $P(x) = x^k P(\frac{1}{x})$, meaning it corresponds to a palindromic point set.

Lemma 2.10 *If $F(x)$ is a non-reciprocal polynomial with integer coefficients, then $F(x)^2$ does not divide any 0-1 polynomial.*

Proof: Let $F(x) = \sum_{k=0}^D a_k x^k$. Then let j be the smallest index such that $a_j \neq a_{D-j}$. By Lemma 2.8, we may assume (by negating F if necessary) that $a_0 = a_D = 1$, hence $j \geq 1$. Let $F(x)^2 = \sum_{k=0}^D A_k x^k$, where $A_0 = a_0^2$, $A_1 = 2a_0a_1$, $A_2 = a_1^2 + 2a_0a_2$, and so on. It is also the case that A_k and A_{2D-k} first differ when $k = j$; specifically $A_j - A_{2D-j} = 2(a_j - a_{D-j})$. Now assume that a polynomial $G(x)$ with integer coefficients exists, such that $F(x)^2 G(x)$ is 0-1. Let $G(x) = \sum_{k=0}^E g_k x^k$, and without loss of generality $g_0 = 1$. Then the coefficients C_k of $F(x)^2 G(x)$ again first differ from C_{E+2D-k} when $k = j$, specifically

$$C_j - C_{E+2D-j} = (A_j - A_{2D-j})g_0 = 2(a_j - a_{D-j})g_0.$$

But the fact that two coefficients differ by ≥ 2 contradicts the assumption that this polynomial is 0-1. ■

Theorem 2.7 now follows from applying Theorem 2.6 in every possible way to the factorization of the generating function of a point set. Lemma 2.10 shows that a power of two generating functions must be obtained in this way; Lemma 2.9 and its corollary show that all generating functions obtained are 0-1, i.e. actually correspond to legitimate point sets.

We will now prove an asymptotic upper bound on $H_1(n)$ [lemk88]. This proof will involve several kinds of polynomial norms, which we will now define. Given a polynomial

$$P(x) = a_0 + a_1 x^1 + \dots + a_k x^k,$$

with integer coefficients, define the L_2 -norm

$$L_2(P) = \left(\sum_{0 \leq i \leq k} a_i^2 \right)^{1/2}$$

and the *Mahler measure* $M(P)$ [mahl76 p. 5]

$$M(P) = |a_k| \prod_{1 \leq i \leq k} \max(|\alpha_i|, 1)$$

where α_i are the roots of P . Two properties of the Mahler measure are multiplicativity: $M(AB) = M(A)M(B)$, and “Specht’s inequality” [spec49] $M(P) \leq L_2(P)$.

An important property of the Mahler measure, due to Smyth [smyt71], is that $M(P) \geq M(x^3 - x - 1) \approx 1.32472$ if P is a non-reciprocal polynomial.

Theorem 2.11 *For all n , $H_1(n) \leq \frac{1}{2}n^\beta$ and $H_1^*(n) \leq \frac{1}{2}n^{2\beta}$, where $\beta = \frac{\ln(2)}{2\ln(\phi)} \approx 1.2324827$, and $\phi \approx 1.32472$ satisfies $\phi^3 - \phi = 1$.*

Proof: Since noncongruent homometric sets are determined by different permutations of the same set of distances, it is clear that for any given n , $H_1(n)$ is finite. By Lemma 2.5 we may assume all the points have integer coordinates, so the generating functions may be assumed to be polynomials with integer coefficients. By Theorem 2.6, the number of point sets homometric to a set with generator polynomial $P(x)$ cannot exceed 2^{F-1} , where $P(x)$ has F irreducible nonreciprocal factors, since we do not count mirror images twice. From Specht’s inequality, Smyth’s lower bound on $M(P)$, and the fact that Mahler measures are multiplicative, $F \leq \ln(L_2(P))/\ln(\phi)$, and we obtain the result since $L_2(P) = \sqrt{n}$ if P is the generator function of a set of n distinct points. Even if the points are not required to be distinct, $L_2(P) \leq n$ so $H_1^*(n) \leq \frac{1}{2}n^{2.4649657}$. ■

It has been conjectured that almost all polynomials with 0-1 coefficients, or *Newman polynomials*, are irreducible. This conjecture is supported by statistical evidence, and is also known to be true for trinomials, since there is an exact formula [fla95,ljun60] for the number of divisors of a Newman trinomial $1 + x^A + x^B$.

If this conjecture is true, then by Theorem 2.6, almost all integer point sets are uniquely determined by their distance multisets, and may be determined by factoring the distance generating function. We conducted an exhaustive search among the n -element subsets of $\{1, 2, \dots, M\}$ for homometric sets. The number $f(M, n)$ of homometric pairs appears in Table 2. Only subsets including 1 and M were considered, and of these, only those that were lexicographically larger than their reflections. Thus if M is even, exactly $\frac{1}{2} \binom{M-2}{n-2} - \frac{1}{2} \binom{(M-2)/2}{(n-2)/2}$ candidate subsets were considered, and among these, $f(M, n)$ homometric pairs were found.

Table 2. A partial census of homometric pairs among n -element subsets of $\{1, 2, \dots, n\}$. The starred entries do not count any pairs arising inside the unique homometric quadruplet with these parameters. For these quadruplets, see Table 1.

M, n	6	7	8	9	10	11	12	13	14	15	16	17
12	1	0	0	1								
13	0	2	0	4	0							
14	1	2	2	7	0	0						
15	1	1	2	11	1	0	0					
16	0	4	6	14	8	4	1	0				
17	0	2	5	25	10	6	7	0				
18	1	3	6	40	16	11	27	2	0	0		
19	1	2	3	44	33	16	45	9	2	5	0	
20	1	3	9	63	38	32	99	15	12	16	0	0
21	1	2	4	78	39	43	148	36	21	50	2	0
22	1	2	11	95	68	78	227	69	65	106	14	2
23	1	2	9	104	62	70	316	88	107	186	27	11
24	2	3	9	144	89	99	541	164	169	405	84	15
25	0	4	5	142	109	123	618	268	313	648	189	61
26	1	4	8	186	109	161	909	364	498	1144	369	154
27	2	4	11	196	112	164	1100	381	681*	1639	601	248
28	1	3	10	232	153	194	1368	529	987	2539	1082	512
29	2	3	10	270	143	220	1720	635*	1070	3702*	1443	886
30	2	4	2	4	13	319	167	247	2246	812	1484	5461

2.2 Homometric Sets for the Beltway Problem

Let $S_d(n)$ be the maximum number of mutually homometric sets on the d -sphere \mathbf{S}^d . Since three points on a circle are uniquely determined up to a dihedral symmetry by their distances, $S_1(n) = 1$ for $n \leq 3$. For $n \geq 4$, the homometric pair given by Patterson [patt44]

$\{0, t, 1+t, 2, 4, \dots, 2(n-3)\}$ and $\{0, t, 2, 4, \dots, 2(n-3), 2n-5+t\} \bmod 2(n-2)$

(where t is a free parameter) proves that $S_1(n) \geq 2$.

The examples in Table 3 prove that $S_1(7) \geq 6$ and $S_1(13) \geq 19$, refuting a claim by Bullough [bull61, pp. 265] that $S_1(n) \leq n - 2$.

We have conducted an exhaustive search for all homometric sets that are n -element subsets of a regular M -gon, $M \leq 31$. The results of our census are in Table 4. By the use of distance generating functions modulo $x^M - 1$, it may be readily shown (see also [buer77] [chie79]) that two subsets of a regular n -gon are homometric if and only if their complements are. For this reason, we have only tabulated n for $4 \leq n \leq M/2$.

As may also be shown using distance generating functions, every n -point subset of the regular $2n$ -gon is homometric to its complement. Of course, it is usually also incongruent to its complement. Hence with probability $\Omega(n^{-1/2})$, a random subset of a regular $2n$ -gon is homometric to at least one other set. For a deeper investigation along these lines, see [rau80]. Thus while homometric turnpike pairs appear exponentially rare, beltway pairs are common.

2.2.1 Constructing Homometric Beltway Sets with Singer Difference Sets

Using Singer difference sets, we can construct beltway instances with a quadratic number of non-isomorphic reconstructions. *Singer difference sets* [sing38]

Table 3. Some n -point h -way-homometric beltway examples on a circle of perimeter M . Examples are included for all record-breaking parameter sets (n, M) , $M \leq 30$, i.e. such that q is larger than for any smaller M with that n .

n	h	M	h sets as binary (hexadecimal) M -bit numbers, each having n ones corresponding to points
4	2	8	E4 D8 (= 11100100 and 11011000 in binary = $\{0,1,2,5\}$ and $\{0,1,3,4\} \bmod 8$ explicitly)
5	2	10	3B0,3C8
5	3	18	34880,34084,32500
6	2	12	F90,F60
6	3	16	F40C,F908,F620
6	4	24	D11800,CC1400,E40220,E22400
6	5	31	68441000,604A2000,68110010,64142000,65020200
7	2	14	3EC0,3F20
7	3	18	3D820,3E410,3CD00
7	6	24	D0C810,D08184,CA4060,D04818,E24840,D18480
8	4	16	F42C,F948,F4C2,ED60
8	6	24	D4C180,D8C140,E46820,E26840,E4C280,E8C240
9	4	18	3EB04,3DD20,3EC14,3ED10
9	8	24	EA48C0,D4C580,E86242,EA4260,E842C8,D584C0,E26A40,E8CA40
10	6	20	F530C,EC31A,F6518,F350C,F3944,EB21C
10	8	24	F51A40,F205A8,EB1680,F50B04,ED02B0,F40B14,F60950,F50348
11	6	22	3F9460,3EE340,3DD380,3F8D10,3F4E20,3E7580
11	8	24	F91298,F584C8,E846D8,EC7242,F33520,ED068C,EE5260,EB4CC0
11	12	30	35982920,3A086494,36904C28,39205348,3A4A6410,36944C20 39A45240,3A486414,35902930,35024B0C,39A05248,36296410
12	12	24	F68722,F90AD8,F60E4A,F684E8,FB12B0,F2B720 FC8D28,F95B10,EC0EB4,F485D8,F72560,F62E12
13	8	26	3DB8740,3F47890,3F8D1A0,3E9D380,3E74B80,3F8B460,3F488F0,3C3BB40
13	19	28	ED60CC8,EB4CC0C,EC58C86,F448D8C,F662650,F9894C8,E606CD4,F46684C,E8D909C F94CC48,ED891C8,ECC5A0C,F623464,EC6CD08,F6468C4,F333520,E656CC0,EC5C90C, ECD1C84
14	16	28	ED05C74,F9CA9A0,FA6AC60,F125770,F4CF054,F6862E8,F2CEAC0,F90ACB8 F115B70,FA468E2,FAB0B30,F60E2CA,F48745C,F44D43C,F537430,FCASCB0 3F929328,3B2DAE40,3DB51B04,3BB60B48,3D9215E4,3EB213A4,3E8266B4, 3D941B46,3D9ACB02,3ACDAD80,3B02B6D8,3EC24AE4,3F532584,3DB21B0A, 3F426534,3E96E224,3D8159B4,3F5249C4,3F6292C8,3ECCA582

[hall67] are n -element subsets of \mathbf{Z}_M , where $M = (q^3 - 1)/(q - 1)$ and $q = n - 1$ is a prime power, with the property that the $n(n - 1) = M - 1$ differences they determine are all of the nonzero elements of \mathbf{Z}_M , each one occurring exactly once. Singer showed that at least one example always exists for each prime power q .

To construct other point sets homometric to Singer sets, we observe that multiplying a Singer set by any element r of \mathbf{Z}_M (r relatively prime to M) preserves the distance set. Unfortunately, as a consequence of Hall's multiplier theorem [hall56] [hall67], if r divides q , then multiplying by r will merely translate, reflect, and/or permute the elements of the difference set.

But it seems likely that upon multiplying a Singer set by any element r of \mathbf{Z}_M such that r is relatively prime to M and $\pm r$ does not divide q , a noncongruent set will be obtained. Let us specifically examine the case when M and q are prime. In this case, $q^3 \equiv 1 \bmod M$, so that q and -1 multiplicatively generate the order-6 dihedral group D_3 (inside \mathbf{Z}_{M-1}) of trivial multiplicative symmetries. Therefore if a Singer set with these parameters exists with *no* further multiplicative symmetries, then it would have $(M - 1)/6$ equivalence classes of noncongruent homometric multiples, each equivalence class of size 6.

By the method outlined in [Hall67 end of Section 11.3], it is easy to generate Singer sets by computer. One may then find all their multiples. After translating and reflecting each such set to reduce it to a canonical (least-lexicographic) form, throwing out degenerate multiples, and sorting to remove redundant sets, one has a large collection of noncongruent homometric

Table 4. A partial census of homometric n -point subsets of the regular M -gon. To explain the entries by example: The $M = 18, n = 9$ entry indicates that there are 512 homometric pairs composed of 9-subsets of the regular 18-gon, 6 such homometric triplets, and 54 such homometric quadruplets.

M, n	4	5	6	7	8	9	10	11
8	1							
9	0							
10	0	3						
11	0	0						
12	1	3	15					
13	1	0	2					
14	0	6	6	48				
15	0	5	25	10				
16	2	10	28,3	40,4	177,0,3			
17	0	0	16	24	52			
18	0	13,1	56,6	118,16	139,11	512,6,54		
19	0	0	21	57	90	156		
20	2	22	96,2	180,11	491,12,32	535,14,16	1973,1,130,0,2	
21	0	0	96	220	276,6	1032,23,7	568,45	
22	0	20	55,5	310,25	540,35	1300,125	1430, 120	6985,5,390,0,10
23	0	0	33	110	429	803,11	1144	1342,33

Table 5. Constructions of n -point h -way-homometric beltway examples on a circle of perimeter $M = \frac{q^3-1}{q-1}$.

n	q	h	M	n	q	h	M
4	3	2	13	8	7	6	57
6	5	5	31	12	11	18	133
18	17	51	307	14	13	20	183
42	41	287	1723	20	19	42	381
60	59	590	3541	24	23	78	553
72	71	852	5113	30	29	132	871
90	89	1335	8011	32	31	110	993

sets which may be used to verify this construction. Explicitly, the first three of these examples are:

$\{0,1,3,9\} \times (1 \text{ or } 2) \bmod 13$ [2 homometric 4-point sets are given]
 $\{0,1,3,8,12,18\} \times (1,2,3,4, \text{ or } 8) \bmod 31$ [5 homometric 6-point sets]
 $\{0,1,3,30,37,50,55,76,98,117,129,133,157,189,199,222,293,299\} \times$
 $(1,13,14,15,16,20,21,103,104,106,108,109,112,113,116,122,125,126,128,$
 $129,131,135,136,137,138,140,141,142,$
 $143,150,154,156,158,159,160,168,175,183,184,186,193,200,202,218,219,220$
 $,222,239,256,269,273) \bmod 307$ [51 homometric 18-point sets]
with each incarnation of each set having the distance (mod M) set $\{1,2,...M\}$.

Table 5 shows this construction yields a $(q + 1)q/6$ homometric sets for every prime q with $q < 100$ and $q^2 + q + 1$ prime. The table gives n,q,h , and M , meaning there are h homometric noncongruent n -point sets on a circle of perimeter M . The right hand column of Table 5 contains some examples of noncongruent homometric multiples of Singer sets for prime q such that $q^2 + q + 1$ is not prime. In these cases, less than $q(q + 1)/6$ sets are obtained, but still, the number is often quite large.

We conjecture that the Singer set construction generates $\geq (q+1)q/6$ homometric $(q+1)$ -point beltway sets for every prime q such that $q^2 + q + 1$ is prime. This would follow from showing that for each such q at least one Singer set with no nontrivial symmetries exists. This seems likely, considering that no Singer set with *any* nontrivial symmetries is known. That there are an infinite number of such q is a standard Hardy-Littlewood conjecture.

2.2.2 A Lower bound on $S_1(n)$

We now present a construction yielding an asymptotically better lower bound on $S_1(n)$. This construction uses the monic cyclotomic polynomial [bate49]

$$\Phi_k(z) \equiv \prod_{\substack{1 \leq m \leq k \\ \gcd(m, k) = 1}} (z - \exp(\frac{2\pi i m}{k})) = \prod_{d|k} (1 - z^d)^{\mu(k/d)} \quad (3)$$

Here $\mu(x)$ is the Moebius function, which is zero if x has a square factor, is $(-1)^c$ if x is the product of c different primes, and $\mu(1) = 1$. $\Phi_k(z)$ is the unique monic irreducible polynomial with integer coefficients whose roots include the k th root of unity. The fact which we will need about the cyclotomic polynomials is that there exist distinct irreducible polynomials $\Phi_1(z), \Phi_2(z), \dots$ with integer coefficients such that

$$z^n - 1 = \prod_{d|n} \Phi_d(z).$$

Theorem 2.12 $S_1(n) \geq \exp(2^{\frac{\ln n}{\ln \ln n} + o(1)})$ for infinitely many n .

Proof: Let p_i denote the i th odd prime, $K \geq 5$, and $Q \equiv 2^K$. Let $R \equiv p_1 p_2 \dots p_K Q$ and $n \equiv (p_1 + 2)(p_2 + 2) \dots (p_K + 2)$. We will construct $\geq 2^Q/4R$ mutually homometric and incongruent n -element subsets of the regular $2R$ -gon.

Throughout the proof, $[K]$ will denote the set $\{1, 2, \dots, K\}$ and U and T will denote subsets of $[K]$. (\bar{T} means $[K] - T$.) The quantities

$$P_T \equiv \prod_{i \in T} p_i, \quad z(T) \equiv \left(\sum_{j \in T} 2^{j-1} \right)$$

each uniquely specify such a set T . (In place of $z(T)$, we could have used any 1-1 mapping between subsets of $[K]$ and the integers m with $0 \leq m \leq 2^K$.) We also define the following polynomials.

$$F_T(x) = \frac{x^{2R} - 1}{x^{2P_T Q} - 1} = 1 + x^{2P_T Q} + x^{4P_T Q} + \dots + x^{2R - 2P_T Q}$$

$$G_T(x) = (x^{P_T Q} + 1) \prod_{i \in T} (x^{2P_T Q/p_i} + 1)$$

$$H_T(x) = (x^{P_T Q} - 1) \prod_{i \in T} (x^{2P_T Q/p_i} - 1)$$

The point sets all correspond to generating polynomials $S(x)$ where

$$S(x) = \sum_T x^{z(T)} F_T(x) \frac{G_T(x) \pm H_T(x)}{2} \quad (4)$$

where all Q of the \pm signs in the sum are independent, hence we have defined 2^Q possible point sets $S(x)$. We will now make a succession of claims, each of which are readily verified by use of the definitions above and the preceding claims, and from which the proof will follow.

1. $F_T(x)H_T(x)$ is divisible by $\Phi_d(x)$ for any d dividing $2R$, except that it is not divisible by $\Phi_d(x)$ with $d = 2P_T Q$.
2. $F_T(x)G_T(x)$ is divisible by $\Phi_d(x)$ for all d which divide $2R$ but do not divide R . In particular, it is divisible by $\Phi_d(x)$ for any d of form $d = 2P_U Q$, $U \subseteq [K]$.
3. Hence any two $S(x)$ are different modulo $\Phi_d(x)$ if $d = 2P_T Q$ and the two $S(x)$ differ in term T . Hence all the point sets $S(x)$ differ modulo $x^{2R} - 1$.
4. The difference generating function $S(x)S(\frac{1}{x})x^{2R}$ is invariant (with respect to the choice of signs in the sum (4) modulo $\Phi_d(x)$ for all d dividing $2R$, hence is invariant modulo $x^{2R} - 1$, hence all these sets $S(x)$ are homometric. This is because all of the “cross terms” of form $H_T(x)G_U(x)F_T(x)F_U(x)$ and $H_U(x)H_T(x)F_T(x)F_U(x)$, the latter with $U \neq T$, are zero modulo $\Phi_d(x)$ for every d dividing $2R$.
5. $G_T(x)$ is a polynomial with 0-1 coefficients, in fact it has exactly $2^{|T|+1}$ 1's. $H_T(x)$ is a polynomial with 0 and ± 1 coefficients, in fact it has exactly $2^{|T|+1}$ nonzero terms all of which coincide with nonzero terms of $G_T(x)$; exactly half of the nonzero coefficients on $H_T(x)$ are -1 's and half are $+1$'s. Hence $[G_T(x) \pm H_T(x)]/2$ has only 0-1 coefficients, and has exactly $2^{|T|}$ 1's. $F_T(x)$ times this has only 0-1 coefficients (since a sum of reciprocals of distinct p_i cannot be an integer, no overlap leading to non 0-1 coefficients is possible) and has exactly $2^{|T|}P_{\bar{T}}$ 1's.
6. If a is such that the coefficient of x^a in term T of the sum is non-zero, then $a = z(T) \bmod Q$. Since all of the $z(T)$'s are distinct and obey $0 \leq z(T)Q$, none of the 1-coefficients in different terms of the sum coincide, hence $S(x)$ has only 0-1 coefficients. Therefore every $S(x)$ really does represent a valid set of distinct points.
7. The total number of points in each set $S(x)$ is, as was claimed earlier

$$n = \sum_T 2^{|T|} P_{\bar{T}} = \prod_{i=1}^K (2 + p_i).$$

8. The total number of distinct, incongruent homometric n -point sets generated (after removal of a possible $4R$ dihedral symmetries) is at least $2^Q/4R$.

Upon applying the Prime Number Theorem $p_i \sim i \ln i$ and approximating various products using Riemann sums of logarithms, we see that

$$\begin{aligned} n &= \left(\frac{(1+o(1))K \ln K}{e} \right)^K = o(R) \\ R &= \left(\frac{2(1+o(1))K \ln K}{e} \right)^K = n^{1+o(1)} \\ K &= \frac{\ln n}{\ln \ln n} + o(1) = \frac{\ln R}{\ln \ln R} + o(1) \end{aligned}$$

from which the theorem follows. ■

As a direct consequence of the fact that $R = n^{1+o(1)}$, there cannot be a polynomial time algorithm to find all beltway reconstructions of a given distance multiset, even if all the input distances are required to be unary integers.

The first time the construction of Theorem 2.12 beats the quadratic construction with Singer sets is with $n = 1,167,075$, $R = 16,336,320$, $K = 6$, when there are at least 282,296,503,645 homometric incongruent n point subsets of the regular $2R$ -gon. The next instance, $K = 7$, involves $n = 24,508,575$, $R = 620,780,160$, and 10^{29} sets. The Singer set construction gives more spectacular and frequent homometric examples for $n < 1,167,075$. It is still unclear what the asymptotic behavior of $S_1(n)$ is, although a $\mathbf{O}(n^{n-2})$ upper bound follows from the algorithm of Theorem 3.5.

2.3 Homometric sets in higher dimensions

In \mathbf{R}^d , $d \geq 2$, at least three noncongruent n point sets exist whenever $n \geq 4$. Gilbert [gilb74] gives the following 2-parameter family of three homometric noncongruent 4-point planar sets. The three sets are $\{A, B, C, D\}$, $\{A, B, C, D'\}$, and $\{A, B, C, D''\}$. Given a general triangle $\triangle ABC$, let a be the midpoint of BC and b be the midpoint of AC . Let K be the line through a perpendicular to the line Aa . Let L be the line through b perpendicular to the line Bb . Then: D is the intersection of L and K . D' is the point of L such that $\text{dist}(D', b) = \text{dist}(D, b)$, and D'' is the point of K such that $\text{dist}(D'', a) = \text{dist}(D, a)$.

Since any sufficiently small perturbation of a regular d -simplex may be reconstructed with all possible assignments of edge lengths,

$$H_d(d+1) \geq \left(\frac{(d+1)d}{2} \right)! / (d+1)! \quad (5)$$

mutually noncongruent homometric $(d+1)$ -point sets exist in \mathbf{R}^d , $d \geq 2$. Thus $3 \leq H_2(4) \leq H_3(4) = 30$.

T. Caelli [cael80] claimed to have a complete characterization of all 4-point homometric pairs in \mathbf{R}^2 , but unfortunately his result is incorrect. Hence the value of $H_2(4)$ is still open within the bounds [3,30].

We will now briefly survey some known results on *few-distance sets*, in order to explain the connections with homometric sets.

Erdős [erd46] defined the function $F_d(n)$ to be the minimum possible number of distinct distances determined by n points in \mathbf{R}^d . One may define the similar function $G_d(n)$ for the case when the n points must lie on the sphere \mathbf{S}^d . It is easy to see that $F_1(n) = n - 1$ and $G_1(n) = \lfloor n/2 \rfloor$.

Any set of n points, all lying on the unit sphere \mathbf{S}^d in $(d+1)$ -dimensional space and determining s distances, must obey [dels77] (see also [koor76])

$$n \leq \binom{d+s}{d} + \binom{d+s-1}{d}, \quad (6)$$

and any set of n points in \mathbf{R}^d determining s distances must obey [bann83]

$$n \leq \binom{d+s}{d}. \quad (7)$$

The d -dimensional integer grid $\{0, 1, 2, \dots, n^{1/d}\}^d$ proves the upper bound $F_d(n) \leq dn^{2/d}$. Thus

$$\left(\frac{d}{e}n\right)^{1/d} - d < F_d(n) < dn^{2/d}, \quad d \geq 3. \quad (8)$$

In two dimensions [soli01]

$$\Omega(n^{6/7}) \leq F_2(n) \leq 0.7044984310 \frac{n}{\sqrt{\ln n}}; \quad (9)$$

the upper bound arises from the points of the equilateral triangle lattice lying inside a circle. P.Erdős offered \$500 for bounds on $F_2(n)$ tight to within a factor of $o(n^\epsilon)$.

In three dimensions [clar90]

$$\Omega(\sqrt{n}(\frac{n}{\lambda_6(n)})^{1/4}) \leq F_3(n) \lesssim \frac{5}{12}(\frac{3}{\pi})^{2/3}n^{2/3}$$

where $\lambda_6(n)/n$ is an extremely slowly growing function related to Davenport-Schinzel sequences [agar89].

Theorem 2.13 *For any fixed d with $d \geq 2$, and for any desired function $g(n)$,*

$$H_d(n) \geq \frac{\binom{g(n)}{n}}{\left(\binom{n}{2} + F_d(g(n)) - 1\right)} \frac{1}{nQ_d(n)}$$

$$S_d(n) \geq \frac{\binom{g(n)}{n}}{\binom{\binom{n}{2} + G_d(g(n)) - 1}{G_d(g(n)) - 1}} \frac{1}{Q_{d+1}(n)}$$

Here $Q_d(n)$ represents the maximum size of any finite group of origin-preserving reflections or rotations preserving the set achieving $H_d(n)$ or $S_{d-1}(n)$; specifically $Q_d(n) \leq 600 d! 2^d n$.

Proof: From the $g(n)$ -point set in \mathbf{R}^d defining the minimum possible number of distances, choose a random n -element subset. The numerator of the bound gives the number of such subsets, while the denominator dominates the number of possible distance multisets. The theorem now follows by the pigeonhole principle; the $nQ_d(n)$ and $Q_{d+1}(n)$ terms remove all possible overcounting due to reflection, rotation, and translation symmetries. ■

Theorem 2.14 For $d \geq 3$ and for any $\epsilon > 0$, there are an infinite number of values of n such that $H_d(n) \geq n^{(d/2-1-\epsilon)n}$.

Proof: Using the grid-set bound $F_d(n) \leq dn^{2/d}$ and choosing $g(n) = n^{d/2-\epsilon}$ in Theorem 2.13 yields

$$H_d(n) \geq \frac{\binom{n^{d/2-\epsilon}}{n}}{\binom{\binom{n}{2} + d o(n)}{d o(n)}} \frac{1}{Q_d n} = e^n n^{(d/2-1-\epsilon)n} n^{d o(n)}.$$

■

Theorem 2.15 For $n > d \geq 2$, $H_d(n) < n^{2dn}$.

Proof: Without loss of generality, we may assume that the point set achieving $H_d(n)$ spans \mathbf{R}^d . Hence, $d+1$ of the point sites must define some nondegenerate d -simplex. We will call some particular $(d-1)$ -simplex defined by d of the point sites the *central simplex*. Consider the following rigid structure defined by $(n-d)d + \binom{d}{2}$ line segments: Every pair of the vertices in the central simplex is joined by a line segment, and each of the remaining $n-d$ sites is joined to each of the d vertices of the central simplex by a line segment. If distances chosen from the $\binom{n}{2}$ available are assigned to each line segment in this structure, and also the $n-d$ points not on the central simplex are assigned \pm signs depending on which side of the hyperplane (defined by the central simplex) they lie on, then clearly the point set is completely defined. The number of possible ways to accomplish these sign and distance assignments, is

$$< \binom{n}{2}^{(n-d)d + \binom{d}{2}} 2^{n-d}.$$

We may remove a factor $d!(n-d)!$ due to automorphism symmetries of the graph structure. $\lfloor d/2 \rfloor$ of the distances in the central simplex (plus one

external edge, if d is odd) may be assumed without loss of generality to be the $\lceil d/2 \rceil$ largest distances available, so that we may remove an additional factor of $\binom{n}{2}^{\lfloor d/2 \rfloor}$ and replace it with a factor

$$(1 + (d-1)d/2)^{\lfloor d/2 \rfloor}.$$

The three reasons that we obtain an upper bound are: firstly, we are counting impossible distance assignments, secondly we are overestimating falling factorials by using powers, and thirdly, if any distances in our set of $\binom{n}{2}$ were equal, then we could have improved our bound even further.

The result follows by combining these terms:

$$H_d(n) \leq \frac{1}{d! (n-d)!} \binom{n}{2}^{(n-d)d + \binom{d}{2} - \lfloor d/2 \rfloor} 2^{n-d} (1 + (d-1)d/2)^{\lfloor d/2 \rfloor} \\ < n^{(2d-1)n - 3d^2/2} e^{n 2^{d^2/2 + (1-d)n}}.$$

Some simpler but weaker bounds are $n^{(2d-1)n} e^n$ and n^{2dn} . ■

In three or higher dimensions, the *logarithm* of the upper and lower bounds of Theorems 2.14 and 2.15 are tight to within a multiplicative factor of 4. Unfortunately in *two* dimensions, the best lower bound we have on $H_2(n)$ is subexponential (from Theorem 2.12), while the best upper bound we have is superexponential. It appears difficult to tighten these bounds, since a significant improvement would resolve Erdős's \$500 problem. Specifically: If there is some positive δ such that n -point planar sets can determine $\mathbf{O}(n^{1-\delta})$ distances, then there are

$$\geq n^{\frac{\delta}{1-\delta}n - o(n)}$$

(a superexponential number) of mutually noncongruent and homometric n -point planar sets. Conversely, if $H_2(n)$ grows more slowly than this, then $H_2(n) = \Omega(n^{1-\delta})$ for all $\delta > 0$. A new attack on the Erdős problem may be possible by using this fact.

A slightly better lower bound may be obtained if we allow overlapping points.

Theorem 2.16 *For all sufficiently large n , there are at least*

$$H_2^*(n) \geq \exp \left\{ \theta(n \exp[-0.7044985\sqrt{\ln n}]) \right\}$$

incongruent homometric sets of n (not necessarily distinct) points in the plane.

Proof: Assume that there is a sequence of n -point planar sets $S[n]$ with $\lesssim \kappa n / \sqrt{\ln n}$ distances. To prove Theorem 2.16, we choose n points at random (only with replacement this time) $S[g(n)]$ where

$$g(n) = \theta(n \exp[-\kappa\sqrt{\ln n}]).$$

Note $g(n) = o(n)$. The number of point sets that may be obtained in this way (after removal of at most $8n$ rotation/reflection equivalences) is at least

$$P = \frac{1}{8n} \binom{n + g(n) - 1}{g(n)}.$$

The number of possible distance multisets that can arise from the random n -point sets is at most

$$D = \binom{\binom{n}{2} + \frac{g(n)}{\sqrt{\ln g(n)}} - 1}{\frac{g(n)}{\sqrt{\ln g(n)}}}.$$

Using Stirling's formula

$$\ln \binom{A}{B} \sim B \ln \frac{A}{B} + B - \frac{1}{2} \ln B + \mathbf{O}(1)$$

to asymptotically analyze $\ln \frac{P}{D}$ gives

$$\ln \frac{P}{D} = \theta(g(n)).$$

Since $H_2^*(n) \geq \frac{P}{D}$ by the pigeonhole principle, using $\kappa \leq c_\Delta$ gives the result. ■

Examples of incongruent sets of points in higher dimensions with stronger kinds of homometricity properties may be obtained from [lind74] [doyv71].

3 Algorithms for Reconstructing Point Sets from their Distance Multiset

In this section we present a pseudo-polynomial time algorithm [gare78] for finding all n -point turnpike reconstructions, as well as an $O(2^n n \log n)$ algorithm which in practice shows excellent behavior, despite having exponential worst-case runtime.

3.1 Turnpike reconstruction by polynomial factoring

The algorithm for turnpike reconstruction which we discuss is due to Rosenblatt and Seymour [ross88], although the analysis showing that it runs in pseudo-polynomial time is due to Lemke and Werman [lemk88]. Given a set of distances d_1, d_2, \dots, d_N , where $N = \binom{n}{2}$, we form the distance generating function

$$Q(x) = n + \sum_i (x^{d_i} + x^{-d_i}).$$

We factor this polynomial into irreducibles over the ring of polynomials with integer coefficients, using a factoring algorithm with runtime polynomial in $\max_i d_i$ [land87] [lens82]. It is convenient to write this factorization as

$$Q(x) = \prod_{i=1}^F P_i(x) P_i\left(\frac{1}{x}\right) R(x) R\left(\frac{1}{x}\right)$$

where the $P_i(x)$'s are monic, irreducible, and non-reciprocal, while all the reciprocal factors are collected in the polynomial $R(x)$. If Q does not have a factorization of this form then no turnpike reconstruction is possible. The point set a_1, a_2, \dots, a_n must have a generating function (1) obeying

$$Q(x) = P(x) P\left(\frac{1}{x}\right).$$

Therefore, we simply try all 2^F possible subsets S of $\{1, \dots, F\}$ as putative factors of $P(x)$, i.e.

$$P_S(x) = \prod_{i \in S} P_i(x) \prod_{i \in \bar{S}} P_i\left(\frac{1}{x}\right) R(x),$$

finding a superset of all the possible point sets $P(x)$ in time polynomial in N and 2^F . To avoid finding reflected solutions we may assume $1 \in S$ and translate the point sets to assure the leftmost point lies at the origin. After eliminating the point sets from generating functions with negative coefficients and sorting to remove possible redundant copies of point sets, we are done.

The runtime claim follows from Theorems 2.6 and 2.11 that if a turnpike reconstruction is possible, then $2^F \leq n^{2.4649657}$. If 2^F does not obey this bound, then no reconstruction is possible, and so we may terminate the algorithm early.

As an example, consider Bloom's distance set $\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 16, 17\}$, whence

$$\begin{aligned} Q(x) &= x^{17} + x^{16} + x^{13} + \dots + x + 6 + x^{-1} + \dots + x^{-13} + x^{-16} + x^{-17} \\ &= P_1(x) P_1\left(\frac{1}{x}\right) P_2(x) P_2\left(\frac{1}{x}\right) \end{aligned}$$

where $P_1(x) = x^6 + x + 1$, $P_2(x) = x^{11} - x^5 + x^4 + 1$ and the point sets $\{0, 1, 4, 10, 12, 17\}$ and $\{0, 1, 8, 11, 13, 17\}$ arise from $P_1(x)P_2(x)$ and $P_1(x)P_2\left(\frac{1}{x}\right)$ respectively.

The factoring algorithm can be made to work even if it is given non-integer real distances by the use of Lemma 2.5, although then "pseudo-polynomial time" will no longer be an applicable notion.

3.2 A Backtracking Algorithm for Turnpike Reconstruction

The turnpike reconstruction problem has a combinatorial flavor which is not reflected by the polynomial factorization algorithm. In this section, we take

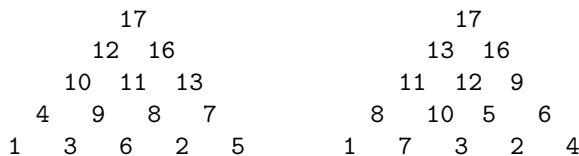


Fig. 1. Distance Pyramids for Bloom's 6-point Homometric Example

a different look at the problem which leads to a practical algorithm for reconstruction.

Let d_{ij} represent the distance between the i th and j th points on a line, ordered from the left. We shall assume each distance to be a real number, $d_{ij} \geq 0$. For convenience, we shall represent these distances by a vector v sorted in increasing order, so that $v[i] \leq v[j]$, $1 \leq ij \leq \binom{n}{2}$. The complete set of distances and their initially unknown assignments can be represented as a triangular matrix or *pyramid* as in Figure 1. The set of distances d_{ij} such that $j - i = l$ defines row l of the pyramid.

As with Pascal's triangle, several identities can be observed on this structure. For example, there is a simple relationship between any four elements forming a "parallelogram" in the pyramid.

Lemma 3.1 $d_{ij} + d_{kl} = d_{il} + d_{kj}$, where $i \leq k \leq l \leq j$.

The relationship between the sums of rows is particularly pretty.

Theorem 3.2 *For any pyramid, the sum of distances on the k th row equals the sum of the distances on the $(n - k)$ th row. Formally:*

$$\sum_{i=1}^{n-k} d_{i(i+k)} = \sum_{i=1}^k d_{i(i+n-k)}.$$

Proof: We will use the triangle equality and count how many times each base distance $d_{i,i+1}$ is added to the row.

$$\begin{aligned} \sum_{i=1}^{n-k} d_{i(i+k)} &= d_{1(k+1)} + d_{2(k+2)} + \dots + d_{(n-k)k} \\ &= \sum_{i=1}^{k-1} i \cdot d_{i,i+1} + k \sum_{i=k}^{n-k} d_{i,i+1} + \sum_{i=1}^{k-1} i \cdot d_{(n-i)(n-i+1)} \\ &= d_{1(n-k+1)} + d_{2(n-k+1)} + \dots + d_{kn} = \sum_{i=1}^k d_{i(i+n-k)} \end{aligned}$$

■

One implication of Theorem 3.2 is that reconstructing a pyramid with an even number of rows (ie. $n - 1 = 2k$) requires solving a partition problem,

since the sum of the $(n^2 - 2n)/8$ distances in the top half of the pyramid equals those of the $(3n^2 - 2n)/8$ in the bottom half.

Based on the triangle equality, we can make certain assignments between distances and pairs of cities. Clearly, the largest distance $v[\binom{n}{2}]$ represents d_{1n} . The reflection of any pyramid is also a pyramid. To eliminate double counting of reflections, we make the convention that a pyramid is in canonical order if it is lexicographically less than its reflection. Thus $d_{2n} = v[\binom{n}{2} - 1]$, which implies $d_{12} = v[\binom{n}{2}] - v[\binom{n}{2} - 1]$. Unfortunately, this represents the extent of our ability to assign distances absolutely.

Theorem 3.3 *There is a $O(2^n n \log n)$ -time algorithm for finding all possible reconstructions of an n -point set from its $\binom{n}{2}$ -element distance multiset.*

Proof: We shall use a backtracking procedure, and repeatedly position the *largest remaining* unpositioned distance. These elements will be filled in from the top of the pyramid, using backtracking to select either the left or right side. Because we are placing the largest available distance in the pyramid, there will be only two possible locations to choose from. The key to making this procedure efficient is to immediately fill in the pyramid any values which are determined by our previous (non-deterministic) selections.

Assume that we have thus far placed l elements along the left-hand side of the pyramid and r on the right, as shown in Figure 2.

All the positions in the shaded regions are determined by the $l + r - 1$ distances d_{1j} and d_{ir} , $(n - l + 1) \leq j \leq n$, $1 \leq i \leq r$. The largest remaining distance must be associated with either $d_{1(n-l)}$ or $d_{(r+1)n}$.

Suppose that the backtrack procedure selects the left side to receive the next element. Thus $d_{1(n-l)}$ is assigned the largest remaining distance. This determines $d_{i(n-l)}$, $2 \leq i \leq r$, since $d_{i(n-l)} = d_{1(n-l)} - d_{1i}$, which are already in the pyramid, as well as $d_{(n-l+1)(n-j)}$, $0 \leq j \leq l - 2$, since $d_{(n-l+1)(n-j)} =$

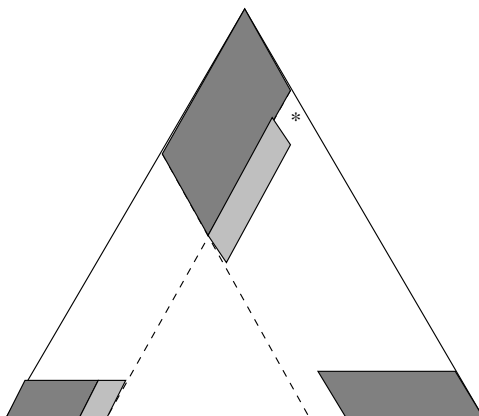


Fig. 2. The Effect of Choosing $d_{1(n-l+1)}$.

$d_{1n} - d_{1(n-l+1)} - d_{(n-j)n}$. If any of these $l + r - 1$ determined values is not a remaining distance, the partial reconstruction cannot be extended into a pyramid and we backup. If they are all remaining distances, we mark them as used and advance to the next level. Since the i th choice determines i values, $n - 1$ choices are sufficient to determine any pyramid. On the i th choice, we will perform i binary searches in a sorted list of the distances to test whether the ones we need are available.

The space complexity of this procedure is optimal - $\mathbf{O}(n^2)$ - since we need maintain only one pyramid if we remove the i elements we positioned when we backtrack from the i th level. (The state information needed at each level is $\mathbf{O}(1)$.) By assigning $d_{1n} = v[\binom{n}{2}]$, $d_{1(n-1)} = v[\binom{n}{2} - 1]$, and $d_{(n-1)n} = d_{1n} - d_{1(n-1)}$, the search is initialized with $l = 2$ and $r = 1$. ■

While this algorithm takes exponential time in the worst case, notice that the exponent is n instead of $\binom{n}{2}$, the total number of distances. Further, it is independent of the magnitude of the distances, unlike the factoring algorithm.

3.3 Performance in Practice

Since the search is pruned whenever one of the i derived values determined on the i th level does not appear in the set of remaining distances, we usually get much more efficient behavior in practice than suggested by the worst-case analysis. For example, if the distance set arises from n real points in general position, then one of the two choices will be pruned immediately with probability 1, so that the procedure will use $\mathbf{O}(n^2 \log n)$ operations.

Another interesting aspect of our algorithm is that if the point set being reconstructed is reflection invariant, then the (unique) solution will be found in $\mathbf{O}(n^2 \log n)$ time, since the entire pyramid is reflection invariant and thus it doesn't matter which side the largest remaining distance is positioned.

The backtrack algorithm was implemented by Chi-long Lin and Yaw-ling Lin to help assess its performance in practice. This implementation includes a simple ordering heuristic which significantly improves performance on randomly generated problems: Since the span of a position in the pyramid overlaps that of its immediate ancestor in all but one base segment, it stands to reason that of the two possible positions for the largest remaining unpositioned distance, the position with the larger parent value is more likely to be correct. Therefore we investigate this choice first.

The algorithm was run on a series of randomly generated examples with n points, such that the distances between neighboring points on the line were independent uniform integer deviates in $[1, m]$. The backtracker usually takes longer when n increases or m decreases, although the run-time for large n and small m tends to fluctuate drastically. Table 6 presents the average size of the backtrack trees arising from 100 random examples for each (n, m) entry - with the exception of starred entries which summarize 25 random examples.

Table 6 shows that so long as the multiplicity of distances is not too excessive, the ordering heuristic almost always makes the correct decision at

Table 6. Mean Sizes of Backtrack Trees Produced by Reconstruction Algorithm.

number of points n	Maximum Neighbor Distance m						
	100	50	25	15	10	5	2
50	49.04	49.10	49.59	50.86	57.26	134.13	3064.21
100	99.18	99.21	99.65	102.43	113.35	269.05	10012.31
150	149.01	149.37	149.95	152.84	171.65	411.89	16730.36
200	199.08	199.46	199.73	206.12	233.74	559.25	7930.00
250	249.02	249.35	249.47	255.15	281.23	1452.03	31272.60*
300	299.10	299.34	299.95	307.65	348.84	826.69	76930.00*
350	349.01	349.12	349.98	358.80	421.56	1563.34	14822.76*
400	399.21	399.23	399.81	408.25	457.21	1172.73	83016.56*

each step, and false steps are quickly detected. These results show that almost all instances with up to several hundred points can be solved in essentially real-time.

We note that the backtracking algorithm can be modified to reconstruct data with experimental errors, by using interval arithmetic [moor66] instead of testing directly for equality. Skiena and Sundaram [skiena94] demonstrated both analytically and experimentally that the algorithm runs in polynomial-time with high probability if the relative errors of the fragment size are bounded by $O(1/n^2)$. Zhang [zhang94] constructed pathological input that causes the exact-data backtrack algorithm to take exponential time.

3.4 Vector differences in higher dimensions

Crystallographic applications require reconstructing sets from the set of pair-wise vector differences. As with the beltway problem, this generates $n(n-1)$ distinct differences.

A simple general method for solving the vector-difference problem uses a subroutine to solve the one-dimensional turnpike problem. With probability 1, projecting the vector differences onto a random direction yields a 1-1 correspondence between the vector differences and the projected distances, i.e. no repeated distances will occur unless the corresponding vector differences are also repeated. (We remove vectors with negative projections.) Therefore, finding all solutions of the projected turnpike problem, and then replacing each distance by its vector equivalent, gives a superset of the possible vector difference reconstructions. Similar ideas also work for the d -vector difference beltway problems modulo a torus.

This leads to an algorithm for vector difference reconstruction running in

$$O(2^n n \log n + dn^2 A)$$

time where A is the number of solutions found for the projected problem, $An^{1.23}$. It is also possible to generalize the backtrack approaches directly to vector-differences in d -space, using the additional restrictions that the

parallelogram identities of Lemma 3.1 must hold as vector identities. The backtrack trees thus generated will never be larger, and indeed usually will be smaller, than for each coordinatewise 1D projected problem.

The Rosenblatt-Seymour polynomial-factoring algorithm may be generalized for the vector-difference problem by using multivariate polynomials. Instead of polynomials $P(\frac{1}{x})$, $P(x)$ and $Q(x)$, we use a d -variate argument

$$\vec{x} = (x_1, x_2, \dots, x_d)$$

and

$$\frac{1''}{\vec{x}} = \left(\frac{1}{x_1}, \frac{1}{x_2}, \dots, \frac{1}{x_d} \right)$$

The exponents are also vector valued – we use the notation

$$\vec{x}^{\vec{e}} = x_1^{e_1} x_2^{e_2} \dots x_d^{e_d}.$$

With these multidimensional conventions, the entire algorithm now works exactly as in one-dimension, using a routine which factors d -variate polynomials with integer coefficients symbolically into irreducibles. Since this is possible in time polynomial in the number of bits in the input coefficient array if the degree d of the polynomial is fixed [lens82], we conclude that there is a pseudo-polynomial algorithm for solving the vector-difference turnpike problem in fixed dimensions.

3.5 A backtracking algorithm for distance reconstruction in d -dimensional space

The “central simplex” method used in the upper bound proof for Theorem 2.15 may be used to construct a backtracking algorithm for finding all n -point sets in d -space with a given distance multiset.

Theorem 3.4 *If $n > d \geq 2$, then all noncongruent n -point sets in \mathbf{R}^d having a given distance multiset may be found in time*

$$\mathbf{O}(n^{(2d-1)n} e^n)$$

Proof: First, construct the “central simplex” in all possible non-isomorphic ways, using at least $\lfloor d/2 \rfloor$ of the largest distances. Let this $(d-1)$ -simplex lie in the hyperplane $x_1 = 0$ without loss of generality. Then for each central simplex, by backtracking consider the addition of one point at a time. Each time a point is added – by selection of the d distances to the central simplex from among the unused distances and the selection of the sign of its x_1 -coordinate – one may prune the backtrack search if the selected distances and the central simplex do not form a valid d -simplex or if the point thus embedded does not assume valid distances to every other currently embedded point. Pruning also occurs if the d -tuple of distances from the new point to the central simplex, is lexicographically larger than some previous d -tuple.

This eliminates $(n-d)!$ automorphisms and is easily implemented by considering candidate d -tuples in lexical order. Each embedding of a point may be accomplished in $\mathbf{O}(d^3)$ time and then all the needed distances may be found in $\mathbf{O}(n \log n)$ time by n binary searches. Every time all points have been successfully embedded, an output occurs. The bound on the run-time of this algorithm now follows from the proof of Theorem 2.15. ■

In fact, if the points lie in general position, then with probability 1 there will be a unique reconstruction and only one branch of the backtracking will not immediately be pruned once the correct central simplex is found. In this case, the algorithm will run in time

$$\mathbf{O}(n^2 \log n + n^{d^2-d-2\lfloor d/2 \rfloor} + n^{2d}).$$

3.6 Backtrack solution of the Beltway Problem

Let d_{ij} be the distance going clockwise from point i to point j . The points are numbered clockwise and the numbers are taken mod n . These satisfy the triangle equality

$$d_{ij} + d_{jk} = d_{ik}$$

and the complement distance identity

$$d_{ij} + d_{ji} = L$$

where L is the perimeter of the beltway loop. The identities

$$kL = \sum_{0 \leq i < n} d_{i(i+k)} \quad \text{and} \quad \binom{n}{2} L = \sum_{0 \leq i, j < n} d_{ij}$$

enable one to determine L .

Theorem 3.5 *There is a $\mathbf{O}(n^n \log n)$ -time algorithm for finding all possible reconstructions of an n -point circular set from its $(n-1)n$ -element difference multiset. This algorithm runs in optimal $\mathbf{O}(n^2)$ space.*

Proof: The $n(n-1)$ distances must be assigned places in a $(n-1) \times n$ cylindrical “tableaux” to solve the problem. The k th row of the “tableaux” consists of the distances $d_{i(i+n-k)}$ for $i = 0, \dots, n-1$. We will use the partial order inequalities

$$d_{(i-1)(j+1)} \geq d_{ij} \geq d_{(i-1)j}$$

which state that a tableaux element is at least as large as the one lying directly below it and also the one below and to the right of it.

Now we describe a $\mathbf{O}(n^n \log n)$ -time reconstruction algorithm. As usual we fill in the distances largest first so that at any time we need only choose which of the n columns to place the next distance in. After each such choice is made, the triangle equalities allow one to fill in various other distances.

In fact, all distances that follow from repeated application of the triangle equality to a new distance may be determined in $\mathbf{O}(\log n)$ time per distance filled in, in the case of Theorem 3.3. We will see that the total number of automatic fill-ins is always $(n - 1)^2$, so that only $n - 1$ decisions need be made by the backtrack algorithm. Of these, the first choice may always be forced to avoid constructing each of n cyclic shifts of a point set, so we really only have to make $n - 2$ decisions. Since each decision involves at most n choices, and each choice may be made (along with all the fill-ins it causes) in $\mathbf{O}(n^2 \log n)$ time, we may conclude that the total run time is $\mathbf{O}(n^n \log n)$. ■

We may also conclude, upon removing mirror symmetries, that $S_1^*(n) \leq \frac{1}{2}n^{n-2}$, as was claimed in the introduction.

The complexity of the beltway problem is at least as great as that of the turnpike problem, as may be shown by a simple reduction; in view of Theorems 2.11 and 2.12, it is probably greater.

4 The complexity of distance embedding problems

The algorithms discussed in Section 3 all had worst-case exponential time behavior, but we have been unable to prove that the problems are intractable. In fact, there is evidence that the turnpike problem is not NP-complete, as a consequence of the fact that the maximum number of solutions is polynomially bounded.

In this section, we prove some related distance embedding problems are intractable.

The *turnpike problem with experimental error bounds*, which one might also call the *distance-interval reconstruction problem*, is as follows.

Instance: A multiset of $\binom{n}{2}$ distance-intervals and a positive number d . These closed intervals are specified by their integer endpoints.

Question: Is there an n -point set in \mathbf{R}^d which satisfies the distance intervals. In other words, is there a 1-1 correspondence between the distances determined by this putative set and the distance-intervals, with each distance lying inside its corresponding interval.

See [papa76] [saxe79] for some similar NP-complete problems.

To prove its NP-completeness, we will use the following “modified 3-Partition problem” (M3PP): Given a set S of $3K$ positive integers A_1, A_2, \dots, A_{3K} , is there a way to partition S into K disjoint triples such that two of the elements in each triple add up to the third. This problem is strongly NP-complete, as is readily shown by a reduction from “numerical matching with target sums,” which is problem SP17 in [gare88].

Theorem 4.1 *The distance-interval reconstruction problem is strongly NP-complete, even if the distance intervals all have arbitrarily small width compared to the value at their midpoints and $d = 1$.*

Proof: Clearly the problem is in NP. Our reduction from M3PP uses the additive constraints to specify K clusters of three points each, with the clusters widely but regularly spaced apart.

Specifically, we claim that if $K \geq 30$, the following set of $(n-1)n/2$ distances arise from a collinear set of $n = 3K$ points if and only if the modified 3-partition problem is solvable

set number	distance interval midpoint	distance interval width	distance multiplicity	total number of distance-intervals
1	$ p - q B/K, 1 \leq pq \leq K$	$5B$	9	$9(K-1)K/2$
2	$A_m, 1 \leq m \leq 3K$	0	1	$3K$
total				$(3K-1)3K/2$

where $B = \sum A_m$. ■

The *arbitrary-dimensional distance reconstruction problem* is as follows.

Instance: a multiset of $(n-1)n/2$ distances and a number $d > 0$.

Question: Is there an n -point set in d -space which realizes the distances.

For the purposes of the proof, we will assume that the *squares* of the distances are to be specified rationals.

Our reduction uses the original version of 3-partition [gare88, problem SP15] where we are given a set S of $3K$ integers A_m , $1 \leq m \leq 3K$, summing to KB and with $B/4 < A_m < B/2$, and we are to determine whether S may be partitioned into K disjoint triples such that each triple sums to B . We will actually merely require that each triple sums to $\leq B$, which of course implies equalities.

Theorem 4.2 *The arbitrary-dimensional distance reconstruction problem is strongly NP-complete, even if the input distances determine a simplex.*

Proof: The fact that the problem is in NP follows from the fact that we may guess the correspondences between distances and point-pairs, and then deducing the embedding may be accomplished via a matrix (Cholesky) factorization, by using a sequence of $\mathbf{O}(n^3)$ rational operations and n extractions of square roots. Each of the coordinates of the points will be expressible in the form $a + \sqrt{b}$ where a and b are rationals, since the Cholesky factorization does not introduce any numbers involving more than one square root. By the usual subdeterminant arguments [papa82] the total number of binary bits in these numbers will be a polynomial, as will also be the case for all the intermediate quantities arising in the computation.

Define $b_g \equiv K^6 + gK^4 + g^2K$, $1 \leq g \leq K$. We claim if $K \geq 2$, then the following set of $(n-1)n/2$ distances, $n = 3K + 2$, is embeddable as a simplex in \mathbf{R}^d , $d = 3K + 1$, if and only if the 3-partition problem below is solvable.

set number	squares of distances	distance multiplicity	total number of distances
1	$(2K^8)^2$	1	1
2a	$1 + (K^8 + b_g)^2, 1 \leq g \leq K$	3	$3K$
2b	$1 + (K^8 - b_g)^2, 1 \leq g \leq K$	3	$3K$
3	$(b_p - b_q)^2 + 2, 1 \leq p < q \leq K$	9	$9(K-1)K/2$
4	$2 - 2\cos(2\pi A_m/B), 1 \leq m \leq 3K$	1	$3K$
total			$(3K+2)(3K+1)/2$

Here $\underline{\cos}(x)$ denotes the best rational approximation p/q to $\cos(x)$ subject to the restrictions that $\underline{\cos}(x) = \cos(y)$ for some y with $0 \leq y < x$, and $|p| \leq q \leq 10B$. (This function may be computed quickly using Taylor series and regular continued fractions [hard79] and has the property that $x - y < 1/(10B)$ if $0 < x \leq \pi$.) Note that if the inputs A_m and B to the 3-partition problem are expressed as unary integers, we may still generate this squared-distance set, and output it in the form of unary rationals, in polynomial time. If both the input and the output are expressed using binary integers, the reduction is even easier.

To prove the claim, regard \mathbf{R}^d as $\mathbf{R} \times \mathbf{R}^{3K}$ for convenience. The largest distance is in set 1, so without loss of generality we may take the two furthest-apart points (the “endpoints”) to lie at $(\pm K^8; \vec{0})$. Now no distance from sets numbered ≥ 3 can be added to any distance from sets numbered ≥ 2 to obtain a distance $\geq 2K^8$, so by the triangle inequality the distances from sets 2a, 2b are precisely the distances involving one endpoint and one other point.

In fact these distances must be paired – each distance from set 2a represents the distance from a point to one endpoint and the corresponding (same m) distance from 2b is the distance to the other endpoint – since any other pairing would entail some distance-pair sum $2K^8$. Therefore we see that the $3K$ non-endpoints have coordinates $(b_m; \vec{x}_{m,i}), 1 \leq m \leq K, 1 \leq i \leq 3$, and each $\vec{x}_{m,i}$ is a $3K$ -vector with unit norm. (The fact that we may use b_m instead of $\pm b_m$ is because none of the distances in sets numbered ≥ 3 have magnitude as large as $\min_{p,q} b_p + b_q$. Now we see the distances in set 3 must be precisely the distances between points with coordinates $(b_p, \vec{x}_{p,i})$ and $(b_q, \vec{x}_{q,j}), p \neq q$, which forces the orthogonality of $\vec{x}_{p,i}$ and $\vec{x}_{q,j}$. We may now assume without loss of generality that $\vec{x}_{p,j}$ is of the form $(0, 0, \dots, 0, e_{p,j}, f_{p,j}, g_{p,j}, 0, 0, \dots, 0)$ where the e, f , and g are in the $3p^{\text{th}}, 3p+1^{\text{th}}$, and $3p+2^{\text{th}}$ positions respectively and $e_{p,j}^2 + f_{p,j}^2 + g_{p,j}^2 = 1$. This is since the three points $\vec{x}_{1,i}$ must lie in some 3-space, without loss of generality [by a rotation of the coordinate system, if necessary] the one spanned by

$$(1, 0, 0, 0, \dots, 0), (0, 1, 0, 0, \dots, 0), \text{ and } (0, 0, 1, 0, \dots, 0),$$

and then the three points $\vec{x}_{2,j}$ must lie in an orthogonal 3-space, without loss of generality the one spanned by

$$(0, 0, 0, 1, 0, 0, 0, \dots, 0), (0, 0, 0, 0, 1, 0, 0, \dots, 0), \text{ and } (0, 0, 0, 0, 0, 1, 0, \dots, 0),$$

and so on. Now, all this is possible if and only if the 3 angles determined by the 3 points $(b_m; \vec{x}_{m,i}), 1 \leq i \leq 3$, in each orthogonal sphere sum to $\leq 2\pi$.

Such an assignment of angles is possible (considering set 4 and the law of cosines) if and only if the 3-partition problem is solvable. ■

5 Conclusions

We have presented several new results concerning algorithms for turnpike reconstruction and the number of homometric sets. In particular, our backtracking algorithm is sufficiently simple and fast in practice that it should be suitable for almost all turnpike problems arising in applications. Several outstanding open problems remain:

1. Improve our bounds on the constant $0.810 < C < 1.233$ defined by

$$C = \limsup_{n \rightarrow \infty} \frac{\ln H_1(n)}{\ln n}$$

2. What is the least value of n such that there are at least 2^r non-congruent homometric n -point sets? Improved bounds for small r can tighten lower bound of problem 1 via Lemma 2.2.
3. Find n -point homometric pairs, $n \geq 7$, with no repeated distances. Piccard [picc39] believed she had shown that no such pair could exist, but Bloom [bloo77] found a counterexample for $n = 6$.
4. What is the asymptotic behavior of $S_1(n)$? It could range from subexponential to superexponential.
5. Find a strongly-polynomial algorithm for turnpike reconstruction. No NP-complete problem is known with the property that the number of solutions is $2^{o(n^{o(1)})}$. Therefore, in light of Theorem 2.4 it seems doubtful that the turnpike problem is NP-hard. Further, find a reasonable algorithm for beltway reconstruction.

References

- [agar89] P. Agarwal, M. Sharir, P. Shor: *Sharp upper and lower bounds on the length of general Davenport-Schinzel sequences*, J. Comb. Theory A, 52 (1989) 228-274.
- [ally88] L. Allison and C. N. Yee: *Restriction Site Mapping is in Separation Theory*, Comput. Appl. Biol. Sci. 4,1 (1988) 97-101
- [bann83] E.Bannai, E.Bannai, D.Stanton: *An upper bound for the cardinality of s -distance subset of Euclidean space II*, Combinatorica 3,2 (1983) 147-152.
- [bate49] P.T. Bateman: *Note on the coefficients of the cyclotomic polynomial*, Bull. AMS 55 (1949) 1180-1181

- [bell88] Bernard Bellon: *Construction of Restriction Maps*, Comput. Appl. Biol. Sci. 4,1 (1988) 111-115
- [bloo77] G.S. Bloom: *A counterexample to a theorem of Piccard*, J. Comb. Theory A22 (1977) 378-379
- [bull61] R.K. Bullough: *On homometric sets*, I. Acta Cryst. 14 (1961) 257-268
II. Acta Cryst. 17 (1964) 295-308
- [buer77] M.J. Buerger: *Exploration of cyclotomic point sets in tautoeikonic complementary pairs*, Z. Kristallogr. 145 (1977) 371-411
- [cael80] T. Caelli: *On generating spatial configurations with identical interpoint distance distributions*, Proc. 7th Australian Conf. Combinatorial Math. Newcastle Australia August 1979 = Springer Lecture Notes in Math. 829 (1980) 69-75
- [chie79] C. Chieh: *Analysis of cyclotomic sets*, Z. Kristallogr. 150 (1979) 261-277
- [clar90] K.L.Clarkson, H. Edelsbrunner, L.J. Guibas, M.Sharir, E.Welzl: *Combinatorial complexity bounds for arrangements of curves and spheres*, Discrete & Comput. Geom. 5,2 (1990) 99-160.
- [daki00] T. Dakic: *On the Turnpike Problem* PhD Thesis, Simon Fraser University, 2000.
- [dels77] P.Delsarte, J.Goethals, J.J.Seidel: *Spherical codes and designs*, Geometriae Dedicata 6,3 (1977) 363-388
- [dixk88] T.I. Dix and D.H. Kieronska: *Errors between sites in restriction site mapping*, Comput. Appl. Biol. Sci. 4,1 (1988) 117-123
- [doyv71] J. Doyen & M. Vandensavel: *Non-isomorphic Steiner quadruple systems*, Bull. Soc. Math. Belgium 23 (1971) 393-410
- [erdo46] P. Erdős: *On sets of distances of n points*, AMM 53 (1946) 248-250
- [fila95] M. Filaseta, *The irreducibility of all but finitely many Bessel polynomials*. Acta Math. 174 (1995), no. 2, 383-397.
- [fila99] M. Filaseta: *On the factorization of polynomials with small Euclidean norm*, pp. 143-163 in *Number theory in progress 1* (volume 2 of 2, e.d K.Györy et al.) W. de Gruyter 1999
- [fink83] A.M. Finkelstein, O. M. Kosheleva, and V. JA. Kreinovic: *On the uniqueness of image reconstruction from the amplitude of radiointerferometric response*. Astro. Sp. Sci. 92 (1983) 31-36
- [gare78] M.R. Garey & D.S. Johnson: *Computers and Intractability: a guide to the theory of NP-completeness*, Freeman 1978
- [gilb74] E.N. Gilbert & L.A. Shepp: *Textures for discrimination experiments*, Bell Laboratories Murray Hill NJ, TM-74-1218-6, TM-74-1215-15 April 15,1974. Filing case 20878.
- [golu83] G.H. Golub & C. Van Loan: *Matrix Computations*, Johns Hopkins University Press 1983
- [hall56] M. Hall Jr.: *A survey of difference sets*, Proc. AMS 7 (1956) 975-986

- [hall67] M. Hall Jr.: *Combinatorial Theory*, Wiley 1967
- [hard79] G.H. Hardy & E.M. Wright: *Introduction to the theory of numbers*, Oxford University Press 1979
- [hose54] R. Hosemann & S.N. Bagchi: *On homometric structures*, Acta. Cryst. 7 (1954) 237-241
- [hose62] R. Hosemann & S.N. Bagchi: *Direct analysis of diffraction by matter*, North-Holland 1962
- [humn01] Special issues on human genome: *Nature* 409 (15 Feb 2001); *Science* 291, 5507 (16 Feb 2001).
- [knut76] D.E. Knuth: *Big omicron and big omega and big theta*, SIGACT News 8,2 (April-June 1976) 18-24
- [koor76] T.M.Koornwinder: *A note on the absolute bound for systems of lines*, Indag. Math. 38,2 (1976) 152-153.
- [land87] Susan Landau: *Factoring polynomials quickly*, Notices Amer. Math. Soc. 34,1 (1987) 3-8
- [lemk88] P. Lemke and M. Werman, *On the complexity of inverting the auto-correlation function of a finite integer sequence, and the problem of locating n points on a line, given the unlabeled distances between them*, manuscript, 1988.
- [lens82] Arjen K. Lenstra, H.W. Lenstra Jr., Laszlo Lovasz: *Factoring polynomials with rational coefficients*, Math. Ann. 261 (1982) 515-534
- [lind74] C.C. Lindner: *On the construction of non-isomorphic Steiner quadruple systems*, Colloq. Math. 29 (1974) 303-306
- [ljun60] W. Ljungren: *On the irreducibility of certain trinomials and quadrinomials* Math. Scandinav. 8 (1960) 65-70; also see H. Tverberg 121-126 same issue.
- [mahl76] K. Mahler: *Lectures on transcendental numbers*, Lecture Notes in Math. 546 Springer-Verlag 1976
- [mcla61] Dan McLachlan, Jr.: *Similarity function for pattern recognition*, J. Appl. Phys. 32 (1961) 1795-1796
- [moor66] R.E. Moore: *Interval Analysis*, Prentice-Hall 1966
- [newb93] L. Newberg and D. Naor, *A Lower Bound on the Number of Solutions to the Probed Partial Digest Problem* Advances in Applied Mathematics 14 (1993).
- [pand01] G. Pandurangan and H. Ramesh, *The Restriction Mapping Problem Revisited*, Journal of Computer and System Sciences (JCSS), to appear
- [papa76] C.H. Papadimitriou: *The NP-completeness of the bandwidth minimization problem*, Computing 16 (1976) 263-270
- [patt35] A. L. Patterson: *A direct method for the determination of the components of interatomic distances in crystals*, Zeitschr. Krist. 90 (1935) 517-542
- [patt44] A.L. Patterson: *Ambiguities in the X-ray analysis of crystal structures*, Phys. Review 65 (1944) 195-201.

- [picc39] Sophie Piccard: *Sur les ensembles de Distances des Ensembles de points d'un Espace Euclidean*, Mem. Univ. Neuchatel 13 (Neuchatel Switzerland 1939)
- [prep85] F. Preparata and M. Shamos, *Computational Geometry*, Springer-Verlag, New York, 1985.
- [rau80] V.G. Rau, L.G. Parkhomov, V.V. Ilyukhin, N.V. Belov: *On the calculation of possible Patterson cyclotomic sets*, I: Doklady Akademii Nauk SSSR 255,4 (1980) 859-861; II: ibid. 255,5 (1980) 1110-1113. (Both in Russian.)
- [rhee88] Gwangsoo Rhee: *DNA restriction mapping from random-clone data*, Technical Report WUCS-88-18, Department of Computer Science, Washington University St. Louis MO 1988
- [ross82] Joseph Rosenblatt & Paul Seymour: *The structure of homometric sets*, SIAM J. Alg. Disc. Methods 3,3 (1982) 343-350
- [saxe79] James B. Saxe: *Embeddability of weighted graphs in k -space is strongly NP-hard*, Proc. 19th Allerton Conference on Computers, Controls, and Communications 19, Urbana, IL (1979) 480-489
- [sham77] M.I. Shamos: *Problems in computational geometry*, Unpublished manuscript, Carnegie Mellon University, Pittsburgh, PA 1977.
- [sing38] James A. Singer: *A theorem in finite projective geometry and some applications to number theory*, Trans. Amer. Math. Soc. 43 (1938) 377-385
- [skie90] S. S. Skiena, W. D. Smith, and P. Lemke: *Reconstructing sets from interpoint distances (extended abstract)* Proc. Sixth ACM Symposium on Computational Geometry (1990) 332-339.
- [skie94] S. S. Skiena and G. Sundaram: *A Partial Digest Approach to Restriction Site Mapping* Bulletin of Mathematical Biology, 56 (1994) 275-294.
- [smyt71] C. J. Smyth: *On the product of the conjugates outside the unit circle of an algebraic integer*, Bull. London Math. Soc. 3 (1971) 169-175
- [soly01] J. Solymosi & Cs. D. Toth: *Distinct distances in the Euclidean plane*, Discr. & Comput. Geom. 25,4 (2001) 629-634.
- [spec49] W. Specht: *Abschätzungen der Wurzeln algebräischer Gleichungen*, Math. Zeit. 52 (1949) 310-321
- [stef78] Mark Stefik: *Inferring DNA structures from segmentation data*, Artificial Intelligence 11 (1978) 85-114
- [tuff88] P. Tuffery, P. Dessen, C. Mugnier, S. Hazout: *Restriction Map Construction Using a 'Complete Sentences Compatibility' Algorithm*, Comput. Appl. Biol. Sci. 4,1 (1988) 103-110
- [zhan94] Z. Zhang: *An exponential example for partial digest mapping algorithm*, J. Computational Biology 1,3 (1994) 235-239.

About Authors

Paul Lemke did this work as a member of the Department of Mathematical Sciences, Troy NY 12180-3590.

Steven Skiena (to whom correspondence should be addressed) is at the Department of Computer Science, SUNY Stony Brook, Stony Brook, NY 11794-4400; *skiena@cs.sunysb.edu*.

Warren D. Smith did most of this work as a member of AT&T Bell Laboratories and the NEC Research Laboratory. His current affiliation is DIMACS, Rutgers University, 96 Frelinghuysen Road, Piscataway, NJ 08854-8018.

Acknowledgments

We thank Bernard Chazelle for first bringing this problem to our attention. A. Blokhuis, Dan Gusfield, A.M. Odlyzko, J.C. Lagarias, N.J.A. Sloane, Pong-Chi Chu, Arch Robison, and Carla Savage provided insights and helpful comments. Discussions with Gene Myers, Webb Miller, and John Turner helped clarify the connection to restriction site mapping.