# Lower envelopes and Davenport-Schinzel sequences

Michiel Smid*

October 30, 2003

## 1   Introduction

Suppose we are given a set $S = \{p_1, p_2, \ldots, p_n\}$ of $n$ points in the plane. Each of these points moves at a constant speed along a straight line. Hence, for each $i$ with $1 \leq i \leq n$, we can write the position of point $p_i$ at time $t$ as

$$p_i(t) = (a_i + v_i t, b_i + w_i t).$$

We want to design a data structure such that we can answer the following type of query: Given any time $t$, find a point of $S$ that, at time $t$, is closest to the origin. That is, we want to find an index $j$ such that

$$|p_j(t)| = \min\{|p_i(t)| : 1 \leq i \leq n\}. \tag{1}$$

Here, $|x|$ denotes the length of the vector $\vec{x}$ or, equivalently, the Euclidean distance between $x$ and the origin.

Here is a solution to this problem. First observe that

$$|p_i(t)| = \sqrt{(a_i + v_i t)^2 + (b_i + w_i t)^2}.$$

It is clear that instead of finding an index $j$ for which (1) holds, we can as well find a $j$ such that

$$|p_j(t)|^2 = \min\{|p_i(t)|^2 : 1 \leq i \leq n\}.$$

---
*School of Computer Science, Carleton University, Ottawa, Ontario, Canada K1S 5B6. E-mail: `michiel@scs.carleton.ca`.

1

For each $i$ with $1 \leq i \leq n$, we define

$$f_i(t) := |p_i(t)|^2 = (a_i^2 + b_i^2) + 2(a_i v_i + b_i w_i)\, t + (v_i^2 + w_i^2)\, t^2.$$

Then, in a query we have to find an index $j$ for which

$$f_j(t) = \min\{f_i(t) : 1 \leq i \leq n\}.$$

The function $f_i(t)$ is a parabola in $t$. Draw the graphs of all these $n$ functions. Given any time $t$, how do we find the point $p_j$ that, at time $t$, is closest to the origin? Here is the answer: Starting at the point $(t, 0)$ on the $x$-axis, walk vertically upwards, until we encounter the first parabola. The index of this parabola gives the answer to the query.

This leads to the following data structure for solving our problem. For each $t \in \mathbb{R}$, let

$$F(t) := \min\{f_i(t) : 1 \leq i \leq n\},$$

i.e., $F$ is the pointwise minimum of the functions $f_1, f_2, \ldots, f_n$. The graph of $F$ is called the *lower envelope* of the $f_i$'s. It consists of parabola segments, each of which is bounded by two vertices. The leftmost (resp. rightmost) segment is bounded by a vertex with $x$-coordinate $-\infty$ (resp. $+\infty$).

Our data structure is just the sequence of all these vertices, except the one with $x$-coordinate $+\infty$, sorted from left to right. For each such vertex $v$, let $s(v)$ be the parabola segment of $F$ that has $v$ as its left endpoint. Then we store with vertex $v$ the index of the parabola that contains the segment $s(v)$.

Given this data structure, a query is answered by performing a binary search with $t$ among the $x$-coordinates of the vertices. If $v$ is the vertex that is immediately to the left of $t$, then the index stored with $v$ is exactly the index of the point that, at time $t$, is closest to the origin.

To analyze the complexity of this data structure, we have to answer the following questions:

1. How many parabola segments does the lower envelope $F$ have?

2. How do we compute the vertices and parabola segments of the lower envelope $F$?

Observe that the query time is logarithmic in the number of segments of $F$.

Consider the first question. Each finite vertex of $F$ is an intersection of two parabolas. Since we have $n$ parabolas, and each two of them intersect at

most twice, the number of vertices of $F$ having a finite $x$-coordinate is less than or equal to $2\binom{n}{2} = n^2 - n$. Hence, $F$ consists of at most $n^2 - n + 1$ parabola segments. Can there be this many segments? The answer is "no": We will prove in Section 2 that $F$ consists of at most $2n-1$ parabola segments.

**Exercise 1** Construct a set of $n$ parabolas whose lower envelope has exactly $2n - 1$ parabola segments.

**Exercise 2** Consider the lower envelope of $n$ non-vertical lines. Prove that it has at most $n$ edges.

## 2    Davenport-Schinzel sequences

**Definition 1** *Let $n \geq 1$, $s \geq 1$ and $m \geq 1$ be integers. A sequence $U = (u_1, u_2, \ldots, u_m)$ is called an $(n, s)$-Davenport-Schinzel sequence, or $DS(n, s)$-sequence for short, if*

1. *$u_i \in \{1, 2, \ldots, n\}$ for all $i$ with $1 \leq i \leq m$,*

2. *$u_i \neq u_{i+1}$ for all $i$ with $1 \leq i < m$,*

3. *there do not exist $s + 2$ indices $1 \leq i_1 < i_2 < \ldots < i_{s+2} \leq m$ such that*
   *$u_{i_1} = u_{i_3} = u_{i_5} = \ldots,$*
   *$u_{i_2} = u_{i_4} = u_{i_6} = \ldots, \;\; and$*
   *$u_{i_1} \neq u_{i_2}.$*

For example, $U$ is a $DS(n, 1)$-sequence, if (i) all symbols of $U$ are integers between 1 and $n$, (ii) adjacent symbols of $U$ are distinct, and (iii) for all $a \neq b$, $1 \leq a \leq n$, $1 \leq b \leq n$, the sequence $U$ does not contain a subsequence of the form $a \ldots b \ldots a$.

Davenport-Schinzel sequences were introduced in 1965 in connection with differential equations. They were re-invented in 1985 by Atallah in a paper on computational geometry for moving objects. (Our example in Section 1 is from Atallah's paper.) These sequences have been used to analyze many problems in discrete and computational geometry.

**Exercise 3** Prove that

1. the lower envelope of a set of $n$ non-vertical lines gives rise to a $DS(n, 1)$-sequence, and

2. the lower envelope $F$ introduced in Section 1 corresponds to a $DS(n, 2)$-sequence.

We already defined the lower envelope of parabolas. We now extend this notion in the obvious way.

**Definition 2** *Let $f_1, f_2, \ldots, f_n$ be a collection of functions. Assume that for each $i$ with $1 \leq i \leq n$, the function $f_i$ is defined on the interval $[l_i, r_i]$, where $-\infty \leq l_i < r_i \leq \infty$. The lower envelope of $f_1, f_2, \ldots, f_n$ is the graph of the function*

$$F(t) := \min\{f_i(t) : 1 \leq i \leq n \text{ and } f_i(t) \text{ is defined}\}.$$

The lower envelope consists of vertices and segments of the $f_i$'s. We refer to these segments as *edges*. If we walk along the lower envelope from left to right, and write down the indices of the $f_i$'s we visit, then we get a sequence whose symbols are integers between 1 and $n$. Clearly, adjacent symbols in this sequence are distinct.

**Theorem 1** *Let $f_1, f_2, \ldots, f_n$ be any collection of continuous functions, each of which is defined on $\mathbb{R}$. Assume that any two of these functions intersect at most $s$ times. Then the lower envelope of the $f_i$'s forms a $DS(n, s)$-sequence.*

*Conversely, let $U$ be any $DS(n, s)$-sequence. There exist $n$ continuous functions, all defined on $\mathbb{R}$, any two of which intersect at most $s$ times, and whose lower envelope is equal to $U$.*

**Proof.** The first claim follows from the fact that each pair of functions intersect at most $s$ times. The proof of the second claim is not given here. ∎

**Example 1** If $f_1, f_2, \ldots, f_n$ is a collection of polynomials, each of degree less than or equal to $s$, then their lower envelope forms a $DS(n, s)$-sequence. In particular, the lower envelope $F$ introduced in Section 1 forms a $DS(n, 2)$-sequence.

In the following two exercises, we generalize Theorem 1 to functions that are defined only on a subinterval of the real numbers.

**Exercise 4** Consider a set of $n$ non-vertical line segments. Prove that their lower envelope forms a $DS(n, 3)$-sequence.

**Exercise 5** Let $f_1, f_2, \ldots, f_n$ be a collection of continuous functions. Assume that for each $i$ with $1 \leq i \leq n$, the function $f_i$ is defined on the interval $[l_i, r_i]$, where $-\infty < l_i < r_i < \infty$. Also, assume that any two of these functions intersect at most $s$ times. Prove that the lower envelope of the $f_i$'s forms a $DS(n, s+2)$-sequence. (The converse of this claim, as in the second claim of Theorem 1, also holds.)

Recall that we wanted to bound the number of edges on the lower envelope $F$ in our problem on moving points. By Theorem 1, it suffices to give an upper bound on the length of any $DS(n, 2)$-sequence. In general, the number of edges on the lower envelope of a collection of $n$ polynomials, each of degree less than or equal to $s$, is bounded from above by the maximum length of any $DS(n, s)$-sequence.

**Definition 3** *Let $n \geq 1$ and $s \geq 1$ be integers. We denote by $\lambda_s(n)$ the maximum length of any $DS(n, s)$-sequence.*

How can we estimate $\lambda_s(n)$? The second claim in Theorem 1 implies that

$$\lambda_s(n) \leq s\binom{n}{2} + 1 = sn(n-1)/2 + 1.$$

Hence, if $s$ is a constant, then $\lambda_s(n) = O(n^2)$. This upper bound is much too crude. (We announced this already for $s = 2$.) The claim is that $\lambda_s(n)$ is linear in $n$ for $s = 1$ and $s = 2$, and *almost* linear in $n$ for any constant $s \geq 3$.

**Theorem 2** *We have*

1. *$\lambda_1(n) = n$, and*

2. *$\lambda_2(n) = 2n - 1$.*

**Proof.** Let $U = (u_1, u_2, \ldots, u_m)$ be any $DS(n, 1)$-sequence. Then, by definition, $U$ does not contain a subsequence of the form $a \ldots b \ldots a$, where $a \neq b$.

We claim that all symbols of $U$ are pairwise distinct. Assume to the contrary that there are indices $i$ and $j$ such that $1 \leq i < j \leq m$ and $u_i = u_j$. Let $a := u_i$ and $b := u_{i+1}$. Since $u_i \neq u_{i+1}$ by the definition of Davenport-Schinzel sequence, we have $i + 1 < j$. Hence, the sequence $U$ contains a subsequence of the form $ab \ldots a$, where $a \neq b$. This is a contradiction.

5

Hence, all symbols of $U$ are pairwise distinct. Since each symbol is an integer between 1 and $n$, it follows that the length of $U$ is less than or equal to $n$. This proves that $\lambda_1(n) \leq n$. It is clear that $(1, 2, 3, \ldots, n)$ is a $DS(n, 1)$-sequence. Therefore, $\lambda_1(n) \geq n$.

To prove the second claim, we first observe that

$$(n, n-1, n-2, \ldots, 3, 2, 1, 2, 3, \ldots, n-2, n-1, n)$$

is a $DS(n, 2)$-sequence. Hence, $\lambda_2(n) \geq 2n - 1$.

It remains to prove that $\lambda_2(n) \leq 2n - 1$. The proof of this claim is by induction on $n$. Since $\lambda_2(1) = 1$, the claim holds for $n = 1$. Now let $n > 1$, and assume that $\lambda_2(i) \leq 2i - 1$ for all $i$ with $1 \leq i < n$.

Let $U = (u_1, u_2, \ldots, u_m)$ be any $DS(n, 2)$-sequence. If we can show that $m \leq 2n - 1$, then it follows that $\lambda_2(n) \leq 2n - 1$.

Consider the first symbol $a := u_1$ of $U$. Let $U'$ be the sequence obtained from $U$ by deleting the first symbol. Hence, we have $U = aU'$. We distinguish two cases.

**Case 1:** $a$ does not occur in $U'$.

Then $U'$ contains at most $n - 1$ different symbols. Hence, $U'$ is a $DS(n - 1, 2)$-sequence and, by the induction hypothesis, its length is less than or equal to $\lambda_2(n - 1) \leq 2(n - 1) - 1 = 2n - 3$. But then, the length of $U$ itself is less than or equal to $1 + (2n - 3) = 2n - 2 \leq 2n - 1$.

**Case 2:** $a$ occurs in $U'$.

We write $U'$ as $U_1 a U_2$, where $a$ does not occur in $U_1$. Hence, we have $U = aU_1aU_2$. Observe that $U_1$ is not empty. Let $i$ be the number of distinct symbols that occur in $U_1$. Then $U_1$ is a $DS(i, 2)$-sequence and, by the induction hypothesis, its length is less than or equal to $2i - 1$.

We claim that no symbol occurs both in $U_1$ and $U_2$. To prove this, assume there is a symbol $b$ that occurs both in $U_1$ and $U_2$. We know already that $a$ does not occur in $U_1$. Hence, $a \neq b$. But then, the sequence $U$ contains a subsequence of the form $a \ldots b \ldots a \ldots b$. This is a contradiction.

It follows that the sequence $aU_2$ contains at most $n - i$ distinct symbols. Hence, it is a $DS(n - i, 2)$-sequence. Since $i > 0$, the induction hypothesis implies that the length of $aU_2$ is less than or equal to $2(n - i) - 1$. This implies the following bound on the length of the entire sequence $U$:

$$|U| = 1 + |U_1| + |aU_2| \leq 1 + (2i - 1) + (2(n - i) - 1) = 2n - 1.$$

This completes the proof. ∎

**Exercise 6** The sequence $U_2$ defined in the proof of Theorem 2 may be empty. Convince yourself that the proof is nevertheless correct.

**Corollary 1** *Let $S$ be a set of $n$ points in the plane, each of which is moving at a constant speed along a straight line. There is a data structure of size $O(n)$ such that for any query value $t$, we can in $O(\log n)$ time find a point of $S$ that—at time $t$—is closest to the origin. This data structure can be built in time $O(n \log n)$.*

**Proof.** The data structure was described in Section 1. Its size is bounded by $\lambda_2(n)$ and its query time is logarithmic in the size. The bound on the building time is left as Exercise 7. ∎

For any constant $s \geq 3$, the value $\lambda_s(n)$ depends on the inverse $\alpha(n)$ of the Ackermann function. This function will be defined in Section 5. We remark here that $\alpha(n)$ is increasing, but *extremely* slowly: Although $\alpha(n)$ goes to infinity as $n$ does, we have $\alpha(n) \leq 4$ for all

$$n \leq \underbrace{2^{2^{2^{\cdot^{\cdot^{\cdot^{2^{2048}}}}}}}}_{2048 \; 2's} \; .$$

Hence, $\alpha(n)$ is *almost* bounded by a constant.

**Theorem 3** *We have*

1. $\lambda_3(n) = \Theta(n \cdot \alpha(n))$,

2. $\lambda_4(n) = \Theta(n \cdot 2^{\alpha(n)})$,

3. *for all constants $s$, $\lambda_{2s}(n) = n \cdot 2^{O(\alpha(n)^{s-1})}$,*

4. *for all constants $s$, $\lambda_{2s+1}(n) = n \cdot \alpha(n)^{O(\alpha(n)^{s-1})}$,*

5. *for all constants $s$, $\lambda_{2s}(n) = n \cdot 2^{\Omega(\alpha(n)^{s-1})}$.*

We will not prove this theorem. (I guess that you, dear student, do not mind!) In short, for any constant $s \geq 3$, the function $\lambda_s(n)$ is *almost* linear in $n$, or, put differently, $\lambda_s(n)$ is only *slightly* superlinear.

**Exercise 7** Let $f_1, f_2, \ldots, f_n$ be polynomials of degree less than or equal to $s$, where $s$ is a constant. Give an algorithm that computes the lower envelope of these functions in $O(\lambda_s(n) \log n)$ time. (*Hint*: Think of merge-sort.)

7

# 3 An example: the lower envelope of line segments

Let $S$ be a set of $n$ non-vertical line segments in the plane. We are interested in the *size* of the lower envelope of these segments. Here, size refers to the number of (finite) vertices. (Clearly, the number of edges is proportional to the size.)

By Exercise 4, the size of the lower envelope of $S$ is bounded from above by the maximum length of any $DS(n, 3)$-sequence, which is $O(n \cdot \alpha(n))$. Moreover, it can be shown that there exists a set of $n$ line segments whose lower envelope has size $\Omega(n \cdot \alpha(n))$. The proofs of both the upper and lower bound are beyond the scope of this course. Instead, we will prove the following weaker result.

**Theorem 4** *Let $S$ be any set of $n$ non-vertical line segments in the plane. The lower envelope of $S$ has size $O(n \log n)$.*

We assume that the segments of $S$ are in *general position*. This means that

1. all $2n$ segment endpoints and all intersections between pairs of segments have different $x$-coordinates,

2. no two segments overlap, i.e., have more than one point in common, and

3. no three segments have a common intersection.

The following beautiful proof is due to Boaz Tagansky. It appeared in the Proceedings of the 11-th Annual ACM Symposium on Computational Geometry, in June 1995.

First some definitions. The lower envelope of $S$ has two types of vertices. First, there are vertices that are defined by endpoints of segments. These will be called *level-0 outer vertices of $S$*. Second, there are vertices that are defined by intersections of segments. These are called *level-0 inner vertices of $S$*. Finally, we define another type of vertex, which is in fact not a vertex of the lower envelope: An intersection $p$ of two segments of $S$ is called a *level-1 inner vertex of $S$*, if there is exactly one segment in $S$ that is strictly below $p$. (See Figure 1.)
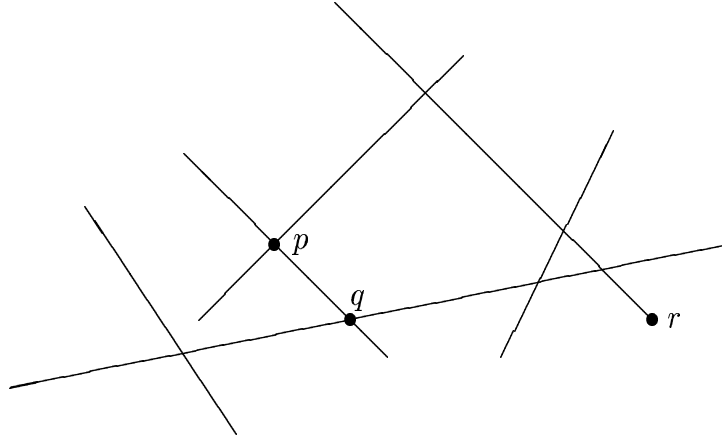
Figure 1: $p$ is a *level-1 inner vertex*, $q$ is a *level-0 inner vertex*, and $r$ is a *level-0 outer vertex*.

---

It is clear that the lower envelope of $S$ has at most $2n$ level-0 outer vertices. Hence, in order to prove Theorem 4, we have to show that there are $O(n \log n)$ level-0 inner vertices.

For $i \in \{0, 1\}$, let $C_i(S)$ be the number of level-$i$ inner vertices of $S$, and let $C_i(n)$ be the maximum value of $C_i(S)$ over all sets $S$ of $n$ non-vertical line segments that are in general position.

Hence, our goal is to prove that $C_0(n) = O(n \log n)$. The proof consists of three steps.

**Step 1:** We give an upper bound for $C_0(S)$ in terms of $C_1(S)$ and an additional factor that can easily be estimated.

Let $p$ be any level-0 inner vertex of $S$. Consider the vertical line through $p$. We move this line to the right, and stop as soon as it

1. encounters a segment endpoint, or

2. encounters a level-1 inner vertex of $S$.

Refer to Figure 2. The main observation is that the line does not encounter any other level-0 inner vertex of $S$ before any of the cases 1. or 2. occurs. Hence, each segment endpoint and each level-1 inner vertex of $S$ is "reached", if at all, from a unique level-0 inner vertex.
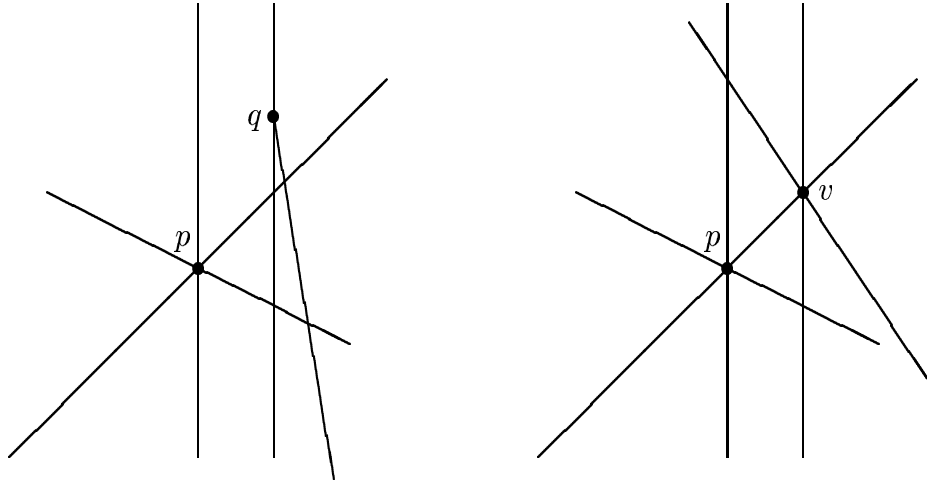
9

Figure 2: *p is a level-0 inner vertex. If we move the vertical line through p to the right, then in the left figure, it first hits at segment endpoint q. In the right figure, the line first hits at level-1 inner vertex v.*

There are $2n$ segment endpoints, and $C_1(S)$ level-1 inner vertices. Therefore,

$$C_0(S) \leq 2n + C_1(S). \qquad (2)$$

**Step 2:** We apply a random sampling analysis.

Take a *random* segment in $S$, and remove it. Let $R$ be the resulting set of $n - 1$ segments. Observe that $C_0(R)$—the number of level-0 inner vertices of $R$—is a random variable. We are interested in the *expected value* of this variable. Observe that a level-0 inner vertex of $R$ is

1. either a level-0 inner vertex of $S$,

2. or a level-1 inner vertex of $S$.

Let $A$ denote the random variable whose value is equal to the number of level-0 inner vertices of $R$ that are also level-0 inner vertex of $S$. Let $B$ denote the random variable whose value is equal to the number of level-0 inner vertices of $R$ that are level-1 inner vertex of $S$. Then

$$C_0(R) = A + B.$$

What is the expected value of $A$? Let $p$ be any level-0 inner vertex of $S$, and let $s_1$ and $s_2$ be the two segments of $S$ that intersect in $p$. Then, $p$ is a level-0

10

inner vertex of $R$ if and only if $s_1$ and $s_2$ are both contained in $R$. Hence, the probability that $p$ is a level-0 inner vertex of $R$ is equal to $(n-2)/n$. This implies that

$$E(A) = \frac{n-2}{n} \cdot C_0(S).$$

In a similar way, we get

$$E(B) = \frac{1}{n} \cdot C_1(S).$$

(Convince yourself that this is true!) Then, by the linearity of expectation, we get

$$E(C_0(R)) = E(A + B) = E(A) + E(B) = \frac{n-2}{n} \cdot C_0(S) + \frac{1}{n} \cdot C_1(S). \quad (3)$$

**Step 3:** We combine (2) and (3):

$$
\begin{aligned}
\frac{n-1}{n} \cdot C_0(S) &= \frac{n-2}{n} \cdot C_0(S) + \frac{1}{n} \cdot C_0(S) \\
&\leq \frac{n-2}{n} \cdot C_0(S) + \frac{1}{n} \left( 2n + C_1(S) \right) \\
&= 2 + \frac{n-2}{n} \cdot C_0(S) + \frac{1}{n} \cdot C_1(S) \\
&= 2 + E(C_0(R)).
\end{aligned}
$$

Since $R$ contains $n-1$ segments, it is clear that $E(C_0(R)) \leq C_0(n-1)$. Therefore,

$$C_0(S) \leq \frac{n}{n-1} \left( 2 + C_0(n-1) \right).$$

This inequality holds for *any* set $S$ of $n$ non-vertical line segments that are in general position. Hence,

$$C_0(n) \leq \frac{2n}{n-1} + \frac{n}{n-1} \cdot C_0(n-1). \quad (4)$$

Observe that $C_0(2) = 1$. We solve the recurrence relation (4) by unfolding it until we "see the solution":

$$
\begin{aligned}
C_0(n) &\leq \frac{2n}{n-1} + \frac{n}{n-1} \cdot C_0(n-1) \\
&\leq \frac{2n}{n-1} + \frac{n}{n-1} \left( \frac{2(n-1)}{n-2} + \frac{n-1}{n-2} \cdot C_0(n-2) \right)
\end{aligned}
$$

11

$$\begin{aligned}
&= & 2\left(\frac{n}{n-1} + \frac{n}{n-2}\right) + \frac{n}{n-2} \cdot C_0(n-2) \\
&\leq & 2\left(\frac{n}{n-1} + \frac{n}{n-2}\right) + \frac{n}{n-2} \cdot \left(\frac{2(n-2)}{n-3} + \frac{n-2}{n-3} \cdot C_0(n-3)\right) \\
&= & 2\left(\frac{n}{n-1} + \frac{n}{n-2} + \frac{n}{n-3}\right) + \frac{n}{n-3} \cdot C_0(n-3) \\
&\leq & \\
&\vdots & \\
&\leq & 2\sum_{i=1}^{n-2}\frac{n}{n-i} + \frac{n}{2} \cdot C_0(2) \\
&= & 2\sum_{j=2}^{n-1}\frac{n}{j} + \frac{n}{2} \\
&= & O(n\log n).
\end{aligned}$$

This proves Theorem 4 (well, at least for segments that are in general position).

**Exercise 8** Prove that $1 + 1/2 + 1/3 + 1/4 + \cdots + 1/n = \Theta(\log n)$. (*Hint:* Estimate the summation by the integral $\int 1/x \cdot dx$.) *Remark:* You may remember from calculus, that $\sum_{i=1}^{n} 1/i - \ln n \to \gamma$ if $n \to \infty$, where $\gamma = 0.5772157\ldots$ is Euler's constant.

# 4 An application of lower envelopes

In this section, we show how lower envelopes can be used to analyze the complexity of an algorithm for solving a simplified version of a problem that comes from the field of neurosurgery: A surgeon wants to remove tissue samples from the brain of a patient for diagnosis purposes. This is done by inserting a probe through a small hole in the skullcap of the patient. In order to minimize the exposure to danger, the point of entry has to be chosen in such a way that the trajectory of the probe stays away from certain brain areas. If we model this trajectory as a ray, and the brain areas we want to avoid by points in three-dimensional space, then we want to find a ray $R$ emanating from the position at which we want to remove the tissue sample such that the minimum distance from the points to $R$ is maximum.
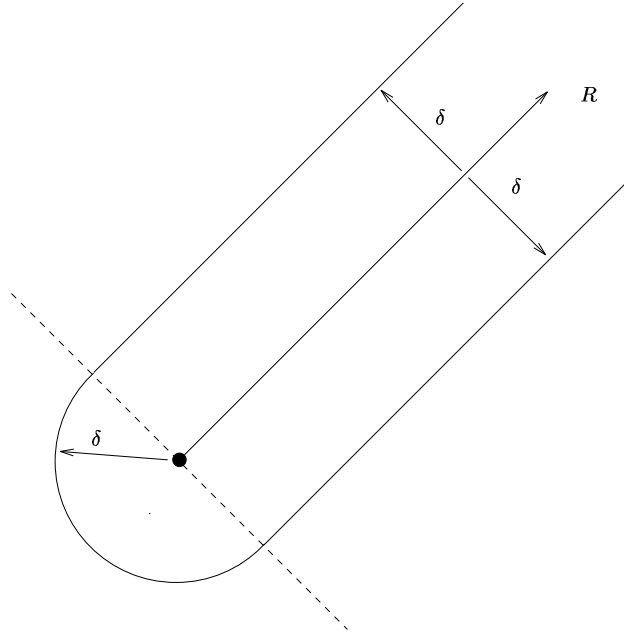
Figure 3: A silo with axis $R$ and radius $\delta$.

We denote the Euclidean distance between a point $p$ and the origin by $|p|$. Moreover, the Euclidean distance between two points $p$ and $q$ is denoted by $d(p, q)$. If $p$ is a point in $\mathbb{R}^d$, and $R$ is a closed subset of $\mathbb{R}^d$, then the distance between $p$ and $R$ is defined as $d(p, R) := \min\{d(p, q) : q \in R\}$. Finally, we define an *anchored ray* as a ray that emanates from the origin.

The above mentioned optimization problem is the three-dimensional version of the following problem.

**Problem 1** *Given a set $S$ of $n$ points in $\mathbb{R}^d$, compute an anchored ray $R$ for which $\min_{p \in S} d(p, R)$ is maximum.*

This problem appeared in the Master's Thesis of Frank Follert at the University of the Saarland in Saarbrücken (Germany).

Let us first give an equivalent formulation of Problem 1: Let $R$ be any ray, and let $\delta \geq 0$. The set of all points in $\mathbb{R}^d$ that are at distance less than or equal to $\delta$ from $R$ is called a *silo* with *axis $R$* and *radius $\delta$*. (See Figure 3.)

**Observation 1** *Problem 1 is equivalent to the following one: Given a set $S$ of $n$ points in $\mathbb{R}^d$, find a silo*

13

1. *whose axis starts in the origin,*

2. *that does not contain any points of $S$, and*

3. *whose radius is maximum.*

We will show that the planar version of Problem 1 can be solved in $O(n \log n)$ time.

Let $S$ be a set of $n$ points in the plane. We want to compute an anchored ray $R$ such that $\min_{p \in S} d(p, R)$ is maximum. Define

$$\delta^* := \max\{\min_{p \in S} d(p, R) : R \text{ is an anchored ray}\}.$$

Let $\delta_l^*$ (resp. $\delta_r^*$) denote the analogous quantity, where we only consider anchored rays that lie on or to the left (resp. right) of the $y$-axis. It is clear that $\delta^* = \max(\delta_l^*, \delta_r^*)$. We show how to compute $\delta_r^*$. The value $\delta_l^*$ can be computed in a symmetric way.

Let $\delta_{min} := \min\{|p| : p \in S\}$. For each $\delta \geq 0$ and each point $p$ of $S$, let $D_p^\delta$ denote the disk with center $p$ and radius $\delta$. For each $\delta$, $0 \leq \delta \leq \delta_{min}$, and each $p \in S$, let $C_p^\delta$ denote the cone consisting of all anchored rays that intersect or touch the disk $D_p^\delta$. (Since $\delta \leq \delta_{min}$, $D_p^\delta$ does not contain the origin. Therefore, $C_p^\delta$ really is a cone.) Observe that $C_p^\delta$ has the origin as its apex. (See Figure 4.)

**Observation 2** *Let $R$ be an anchored ray, let $\delta \geq 0$, let $s$ be the silo with axis $R$ and radius $\delta$, and let $p$ be a point in the plane. Then $p$ is contained in $s$ if and only if $R$ intersects the disk $D_p^\delta$.*

This immediately leads to:

**Observation 3** *We have*

1. $0 \leq \delta_r^* \leq \delta_{min}$.

2. *$\delta_r^*$ is the maximum value of $\delta$, $0 \leq \delta \leq \delta_{min}$, such that there is an anchored ray in the halfplane $x \geq 0$ that does not intersect the interior of any disk $D_p^\delta$, $p \in S$.*

3. *$\delta_r^*$ is the minimum of $\delta_{min}$ and the minimum value of $\delta$, $0 \leq \delta \leq \delta_{min}$, such that the cones $C_p^\delta$, $p \in S$, cover the halfplane $x \geq 0$.*
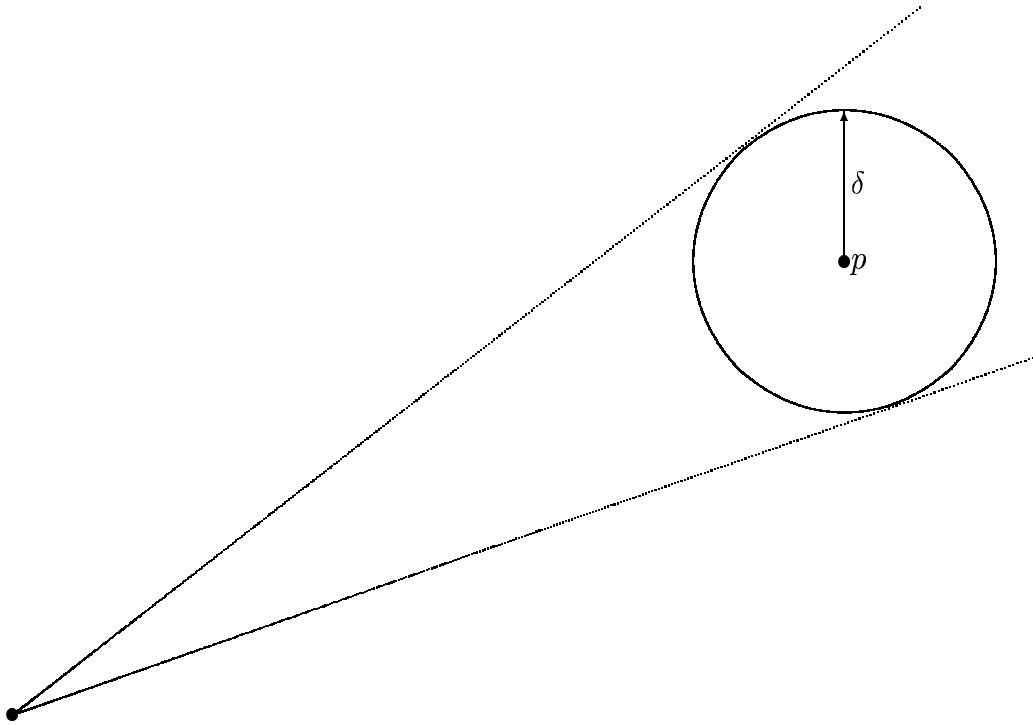
14

Figure 4: $D_p^\delta$ is the disk with center $p$ and radius $\delta$. The cone $C_p^\delta$ is bounded by the two rays emanating from the origin.

It turns out that the third property is easiest to work with.

Let $\delta$ be such that $0 \leq \delta \leq \delta_{min}$, and let $p \in S$. Consider the intersection of the cone $C_p^\delta$ with the halfplane $x \geq 0$. Let $I_p(\delta)$ be the interval of slopes spanned by all anchored rays that lie in this intersection. We represent each slope by the angle between the ray and the positive $x$-axis. Hence, $I_p(\delta) \subseteq [-\pi/2, \pi/2]$.

It is clear that the cones $C_p^\delta$, $p \in S$, cover the halfplane $x \geq 0$, if and only if the intervals $I_p(\delta)$, $p \in S$, cover $[-\pi/2, \pi/2]$. Hence:

**Observation 4** $\delta_r^*$ *is the minimum of*

1. $\delta_{min}$, *and*

2. *the minimum value of $\delta$, $0 \leq \delta \leq \delta_{min}$, such that the intervals $I_p(\delta)$, $p \in S$, cover $[-\pi/2, \pi/2]$.*

Let us look at the intervals $I_p(\delta)$ more closely. We can easily write down such an interval explicitly:
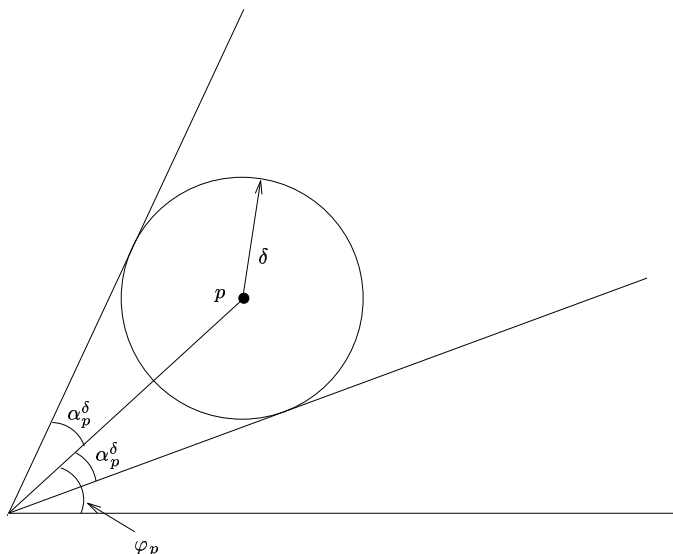
15

Figure 5: Illustration of the angles $\varphi_p$ and $\alpha_p^\delta$.

Let $p$ have coordinates $(p_1, p_2)$, and let $\varphi_p$, $-\pi < \varphi_p \leq \pi$, be the angle between the vector $\vec{p}$ and the positive $x$-axis. Then, $\sin \varphi_p = p_2/|p|$. Also, for each $\delta$, $0 \leq \delta \leq \delta_{min}$, let $\alpha_p^\delta$ be the angle between $\vec{p}$ and an anchored ray that is tangent to the disk $D_p^\delta$. (There are two such tangents, but both define the same angle.) Then, $0 \leq \alpha_p^\delta \leq \pi/2$ and $\sin \alpha_p^\delta = \delta/|p|$. (See Figure 5.) If $p_1 \geq 0$, then

$$I_p(\delta) = \begin{cases} [\varphi_p - \alpha_p^\delta, \varphi_p + \alpha_p^\delta] & \text{if } 0 \leq \delta \leq \delta_{min} \text{ and } \delta \leq p_1, \\ [\varphi_p - \alpha_p^\delta, \pi/2] & \text{if } p_1 \leq \delta \leq \delta_{min} \text{ and } p_2 \geq 0, \\ [-\pi/2, \varphi_p + \alpha_p^\delta] & \text{if } p_1 \leq \delta \leq \delta_{min} \text{ and } p_2 \leq 0. \end{cases}$$

If $p_1 \leq 0$, then

$$I_p(\delta) = \begin{cases} \emptyset & \text{if } 0 \leq \delta \leq \delta_{min} \text{ and } \delta \leq -p_1, \\ [\varphi_p - \alpha_p^\delta, \pi/2] & \text{if } -p_1 \leq \delta \leq \delta_{min} \text{ and } p_2 \geq 0, \\ [-\pi/2, \varphi_p + \alpha_p^\delta] & \text{if } -p_1 \leq \delta \leq \delta_{min} \text{ and } p_2 \leq 0. \end{cases}$$

Using the intervals $I_p(\delta)$ has the disadvantage that we need non-algebraic functions. In order to stay within the algebraic computation tree model—whose operations are much easier to implement—our algorithm works with

16

the intervals
$$J_p(\delta) := \sin\left(I_p(\delta)\right) = \{\sin\gamma : \gamma \in I_p(\delta)\}.$$

Observe that $I_p(\delta) \subseteq [-\pi/2, \pi/2]$ and that the function $\sin(\cdot)$ is increasing on $[-\pi/2, \pi/2]$. Therefore, if $I_p(\delta) = [\ell, r]$, then $J_p(\delta) = [\sin\ell, \sin r]$.

Using the relations $\sin\varphi_p = p_2/|p|$, $\cos\varphi_p = p_1/|p|$, $\sin\alpha_p^\delta = \delta/|p|$, $\cos\alpha_p^\delta = \sqrt{p_1^2 + p_2^2 - \delta^2}/|p|$, and $\sin(x+y) = \sin x \cos y + \cos x \sin y$, we get the following expressions for $J_p(\delta)$. If $p_1 \geq 0$, then

$$J_p(\delta) = \begin{cases} \left[\dfrac{p_2\sqrt{p_1^2+p_2^2-\delta^2}-p_1\delta}{|p|^2}, \dfrac{p_2\sqrt{p_1^2+p_2^2-\delta^2}+p_1\delta}{|p|^2}\right] & \text{if } 0 \leq \delta \leq \delta_{min} \text{ and } \delta \leq p_1, \\[4mm] \left[\dfrac{p_2\sqrt{p_1^2+p_2^2-\delta^2}-p_1\delta}{|p|^2}, 1\right] & \text{if } p_1 \leq \delta \leq \delta_{min} \text{ and } p_2 \geq 0, \\[4mm] \left[-1, \dfrac{p_2\sqrt{p_1^2+p_2^2-\delta^2}+p_1\delta}{|p|^2}\right] & \text{if } p_1 \leq \delta \leq \delta_{min} \text{ and } p_2 \leq 0. \end{cases}$$

If $p_1 \leq 0$, then

$$J_p(\delta) = \begin{cases} \emptyset & \text{if } 0 \leq \delta \leq \delta_{min} \text{ and } \delta \leq -p_1, \\[4mm] \left[\dfrac{p_2\sqrt{p_1^2+p_2^2-\delta^2}-p_1\delta}{|p|^2}, 1\right] & \text{if } -p_1 \leq \delta \leq \delta_{min} \text{ and } p_2 \geq 0, \\[4mm] \left[-1, \dfrac{p_2\sqrt{p_1^2+p_2^2-\delta^2}+p_1\delta}{|p|^2}\right] & \text{if } -p_1 \leq \delta \leq \delta_{min} \text{ and } p_2 \leq 0. \end{cases}$$

The value of $\delta_r^*$ is equal to the minimum of $\delta_{min}$ and the minimum value of $\delta$, $0 \leq \delta \leq \delta_{min}$, such that the intervals $J_p(\delta)$, $p \in S$, cover $[-1, 1]$.

For $p \in S$, let

$$R_p := \{(x, \delta) : 0 \leq \delta \leq \delta_{min}, x \in J_p(\delta)\}.$$

The region $R_p$ is contained in the rectangle $[-1, 1] \times [0, \delta_{min}]$.

**Observation 5** $\delta_r^*$ *is the minimum of*

1. $\delta_{min}$, *and*

2. *the minimum value of* $\delta$, $0 \leq \delta \leq \delta_{min}$, *such that the horizontal segment with endpoints* $(-1, \delta)$ *and* $(1, \delta)$ *is completely contained in* $\bigcup_{p \in S} R_p$.

Let $l_p$ be the lower envelope of $R_p$. Then, $l_p$ is the graph of a continuous function on a subinterval of $[-1, 1]$. Finally, let $L$ be the lower envelope of the graphs $l_p$, $p \in S$, and the line segment with endpoints $(-1, \delta_{min})$ and $(1, \delta_{min})$.

**Observation 6** $\delta_r^*$ *is the y-coordinate of a highest vertex of* $L$.

Observation 6 leads to the following simple algorithm for computing $\delta_r^*$ and a corresponding ray.

**Step 1:** Compute the functions $l_p$, $p \in S$.

**Step 2:** Compute the lower envelope $L$ of the functions $l_p$, $p \in S$, and the horizontal segment with endpoints $(-1, \delta_{min})$ and $(1, \delta_{min})$.

**Step 3:** Walk along $L$ and find a highest vertex on it. Let this vertex have coordinates $(a, \delta)$.

**Step 4:** Output $\delta$ and the anchored ray $R := \{(x, ax/\sqrt{1 - a^2}) : x \geq 0\}$.

To prove the correctness of this algorithm, consider the vertex $(a, \delta)$ that is found in Step 3. Observation 6 implies that $\delta = \delta_r^*$. Let $\varphi$ be the angle such that $-\pi/2 \leq \varphi \leq \pi/2$ and $\sin \varphi = a$. Let $R^*$ be the anchored ray that makes an angle of $\varphi$ with the positive $x$-axis. Then $\delta = \min_{p \in S} d(p, R^*)$. It is easy to see that $R = R^*$.

It is clear that the running time of our algorithm depends on the complexity of the lower envelope $L$.

**Lemma 1** *The names of the points that correspond to the edges of* $L$*, when we traverse* $L$ *from left to right, form a* $DS(n + 1, 2)$*-sequence. Hence,* $L$ *has size* $O(n)$.

Suppose we have proved this lemma. Then we can easily prove the main result of this section:

**Theorem 5** *Let* $S$ *be a set of* $n$ *points in the plane. In* $O(n \log n)$ *time, we can compute an anchored ray* $R^*$ *for which* $\min_{p \in S} d(p, R^*)$ *is maximum.*

**Proof.** Consider the algorithm given above. Step 1 takes $O(n)$ time. By Lemma 1 and Exercise 7, Step 2 takes $O(n \log n)$ time. Step 3 takes $O(n)$ time and, finally, Step 4 takes $O(1)$ time. Hence, it takes $O(n \log n)$ time to compute $\delta_r^*$. In the same amount of time, we can compute $\delta_l^*$. ∎

**Remark 1** The result of Theorem 5 is optimal in the algebraic computation tree model. (Observe that our algorithm works into this model.)

It remains to prove Lemma 1. The proof follows from a careful analysis of the lower envelope $L$.

Let $B_l$, $B_r$, $B_t$ and $B_b$ be the left, right, top and bottom sides of the rectangle $[-1, 1] \times [0, \delta_{min}]$, respectively.

Let $p = (p_1, p_2)$ be a point of $S$, and consider the graph $l_p$. If $p_1 \geq 0$, then $l_p$ consists of a decreasing part $l_p^-$ that has $(p_2/|p|, 0)$ as its lowest and rightmost endpoint, and an increasing part $l_p^+$ that has $(p_2/|p|, 0)$ as its lowest and leftmost endpoint. Moreover, $l_p^-$ (resp. $l_p^+$) has its leftmost (resp. rightmost) endpoint on $B_l$ or $B_t$ (resp. $B_r$ or $B_t$). If $p_1 \leq 0$ and $p_2 \geq 0$, then $l_p$ is decreasing from some point on $B_t$ to some point on $B_r$. Finally, if $p_1 \leq 0$ and $p_2 \leq 0$, then $l_p$ is increasing from some point on $B_l$ to some point on $B_t$.

Let $p = (p_1, p_2)$ and $q = (q_1, q_2)$ be two distinct points of $S$. We claim that the graphs $l_p$ and $l_q$ intersect at most twice. First, we give a geometric explanation for this claim. Then, in Lemma 2 below, we give a rigorous proof.

For the intuitive explanation, assume that $p_1$ and $q_1$ are both positive and that $\varphi_q > \varphi_p$. For each $\delta$, $0 \leq \delta \leq \delta_{min}$, let $U_p(\delta)$ (resp. $L_p(\delta)$) be the anchored ray that is upper (resp. lower) tangent to the disk $D_p^\delta$. Define $U_q(\delta)$ and $L_q(\delta)$ analogously.

Intersections of $l_p$ and $l_q$ are in one-to-one correspondence with values of $\delta$ such that $\{U_p(\delta), L_p(\delta)\} \cap \{U_q(\delta), L_q(\delta)\} \neq \emptyset$.

Consider what happens when we grow $\delta$ from 0 to $\delta_{min}$. More precisely, assume we let $\delta$ grow at a constant "speed" of, say, one meter per second. Initially, $U_p(\delta) = L_p(\delta)$ and $U_q(\delta) = L_q(\delta)$. If $\delta$ grows, then the tangents $U_p(\delta)$ and $L_p(\delta)$ move in opposite directions. Similarly, the tangents $U_q(\delta)$ and $L_q(\delta)$ move in opposite directions. (See Figure 6.) Clearly, there is exactly one $\delta_0$ such that $L_q(\delta_0) = U_p(\delta_0)$. This corresponds to an intersection between $l_q^-$ and $l_p^+$. Also, for $\delta < \delta_0$, there are no intersections between $l_p$ and $l_q$. Now we grow $\delta$ further, from $\delta_0$ to the next "time" $\delta_1$ at which $\{U_p(\delta_1), L_p(\delta_1)\} \cap \{U_q(\delta_1), L_q(\delta_1)\} \neq \emptyset$. (If there is no such "time", then the graphs $l_p$ and $l_q$ intersect exactly once, and we are done.) At "time" $\delta_1$, we have $U_p(\delta_1) = U_q(\delta_1)$ or $L_p(\delta_1) = L_q(\delta_1)$. Assume w.l.o.g. that $U_p(\delta_1) = U_q(\delta_1)$. This corresponds to the second intersection between $l_p$ and $l_q$; more precisely, an intersection between $l_p^+$ and $l_q^+$. Observe that then $U_p(\delta)$ must move faster than $U_q(\delta)$. Hence, for $\delta > \delta_1$, these two tangents never coincide any more. That is, $l_p^+$ and $l_q^+$ intersect only once. Now look at $L_p(\delta)$ and $L_q(\delta)$: Since $L_p(\delta)$ and $U_p(\delta)$ (resp. $L_q(\delta)$ and $U_q(\delta)$) move at the same, but opposite speeds, $L_q(\delta)$ will never overtake $L_p(\delta)$. That is, $l_p^-$ and $l_q^-$ do not
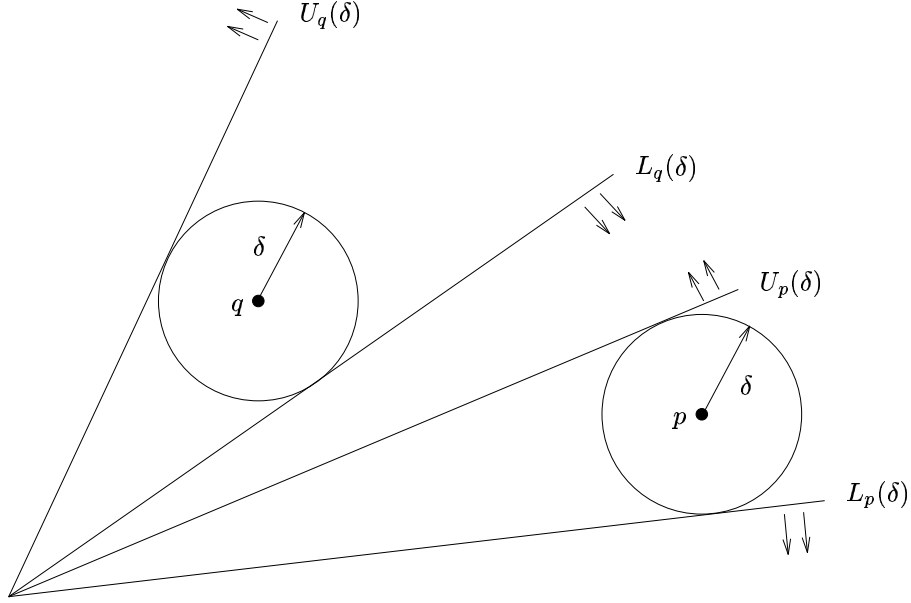
Figure 6: Growing $\delta$ from 0 to $\delta_{min}$.

intersect.

This concludes the intuition why the graphs $l_p$ and $l_q$ intersect at most twice. We now prove this rigorously.

**Lemma 2** *Let $p$ and $q$ be two distinct points of $S$. Then, the graphs $l_p$ and $l_q$ intersect at most twice.*

**Proof.** Assume first that $\varphi_p = \varphi_q$. Then, $|p| \neq |q|$. This implies that for all $\delta$, $0 \leq \delta \leq \delta_{min}$, the cone $C_p^\delta$ is completely contained inside $C_q^\delta$, or vice versa. As a result, $l_p$ and $l_q$ have only one intersection point, with $y$-coordinate zero.

Assume from now on that $\varphi_p \neq \varphi_q$. Furthermore, assume w.l.o.g. that $\varphi_q > \varphi_p$. Let $I$ be the interval of all values $\delta$ such that $J_p(\delta)$ and $J_q(\delta)$ are both non-empty. Consider the function

$$f(\delta) := \varphi_q - \varphi_p - \alpha_q^\delta - \alpha_p^\delta,$$

for $\delta \in I$. Then $f(\delta) = 0$ if and only if the increasing part of $l_p$ and the decreasing part of $l_q$ have an intersection point with $y$-coordinate $\delta$. Observe that

$$f(\delta) = \varphi_q - \varphi_p - \arcsin(\delta/|q|) - \arcsin(\delta/|p|).$$

20

The derivative of $f$ is equal to

$$f'(\delta) = \frac{-1}{\sqrt{|q|^2 - \delta^2}} + \frac{-1}{\sqrt{|p|^2 - \delta^2}}.$$

Hence, $f'$ is strictly negative, which implies that $f$ has at most one root.

Next let

$$g(\delta) := \varphi_q - \varphi_p + \alpha_q^\delta - \alpha_p^\delta,$$

for $\delta \in I$. The roots of $g$ are in one-to-one correspondence with the intersections between the increasing parts of $l_p$ and $l_q$. We have

$$g'(\delta) = \frac{1}{\sqrt{|q|^2 - \delta^2}} - \frac{1}{\sqrt{|p|^2 - \delta^2}}.$$

If $|p| = |q|$, then $g(\delta) = \varphi_q - \varphi_p$, which is never zero. If $|p| \neq |q|$, then $g'$ is either strictly positive or strictly negative for all $\delta \in I$. Hence, the function $g$ has at most one root.

In a completely symmetric way, it follows that the function

$$h(\delta) := \varphi_q - \varphi_p - \alpha_q^\delta + \alpha_p^\delta,$$

for $\delta \in I$, has at most one root. That is, the decreasing parts of $l_p$ and $l_q$ intersect at most once.

Now we can prove the lemma. First assume that $p_1 \leq 0$. Then $l_p$ consists only of a decreasing part, or only of an increasing part. If $q_1 \leq 0$, then $l_q$ consists of one monotone part. The above analysis shows that in this case, $l_p$ and $l_q$ intersect at most once. If $q_1 > 0$, then $l_q$ consists of two monotone parts, each of which intersects $l_p$ at most once. Hence, in this case, $l_p$ and $l_q$ intersect at most twice.

If $p_1 \geq 0$ and $q_1 \leq 0$, then a symmetric argument shows that $l_p$ and $l_q$ intersect at most twice.

It remains to consider the case when $p_1 > 0$ and $q_1 > 0$. We proved above that the increasing part $l_p^+$ of $l_p$ and the decreasing part $l_q^-$ of $l_q$ intersect at most once.

Assume that $l_p^+$ and $l_q^+$ intersect. Then the analysis above shows that they intersect exactly once. Since the function $g$ is monotone, $g(0) = \varphi_q - \varphi_p > 0$, and $g$ has a root, this function is decreasing. But this implies that $h$ is increasing. Since $h(0) > 0$, the function $h$ does not have any root, which proves that $l_p^-$ and $l_q^-$ do not intersect.

If $l_p^-$ and $l_q^-$ intersect, then it follows in a completely symmetric way that $l_p^+$ and $l_q^+$ do not intersect.

This proves that $l_p$ and $l_q$ intersect at most twice. ∎

**Proof of Lemma 1:** First observe that the lower envelope $L$ is defined by the $n$ graphs $l_p$, $p \in S$, and the horizontal line segment $h$ with endpoints $(-1, \delta_{min})$ and $(1, \delta_{min})$. Therefore, the sequence $U$ corresponding to $L$ is over an alphabet of size $n + 1$. We have to show that $U$ does not contain a subsequence of the form $p \ldots q \ldots p \ldots q$. Clearly, the segment $h$ does not lead to such a subsequence. Also, for any two points $p, q \in S$, $p \neq q$, the graphs $l_p$ and $l_q$ do not give such a subsequence: this follows from the fact that $l_p$ and $l_q$ intersect at most twice, and from the restrictions on the endpoints of these graphs. ∎

# 5 The Ackermann function and its inverse

In this section, we define the extremely slowly growing inverse of the Ackermann function. First, we define the Ackermann function itself.

We will use the following notation. If $f$ is a function and $i$ is a nonnegative integer, then $f^{(i)}$ denotes the $i$-th iterate of $f$. That is, $f^{(0)}$ is the identity function and for $i \geq 0$, $f^{(i+1)}$ is defined by $f^{(i+1)}(x) := f(f^{(i)}(x))$ for all $x$.

For any $k \geq 0$, we define the function $A_k : \mathbb{N} \longrightarrow \mathbb{N}$ recursively, as follows:

1. For all $x \in \mathbb{N}$, $A_0(x) := x + 1$.

2. For $k \geq 0$ and $x \in \mathbb{N}$, $A_{k+1}(x) := A_k^{(x)}(x)$.

To get an idea of the behavior of these functions, we consider a few of them.

For $x = 0$, we have $A_0(0) = 1$ and $A_{k+1}(0) = A_k^{(0)}(0) = 0$ for all $k \geq 0$. For $x = 1$, we have $A_0(1) = 2$ and

$$A_{k+1}(1) = A_k^{(1)}(1) = A_k(1) = \ldots = A_0(1) = 2$$

for all $k \geq 0$.

Let $x \geq 2$. Then $A_0(x) = x + 1$ and

$$\begin{aligned} A_1(x) &= A_0^{(x)}(x) = A_0(A_0^{(x-1)}(x)) = A_0^{(x-1)}(x) + 1 \\ &= A_0(A_0^{(x-2)}(x)) + 1 = A_0^{(x-2)}(x) + 2. \end{aligned}$$

Continuing in this way, we get

$$A_1(x) = A_0^{(1)}(x) + x - 1 = A_0(x) + x - 1 = 2x.$$

For $k = 2$, we get

$$\begin{aligned}
A_2(x) &= A_1^{(x)}(x) = A_1(A_1^{(x-1)}(x)) = 2 \cdot A_1^{(x-1)}(x) \\
&= 2 \cdot A_1(A_1^{(x-2)}(x)) = 2^2 \cdot A_1^{(x-2)}(x) \\
&= \ldots = 2^{x-1} \cdot A_1^{(1)}(x) = 2^{x-1} \cdot A_1(x) = x \cdot 2^x.
\end{aligned}$$

In particular, we have

$$A_2(x) \geq 2^x.$$

Next we consider $A_3$:

$$A_3(x) = A_2^{(x)}(x) = A_2(A_2^{(x-1)}(x)) \geq 2^{A_2^{(x-1)}(x)},$$

which implies

$$A_3(x) \geq \underbrace{2^{2^{2^{\cdot^{\cdot^{2^x}}}}}}_{x\ 2's}.$$

The function $A_4$ grows so fast that we only consider $A_4(2)$:

$$A_4(2) = A_3^{(2)}(2) = A_3(A_3(2)).$$

Since

$$A_3(2) = A_2^{(2)}(2) = A_2(A_2(2)) = A_2(8) = 2048,$$

we get

$$A_4(2) = A_3(2048) \geq \underbrace{2^{2^{2^{\cdot^{\cdot^{2^{2048}}}}}}}_{2048\ 2's}.$$

Now we can define our extremely rapidly growing function $A : \mathbb{N} \longrightarrow \mathbb{N}$:

$$A(k) := A_k(2) \text{ for } k \geq 0.$$

This function is called the Ackermann function. Observe that

$$A(0) = 3, A(1) = 4, A(2) = 8, A(3) = 2048,$$

23

whereas

$$A(4) \geq \underbrace{2^{2^{2^{\cdot^{\cdot^{\cdot^{2^{2^{2048}}}}}}}}}_{2048 \ 2's} .$$

The function we are actually interested in, is its inverse $\alpha : \mathbb{N} \longrightarrow \mathbb{N}$, defined by

$$\alpha(n) := \min\{k \geq 0 : A(k) \geq n\}.$$

We claim that for all practical purposes, $\alpha(n)$ is at most 4. We have

$$\alpha(0) = \alpha(1) = \alpha(2) = \alpha(3) = 0,$$

$$\alpha(4) = 1,$$
$$\alpha(5) = \alpha(6) = \alpha(7) = \alpha(8) = 2,$$
$$\alpha(9) = \alpha(10) = \ldots \alpha(2048) = 3,$$
$$\alpha(2049) = 4.$$

Let $n$ be such that $\alpha(n) \geq 5$. Then $A(k) < n$ for $0 \leq k \leq 4$. In particular, $n > A(4)$. We have seen, however, that $A(4)$ is a number beyond comprehension.

**Exercise 9** Prove that $\alpha$ is well-defined. That is, prove that for each $n \geq 0$, there is a $k \geq 0$ such that $A(k) \geq n$. Also, prove that $\alpha(n) \to \infty$ for $n \to \infty$.

**Remark 2** Ackermann defined "his" function in 1928. You may remember it from the course Theoretical Computer Science, as a function that is recursive (i.e., computable by a Turing machine), but not primitive recursive. The results of the preceding sections show that this function, which looks completely artificial, in fact occurs in nature (well, in Euclidean nature). In the literature, different definitions of the Ackermann function appear. All these functions grow at roughly the same rate.

# 6   Remarks

The book

> *Davenport-Schinzel sequences and their geometric applications*, by Micha Sharir and Pankaj Agarwal, Cambridge University Press, 1995

is entirely devoted to lower envelopes and Davenport-Schinzel sequences.

The three-dimensional version of Problem 1 in Section 4 can be solved using the *parametric search technique*, in $O(n \log^4 n)$ time. See

Follert *et al.*, *Computing a largest empty anchored cylinder, and related problems*, International Journal of Computational Geometry & Applications **7** (1997), pp. 563-580.