

# Output-Sensitive Algorithms for Tukey Depth and Related Problems\*

David Bremner                      Dan Chen  
University of New Brunswick      Carleton University

John Iacono                      Stefan Langerman  
Polytechnic University          Université Libre de Bruxelles

Pat Morin  
Carleton University

## Abstract

The *Tukey depth* (Tukey 1975) of a point  $p$  with respect to a finite set  $S$  of points is the minimum number of elements of  $S$  contained in any closed halfspace that contains  $p$ . Algorithms for computing the Tukey depth of a point in various dimensions are considered. The running times of these algorithms depend on the value of the output, making them suited to situations, such as outlier removal, where the value of the output is typically small.

**Keywords:** Tukey depth, halfspace depth, algorithms, computational statistics, computational geometry, fixed-parameter tractability

## 1 Introduction

Let  $S$  be a set of  $n$  points in  $\mathbb{R}^d$ . The *Tukey depth*, or *halfspace depth* of a point  $p \in \mathbb{R}^d$  with respect to  $S$  can be defined in several equivalent ways [32]:

$$\text{depth}(p, S) = \min\{|h \cap S| : h \text{ is a closed halfspace containing } p\} \quad (1)$$

$$= \min\{|h \cap S| : h \text{ is a closed halfspace with } p \text{ on its boundary}\} \quad (2)$$

$$= \min\{|S'| : p \text{ is outside the convex hull of } S \setminus S'\} \quad (3)$$

A point of maximum Tukey depth serves as a  $d$ -dimensional generalization of the (1-dimensional) median that has many nice statistical properties including being robust against

---

\*This research was partly funded by the NSERC Canada.

outliers, invariant under affine transformations, and monotone. The contours of the Tukey depth function<sup>1</sup> are generalizations of 1-dimensional percentiles that also have many nice properties including convexity, robustness, and monotonicity [27, 28, 31]. Algorithms for computing a point  $p \in \mathbb{R}^d$  of maximum Tukey depth have a rich history [15, 14, 5] that has recently culminated (from a theoretical point of view) in Chan’s  $O(n \log n + n^{d-1})$  expected time algorithm.

In this paper we consider the simpler, but still very difficult, problem of computing the Tukey depth of a given point  $p$  with respect to a set  $S$ . If the dimension  $d$  of the problem is part of the input, then this problem is NP-hard [13], and is even APX-hard [1]. The current paper presents algorithms whose running times are dependent on the dimension  $d$  and the value,  $k$ , of the output. In some applications, such as outlier removal, the goal is to identify quickly the data points of small depth (so they can be removed). Our algorithms are particularly well-suited to such applications since they run quickly when the depth of  $p$  is small. Specifically, we present the following results:

1. A simple  $O(n + k \log k)$  time algorithm to compute the Tukey depth of a point in  $\mathbb{R}^2$  (Section 2). The most complicated data structure used in this algorithm is a binary heap.
2. An  $O(n + (n - k) \log(n - k))$  time algorithm to find the largest clique in an interval graph, where  $k$  is the size of the clique found (Section 3). This problem is a generalization of the Tukey depth problem in  $\mathbb{R}^2$ .
3. An  $O(n \log n + k^2 \log n)$  time algorithm to compute the Tukey depth of a point in  $\mathbb{R}^3$  and an  $O(n + k^{11/4} n^{1/4} \log^{O(1)} n)$  time algorithm to compute the Tukey depth of a point in  $\mathbb{R}^4$  (Section 4). These algorithms rely on results of Chan on linear programming with violated constraints [6] which in turn rely on sophisticated range searching data structures [16, 23] and/or dynamic convex hull data structures [4].
4. A simple  $O(d^k \text{LP}(n, d - 1))$  time algorithm to compute the Tukey depth of a point in  $\mathbb{R}^d$ , where  $\text{LP}(n, d)$  denotes the time required to determine the feasibility of a linear program having  $n$  constraints and  $d$  variables (Section 5). Not surprisingly, this algorithm is also based on linear programming with violated constraints and is obtained by presenting a fixed-parameter tractable algorithm for a parameterization of the NP-hard MAXIMUMFEASIBLESUBSYSTEM problem.

For the remainder of this paper we use the following notations: For points  $p, q \in \mathbb{R}^d$ ,  $p_i$  denotes the  $i$ th coordinate of  $p$ ,  $\|p\| = (\sum_{i=1}^d p_i^2)^{1/2}$ , and  $p \cdot q = \sum_{i=1}^d p_i q_i$ . The unit sphere in  $\mathbb{R}^{d+1}$  is denoted by  $\mathbb{S}^d = \{p \in \mathbb{R}^{d+1} : \|p\| = 1\}$ . The top side of this sphere is denoted by  $\mathbb{S}_+^d = \{p \in \mathbb{S}^d : p_{d+1} > 0\}$ , the bottom side is denoted by  $\mathbb{S}_-^d = \{p \in \mathbb{S}^d : p_{d+1} < 0\}$  and the equator is denoted by  $\mathbb{S}_0^d = \{p \in \mathbb{S}^d : p_{d+1} = 0\}$ .

---

<sup>1</sup>The  $\ell$ -contour of the tukey depth function is defined as  $\Gamma_\ell(S) = \{p \in \mathbb{R}^d : \text{depth}(p, S) \geq \ell\}$ .

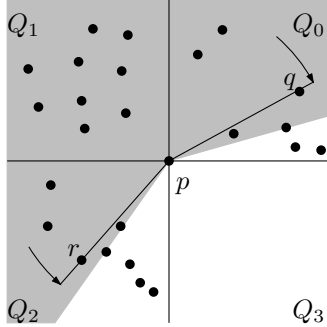


Figure 1: Computing the quantity  $\text{depth}_1(p, S)$ .

## 2 An Algorithm for Points in $\mathbb{R}^2$

In this section we give a simple  $O(n + k \log k)$  time algorithm to compute the Tukey depth of a point  $p \in \mathbb{R}^2$  with respect to a set  $S$  of  $n$  points in  $\mathbb{R}^2$ . We first note that an  $O(n \log n)$  time *sort-and-scan* algorithm is easily obtained by sorting the points of  $S$  radially about  $p$  and then scanning the resulting sorted list using two pointers [15]. The main idea behind our algorithm is to reduce the problem to a *kernel* of size  $O(k)$  on which we can apply this sort-and-scan algorithm.

The algorithm begins by partitioning  $\mathbb{R}^2$  into 4 quadrants around  $p$  that, in counterclockwise order, we denote by  $Q_0, \dots, Q_3$ . The algorithm then simultaneously begins computing the 4 quantities  $\text{depth}_0(p, S), \dots, \text{depth}_3(p, S)$  where

$$\text{depth}_i(p, S) = \min\{|h \cap S| : h \text{ is a closed halfspace containing } Q_i\} . \quad (4)$$

Clearly,  $\text{depth}(p, S) = \min\{\text{depth}_i(p, S) : 0 \leq i \leq 3\}$  since any closed halfspace containing  $p$  contains at least one of the four quadrants. In the remainder of this section we will describe how to compute  $k_i = \text{depth}_i(p, S)$  in  $O(n + k_i \log k_i)$  time. Since the computation can stop once  $\text{depth}_i(p, S)$  has been computed for the index  $i$  that minimizes (4), running the computation of  $k_0, \dots, k_3$  in parallel yields an  $O(n + k \log k)$  time algorithm, where  $k = \text{depth}(p, S)$ .

Let  $S_i = S \cap Q_i$ . To compute  $\text{depth}_i(p, S)$  we create two binary heaps  $H_{i-1}$  and  $H_{i+1}$  that store the elements of  $S_{i-1}$ , respectively  $S_{i+1}$ , in clockwise, respectively, counterclockwise, order around  $p$ .<sup>2</sup> Creating these two heaps takes  $O(n)$  time using the standard bottom-up algorithm to construct a binary heap [9, Chapter 6]. Next we extract elements one at a time from each of  $H_{i-1}$  and  $H_{i+1}$  until either (a) one of the heaps is empty or (b) we extract two elements  $q$  from  $H_{i-1}$  and  $r$  from  $H_{i+1}$  such that the angle  $\angle qpr > \pi$ . Suppose we have extracted  $\ell$  elements each from  $H_{i-1}$  and  $H_{i+1}$  when this occurs. Then, any closed halfspace containing  $Q_i$  contains  $\ell - 1$  elements extracted from  $H_{i-1}$  or it contains  $\ell - 1$  elements

<sup>2</sup>Here and in the remainder of this section  $S_i$  is treated implicitly as  $S_{(i \bmod 4)}$ .

extracted from  $H_{i+1}$ . Therefore,

$$|S_i| + \ell - 1 \leq \text{depth}_i(p, S) .$$

On the other hand, the closed halfspace containing  $Q_i$  and bounded by the line through  $p$  and the last element extracted from  $H_{i-1}$  contains at most  $|S_i| + 2\ell - 1$  points of  $S$ , so

$$\text{depth}_i(p, S) \leq |S_i| + 2\ell - 1 .$$

Next, we continue to extract as many elements as possible from each of  $H_{i-1}$  and  $H_{i+1}$  up to a maximum of an additional  $\ell - 1$  elements each. The total time required to extract these at most  $4\ell - 2$  elements from the two heaps is  $O(\ell \log n)$ . By sorting and scanning all the elements extracted from the heaps plus the elements of  $S_i$  we can then compute  $\text{depth}_i(p, S)$  in an additional

$$O((|S_i| + \ell) \log n) = O(k_i \log n)$$

time. This yields an a total running time of

$$O(n + k_i \log n) = O(n + k_i \log k_i) ,$$

as required.<sup>3</sup> This completes the proof of:

**Theorem 1.** *The Tukey depth of a point  $p$  with respect to a set  $S$  of  $n$  points in  $\mathbb{R}^2$  can be computed in  $O(n + k \log k)$  time, where  $k$  is the value of the output.*

### 3 An Algorithm for Max-Clique in Interval Graphs

The problem of computing Tukey depth in  $\mathbb{R}^2$  can be viewed as a problem on a set of circular arcs. By the second definition of Tukey depth (Equation (2)), computing the Tukey depth of  $p$  is equivalent to finding a unit normal vector  $v$  such that the halfspace with  $p$  on its boundary and having inner normal  $v$  contains as few points of  $S$  as possible. Note that the set of unit normals in  $\mathbb{R}^2$  is homeomorphic to the unit circle  $\mathbb{S}^1$  and that each point  $q \in S$  defines an open circular arc of  $\mathbb{S}^1$  such that all choices of  $v$  in this circular arc yield a closed halfspace with  $p$  on its boundary that does not contain  $q$ . In fact, all the open circular arcs obtained this way are half-circles. Thus, the Tukey depth problem reduces to the problem of finding a vector  $v$  that is contained in the largest number of half-circles. Although it is not immediately apparent, the correctness of the algorithm in Theorem 1 relies crucially on the fact that the arcs of  $\mathbb{S}^1$  defined by the Tukey depth problem are all half-circles and not arcs of arbitrary length.

An obvious generalization of the Tukey depth problem is that of, given a set of  $n$  circular arcs of  $\mathbb{S}^1$ , finding a point  $p \in \mathbb{S}^1$  contained in the largest number of arcs. Note that, in

---

<sup>3</sup>This equation follows from the analysis of two cases. If  $k_i \leq n/\log n$  then  $k_i \log k_i \leq n$ . On the other hand, if  $k_i > n/\log n$  then  $\log k_i > \log n - \log \log n$  so  $k_i \log n < 2k_i \log k_i$  for all  $n \geq 16$ .

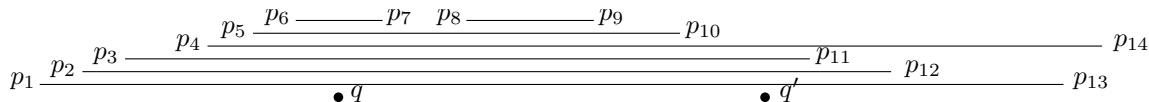


Figure 2: An illustration of Lemma 1 and Lemma 2. The point  $q \in [p_6, p_7]$  is contained in 6 intervals. Therefore, by Lemma 1, the point  $q' \in [p_{10}, p_{11}]$  is contained in at least  $6 - 4 = 2$  intervals and by Lemma 2  $q'$  is contained in at most  $14 - 6 - 4 = 4$  intervals.

this generalization, the arcs can have arbitrary and different lengths. Like the Tukey depth problem, this problem is easily solved in  $O(n \log n)$  time by the sort-and-scan algorithm.

Unfortunately, for the problem of finding a point  $p$  contained in the largest number of circular there can be no algorithm whose running time depends on the number  $k$  of arcs containing  $p$  or even on the number  $(n - k)$  of arcs not containing  $p$ . There are two reasons for this. The problem of testing if a set of circular arcs are all pairwise-disjoint has an  $\Omega(n \log n)$  lower bound [3]. This problem can be solved by finding a point  $p$  contained in the largest number of arcs and then determining, in  $O(n)$  time, how many arcs contain  $p$ . Thus, even in the cases where  $k \in \{1, 2\}$ , the  $\Omega(n \log n)$  lower bound holds. At the other end of the scale, testing whether a set of  $n$  circular arcs covers  $\mathbb{S}^1$  has an  $\Omega(n \log n)$  lower-bound [3]. This problem is equivalent, by taking the complement of each arc, to the problem of finding the point contained in the maximum number of arcs and checking if this point is contained in  $n$  of the complementary arcs. In particular, the original set of arcs do not cover  $\mathbb{S}^1$  if and only if there is a point  $p$  contained in every complementary arc. Thus, the  $\Omega(n \log n)$  lower bound holds even when  $k \in \{n, n - 1\}$ .

Since we can not hope to find an output-sensitive algorithm for circular arcs of  $\mathbb{S}^1$ , we settle for the next best thing. Let  $I$  be a set of real intervals. Here we describe an  $O(n + (n - k) \log(n - k))$  time algorithm to find a point  $p \in \mathbb{R}$  that is contained in the largest number of intervals in  $I$ . Here  $k$  is the number of intervals in  $I$  that contain  $p$ .

Let  $p_1, \dots, p_{2n}$  denote the  $2n$  endpoints of the intervals in  $I$ , in increasing order. For convenience we use the convention that  $p_i = -\infty$  for  $i \leq 0$  and  $p_i = +\infty$  for  $i > 2n$ . Note that the  $p_1, \dots, p_{2n}$  notation is for convenience only and the algorithm we describe does not require that the endpoints be given in sorted order, nor will it sort them into this order. Together, the following two observations imply that all the points contained in many intervals are clustered together. (Refer to Figure 2.)

**Lemma 1.** *Let  $q \in [p_i, p_{i+1}]$  be a point contained in  $k$  intervals of  $I$ . Then, for any  $0 \leq r \leq n$ , every point  $q' \in [p_{i-r}, p_{i+r+1}]$  is contained in at least  $k - r$  intervals of  $I$ .*

*Proof.* Without loss of generality, assume that  $q' \in [q, p_{i+r+1}]$ . There are at most  $r$  endpoints of intervals in  $I$  contained in the interval  $[q, q']$ . Therefore there are at most  $r$  intervals that contain  $q$  but not  $q'$ .  $\square$

**Lemma 2.** *Let  $q \in [p_i, p_{i+1}]$  be a point contained in  $k$  intervals of  $I$ . Then, for any  $n - k \leq r \leq n$ , every point  $q' \notin [p_{i-r}, p_{i+r+1}]$  is contained in at most  $2n - k - r$  arcs of  $C$ .*

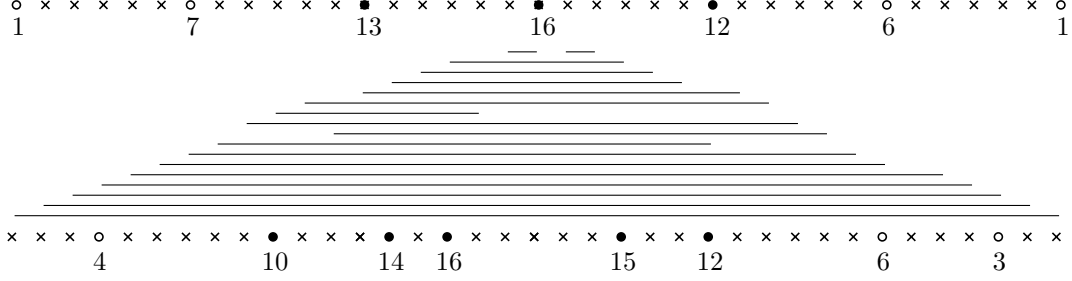


Figure 3: An illustration of the algorithm for interval graphs using sampling at regular intervals (above) and using random samples (below). The rows of dots show the endpoints of the intervals with samples drawn as circles labelled with the number of intervals that contain them. High-depth samples are indicated as filled circles.

*Proof.* Without loss of generality, assume that  $q' > p_{i+r+1}$ . Then, as we walk from  $q$  to  $q'$  we encounter at least  $r$  endpoints of intervals in  $I$ . At most  $n - k$  of these endpoints are left endpoints of intervals and at least  $r - (n - k)$  of these are right endpoints. Thus, the number of intervals that contain  $q'$  is at most

$$k + (n - k) - (r - (n - k)) = 2n - k - r \ ,$$

as required.  $\square$

We first explain our algorithm at a high level in which we deliberately ignore several important implementation details that are discussed later. Refer to Figure 3. Suppose we are given a value  $k$  and only wish to find a value  $p \in \mathbb{R}$  contained in at least  $k$  intervals of  $I$ . We begin by taking a regular sample  $s_1, \dots, s_{2t}$  of  $p_1, \dots, p_{2n}$  so that any interval  $[s_i, s_{i+1}]$  between two consecutive sample points contains at most  $n/t$  points of  $p_1, \dots, p_{2n}$ . We then compute, for each sample point  $s_i$  the number of intervals in  $I$  that contain  $s_i$ . By Lemma 1, if there exists any point  $p \in \mathbb{R}$  contained in  $k$  intervals of  $I$  then the two sample points  $s_j$  and  $s_{j+1}$  on either side of  $p$  are *high-depth samples* that are each contained in at least  $k - n/t$  intervals of  $I$ . Furthermore, by Lemma 2, the only high-depth samples are contained in the interval  $[p_{i-r}, p_{i+r}]$  for  $r = 2(n - k) + n/t$ .

If we choose  $t = \sqrt{n}$  then  $r = O(n - k + \sqrt{n})$ . Thus, by computing an interval  $[p_a, p_b]$  that contains all high-depth samples we can find the point  $p$  contained in the largest number of intervals of  $C$  by applying the standard sort-and-scan algorithm on the  $O(n - k + \sqrt{n})$  endpoints of the intervals of  $C$  that fall in the interval  $[p_a, p_b]$ . The running time of the sort-and-scan algorithm is  $O(m \log m)$  where  $m$  is the number of points to be scanned. In this case  $m = O(n - k + \sqrt{n})$  for a running time of

$$O((n - k + \sqrt{n}) \log(n - k + \sqrt{n})) = O(n + (n - k) \log(n - k)) \ ,$$

as required. Note that the value of  $k$  has very little effect on the execution of the algorithm other than to bound the number of endpoints in the interval  $[p_a, p_b]$  on which the sort-and-scan algorithm is run. In implementing the above ideas, several complications arise:

1. The value of  $k$  is not known in advance. However, we do not need the exactly value of  $k$  and the value of  $k$  can be estimated using Lemma 1. In particular, since the number of endpoints between any two consecutive samples is at most  $\sqrt{n}$  the number of intervals containing  $s_i$  estimates, to within an additive error of  $\sqrt{n}$ , the number of intervals containing any point  $q \in [s_{i-1}, s_{i+1}]$ . In particular, we can estimate the value of  $k$  by computing, for each sample point  $s_i$ , the number of intervals of  $I$  that contain  $s_i$  (see Issue 3, below) and use the maximum of these values as an estimate for  $k$ .
2. We can not obtain a perfectly regular sample  $s_1, \dots, s_{2\sqrt{n}}$  of  $p_1, \dots, p_{2n}$  in  $O(n)$  time. However, we do not require a perfectly regular sample. By taking a random sample of size  $c\sqrt{n} \log n$  for an appropriate constant  $c$  we obtain a set of samples  $s_1, \dots, s_{c\sqrt{n} \log n}$  such that, with probability  $1 - n^{-\Omega(c)}$ , no interval  $[s_i, s_{i+1}]$  contains more than  $\sqrt{n}$  endpoints of intervals of  $I$  [21]. Such a sample still guarantees that the number of intervals containing a point  $q \in [s_i, s_{i+1}]$  is estimated, to within an additive error of  $\sqrt{n}$ , by the number of intervals containing  $s_i$ .
3. We can not compute, in  $O(n)$  time, for each sample point  $s_i$ , the number of intervals of  $I$  that contain  $s_i$ . However, random sampling helps again here. Let  $d(s_i)$  denote the number of elements of  $I$  that contain  $s_i$ . By taking a random sample  $I' \subseteq I$ ,  $|I'| = c\sqrt{n}$  we can determine for each  $s_i$  a number  $d_i$  such that, with probability  $1 - n^{-\Omega(c)}$

$$d(s_i) - O(n^{4/5}) \leq d_i \leq d(s_i) + O(n^{4/5}) .$$

By storing the  $\sqrt{n}$  elements of  $I'$  in an interval tree [22] and then querying this interval tree with the  $c\sqrt{n} \log n$  sample elements the numbers  $d_1, \dots, d_{c\sqrt{n} \log n}$  can be computed in  $O(\sqrt{n} \log^2 n)$  time.

Note that each of the above steps can be accomplished in  $o(n)$  time, determining the endpoints contained in  $[p_a, p_b]$  can be done in  $O(n)$  time, and the final sort-and-scan step takes  $O(n + (n - k) \log(n - k))$  time. The correctness of the resulting output depends on the success of the random sampling steps described in points 2 and 3, above. However, Lemma 2 implies that this final sort-and-scan step allows us to check the correctness of the output (by counting the number of intervals containing  $p_a$  and  $p_b$  and comparing this to the computed value of  $k$ ) and restarting the algorithm from scratch if necessary. This yields:

**Theorem 2.** *There exists a randomized algorithm that, given a set  $I$  of  $n$  real intervals, finds a value  $p \in \mathbb{R}$  contained in the largest number of intervals of  $I$  and that runs in  $O(n + (n - k) \log(n - k))$  expected time, where  $k$  is the number of intervals containing  $p$ .*

## 4 Algorithms for Points in $\mathbb{R}^3$ and $\mathbb{R}^4$

The previous section showed how the problem of computing the Tukey depth of a point in  $\mathbb{R}^2$  is equivalent to the problem of finding a point contained in the largest number of halfcircles on the unit circle  $\mathbb{S}^1$ . Our main goal in this section is to present a reduction from Tukey

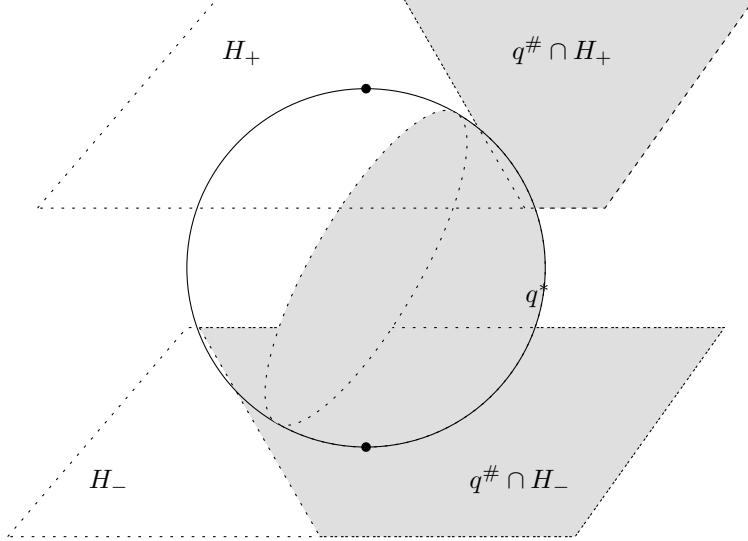


Figure 4: Computing the Tukey depth of a point in  $\mathbb{R}^d$  reduces to two MAXIMUMFEASIBLESUBSYSTEM problems in  $\mathbb{R}^{d-1}$ .

depth to the MAXIMUMFEASIBLESUBSYSTEM (MFS) Problem. This will allow us to apply known algorithms for the MFS problem to give theoretically efficient algorithms for Tukey depth in  $\mathbb{R}^3$  and  $R^4$ . The reduction to MFS will also prove useful in the next section when we discuss the general dimensional case.

The reduction in  $\mathbb{R}^d$  proceeds as follows: Each point  $q \in S$  defines an open halfsphere  $q^* = \{v \in \mathbb{S}^{d-1} : v \cdot q < 0\}$ . That is, all vectors in  $q^*$  are the inner normals of halfspaces that contain  $p$  but do not contain  $q$ . Thus, the problem of determining the Tukey depth of  $p$  reduces to the problem of finding the point contained in the largest number of halfspheres in  $S^* = \{q^* : q \in S\}$ .

We observe that this problem can be solved by solving two MAXIMUMFEASIBLESUBSYSTEM problems in  $\mathbb{R}^{d-1}$ . Refer to Figure 4. Each open halfsphere  $q^* \in S^*$  is the intersection of an open halfspace  $q^\#$  with  $\mathbb{S}^{d-1}$ . Consider the intersection of  $q^\#$  with the hyperplane  $H_+ = \{(x_1, \dots, x_d) : x_d = 1\}$ . By central projection, there is a bijection between points in  $\mathbb{S}_+^{d-1}$  and  $H_+$  and this bijection has the property that  $r \in \mathbb{S}_+^{d-1}$  is in  $q^*$  if and only if the projection of  $r$  is in  $q^\# \cap H_+$ . Thus, finding the point in  $\mathbb{S}_+^{d-1}$  contained in the largest number of halfspheres is equivalent to finding a point in  $H_+$  contained in the largest number of halfspaces. A similar statement holds regarding  $\mathbb{S}_-^{d-1}$  using the hyperplane  $H_- = \{(x_1, \dots, x_d) : x_d = -1\}$ . Finally, we observe that we do not need to consider solutions on the equator  $\mathbb{S}_0^{d-1}$  because our input consists of *open* halfspheres.<sup>4</sup>

<sup>4</sup>A solution to the Tukey depth problem is obtained when we find a set  $n - k$  open halfspheres that have a non-empty common intersection. The intersection of a finite set of open halfspheres is either an empty set or a set with positive measure. Since  $\mathbb{S}_0^{d-1}$  has measure 0 any solution to the Tukey depth problem contains points not in  $\mathbb{S}_0^{d-1}$ .



The above discussion shows that computing the Tukey depth of a point in  $\mathbb{R}^d$  reduces to two instances of the problem MAXIMUMFEASIBLESUBSYSTEM problem in  $\mathbb{R}^{d-1}$ : Given a set  $K$  of  $n$  halfspaces in  $\mathbb{R}^{d-1}$ , find the subset  $K'$  of  $K$  of minimum cardinality such that  $\cap(K \setminus K')$  is non-empty. (The set  $K' \setminus K$  is called a *maximum feasible subsystem*.) The current best results for MAXIMUMFEASIBLESUBSYSTEM in small dimensions are due to Chan [6]. Using two instances of his algorithm for MAXIMUMFEASIBLESUBSYSTEM in  $\mathbb{R}^2$ , respectively,  $\mathbb{R}^3$ , and running them in parallel gives:

**Theorem 3.** *The Tukey depth of a point  $p$  with respect to a set  $S$  of  $n$  points in  $\mathbb{R}^3$  can be computed in  $O(n \log n + k^2 \log n)$  time, where  $k$  is the value of the output.*

**Theorem 4.** *The Tukey depth of a point  $p$  with respect to a set  $S$  of  $n$  points in  $\mathbb{R}^4$  can be computed in  $O(n \log n + k^{11/4} n^{1/4} \log^{O(1)} n)$  time, where  $k$  is the value of the output.*

## 5 An Algorithm for Points in $\mathbb{R}^d$

Finally, we consider the general case of point sets in  $\mathbb{R}^d$ . In the previous section we showed that computing the Tukey depth of a point  $p$  with respect to a set  $S$  of  $n$  points in  $\mathbb{R}^d$  can be reduced to two instances of MAXIMUMFEASIBLESUBSYSTEM in  $\mathbb{R}^{d-1}$ . In this section we give a fixed-parameter tractable [10] algorithm for MAXIMUMFEASIBLESUBSYSTEM that, for any fixed value of the output,  $k$ , is polynomial in the dimension,  $d$ .

The algorithm uses linear programming as a subroutine in the following way: Given a collection  $K$  of halfspaces in  $\mathbb{R}^{d-1}$ , an algorithm for linear programming can be used to either

1. Determine a point  $p \in \cap K$  if such a point exists or,
2. report a subset  $B \subseteq K$ ,  $|B| \leq d$ , such that  $\cap B = \emptyset$ .

The set  $B$  reported in the latter case is called a *basic infeasible subsystem*. Standard combinatorial algorithms for linear programming, including algorithms for linear programming in small dimensions [8, 11, 18, 19, 29, 30] as well as the simplex method (cf., [7]), can generate a basic infeasible subsystem given an infeasible linear program. The details of finding a basic infeasible subsystem using linear programming are discussed further in Appendix A.

Let  $\text{BIS}(K)$  denote a routine that outputs a basic infeasible subsystem of  $K$  if  $K$  is infeasible, and that outputs the empty set otherwise. The following algorithm solves the MAXIMUMFEASIBLESUBSYSTEM decision problem:

MFS( $K, k$ )

- 1: { $\star$  determine if there exists  $K' \subseteq K$ ,  $|K'| \leq k$ , such that  $\cap(K \setminus K') \neq \emptyset \star$ }
  - 2:  $B \leftarrow \text{BIS}(K)$
  - 3: **if**  $B = \emptyset$  **then**
  - 4:     **return** true
  - 5: **if**  $k = 0$  **then**
  - 6:     **return** false
  - 7: **for** each  $h \in B$  **do**
  - 8:     **if** MFS( $K \setminus \{h\}, k - 1$ ) = true **then**
  - 9:         **return** true
  - 10: **return** false
- The correctness of the above algorithm is easily established by induction on the value of  $k$ . The running time of the algorithm is given by the recurrence

$$T(n, k) \leq \text{LP}(n, d - 1) + dT(n - 1, k - 1) ,$$

where  $\text{LP}(n, d)$  denotes the running time of an algorithm for finding a basic infeasible system in a linear program with  $n$  constraints and  $d$  variables. This recurrence readily resolves to  $O(d^k \text{LP}(n, d - 1))$ . Using this as a subroutine for Tukey depth computation we obtain our final result:

**Theorem 5.** *The Tukey depth of a point  $p$  with respect to a set  $S$  of  $n$  points in  $\mathbb{R}^d$  can be computed in  $O(d^k \text{LP}(n, d - 1))$  time, where  $k$  is the value of the output and  $\text{LP}(n, d)$  is the time to solve a linear program with  $n$  constraints and  $d$  variables.*

**Remark:** The algorithm described above is closely related to Matoušek’s algorithm for MAXIMUMFEASIBLESUBSYSTEM which, in our setting, has a running time of  $O(k^d \text{LP}(n, d))$  [17]. In the language of fixed-parameter tractability, the primary difference between the two algorithms is that Matoušek’s algorithm explores the search tree in breadth-first order and uses a dictionary to ensure that identical subtrees are not explored, whereas the current algorithm explores the search tree in depth-first order. The two algorithms can, of course, be combined to obtain an algorithm with running time  $O(\min\{k^d + d^k\} \text{LP}(n, d))$ .

## Acknowledgement

Dan Chen, David Bremner, and Pat Morin are grateful to Diane Souvaine and the Radcliffe Institute for hosting the *Workshop on Computational Aspects of Statistical Data Depth Analysis*, July 7–10, 2006. In particular the work in Section 5 is the result of discussions that began at the workshop.

## References

- [1] E. Amaldi and V. Kann. The complexity and approximability of finding maximum feasible subsystems of linear relations. *Theoretical Computer Science*, 147:181–210, 1995.

- [2] P. A. Beling and N. Megiddo. Using fast matrix multiplication to find basic solutions. *Theoretical Computer Science*, 205(1-2):307–316, September 1998.
- [3] M. Ben-Or. Lower bounds for algebraic computation trees (preliminary report). In *Proceedings of the 15th Annual ACM Symposium on Theory of Computing (STOC'83)*, pages 80–86, 1983.
- [4] G. S. Brodal and R. Jacob. Dynamic planar convex hull. In *Proceedings of the 43rd IEEE Symposium on Foundations of Computer Science (FOCS 2002)*, pages 617–626, 2002.
- [5] T. M. Chan. An optimal randomized algorithm for maximum Tukey depth. In *Proceedings of the 15th ACM-SIAM Symposium on Discrete Algorithms*, pages 423–429, 2004.
- [6] T. M. Chan. Low-dimensional linear programming with violations. *SIAM Journal on Computing*, 34:879–893, 2005.
- [7] V. Chvátal. *Linear programming*. A Series of Books in the Mathematical Sciences. W. H. Freeman and Company, New York, 1983.
- [8] K. L. Clarkson. Las Vegas algorithms for linear and integer programming when the dimension is small. *Journal of the ACM*, 42:488–499, 1995.
- [9] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms*. McGraw Hill, 2nd edition, 2001.
- [10] R. G. Downey and M. R. Fellows. *Parameterized Complexity*. Monographs in Computer Science. Springer, 1998.
- [11] M. E. Dyer. Linear time algorithms for two- and three-variable linear programs. *SIAM Journal on Computing*, pages 31–45, 1984.
- [12] C. C. Gonzaga. An algorithm for solving linear programming problems in  $O(n^3L)$  operations. In *Progress in mathematical programming (Pacific Grove, CA, 1987)*, pages 1–28. Springer, New York, 1989.
- [13] D. S. Johnson and F. P. Preparata. The densest hemisphere problem. *Theoretical Computer Science*, 6:93–107, 1978.
- [14] S. Langerman and W. Steiger. Optimization in arrangements. In *Proceedings of the 20th Symposium on Theoretical Aspects of Computer Science*, volume 2607 of *Lecture Notes in Computer Science*, pages 50–61. Springer-Verlag, 2003.
- [15] J. Matoušek. Computing the center of planar point sets. In J. E. Goodman, R. Pollack, and W. Steiger, editors, *Computational Geometry: Papers from the Special Year*, pages 221–230. AMS, Providence, 1991.

- [16] J. Matoušek. Reporting points in halfspaces. *Computational Geometry: Theory and Applications*, 2:169–186, 1992.
- [17] J. Matoušek. On geometric optimization with few violated constraints. *Discrete & Computational Geometry*, 14:365–384, 1995.
- [18] N. Megiddo. Linear time algorithms for linear programming in  $\mathbb{R}^3$  and related problems. *SIAM Journal on Computing*, 12:759–776, 1983.
- [19] N. Megiddo. Linear programming in linear time when the dimension is fixed. *Journal of the ACM*, 31:114–127, 1984.
- [20] N. Megiddo. On finding primal- and dual-optimal bases. *ORSA J. Comput.*, 3(1):63–65, 1991.
- [21] K. Mulmuley. *Computational Geometry: An Introduction Through Randomized Algorithms*. Prentice-Hall, 1998.
- [22] F. P. Preparata and M. I. Shamos. *Computational Geometry: An Introduction*. Springer-Verlag, New-York, 1985.
- [23] E. Ramos. On range reporting, ray shooting, and  $k$ -level construction. In *Proceedings of the 15th ACM Symposium on Computational Geometry (SoCG 1999)*, pages 390–399, 1999.
- [24] J. Renegar. A polynomial-time algorithm, based on Newton’s method, for linear programming. *Math. Programming*, 40(1, (Ser. A)):59–93, 1988.
- [25] C. Roos. An  $O(n^3L)$  approximate center method for linear programming. In *Optimization (Varetz, 1988)*, volume 1405 of *Lecture Notes in Math.*, pages 147–158. Springer, Berlin, 1989.
- [26] C. Roos and J.-Ph. Vial. A polynomial method of approximate centers for linear programming. *Math. Programming*, 54(3, Ser. A):295–305, 1992.
- [27] P. J. Rousseeuw and I. Ruts. Constructing the bivariate Tukey median. *Statistica Sinica*, 8:827–839, 1998.
- [28] I. Ruts and P. J. Rousseeuw. Computing depth contours of bivariate point clouds. *Computational Statistics and Data Analysis*, 23:153–168, 1996.
- [29] R. Seidel. Small-dimensional linear programming and convex hulls made easy. *Discrete & Computational Geometry*, 6:423–434, 1991.
- [30] M. Sharir and E. Welzl. A combinatorial bound for linear programming and related problems. In *Proceedings of the 9th Symposium on Theoretical Aspects of Computer Science*, volume 5777 of *Lecture Notes in Computer Science*, pages 569–579. Springer-Verlag, 1992.

- [31] C. G. Small. A survey on multidimensional medians. *International Statistics Review*, 58:263–277, 1990.
- [32] J. W. Tukey. Mathematics and the picturing of data. In *Proceedings of the International Congress of Mathematicians*, volume 2, pages 523–531, 1975.
- [33] P. M. Vaidya. An algorithm for linear programming which requires  $O(((m+n)n^2 + (m+n)^{1.5}n)L)$  arithmetic operations. *Math. Programming*, 47(2, (Ser. A)):175–201, 1990.
- [34] S. A. Vavasis and Y. Ye. Identifying an optimal basis in linear programming. *Ann. Oper. Res.*, 62:565–572, 1996. Interior point methods in mathematical programming.

## A Computing a Basic Infeasible Subsystem

This appendix explains how, given an infeasible linear program, to find a basic infeasible subsystem of that linear program. This routine is required as part of the algorithm for solving MAXIMUMFEASIBLESUBSYSTEM described in Section 5.

For any matrix  $M$ , let  $M_J$  denote the set of rows indexed by  $J$ . Given a system of linear inequalities  $Mx \geq b$ ,  $M \in \mathbb{R}^{m \times d}$ , a *basic infeasible subsystem* is a subset of  $\{1 \dots m\}$  such that the system  $M_I x \geq b_I$  is infeasible, and  $|I| \leq d + 1$ . We consider the standard first stage simplex problem (see e.g. [7], p. 39). Let  $e$  denote the  $m$ -vector of all ones,  $c$  the length  $d + 1$  binary vector with exactly one one in the last position and let  $A = [Me]$ . We can write the first stage LP for our system as

$$\begin{aligned} \min c^T x &= x_{d+1} \\ \text{subject to} & \\ Ax &\geq b \end{aligned} \tag{P}$$

In the case of an infeasible system, the optimal value of this LP will be strictly positive. The dual LP of (P) is

$$\begin{aligned} \max b^T y & \\ \text{subject to} & \\ yA &= c \\ y &\geq 0 \end{aligned} \tag{D}$$

In what follows, we generally follow the notation of [20], except that we interchange the definitions of the primal and dual LPs. Define a *basic partition* (or just *basis*)  $(\beta, \eta)$  as a partition of the row indices of  $A$  such that  $A_\beta$  is nonsingular. For each basic partition, we define a *primal basic solution*

$$x^* = A_\beta^{-1} b_\beta$$

and a *dual basic solution*

$$y^* = c A_\beta^{-1}$$

We say that a basis is *primal feasible* (resp. *dual feasible*) if  $x^*$  is feasible for (P) (respectively  $y^*$  is feasible for (D)). It is a standard result of linear programming duality that a basis which is both primal and dual feasible defines a pair  $(x^*, y^*)$  of optimal solutions to the primal and dual LP's; such a partition is called an *optimal basis partition*

In general LP algorithms (either directly in the case of Simplex type algorithms, or via postprocessing using e.g. [20, 34, 2]) provide an optimal basis partition  $(\beta, \eta)$ . Consider the relaxed LP

$$\begin{aligned} & \min c^T x \\ & \text{subject to} \\ & A_\beta x \geq b_\beta \end{aligned} \tag{R}$$

It is easy to verify that an optimal basis partition for (P) is also primal and dual feasible for (R). This implies that the system  $M_\beta x \geq b_\beta$  is infeasible, and provides a basic infeasible system. Using interior point algorithms (see [12, 24, 33, 25, 26]), a solution to the first stage LP can be found  $O(m^3 L)$  time, where  $L$  is the number of bits in the input. Using the algorithm of Beling and Megiddo [2], an optimal basis partition can be computed in  $O(m^{1.594} d)$  time (where the bound is based on the best known methods for fast matrix multiplication).