

# ASSIGNMENT 2 – COMP 252

Alexandre St-Aubin & Jonathan Campana

January 30, 2024

1. **ALGORITHM DESIGN.** You are given  $n$  vectors  $x_1, \dots, x_n$  in  $\mathbb{Z}^n$ . Design an efficient algorithm in the ram model for computing for each  $x_i$  one of its nearest neighbors among the other points, using the standard Euclidean metric to measure distances. You can't use real numbers, and operations like square root are not available. Nevertheless, show how this can be done in  $o(n^3)$  worst-case time.

*Solution:*

We first note that minimizing the Euclidean distance between 2 vectors of length  $n$  is equivalent to minimizing the sum of squared differences between each individual components. The reason for this is that the square root is a monotone increasing function. Let  $x_i = (x_{i,1}, \dots, x_{i,n})$ ,  $x_j = (x_{j,1}, \dots, x_{j,n})$  both in  $\mathbb{Z}^n$ , and define

$$d_2^2(x_i, x_j) := \sum_{k=1}^n (x_{i,k} - x_{j,k})^2 = \sum_{k=1}^n (x_{i,k}^2 - 2x_{i,k}x_{j,k} + x_{j,k}^2) = \langle x_i, x_i \rangle + 2\langle x_i, x_j \rangle + \langle x_j, x_j \rangle \quad (1)$$

It is obvious that  $d_2^2 \sim O(n)$  in the RAM model, and that no real numbers, nor square roots were used to compute it. Now, we notice that  $\frac{n(n-1)}{2}$  distances need to be computed, and in view of dynamic programming, we find a way to compute everything at once in order to reduce the complexity that would occur if we were to get each  $d_2^2$  individually, namely,  $O(n^3)$ . Construct the following matrix,

$$X = \begin{pmatrix} x_{1,1} & x_{1,2} & x_{1,3} & \dots & x_{1,n} \\ x_{2,1} & x_{2,2} & x_{2,3} & \dots & x_{2,n} \\ x_{3,1} & x_{3,2} & x_{3,3} & \dots & x_{3,n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & x_{n,3} & \dots & x_{n,n} \end{pmatrix}$$

Then,

$$X \cdot X^T = \begin{pmatrix} \sum_{k=1}^n x_{1,k}x_{1,k} & \sum_{k=1}^n x_{1,k}x_{2,k} & \dots & \sum_{k=1}^n x_{1,k}x_{n,k} \\ \sum_{k=1}^n x_{2,k}x_{1,k} & \sum_{k=1}^n x_{2,k}x_{2,k} & \dots & \sum_{k=1}^n x_{2,k}x_{n,k} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{k=1}^n x_{n,k}x_{1,k} & \sum_{k=1}^n x_{n,k}x_{2,k} & \dots & \sum_{k=1}^n x_{n,k}x_{n,k} \end{pmatrix} = \begin{pmatrix} \langle x_1, x_1 \rangle & \langle x_1, x_2 \rangle & \dots & \langle x_1, x_n \rangle \\ \langle x_2, x_1 \rangle & \langle x_2, x_2 \rangle & \dots & \langle x_2, x_n \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle x_n, x_1 \rangle & \langle x_n, x_2 \rangle & \dots & \langle x_n, x_n \rangle \end{pmatrix}$$

where the entries of  $X \cdot X^T$  are exactly the dot products needed in (1). Therefore, by employing STRASSEN'S algorithm, we can efficiently compute the product  $X \cdot X^T$  with a time complexity of  $O(n^{2.807})$ . Subsequently, each of the  $\frac{n(n-1)}{2}$  distances can be computed in  $O(n^2)$  time, as computing one distance will take constant time by accessing the dot products in the matrix previously computed. Simultaneously, these distances can be added to an ordered list (corresponding to each individual vector) in constant time. Upon completion of this process, we will have the closest vector to any given vector readily accessible. The algorithm outlined above has a complexity of  $O(n^{2.807})$ , demonstrating that it can be accomplished in worst-case time less than  $o(n^3)$ .

## 2. DYNAMIC PROGRAMMING: COMPUTING THE OPTIMAL STAR.

---

**Algorithm 1:** Algorithm to compute the optimal star.

---

**Input:** A matrix  $D[i, j]$  of pairwise distances between cities  $1 \leq i, j \leq n$ .**Output:** The cost of the optimal star that links the given cities.

```

// initialize M[v, S] for sets of size 1
1 for all  $i \in [1, n]$  do
2   for all  $j \neq i$  do
3      $M[i, \{j\}] \leftarrow D[i, j]$ ;

// loop through each central city
4 for all  $v \in [1, n]$  do
  // loop through each subset size
5   for all  $k \in [1, n-1]$  do
    // loop through each subset of size k excluding the central city
6     for all  $S$  with  $|S| = k$ ,  $S \subseteq \{1, \dots, n\} \setminus \{v\}$  do
7        $M[v, S] \leftarrow \min_{l \in S} \{M[v, S \setminus \{l\}] + \min\{d(l, x) : x = v \text{ or } \deg(x) = 1\}\};$ 

8 return  $\min_{v \in [1, n]} \{M[v, S] : |S| = n\}$ 

```

---

*Remark.* The above is heavily inspired from the *Held–Karp algorithm*, which was presented in class.

**Analysis of the algorithm**

Allow  $\Omega$  to represent the set of all vertices (cities). The outer loop iterates through each potential central vertex  $v$  of the star, while the subsequent loops systematically explore subsets of  $\Omega \setminus v$  in a bottom-up fashion. Starting with subsets of size 1 (i.e.,  $S = s$ ), we know  $M[v, S] = d(v, s)$ . We then progress through sets of increasing size, using our knowledge from preceding, smaller sets to obtain an optimal set length, until we reach  $M[v, \Omega \setminus v]$ , thereby identifying the optimal star configuration with the central vertex  $v$ .

Let's delve into the time complexity analysis. Line 4 executes  $n$  times. Lines 5 and 6 represent all possible subsets of a set with  $n-1$  vertices, resulting in  $2^{n-1}$  iterations. We can safely assume line 7 to be  $O(n^2)$ , as we are looping through each  $l \in \Omega$ , and for each  $l$ , we loop through each endpoint, of which we can assume there are at most  $n$ . Combining these components, the overall time complexity is given as  $O(n^3 \cdot 2^n)$ . We note that lines 1 & 8 run in  $O(n)$  time, so they play no significant role in the overall time complexity of the algorithm.

3. INDUCTION. We are given the recurrence

$$T_n = 2T_{\frac{n}{a}} + 7T_{\frac{n}{a^2}} + 1,$$

where  $a \geq 2$  is a given integer, and  $n$  is restricted to be a power of  $a$ . We also know that  $T_1 = T_a = 1$ .

(a)  $T_n = \Omega(n^c)$  for some constant  $c$ .

*Proof. Base Case:* We consider the base case where  $n = a^2$ , then,

$$\begin{aligned} T_{a^2} &= 2T_{\frac{a^2}{a}} + 7T_{\frac{a^2}{a^2}} + 1 \\ &= 2T_a + 7T_1 + 1 \\ &= 2(1) + 7(1) + 1 \\ &= 10 \end{aligned}$$

We will denote the following to get our lower bound,

$$T_{a^2} = 10 \geq \alpha(a^2)^c$$

Thus we will let  $\alpha = \frac{1}{(a^2)^c}$ , so the inequality holds.

**Induction Step:** Here we will be supposing that  $T_n \geq \alpha n^c$ .

$$\begin{aligned} T_n &= 2T_{\frac{n}{a}} + 7T_{\frac{n}{a^2}} + 1 \\ &\geq 2\alpha\left(\frac{n}{a}\right)^c + 7\alpha\left(\frac{n}{a^2}\right)^c + 1 \\ &\geq 2\alpha\left(\frac{n}{a}\right)^c + 7\alpha\left(\frac{n}{a}\right)^c \\ &= \alpha n^c \left( \frac{2}{a^c} + \frac{7}{a^{2c}} \right) \end{aligned}$$

Thus, we must find  $c$  such that

$$\begin{aligned} \frac{2}{a^c} + \frac{7}{a^{2c}} &\leq 1 \\ \frac{1}{a^c} \left( 2 + \frac{7}{a^c} \right) &\leq 1 \end{aligned}$$

Let  $x = a^c$ ,

$$\begin{aligned}\frac{1}{x} \left( 2 + \frac{7}{x} \right) &\leq 1 \\ \left( 2 + \frac{7}{x} \right) &\leq x \\ x - \frac{7}{x} - 2 &\geq 0 \\ x^2 - 2x - 7 &\geq 0\end{aligned}$$

By quadratic formula,

$$\begin{aligned}x &\geq \frac{2 \pm \sqrt{(-2)^2 - 4(1)(-7)}}{2} \\ &= \frac{2 \pm \sqrt{32}}{2} \\ &= \frac{2 \pm 4\sqrt{2}}{2} \\ &= 1 \pm 2\sqrt{2}\end{aligned}$$

Since we know that  $x > 0$ , we take the positive  $x$  from the quadratic formula. Therefore,

$$\begin{aligned}x = a^c &= 1 + 2\sqrt{2} \\ \log_a a^c &= \log_a (1 + 2\sqrt{2}) \\ c &= \log_a (1 + 2\sqrt{2})\end{aligned}$$

So we have found our  $c$ .

□

(a) Show by induction that  $T_n = O(n^c)$  for the same maximal  $c$  found above.

*Proof. Base Case:* We consider the base cases for the upper bound, we know that

$$\begin{aligned}T_1 &= 1 \\ T_a &= 1\end{aligned}$$

and we want to prove that  $T_n = O(n^c)$  thus we show  $\exists \alpha \in^+$  such that

$$T_1 \leq \alpha 1^c - 1 = \alpha 1^{\log_a(1+2\sqrt{2})} - 1 = \alpha(1) - 1 = \alpha - 1 \quad \text{where } \alpha \geq 2$$

and similarly,

$$\begin{aligned} T_a &\leq \alpha a^{\log_a(1+2\sqrt{2})} - 1 = \alpha(1+2\sqrt{2}) - 1 \\ &= \alpha + 2\alpha\sqrt{2} - 1 \\ &= (\alpha - 1) + 2\alpha\sqrt{2} \quad \text{where } \alpha \geq 2 \end{aligned}$$

Which the inequality will hold  $\forall n$  through induction.

**Induction step:** Since  $n$  is a power of  $a$ , we can express it as  $x = a^k$ ,  $k \in \mathbb{N}$ . Let  $y \in \mathbb{N}$  s.t  $y \geq 2$ , it follows that  $\forall k \in \{0, 1, \dots, y-1\}$ ,

$$\begin{aligned} T_n = T_{a^k} &\leq \alpha(a^k)^c - 1 = \alpha(a^c)^k - 1 \\ &= \alpha(a^{\log_a 1+2\sqrt{2}})^k - 1 \\ &= \alpha(1+2\sqrt{2})^k - 1 \end{aligned}$$

where  $\alpha \geq 2$  from our base case. Therefore assuming the induction hypothesis and letting  $n = a^y$ ,

$$\begin{aligned} T_n = T_{a^y} &= 2T_{a^{y-1}} + 7T_{a^{y-2}} + 1 \\ &\leq 2\alpha(a^{y-1})^c - 2 + 7\alpha(a^{y-2})^c - 7 + 1 \quad \text{by induction hypothesis} \\ &= 2\alpha(1+2\sqrt{2})^{y-1} + 7\alpha(1+2\sqrt{2})^{y-2} - 8 \\ &= \alpha(1+2\sqrt{2})^{y-1} \left( 2 + \frac{7}{(1+2\sqrt{2})} \right) - 8 \\ &= \alpha(1+2\sqrt{2})^y - 8 \quad \text{*see below for equality} \\ &\leq \alpha(1+2\sqrt{2})^y - 1 \end{aligned}$$

Thus we are done.

\* In the above part of the induction is used the fact that,

$$2 + \frac{7}{1+2\sqrt{2}} = 1 + 2\sqrt{2}$$

We show this is the case,

$$\begin{aligned} 2 + \frac{7}{1+2\sqrt{2}} &= \frac{2(1+2\sqrt{2}) + 7}{1+2\sqrt{2}} \\ &= \frac{2+4\sqrt{2}+7}{1+2\sqrt{2}} \\ &= \frac{9+4\sqrt{2}}{1+2\sqrt{2}} \left( \frac{1+2\sqrt{2}}{1+2\sqrt{2}} \right) \\ &= \frac{(9+4\sqrt{2})(1+2\sqrt{2})}{9+4\sqrt{2}} \\ &= 1+2\sqrt{2} \end{aligned}$$

□

## 4. SORTING WITH DUPLICATES.

- (i) Give a divide-and-conquer algorithm that sorts  $n$  numbers with  $1 \leq k \leq n$  uniques. This can be done in time  $O(n \log_2(k+1))$  with a ternary comparison oracle.

See Algorithm 2 on next page.

- (ii) Prove the complexity claim.

*Proof.* For the sake of simplicity, we assume that  $k = 2^i$  for some  $i \in \mathbb{N}$ . To prove that the complexity of our algorithm is  $O(n \log(k+1))$ , we count the number of times that it uses the ternary oracle (number of comparisons made), as the rest costs 0, by the oracle model of complexity. We notice that the comparisons occur only at the *merge* step of the algorithm, so we'll start counting at the first merge, namely when we have  $n$  subarrays of size 1. We assume the worst case, i.e. that no duplicates are found in the merges until level  $\log_2 k$  is reached (see Figure 1). When this level is reached, we have  $\frac{n}{k}$  sets of length at most  $k$ , and the next step will surely find duplicates, thereby reducing the running time compared to plain old merge sort.

Let's count the number of comparison needed to get to level  $\log_2 k$  of the merging. In order to get one subarray of length  $k$ ,  $k$  subarrays of length 1 must be merged into 1. This process will take  $\log_2 k$  merges, each one costing  $\frac{(2^i-1)k}{2^i}$  comparisons, that is,

$$\sum_{i=1}^{\log_2 k} \frac{(2^i-1)k}{2^i} \leq \sum_{i=1}^{\log_2 k} k = n \log_2 k \text{ comparisons.}$$

And, since there are  $\frac{n}{k}$  such subarrays, it follows that the complexity of the above process is  $O(n \log_2 k)$ .

We now consider the second and last part of our computation. At this step, we have  $\frac{n}{k}$  sets of size  $k$ . In the normal merge sort algorithm, the next step would generate half as many sets of double the size ( $2k$ ). However, our algorithm combines duplicate terms into a single tuple element, and since there are only  $k$  elements that are different in our set,  $k$  is an upper bound for the size of a subarray. Hence, we instead obtain half the amount of sets, but still of size at most  $k$ . This is also true for all the following merges.

Let's compute the time complexity of merging  $\frac{n}{k}$  subarrays until we have 1 array of length  $k$ . Well, it is clear that each merge will cost at most  $k$  comparisons, as in the worst case, we are merging two arrays of length  $k$ , which must therefore be identical, as they are ordered and contain the same  $k$  elements. Thus, the oracle will yield  $k$  equalities, which suffices to merge the arrays. So, the total number of comparisons is given by

$$k \sum_{i=0}^{\log_2 \frac{n}{k}} 2^i = k(2^{\log_2 \frac{n}{k}} - 1) = k\left(\frac{n}{k} - 1\right) = n - k$$

Combining both processes, it is easy to see that the total time complexity is

$$O(n - k + n \log_2(k)) = O(n \log_2(k+1))$$

*Remark.* The reason why the complexity is  $O(n \log_2(k+1))$ , and not  $O(n \log_2(k))$  is that in the case where  $k = 1$ , if we need the complexity to be  $O(n \log_2(1))$ , not  $O(0)$ .

□

---

**Algorithm 2:** DC algorithm to sort a list of  $n$  elements with  $1 \leq k \leq n$  duplicates.

---

**Input:** An arbitrary array of integers.**Output:** The sorted array.

```

// Initialize the array with Tuples
1 for  $i = 0$  to  $l(\text{array})$  do
2    $\text{array}[i] \leftarrow \text{new Tuple}(i, 1);$ 

// A function to compare tuples, each use costs 1. Anything else is free.
3 Function TernaryOracle( $\text{tuple1}, \text{tuple2}$ ):
4   if  $\text{tuple1}[0] = \text{tuple2}[0]$  then
5     return 0;
6   else if  $\text{tuple1}[0] < \text{tuple2}[0]$  then
7     return -1;
8   else
9     return 1;

10 Function MergeSort( $\text{array}$ ):
11   if  $l(\text{array}) < 2$  then
12     return  $\text{array}$ ;
13   else
14      $\text{mid} \leftarrow \lfloor l(\text{array})/2 \rfloor;$ 
15     return Merge (MergeSort ( $\text{array} [: \text{mid}]$ ), MergeSort ( $\text{array} [\text{mid}:]$ ));

16 Function Merge( $\text{Left}, \text{Right}$ ):
17    $i, j \leftarrow 0;$ 
18    $\text{array} \leftarrow \text{new array};$ 
19   while  $i \leq l(\text{Left}) \ \&\& \ j \leq l(\text{Right})$  do
20     switch TernaryOracle( $\text{Left}[i], \text{Right}[j]$ ) do
21       case 0 do
22          $\text{array}[i + j] \leftarrow \text{new Tuple}(\text{Left}[i][0], \text{Left}[i][1] + \text{Right}[j][1]);$ 
23          $i ++, j ++;$ 
24       case 1 do
25          $\text{array}[i + j] \leftarrow \text{Left}[i];$ 
26          $i ++;$ 
27       case -1 do
28          $\text{array}[i + j] \leftarrow \text{Right}[j];$ 
29          $j ++;$ 
30   // then add the rest of the array that we didn't reach the end of to
   the current array:
    $\text{array} \leftarrow \text{array} \& (i \leq l(\text{Left})) * \text{Left}[i : l(\text{Left})] \& (j \leq l(\text{Right})) * \text{Right}[j : l(\text{Right})]$ 

```

---

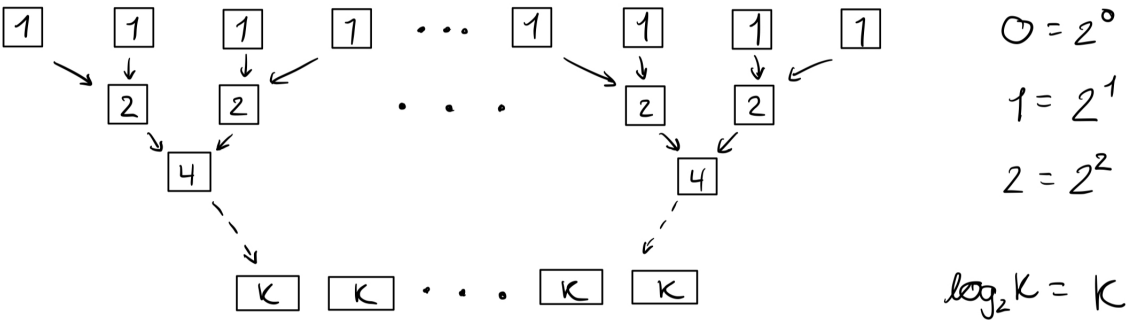


Figure 1: The merging of subarrays.