

A Lecture on Divide-and-Conquer Algorithms and the Select Problem

Olivia Shi, Yutong Zhang

February 27, 2018

This is the augmented transcript of a lecture given by Luc Devroye on January 18th, 2018 for the Honours Algorithms and Data Structures class (COMP 252, McGill University). The subject was Euclid's algorithm, ham-sandwich and pancake theorems, half space counting and the linear time selection problem.

Euclid's Algorithm

Euclid's algorithm efficiently computes the greatest common divisor (GCD) of two nonnegative numbers m, n ($m < n$), the largest number that divides both of them without leaving a remainder.

$\text{GCD}(n, m)$

- 1 **if** $m == 0$
- 2 **return** $m = 0$
- 3 **else**
- 4 **return** $\text{GCD}(m, n \bmod m)$

Denote T_n as the worst case "time" for any $\text{GCD}(k, l)$ with $n \geq k \geq l > 0$. We have the complexity $T_n \leq 2 \log_2 n$, since n gets halved every second iteration, and a number n can get halved at most $\log_2 n$ times.

Exercise 1. Using the bit model, show that the complexity of $\text{GCD}(n, m)$ is $O(\log_2 n \cdot \log_2 m)$.

Ham-Sandwich Theorem (3D)

The Ham-Sandwich Theorem¹ describes the following: take a sandwich made of a slice of ham and two slices of bread. As long as one's knife is long enough, one can cut all three pieces in half in only one pass. The precise mathematical statement of the theorem, generalized to n dimensions, is that given n compact sets in \mathbb{R}^n , there is a hyperplane that bisects each set so that the two halves of both sets have equal measure.

CAN WE CUT A HAMBURGER EXACTLY IN HALF SO THAT EACH HALF HAS EXACTLY 50% OF THE CHEESE, 50% OF THE MEAT, AND 50% OF THE BREAD? Yes.

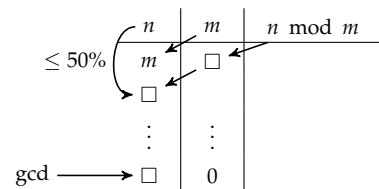


Figure 1: Illustration of Euclid's method. The number written two rows below n is at most 50% of m — we say that n is "halved".

¹ Libgober [2008]

RED AND BLACK POINTS

In 2D, this theorem is known as the pancake theorem. It implies that we can perfectly bisect n red points and n black points in the plane, so that each side has $n/2$ red and $n/2$ black points.

Application: We can divide m points into four sets of $n/4$ points each by two lines (Figure 2). First draw a horizontal line that divides the points in half. Color all points below the line red, and above the line black. By the pancake theorem, we can find a second line that evenly divides red and black, yielding $n/4$ points in each of the four sets of the partition.

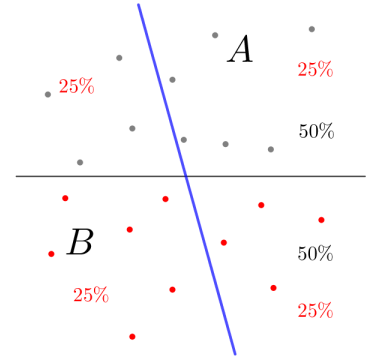


Figure 2: Illustration of the application.

Half Space Counting Problem

Let x_1, \dots, x_n be drawn from \mathbb{R}^2 . Make a data structure such that one can “efficiently” answer queries that take as input a line given by the user, and outputs the number of points on one side of the line.

DATA STRUCTURE (DIVIDE-AND-CONQUER) $COUNT(\ell, S)$ is a divide-and-conquer algorithm for computing the number of points of S on one side of ℓ , say the side that contains the origin. Assume that we partitioned our space recursively using the 25%-trick suggested by the pancake theorem.

we return the count to the set S .

$COUNT(\ell, S)$

```

1  if  $|S| \leq 10$ 
2      do it manually
3  else
4      determine the three sets cut by  $\ell$ , say  $A, B, C$ 
5      if the fourth set is on the good side of  $\ell$ 
6          return  $|S|/4 + COUNT(\ell, A) + COUNT(\ell, B) + COUNT(\ell, C)$ 
7      else
8          return  $COUNT(\ell, A) + COUNT(\ell, B) + COUNT(\ell, C)$ 
    
```

We can indeed determine A, B, C in constant time. In the RAM model, $T_n = 1 + 3T_{n/4}$, which yields $T_n = \Theta(n^{\log_4 3})$ by the master theorem.

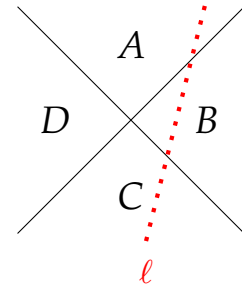


Figure 3: Illustration of the $COUNT(\ell, S)$ algorithm.

Linear time selection: finding the k^{th} smallest element

We present a selection algorithm invented by Blum, Floyd, Pratt, Rivest and Tarjan² for finding the k^{th} smallest number in a list or array; such a number is called the k^{th} order statistic. This includes the cases of finding the minimum, maximum, and median elements.

² Blum et al. [1973]

In particular, we want to find the k^{th} smallest number in an un-ordered set $S = \{x_1, \dots, x_n\}$, and can use a comparison oracle. By sorting we would have time complexity $O(n \log_2 n)$.

We will give an $O(n)$ complexity solution, called $\text{SELECT}(k, S)$, where S is the collection of elements, and $1 \leq k \leq |S|$.

$\text{SELECT}(k, S)$

```

1  if  $|S| \leq 5$ 
2      SORT( $S$ )
3      return the  $k^{\text{th}}$  smallest element of  $S$ 
4  else
5      group all elements in groups of 5, and find the median in
        each group and call the set of medians  $M$ 
        // This costs 6 comparisons.
6      let  $m = \text{SELECT}(|M|/2, M)$ 
        //  $m$  is the median of the medians
7      compare all elements of  $S$  with  $m$ , forming the sets  $L$  and  $R$ 
        of smaller and larger elements
8      if  $|L| == k - 1$ 
9          return  $m$ 
10     elseif  $|L| \geq k$ 
11         return  $\text{SELECT}(k, L)$ 
12     else
13         return  $\text{SELECT}(k - |L| - 1, R)$ 

```

The bound on the cost for every critical step of the $\text{SELECT}(k, S)$ algorithm is listed in the table below:

Line number	Bound on the cost
1 - 3	≤ 7
5	$\leq 6n/5$
6	$= T_{n/5}$
7	$= n$
11	$= T_{ L } \leq T_{7n/10}$
13	$= T_{ R } \leq T_{7n/10}$

A POSSIBLE IMPROVEMENT. As shown in Figure 5, L and R each have at least $3|S|/10$ and at most $7|S|/10$ elements. Hence, to identify L and R , if we program correctly, we only require at most $4n/10$ new comparisons. Therefore, the recurrence that calculates the complexity can be reduced to $T_n \leq 8n/5 + T_{n/5} + T_{7n/10}$. We will prove by induction that $T_n \leq Cn$ for all $n \in \mathbb{N}$.

Proof. Base case: If $n \leq 5$, then $C \geq 7$ is a safe choice.

Induction hypothesis: Assume $T_k \leq Ck$ for all $k < n$.

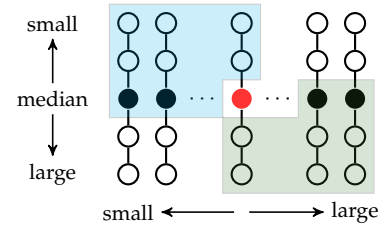


Figure 4: Medians of the median.

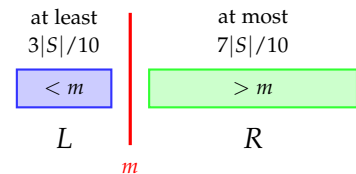


Figure 5: Illustration of L and R .

Inductive step: $T_n \leq 8n/5 + T_{n/5} + T_{7n/10} \leq 8n/5 + C \cdot 9n/10$.
 This should be $\leq Cn$, so we have the requirement $8n/5 \leq C \cdot n/10$, or $C \geq 16$. We have thus shown that $T_n \leq 16n$ for all $n \in \mathbb{N}$. \square

MEDIAN-OF-3 RECURRENCE. When we replace the median-of-5 step by a median-of-3 step, then one can see that $T_n = \Theta(n) + T_{n/3} + T_{2n/3}$. This equation is analyzed via a recursion tree. Note that in each level the work adds up to n . We have $(\frac{2}{3})^k n = 1$, where k is the height of the recursion tree, so $k = \log_2 n / \log_2(3/2)$. Therefore, $T_n = \Theta(n \log_2 n)$.

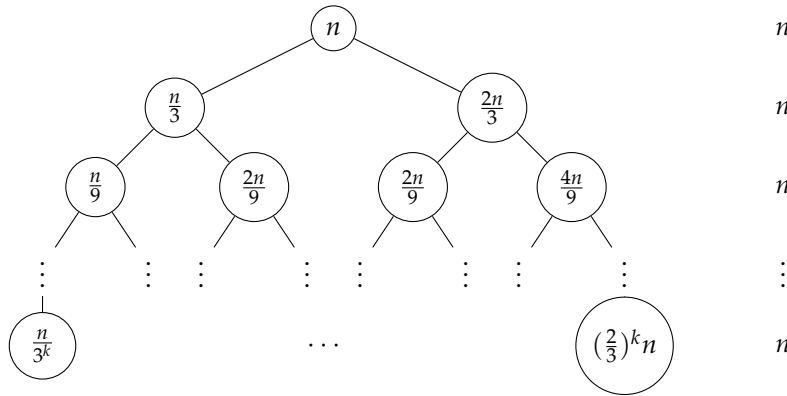


Figure 6: Recursion tree.

Observe that the recursion tree is not balanced. Nevertheless, the work at each level is precisely n .

References

M. Blum, R. W. Floyd, V. R. Pratt, R. L. Rivest, and R. E. Tarjan. Time bounds for selection. *Journal of Computer and System Sciences*, pages 448–461, Aug 1973. DOI: 10.1016/S0022-0000(73)80033-9. URL <http://people.csail.mit.edu/rivest/pubs/BFPRT73.pdf>.

Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms*. The MIT Press, 3rd edition, 2009. ISBN 978-0-262-53305-8.

Brian Libgober. The Borsuk-Ulam and ham-sandwich theorems. May 2008. URL <http://www.math.uchicago.edu/~may/VIGRE/VIGRE2008/REUPapers/Libgober.pdf>.