

INF1008
DATA STRUCTURES & ALGORITHMS
PART 1 ASSIGNMENT

Group 10

Name	Student ID
Wong Yok Hung	2202391
Christopher Kok	2203503
Chong Hou Wei	2201783
Vivian Ng	2203557
Methinee Ang	2202781

Contents

Task 1	Page 3
1.1 User Guide	3
1.2 Structure/Design	3
1.3 Limitations	3
1.4 Testing	3
Task 2	Page 5
2.1 User Guide	5
2.2 Structure/Design	5
2.3 Limitations	5
2.4 Testing	5
Task 3	Page 7
Bibliography	

List of Figures

1.1	Quick Sort vs Quick Select time efficiency	4
2.1	Quicksort implementation vs BTreeMap implementation time efficiency	6

Task 1

1.1 User Guide

To compile the executable, the user must have an updated version of rustlang's cargo tool installed (if they would like to conduct tests on the algorithm as well as run it). The user then just needs to input into the command line: `cargo run <phone number> <phone number> <phone number>`. Alternatively, they may run the following on the command line should the program already be compiled: `qn_1.exe <phone number> <phone number> <phone number>`. The program will then output the median phone number from the given input phone numbers. Note that phone number arguments with whitespace in them must be delimited by quotation marks.

Example `cargo run 123-456-7890 "(323) 456-7890" +1 223-456-7890 1-322-345-7890 "322 555 0000"`

1.2 Structure/Design

The program takes in an input vector of strings, after which it will clean the input string into phone numbers of length 10. This process is $O(n)$, as it must iterate across the entire vector once. Next, the vector is passed into a function to find the median values. If the length of the array is an even number, the algorithm runs the quickselect algorithm twice, otherwise once.

The quickselect algorithm is based on the quicksort algorithm, where it only recurses on the slice of the vector where the desired k-th value is. This leads to quickselect having an average case time complexity of $O(n)$, with a worst case of $O(n^2)$, compared to quicksort's average and worst case time complexity of $O(n \log n)$ and $O(n^2)$ respectively. Due to this $O(n)$ average case, we opted to use a vector as the data structure over other data structures like linked lists or tree-based structures.

1.3 Limitations

Despite the $O(n)$ average case time complexity, it shares the same $O(n^2)$ worst case time complexity due to its similarities with quicksort, where a non-optimally chosen pivot leads to $O(n^2)$. This may be alleviated by using hybrid algorithms such as introselect, which uses both the quickselect and median of medians algorithm, depending on which algorithm is expected to perform better for the given input slice.

Due to our implementation of quickselect's inability to select more than one k, the quickselect algorithm must be run twice if the input array's length is an even number, doubling the runtime of the algorithm.

1.4 Testing

Four main tests are conducted to verify the correctness of the algorithm and its performance. The program passes all tests.

- Testing to find the median value for an unsorted array against the value of the median index of a sorted copy of the unsorted array. This array has a random length and randomly chosen integer values. This test is run 100 times.
- Testing that quickselect finds the correct value in the worst case scenario (a sorted array)

- Testing that quickselect finds the median value for an unsorted array of randomly generated 10 digit phone numbers. This test is run 100 times.
- Comparing the efficiency of the algorithm against an implementation of quicksort. Like test 2, the array consists of randomly generated phone numbers, and the comparisons are conducted from an array length of 1 to 100,000. To aid in completing the test faster, the comparisons for each array length are split among all logical processors of the user's CPU. For slices of the arrays that are the same length, their comparisons are run on the same thread to ensure there is no variance between the conditions of the test for both algorithms. Each test is then run 5 times and a mean of the time taken for both algorithms to find the k -th smallest value is recorded to a file and is shown below:

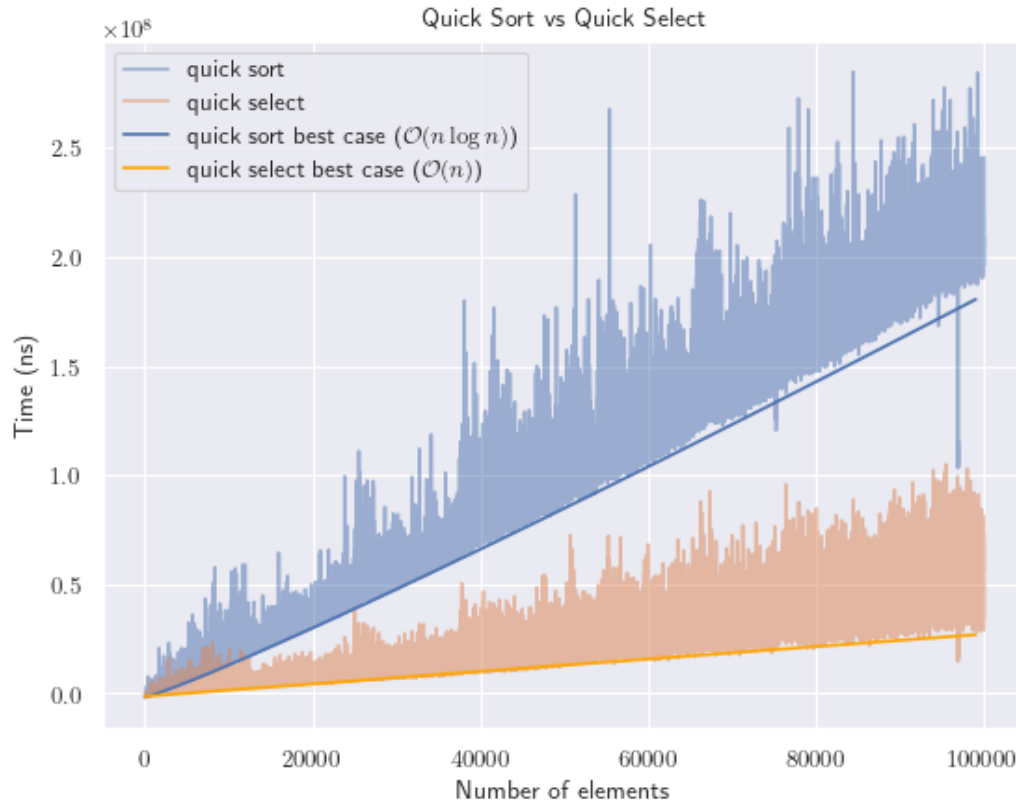


Figure 1.1: Quick Sort vs Quick Select time efficiency

Task 2

2.1 User Guide

The program takes 3 positional arguments:

1. A filename of a file which contains the list of phone numbers.
2. A target phone number that the user would like to find the nearest K numbers for. Note that if the phone number contains whitespace, it must be delimited by quotation marks.
3. A number K for which the program will find the nearest K phone numbers to argument 2. If $K \leq 0$, the program panics with the error message, “ k must be greater than 0”.

Example A user may run the program in one of the following manners:

- `cargo run <filename> <target number> <K>`
- `Question_2_rust.exe <filename> <target number> <K>`

2.2 Structure/Design

The phone numbers are first stored in a vector, then are sanitized one by one to ensure only digits remain, then parsed as an integer. This step is $O(n)$, as it must iterate through the entire input list. The numbers are then stored in a hashmap as keys, where the values are the number of times the number has appeared in the input vector. This step is also $O(n)$, as it iterates through the vector which contains the phone numbers.

The absolute difference between each number and the target number is then stored in a BTreeMap, with the values being the number themselves. The BTreeMap was chosen as it has two properties that benefit us. Firstly, due to it being a form of a B-Tree, the keys are sorted by value. Secondly, the B-Tree makes each node contain $B - 1$ to $2B - 1$ elements in a contiguous array for some choice of B , which improves cache efficiency [1]. Construction of a BTreeMap is $O(n \log n)$.

We then iterate through K phone numbers of the BTreeMap which has the worst case of $O(K \log n)$, where there are K unique closest numbers to the target number. If $K > n$, then the worst case is $O(n \log n)$. We then print each of the closest K phone numbers as many times as they appeared in the input vector by checking their values in the hashmap, which takes $O(1)$ time. Overall, the time complexity of the program is $O(n \log n)$.

2.3 Limitations

As the time complexity for the algorithm is $O(n \log n)$, it performs similarly to just running a simple quicksort on the input, then printing K unique phone numbers closest to the target number, which would be far less complex to implement compared to having a BTreeMap in the implementation.

2.4 Testing

The following tests are conducted, and the program passes all.

1. Testing that the algorithm can select the closest k values from a predefined array.
2. Testing that the algorithm will return the same values for a randomly generated array of random length, with random 10-digit phone numbers as the elements.
3. Comparing the efficiency of the algorithm against an implementation of quicksort. The array consists of randomly generated phone numbers, and the comparisons are conducted from an array length of 1 to 10,000. To aid in completing the test faster, the comparisons for each array length are split amongst all logical processors of the user's CPU. For slices of the arrays that are the same length, their comparisons are run on the same thread to ensure there is no variance between the conditions of the test for both algorithms. Each test is then run 5 times and a mean of the time taken for both algorithms to find the k -th closest phone numbers is recorded to a file and is shown below:

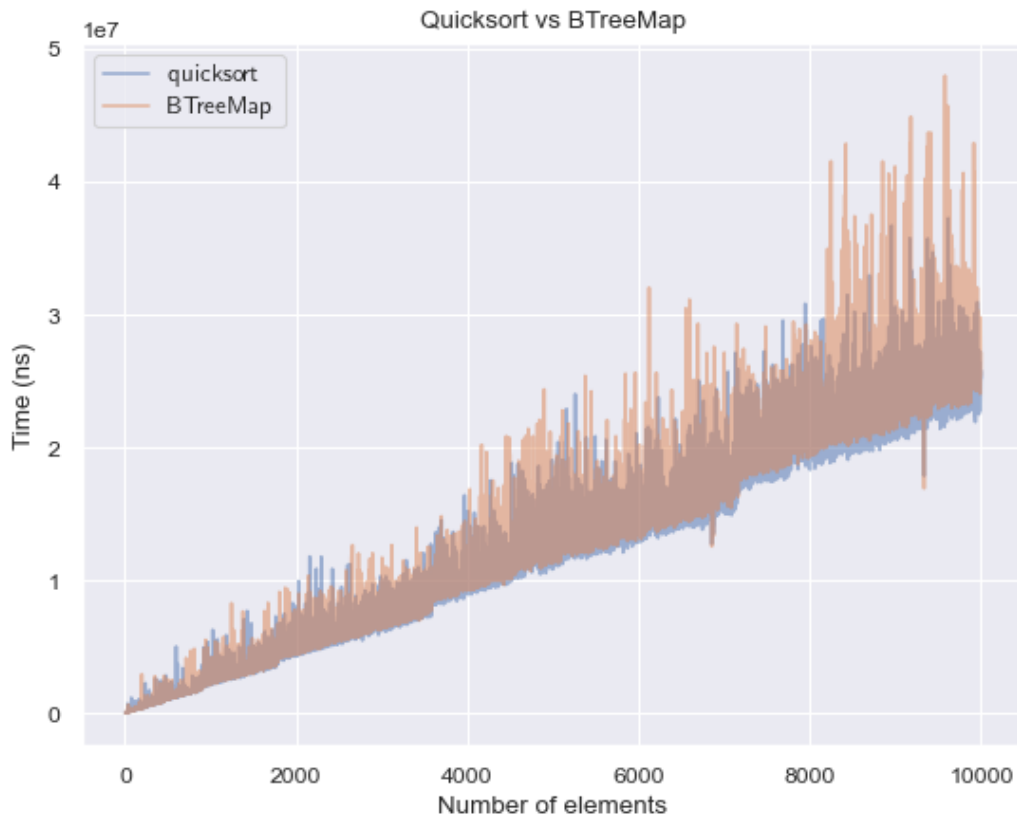


Figure 2.1: Quicksort implementation vs BTreeMap implementation time efficiency

3

Bibliography

- [1] RustLang. “BTreeMap in std::collections - rust,” std - Rust. (May 15, 2015), [Online]. Available: <https://doc.rust-lang.org/stable/std/collections/struct.BTreeMap.html> (visited on 03/18/2023).