

Algorithms: Selection in Linear Time

2017-18570 Sungchan Yi

Dept. of Computer Science and Engineering

April 6th, 2019

Contents

1	Introduction	1
2	Implementation	2
2.1	Randomized Select	2
2.2	Deterministic Select	3
2.3	Checker	5
3	Experiments	5
3.1	Randomized Select	6
3.2	Deterministic Select	6
4	Conclusion	7

1 Introduction

In this homework, we implement the randomized and deterministic selection algorithms with C++ language. The deterministic selection algorithm must run in $O(n)$ time. Given n elements in an array and an integer k ($0 \leq k < n$), the algorithm must find the k -th smallest element in the array.¹ Additionally, we implement a checker program for the correctness for the implemented algorithms. The following are required:

- Implement the randomized and deterministic selection algorithms.
- Run each algorithm for given inputs, measure the running time, print the k -th smallest element.
- Check the correctness of each algorithm by using a checker program that runs in $O(n)$ time. Print the result of checking.

¹For simplicity in programming, we will count from 0.

- After measuring the running time for inputs of various sizes, compute the ratios of the constant hidden in the asymptotic complexities of each algorithm.

The following are the restrictions on given inputs.

- n, k are integers, with $1 \leq n \leq 1,000,000,000$, $0 \leq k < n$.
- All the elements given in the array are distinct.²

2 Implementation

In this section, we will cover the implementation of each algorithm. The basic idea is divide and conquer. Select a pivot and partition the given array into subproblems. The two algorithms discussed here differs only on the method of selecting the pivot.

For the codes, assume the necessary header files are already included.

2.1 Randomized Select

```

1  int randomized_select(int a[], int n, int k) {
2      return randkth(a, 0, n - 1, k + 1);
3  }
```

The RANDOMIZED-SELECT algorithm chooses a pivot randomly from the array. The algorithm is implemented in `randkth` method. `randkth(arr, st, ed, k)` will return the k -th smallest element from `arr[st..ed]`.³

```

1  int randkth(int arr[], int st, int ed, int k) {
2      if(st == ed) // base case for recursion
3          return arr[st];
4      int pivot = randomized_partition(arr, st, ed); // index of pivot
5
6      // index of pivot from strating from st
7      int idx = pivot - st + 1;
8
9      // Divide and Conquer
10     if(k < idx) // search left
11         return randkth(arr, st, pivot - 1, k);
12     else if(k > idx) // search right
13         return randkth(arr, pivot + 1, ed, k - idx);
14     else // element is found
15         return arr[pivot];
16 }
```

²This condition is required for the deterministic selection algorithm to run in $O(n)$ time.

³ $A[i..j]$ denotes the elements from $A[i]$, ..., $A[j]$.

For `randomized_partition`, select a random index from `st..ed` and use it as a pivot.

```
1  int randomized_partition(int arr[], int st, int ed) {
2      srand(time(NULL));
3      int pIdx = st + rand() % (ed - st + 1); // random index
4      swap(arr, ed, pIdx); // put it at the end of array
5      int x = arr[ed]; // pivot element
6      int i = st - 1; // track the correct place of pivot
7      for(int j = st; j < ed; ++j) {
8          if(arr[j] <= x)
9              swap(arr, ++i, j); // increment i then swap
10     }
11     swap(arr, ++i, ed); // put pivot at correct place
12     return i;
13 }
```

`swap(A, i, j)` will swap `A[i]` and `A[j]`. Its implementation is trivial, so we omit it here.

According to the CLRS textbook, the above algorithm will run in expected $O(n)$ time. But the worst case complexity is $O(n^2)$, since we get extremely unlucky by the random function choosing the smallest/largest element from `A[st..ed]`. That choice of pivot will lead to an unbalanced partition of input array.

2.2 Deterministic Select

```
1  int detkth(int a[], int n, int k) {
2      return detkth(a, 0, n - 1, k + 1);
3  }
```

The deterministic selection algorithm chooses a pivot by using “median of medians”. First divide the array into $\lceil n/5 \rceil$ sub-arrays, each with 5 elements, except for the last sub-array. Then choose a median from each sub-array and from the chosen medians, choose a median and use it as the pivot for partition. `detkth(arr, st, ed, k)` will return the k -th smallest element from `arr[st..ed]`.

```
1  int detkth(int arr[], int st, int ed, int k) {
2      if(st == ed) // base case
3          return arr[st];
4      int n = ed - st + 1; // number of elements
5
6      // array for storing medians of each sub-array
7      int* med = (int*) malloc((n + 4) / 5 * sizeof(int));
8      // find medians of each subarray
9      int i;
```

```

10     for(i = 0; i < n / 5; ++i)
11         med[i] = median(arr + st + 5 * i, 5);
12     if(5 * i < n) {
13         med[i] = median(arr + st + 5 * i, n % 5);
14         i++;
15     }
16
17     // choose median of medians
18     int medOfMed = (i == 1) ? med[i-1] : detkth(med, 0, i-1, i/2);
19
20     // partition the array with median of medians
21     int pivot = partition(arr, st, ed, medOfMed);
22     int idx = pivot - st + 1;
23
24     // divide and conquer
25     if(idx == k)
26         return arr[pivot];
27     else if(idx > k)
28         return detkth(arr, st, pivot - 1, k);
29     else
30         return detkth(arr, pivot + 1, ed, k - idx);
31 }

```

`i == 1` in line 18 was to handle the case with the input array size less than 6. To find the median for each sub-array, `median(arr, n)` method was called, and it will find the median of n elements starting from `arr` (pointer) by insertion sort.

```

1  int median(int arr[], int n) {
2      for(int i = 1; i < n; i++) {
3          for(int j = i; j > 0; j--) {
4              if(arr[j] < arr[j - 1]) swap(arr, j, j - 1);
5              else break;
6          }
7      }
8      return arr[n / 2];
9  }

```

partition process is done similarly.

```

1  int partition(int arr[], int st, int ed, int x) {
2      // search for x in arr, and move it to the end
3      int i;
4      for(i = st; i < ed; ++i) {
5          if(arr[i] == x) break;
6      }
7      swap(arr, i, ed);

```

```

8      i = st - 1;
9      for(int j = st; j < ed; ++j) {
10         if(arr[j] <= x) // increment i then swap arr[i] and arr[j]
11             swap(arr, ++i, j);
12     }
13     swap(arr, ++i, ed);
14     return i;
15 }

```

According to the textbook, the above algorithm will run in $O(n)$ time.

2.3 Checker

To check the correctness, we implement a checker program. The basic idea is to count the number of elements less than or equal to the return value of above algorithms.

```

1  bool checker(int a[], int n, int k, int ans){
2      int cnt = 0;
3      for(int i = 0; i < n; ++i) {
4          if(ans >= a[i]) cnt++; // count elements
5      }
6      if(cnt == k) // cnt must be k to be true
7          return true;
8      else return false;
9  }

```

This checker runs in $O(n)$ time because it sweeps the given array and compares each element with `ans`.

3 Experiments

Here is how the test was done.

- Test environment
 - OS: Ubuntu 18.04.2 LTS
 - CPU: Intel Core i5-6200U CPU @ 2.30GHz \times 4
 - RAM: 8GB

- Checker Code

`chrono` header file was used to check the wall clock.

```

1  #include <chrono>
2  using namespace std::chrono;

```

```

3  ... // get input and etc. omitted here
4  auto start = high_resolution_clock::now(); // start time
5  int ans1 = randomized_select(arr, n, k);
6  auto stop = high_resolution_clock::now(); // end time
7
8  // check correctness
9  if(checker(arr, n, k, ans1)) printf("%s", "correct");
10 else printf("%s", "incorrect");
11
12 auto duration = duration_cast<microseconds>(stop - start);
13 printf(", Execution time ");
14 std::cout << duration.count() / 1000.0 << " ms\n";

```

Similar code was also run for `deterministic_select` algorithm.

- The checker code was saved to `checker.cpp`.
 - Compile: `g++ check.cpp -o check`
 - Run: `./check.cpp < 1.txt`

There were no reports of incorrectness during the following experiments.

3.1 Randomized Select

Since the running time of randomized algorithms differ for every execution, we ran the algorithm 10,000 times and take the average execution time.

Array Size (n)	20,000	60,000	200,000	600,000
Average Time (ms)	0.363094	0.926259	1.955946	9.627423

Table 1: Average running time of RANDOMIZED-SELECT algorithm

The estimated complexity is $T(n) = 0.01615n - 335.4$ (μs), with $R^2 = 0.9789$. The R^2 value is less than the deterministic select due to the randomness of the algorithm.

3.2 Deterministic Select

The running time for deterministic algorithm does not differ very much. We ran it once here.

Array Size (n)	20,000	60,000	200,000	600,000
Time (ms)	1.922	6.265	17.785	54.067

Table 2: Running time of deterministic select algorithm

The estimated complexity is $T(n) = 0.09034n + 269.3$ (μs), with $R^2 = 0.9996$. Just for curiosity, we ran the program 1,000 times to see the average.

Array Size (n)	20,000	60,000	200,000	600,000
Average Time (ms)	0.902151	2.726879	10.423119	31.067172

Table 3: Average running time of deterministic select algorithm

It can be easily seen that the deterministic algorithm show linearity.

4 Conclusion

We calculated the ratio of running times by dividing the results of Table 2 by Table 1, for each n .

Array Size (n)	20,000	60,000	200,000	600,000
Ratio	5.293	6.764	9.093	5.616

Table 4: Ratio of running times for each n

With only the leading coefficients of the complexities for each algorithm, we can see that the ratio is $0.09034/0.01615 = 5.593808$. We conclude that the deterministic selection algorithm is around 5.5 times slower than the randomized selection algorithm.