# Random Algorithms (RA-MIRI)
# Assignment 2

Daniel Benedí

November 2022

## 1 QuickSelect

In this section, the algorithm QuickSelect designed by Hoare will be implemented and the number of comparisons made between array elements,$C_n^{(\alpha n)}$, will be counted , so we can measure the average number of comparisons $\bar{C}_n^{(\alpha n)}$ and its standard deviation $\sigma_n^{(\alpha n)}$.

### 1.1 Implementation

From the `find` algorithm given by Hoare [2], we can define the following QuickSelect implementation:

The implementation provided uses the standard library of C++ (further detail in appendix A). To choose the pivot, we use a pseudo-random numbers engine provided by the standard library called `std::mt19937` (for implementation details check Makoto and Takuji paper [3]). We have chosen this engine instead of a linear congruential because with the right parameters it has the longest non-repeating sequence and we want to generate a lot of numbers, although it is slower and has greater state storage requirements. It makes use of `std::random_engine` which is a uniformly-distributed integer random number generator that produces non-deterministic random numbers or pseudo-random numbers if a non-deterministic source is not available. After generating the random (or pseudo-random) bits, they are passed to `std::uniform_int_distribution` so it can generate the index of the pivot that is going to be used. Because initializing the engine and the random device is highly cost task, we have used the `static` keyword, so it is only initialized once over all the experiments.

In order to get robust results from the experiments, the experiment has been realized 100 times shuffling the input vector every time so we could get a different permutation every time. Also, for each input size we have tried $\alpha$ values from 0 to 1 with a step of 0.01, so we could also get a fine-grained resolution of the costs. We have tried different input sizes in the range of $2 * 10^4$ to $2 * 10^6$, but not all the values in the range, just some of them. In order to use the whole

**Algorithm 1** QuickSelect

---

1: **procedure** QUICKSELECT$(A, j)$      ▷ QuickSelect method to find the $j$-th element $(\alpha = \frac{j}{n})$
2:      min $\leftarrow 1$                               ▷ First index of array
3:      max $\leftarrow |A|$                            ▷ Last index of array
4:      $C_n \leftarrow 0$                       ▷ Number of comparisons made
5:      $k$
6:      **do**
7:          $a \leftarrow$ pick u.a.r a element from $A$
8:          $l \leftarrow$ min, $r \leftarrow$ max                  ▷ Partially ordered elements
9:          **while** $l < r$ **do**      ▷ Order all the elements according to the pivot $a$
10:              **while** $l \leq$ max $\wedge A[l] < a \wedge l < r$ **do**
11:                  $l \leftarrow l + 1, C_n \leftarrow C_n + 1$
12:              **end while**
13:              $C_n \leftarrow C_n + 1$
14:              **while** $r \geq$ min $\wedge A[r] > a \wedge l < r$ **do**
15:                  $r \leftarrow r - 1, C_n \leftarrow C_n + 1$
16:              **end while**
17:              $C_n \leftarrow C_n + 1$
18:              SWAP$(A[l], A[r])$
19:          **end while**
20:          $k = l$
21:          **if** $i < k$ **then**      ▷ If pivot is to the right of the looked for element
22:              max $\leftarrow k - 1$ $i > k$      ▷ If pivot is to the left of the looked for element
23:              min $\leftarrow k + 1$
24:          **end if**
25:      **while** $i \neq r$
26: **end procedure**

---

CPU and decrease the time needed, we have made use of the library OpenMP to parallelize the experiments for each $n$.

We have also implemented the option to take $2 * t + 1$ elements and choose the pivot that falls closer to the $j$-th position in that sample. It is known that the expected number of comparisons is lower than the number of comparison of the previous method, but because of a lack of time we haven't been able to do the proper experimentation to add it in this report.

## 1.2 Results

In order to visualize our results, we have to take into account that we have two input variables and one output, the average cost $\bar{C}_n^{(\alpha n)}$. Theoretically, we know that $\frac{\bar{C}_n^{(\alpha n)}}{n} \approx 2 - 2(\alpha \ln \alpha + (1-\alpha) \ln (1 - \alpha))$. The figure 1 shows this proportion depending on the input size, $n$ and $\alpha$
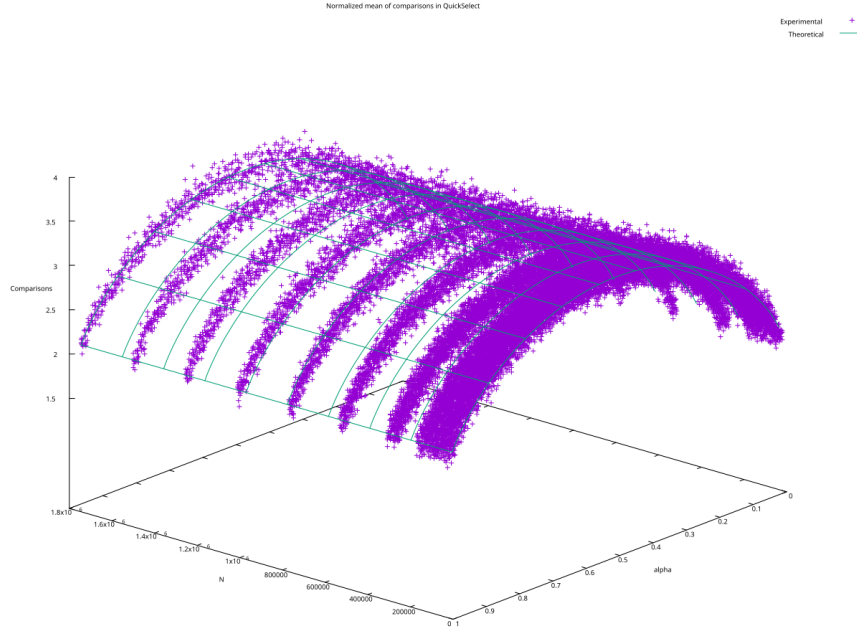


Figure 1: Average comparisons divided by the case, depending on the input size and $\alpha$

We can observe that apparently the cost experimentally fits the cost theoretically. When we collapse the input size, $n$, dimmension, we obtain the figure 2. We can observe in it that it seems to follow the theoretical approximation quite nice, but moving around it.
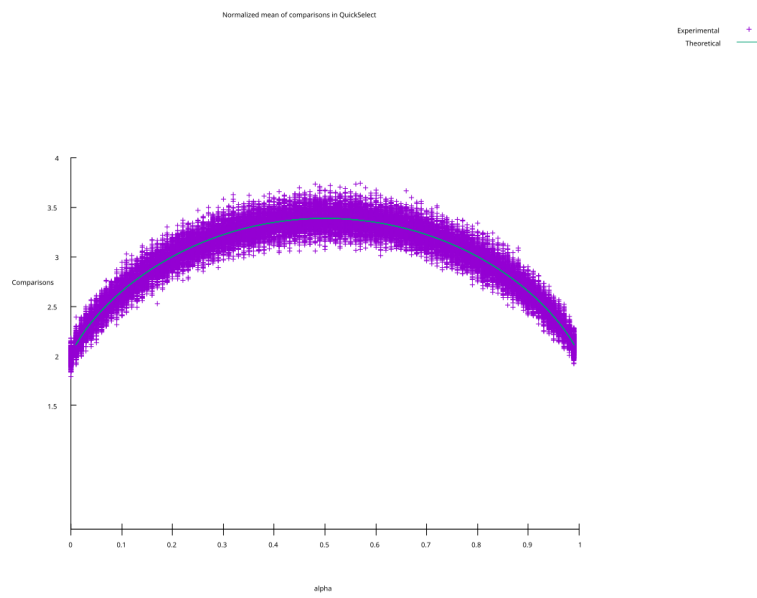
Normalized mean of comparisons in QuickSelect

Experimental  +
Theoretical  —

Comparisons

alpha

Figure 2: Average comparisons divided by the case, only depending on $\alpha$

4

(a) Standard error of mean comparisons divided by input size over both dimensions

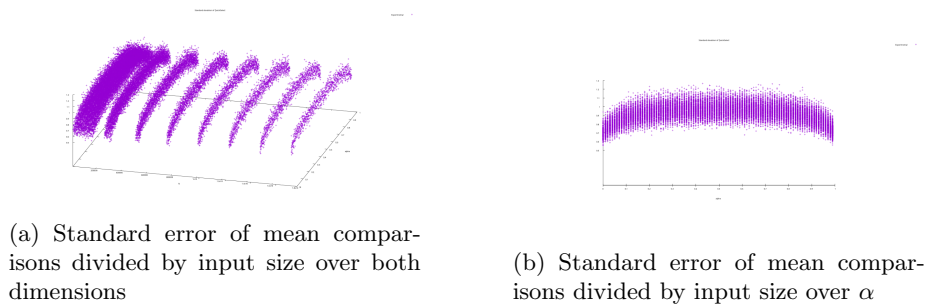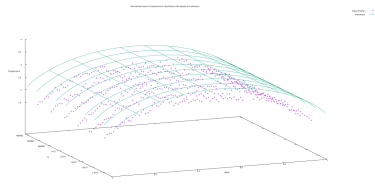(b) Standard error of mean comparisons divided by input size over $\alpha$

Figure 3: Standard error of QuickSelect

Previous figures showed us that the number of comparisons are quite close to the theoretical approximation. In figure 3 we analyzed the standard deviation as $\frac{\sigma_n^{(\alpha n)}}{n}$. We observe that it has a high standard deviation. This means that although it is not exactly the theoretical approximation, we know that the obtained values could be the theoretical value because the high standard deviation (We should do a statistical test, t-test, to determine if the value could be the theoretical or not).
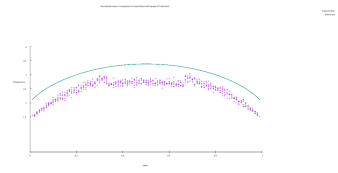
## 1.3 Pivot with sample of $(2t + 1)$ elements

Instead of choosing randomly the pivot, we also provide an implementation which sample $2t+1$ elements u.a.r and choose the element $\alpha$ in the sample. The key with this method is that $t$ is given in compilation time so it can be implemented in constant time and it improves the number of comparisons needed.

In figure 4, we show the performance for $t = 1$ for this algorithm. We have done not too many experiments because of a lack of time, but we can still observe some significant results. Further experimentation should be realized to analyze how the parameter $t$ affects the number of comparisons. For instance, the subfigure 4a and subfigure 4b show that it have significantly less comparisons than QuickSelect (figure 1 and figure 2). Also it is noticeable that the standard error is also lower with the sample than the other one.

(a) Mean comparisons divided by input size over both dimensions. Compared with the theoretical cost of the simple QuickSelect.



(b) Mean comparisons divided by input size over $\alpha$. Compared with the theoretical cost of the simple QuickSelect.



(c) Standard error of mean comparisons divided by input size over both dimensions



(d) Standard error of mean comparisons divided by input size over $\alpha$

Figure 4: Standard error of QuickSelect

# 2    Randomized Selection

In this section, the algorithm Randomized Selection designned by Floy and Rivest [1] with the version studied in the lectures. This version will be used to obtain the median, but can be easily modified to get any element.

## 2.1 Implementation

The algorithm showed the lectures is:

---
**Algorithm 2** Randomized Selection

---
1: **procedure** RANSELECT($A$)     ▷ Randomized Selection method to find the median
2:     $R \leftarrow$ Sample $\lceil n^{\frac{3}{4}} \rceil$ u.a.r with replacement from $A$
3:     SORT(R)
4:     $d \leftarrow R[\lfloor \frac{1}{2} n^{\frac{3}{4}} - \sqrt{n} \rfloor]$
5:     $u \leftarrow R[\lfloor \frac{1}{2} n^{\frac{3}{4}} + \sqrt{n} \rfloor]$
6:     $C \leftarrow x \in A | d <= x <= u$
7:     $l_d \leftarrow |x \in A | x < d|$
8:     $l_u \leftarrow |x \in A | u < x|$
9:     **if** $l_d > \frac{n}{2} \wedge l_u > \frac{n}{2}$ **then**
10:         **FAIL**
11:     **end if**
12:     **if** $|C| > 4n^{\frac{3}{4}}$ **then**
13:         **FAIL**
14:     **end if**
15:     SORT($C$)
16:         **return** $C[\lceil \frac{n}{2} \rceil - l_d + 1]$
17: **end procedure**

---

If instead of choosing the median, we want a $\alpha$ element then we need to change in line 4 and 6 the $\frac{1}{2}$ with $\alpha$, the line 9 by $l_d > \alpha n \wedge l_u > n(1 - \alpha)$ and inn the line 16 by $C[\lceil \alpha n \rceil - l_d + 1]$.

The implementation provided uses the standard library of C++ (further detail in appendix A). To sample uniformly at random with replacement, we use a pseudo-random numbers engine provided by the standard library called `std::mt19937` (for implementation details check Makoto and Takuji paper [3]). We have chosen this engine instead of a linear congruential because with the right parameters it has the longest non-repeating sequence and we want to generate a lot of numbers, although it is slower and has greater state storage requirements. It makes use of `std::random_engine` which is a uniformly-distributed integer random number generator that produces non-deterministic random numbers or pseudo-random numbers if a non-deterministic source is not available. After generating the random (or pseudo-random) bits, they are passed to `std::uniform_int_distribution` so it can generate the index of the sampled element that is going to be used. Because initializing the engine and the random device is a highly cost task, we have used the `static` keyword, so it is only initialized once over all the experiments.

## 2.2 Results and comparison

Similarly, to the previous algorithm we have realized 100 times for each input size and shuffling the input vector each time so we could get a different permutation. We have tried different input size in the range of $2 * 10^4$ to $2 * 10^6$ with a step of 1000. In this case, we have been able to try all of them because we weren't trying different $\alpha$, only the median. To be able to compare this algorithm with QuickSelect, we had to slice the results in the dimension with $\alpha = 0.5$. To have a comparable metric with QuickSelect, we have also divided the average number of comparisons by the input size.



Figure 5: Comparison of mean comparisons in QuickSelect and Randomized Selection

In the figure 5, we observe the proportion number of comparisons decrease with the input size. Compared with QuickSelect this means Randomized Selection performs much better than QuickSelect. We have fitted the theoretical curves and we see that for QuickSelect it fits nicely, for Randomized Selection we observe that the theoretical curve is an upper-bound but not tight which fits with the small-o notation used in the lectures.

The figure 6 shows that not only Randomized Selection has a lower cost in terms of comparisons, but also has a lower standard error which means that the total number of comparisons is much more stable. A drawback of Randomized Selection is that it can fail which it doesn't happen in QuickSelect, but the probability of fail is really low, $P[\textbf{FAIL}] \leq \frac{1}{n^{1/4}}$, and in all the experiments realized we never got a failure from Randomized Selection.

Figure 6: Comparison of standard deviation in comparisons in QuickSelect and Randomized Selection

# Appendices

## A  C++ Implementation

```
1   //*********************************************************
2   //  File:    assignment2.cpp
3   //  Author:  Daniel Bened  Garc a
4   //  Date:
5   //  Coms:
6   //*********************************************************
7
8   #include <iostream>
9   #include <random>
10  #include <fstream>
11  #include <numeric>
12  #include <vector>
13  #include <algorithm>
14  #include <sstream>
15  #include <optional>
16  #include <omp.h>
17
18  /*
```

```
19    *  Recursive  implementation
20    *        int  r = RanPartition(A, p, q);
21
22    *        if(r == j){
23    *                return A[r];
24    *        } else if( i < r){
25    *                return  QuickSelect(A, i, min, r-1);
26    *        } else {
27    *                return  QuickSelect(A, i-r, r+1, max);
28    *        }
29    *
30    *
31    * Note:  It  returns  only  the  number  of  comparisons,  but  it  can  also  be  return
32    *        the  i-th  element
33    */
34  int QuickSelect(int* A, int i, int min, int max){
35          // Random device and Mersenne twistter are static so they are stored
36          // generated only once
37          static std::random_device rd;
38          static std::mt19937_64 g(rd());
39
40          std::shuffle(A+min, A+max, g);
41
42          max--;
43          int r;
44          int Cj = 0;
45          do{
46                  // Create distribution for range min, max in that iteration
47                  std::uniform_int_distribution<int> distrib(min,max);
48
49                  // Choose an element a in A uniformly at random
50                  int a = A[distrib(g)];
51
52                  // Order according the pivot element a
53                  int left = min, right = max;
54                  while(left < right){
55                          while(left <= max && A[left] < a && left < right){
56                                  left++;
57                                  Cj++;
58                          }
59                          Cj++;
60                          while(right >= min && A[right] > a && left < right){
61                                  right--;
62                                  Cj++;
63                          }
64                          Cj++;
```

```
65                          int temp = A[left];
66                          A[left] = A[right];
67                          A[right] = temp;
68                          left++; right--;
69                      }
70
71                  r = left;
72                  if(i < r){
73                          max = r - 1;
74                  }else if(i > r){
75                          min = r + 1;
76                  }
77          }while(i != r);
78
79          return Cj;
80  }
81
82  void task1(){
83          const int MIN = 20000, MAX = 2000000;
84          std::ofstream comparisons("quickselect_comparisons.csv");
85
86          #pragma omp parallel for
87          for(int N = MIN; N < MAX; N+=1000){
88                  int* vec = new int[N];
89                  std::iota(vec, vec+N,1);
90
91                  std::stringstream buff;
92                  for(double p = 0.00; p <= 1; p += 0.01){
93                          int j = int(p*(N-1));
94                          buff << N << "," << j;
95                          for(int attempt = 0; attempt < 100; attempt++){ //Do
96                                  int Cj = QuickSelect(vec, j, 0, N);
97                                  buff << "," << Cj;
98                          }
99                          buff << std::endl;
100                 }
101                 comparisons << buff.str();
102                 delete[] vec;
103         }
104 }
105
106 /*
107  * Same implementation of QuickSelect but now the pivot will depend on the
108  * sample u.a.r of (2t+1) elements
109  */
110 template<int t>
```

```
111   int QuickSelect_t(int* A, int i, int min, int max){
112           // Random device and Mersenne twistter are static so they are stored
113           // generated only once
114           static std::random_device rd;
115           static std::mt19937_64 g(rd());
116
117           std::shuffle(A+min, A+max,g);
118
119         max--;
120         int r;
121         int Cj = 0;
122         do{
123                   // Create distribution for range min, max in that iteration
124                   std::uniform_int_distribution<int> distrib(min,max);
125
126                   // Choose 2t+1 elements in A uniformly at random
127                   int a;
128                   if( (max+1-min) > (2*t+1) ){
129                           int S[2*t+1];
130                           for(int k = 0; k < 2*t+1; k++)
131                                   S[k] = A[distrib(g)];
132                           std::sort(S,S+2*t+1);
133                           // Choose pivot depending on the sample
134                           a = S[(i-min)/((max+1-min)/(2*t+1))];
135                   }else{
136                           a = A[distrib(g)];
137                   }
138
139                   // Order according the pivot element a
140                   int left = min, right = max;
141                   while(left < right){
142                           while(left <= max && A[left] < a && left < right){
143                                   left++;
144                                   Cj++;
145                           }
146                           Cj++;
147                           while(right >= min && A[right] > a && left < right){
148                                   right--;
149                                   Cj++;
150                           }
151                           Cj++;
152                           int temp = A[left];
153                           A[left] = A[right];
154                           A[right] = temp;
155                           left++; right--;
156                   }
```

```
157
158                       r = left;
159                       if(i < r){
160                               max = r − 1;
161                       }else if(i > r){
162                               min = r + 1;
163                       }
164               }while(i != r);
165
166               return Cj;
167    }
168
169    void task1_extra(){
170               const int MIN = 20000, MAX = 2000000;
171               std::ofstream comparisons("quickselect_sample_comparisons.csv");
172
173               #pragma omp parallel for
174               for(int N = MIN; N < MAX; N+=1000){
175                       int* vec = new int[N];
176                       std::iota(vec, vec+N,1);
177
178                       std::stringstream buff;
179                       for(double p = 0.00; p <= 1; p += 0.01){
180                               int j = int(p*(N−1));
181                               buff << N << ",_" << j;
182                               for(int attempt = 0; attempt < 100; attempt++){ //Do
183                                       int Cj = QuickSelect_t<1>(vec, j, 0, N);
184                                       buff << ",_" << Cj;
185                               }
186                               buff << std::endl;
187                       }
188                       comparisons << buff.str();
189                       delete[] vec;
190               }
191    }
192
193
194
195    /*
196     *        @return The median element of A[min..max]
197     *
198     *        R := Sample ceil(n**0.75) from A uniform and w/replacement
199     *        Sort(R)
200     *        d := floor(0.5*n**0.75−n**0.5) element in R
201     *        u := floor(0.5*n**0.75+n**0.5) element in R
202     *
```

13

```
203    *       C    := elements between d and u
204    *       l_d := num elements lower than d
205    *       l_u := num elements bigger than u
206    *
207    *       if l_d > n/2 or l_u > n/2
208    *               FAIL
209    *
210    *       if |C| > 4n**0.75
211    *               FAIL
212    *
213    *       Sort(C)
214    *       return floor(0.5*n) - l_d + 1 element in C
215    */
216   std::optional<std::pair<int,int>> RanSelect(std::vector<int> A){
217           static std::random_device rd;
218           static std::mt19937_64 g(rd());
219
220           std::shuffle(A.begin(), A.end(),g);
221
222           int C_n = 0;
223           auto comparator = [&C_n](int a, int b){C_n++; return a < b;};
224
225           int n = A.size();
226           double n_3_4 = std::pow(n, 0.75);
227           double n_1_2 = std::sqrt(n);
228
229           std::vector<int> R(std::ceil(n_3_4));
230
231           std::sample(A.begin(), A.end(),
232                                       R.begin(),
233                                       int(std::ceil(n_3_4)), g);
234
235           std::sort(R.begin(), R.end(), comparator);
236
237           int d = R[std::floor(0.5*n_3_4-n_1_2)],
238                   u = R[std::floor(0.5*n_3_4+n_1_2)];
239
240           std::vector<int> C;
241           int l_d = 0, l_u = 0;
242           for(int elem : A){
243               if(elem < d) l_d++;
244               else if(elem > u) l_u++;
245               else C.push_back(elem);
246               C_n++;
247           }
248
```

```
249              if(l_d > A.size()/2  ||  l_u > A.size()/2){
250                      std::cerr << "FAIL:_l_d_or_l_u_too_big" << std::endl;
251                      return std::nullopt;
252              }
253
254              if(C.size() > 4*n_3_4){
255                      std::cerr << "FAIL:_C_is_too_big" << std::endl;
256                  return std::nullopt;
257              }
258
259              std::sort(C.begin(), C.end(), comparator);
260
261              return std::optional<std::pair<int,int>>(std::make_pair(C[n/2 - l_d +
262  }
263
264  void task2(){
265              const int MIN = 20000, MAX = 2000000;
266              std::ofstream comparisons("randselect_comparisons.csv");
267
268              #pragma omp parallel for
269              for(int N = MIN; N < MAX; N+=1000){
270                      std::vector<int> vec(N);
271                      std::iota(vec.begin(), vec.end(),1);
272
273
274                      std::stringstream buff;
275                      int tries = 0;
276                      for(int attempt = 0; attempt < 100; attempt++){ //Do 100 rand
277                              std::optional<std::pair<int,int>> res = std::nullopt;
278                              while(!res){
279                                      tries++;
280                                      res = RanSelect(vec);
281                              }
282                              buff << ",_" << res.value().second;
283                      }
284
285                      std::stringstream temp;
286                      temp << N << ",_" << tries - 100 << buff.rdbuf() << std::endl
287                      buff = std::move(temp);
288                      comparisons << buff.str();
289              }
290  }
291  int main(int argc, char* argv[]){
292
293              //task1();
294              task1_extra();
```

```
295              //task2();
296    }
```

# References

[1] Robert W. Floyd and Ronald L. Rivest. "Algorithm 489: The Algorithm SELECT—for Finding the Ith Smallest of n Elements [M1]". In: *Commun. ACM* 18.3 (Mar. 1975), p. 173. ISSN: 0001-0782. DOI: 10.1145/360680.360694. URL: https://doi.org/10.1145/360680.360694.

[2] C. A. R. Hoare. "Algorithm 65: Find". In: *Commun. ACM* 4.7 (July 1961), pp. 321–322. ISSN: 0001-0782. DOI: 10.1145/366622.366647. URL: https://doi.org/10.1145/366622.366647.

[3] Makoto Matsumoto and Takuji Nishimura. "Mersenne Twister: A 623-Dimensionally Equidistributed Uniform Pseudo-Random Number Generator". In: *ACM Trans. Model. Comput. Simul.* 8.1 (Jan. 1998), pp. 3–30. ISSN: 1049-3301. DOI: 10.1145/272991.272995. URL: https://doi.org/10.1145/272991.272995.