

CAB302 - Assignment 2

Luke Josh, Jason Queen

May 24, 2016

Contents

1 Summary

The purpose of this report is to analyse and compare the average complexity of the selection median algorithm against a brute force solution. This report uses algorithm analysis techniques to experimentally determine the average case efficiency of both the selection median and brute force median algorithm. The expected theoretical efficiencies of the algorithms are evaluated using mathematical analysis and are contrasted against the computed experimental results. The results for each algorithm are also compared to each other and the most efficient algorithm is determined.

2 Description of the Algorithms

2.1 Brute Force

The brute force median algorithm works by manually checking each value of the array and determining if the element's position in a sorted array is in the median position of $\frac{n}{2}$. It does this by checking each element of the array against every other element of the array, and counting the number of elements that are less than the value, and the number of elements that are greater than the value.

Recalling that the median value of a list of numbers is the element that occurs in center position (rounded up for the purpose of this algorithm), the algorithm tests if a value is the median value by testing if half the elements in the array are less than the value. This can be seen in the pseudocode included in the next section.

2.1.1 The Algorithm

Algorithm 1 Brute Force Median

```
1: function BRUTEFORCEMEDIAN( $A[0..n-1]$ )
2:    $k \leftarrow \lfloor n/2 \rfloor$ 
3:   for  $i \leftarrow 1$  to  $n-1$  do
4:      $numsmaller \leftarrow 0$ 
5:      $numequal \leftarrow 0$ 
6:     for  $j \leftarrow 0$  to  $n-1$  do
7:       if  $A[j] < A[i]$  then
8:          $numsmaller \leftarrow numsmaller + 1$ 
9:       else
10:        if  $A[j] = A[i]$  then
11:           $numequal \leftarrow numequal + 1$ 
12:        end if
13:      end if
14:    end for
15:    if  $numsmaller < k$  and  $k \leq (numsmaller + numequal)$  then
16:      return  $A[i]$ 
17:    end if
18:  end for
19: end function
```

2.2 Median Selection Algorithm

The selection algorithm for finding the median value of the array derives much of its logic from the QuickSelect algorithm. QuickSelect implements a random selection of a pivot element, which is used to partially sort the array around that pivot - which it then recursively continues on the subarray containing the median. QuickSelect has been quoted to have an average running complexity of $\Theta(n)$ (?) - which is a value we can expect to get for the median selection algorithm.

The selection algorithm uses a recursive method to determine the median of the array. The recursive function first performs a partition on the array where all elements less than a pivot element (starting with element 0 as the first pivot) are swapped until the pivot element swapped into its correct position. All elements with an index lower than that of the pivot, now also have a value which is less than that of the pivot. If the new position of the pivot is the middle of the array then the value of the pivot is returned as it is the median of the array. This check is the base case of the recursive function. Otherwise, if the position of the pivot is less than or greater than that of the middle index then the recursive function is called again and then performs the partition on either the lower or upper partition of the array respectively. This is performed until the base case is met and the element in the middle of the array is returned.

It should be noted that there is no guarantee that the whole array will be completely sorted by this process. However, all elements at a lower index than that of the pivot will have a value that is lower than that of the pivot element and all elements with an index greater than that of the pivot element will have values that are also greater than the value of the median. Furthermore, in contract to the brute force algorithm, when this algorithm is performed on an even array size, returns the value which is towards the lower section of the sorted array. An example of this process is as follows, on an array $A = [3, 5, 4, 2, 1]$:

$A = [3, 5, 4, 2, 1] \rightarrow 5$ and 4 are both greater than 3 , do nothing

$A = [3, 2, 4, 5, 1] \rightarrow 2 < 3$, swap 5 with 2

$A = [3, 2, 1, 5, 4] \rightarrow 1 < 3$, swap 4 with 1

$A = [1, 2, 3, 5, 4] \rightarrow$ Finally, swap 3 with 1

If the first value happens to be the median value of the array, the algorithm stops. However, if the index returned is not the median index, there are two possibilities:

1. If the index is less than the median index, the same process as above is performed, ignoring the first element of the array.
2. If the index is greater than the median index, the same process as above is performed, ignoring the last element of the array.

This process continues until the value with an index equal to the median index is found.

2.2.1 The Algorithm

Algorithm 2 Selecion Median

```

1: function MEDIAN( $A[0..n - 1]$ )
2:   if  $n = 1$  then
3:     return  $A[0]$ 
4:   else
5:      $\text{Select}(A, 0, \lfloor n/2 \rfloor, n - 1)$ 
6:   end if
7: end function
8:
9: function SELECT( $A[0..n - 1], l, m, h$ )
10:   $pos \leftarrow \text{Partition}(A, l, h)$ 
11:  if  $pos = m$  then
12:    return  $A[pos]$ 
13:  end if
14:  if  $pos > m$  then
15:    return SELECT( $A, l, m, pos - 1$ )
16:  end if
17:  if  $pos < m$  then
18:    return SELECT( $A, pos + 1, m, h$ )
19:  end if
20: end function
21:
22: function PARTITION( $A[0..n - 1], l, h$ )
23:   $pivotval \leftarrow A[l]$ 
24:   $pivotloc \leftarrow l$ 
25:  for  $j \leftarrow l + 1$  to  $h$  do
26:    if  $A[j] < pivotval$  then
27:       $pivotloc \leftarrow pivotloc + 1$ 
28:      swap( $A[pivotloc], A[j]$ )
29:    end if
30:  swap( $A[l], A[pivotloc]$ )
31: end for
32:  return  $pivotloc$ 
33: end function

```

3 Theoretical Analysis of the Algorithms

3.1 Choice of Problem Size

The algorithms are tested using randomly generated arrays of sizes 1 through to 999 in steps of 3. Additionally, the test for each array size is repeated 2000 times in order to normalize the results.

The upper limit was chosen as array sizes greater than 1000 have brute force time complexities which are significantly large. This limits the number of averages that can be taken for each array size as the total time to run all tests increases to a value that is impractical to perform repeatedly.

The odd step size of 3 is chosen so that tests are performed for both odd and even array sizes.

3.2 Brute Force Median

3.2.1 Choice of Basic Operations

The operation that best defines the complexity and running time of the brute force median algorithm is the comparison $A[j] < A[i]$. This comparison operation is performed more than any other operation in the algorithm - a minimum of $n - 1$ times, and a maximum of $(n - 1)^2$ times.

3.2.2 Average Case Efficiency

The average case efficiency of the algorithm can be derived by considering then number of operations required to determine the median of an arbitrarily sized array. Consider an array $A = [a_1, a_2, \dots, a_n]$ with it's median value placed at position k . The algorithm must check each element in the array up to and including $A[k]$ against each element in the array, including itself. Thus, the algorithm must perform $k \cdot n$ comparisons to determine that $A[k]$ is the median value of the array.

To determine the average number of operations for an array of size n , we must consider the pairity of the size of the array, as the algorithm chooses the value to the left of the midpoint when there are an even number of elements. We will first assume that each value in the array has an equivalent chance of being the median. In this case, the average case number of operations to determine the median is average number of operations for each $M = A[k] \forall k \in [0, 1, \dots, n]$:

$$c_{average} = \frac{\sum_{j=1}^n j \cdot n}{n} \quad (1)$$

$$c_{average} = \sum_{j=1}^n j \quad (2)$$

$$c_{average} = \frac{n^2 + n}{2} \quad (3)$$

To account for the issue of pairty, consider two sorted arrays, $A = [a_1, a_2, \dots, a_n]$ and $B = [a_1, a_2, \dots, a_{n+1}]$, where n is odd. In each of these arrays, the median value will be the same, $M = a_{\lfloor \frac{n}{2} \rfloor}$, and will have to check the same number of elements to find it, however, in B , each element is checked against $n + 1$ other elements, as opposed to A 's n other elements. We can use this to restate the above value for the average number of comparisons as $c_{oddaverage} = \frac{n^2 + n}{2}$, and conversely, $c_{evenaverage} = \frac{n^2}{2}$. These two statements can be combined using a modulo operator to give the true average number of comparisons for an arbitrarily sized array as:

$$c_{average} = \frac{n^2 + (n \bmod 2) \cdot n}{2} \quad (4)$$

Which gives the algorithm an average case complexity of $\Theta(n^2)$.

The above caluculations are assuming the median value for the array only occurs once. In other cases, an array may take the form $[1, 2, 3, 3, 3, 4, 5]$, which will skew the results for the average number of comparisons. For example, an array $A = [1, 2, 4, 5, 3, 3, 3]$, the median value is 3, but 3 appears more than once. This means the algorithm will find the median at the 5th position, but it must check each value against 7 other values.

3.3 Selection Median

3.3.1 Choice of Basic Operations

The operation that best defines the complexity and running time of the selection median algorithm is the comparison $A[j] < pivotval$ which is performed by the partition sort logic borrowed from the Quicksort algorithm. This comparison operation is performed more than any other operation in the algorithm - a minimum of $n - 1$ times if the median value is in the first position of the array, and a maximum of $(n - 1)^2$ times, if the median value is in its correct position at the centre of the array.

3.3.2 Average Case Efficiency

The average case number of comparisons performed in the selection algorithm is considerably more complex than that of the brute force implementation. It has been stated that the selection algorithm is extremely similar to the QuickSort algorithm - and as such, a similar computational complexity is expected. Berman and Paul state that the number of comparisons to be expected for an array of size n is less than $4 \cdot n$ (? , ?).

4 Methodology, Tools and Techniques

4.1 Programming Environment

All testing and computation was executed on a Microsoft Surface running Windows 8.1 64 bit, Intel Core i5 Processor (1.9GHz), 4GB DDR4 RAM.

4.2 Implementation of the Algorithms

The algorithms and their tests were both implemented in C++, compiled using the GCC GNU Compiler in the Code Blocks IDE. The source code for these algorithms is included in Appendix 1.

In general, operation counts were performed by running the algorithm a large number of times on randomly generated arrays, appending the number of operations for each iteration to a pointer to a integer running total, and then calculating the average. The computation time for algorithms was generated in much the same way - using clock types to find the time between two points in the code (namely, the start and end of the algorithm), and then repeating it and taking the average.

4.3 Generating Test Data and Running the Experiments

I think I've already covered this, read again when you're not tired

5 Experimental Results

All graphs referenced in this section have been included in section 2 of the appendix.

5.1 Functional Testing

To ensure the functional correctness of the implemented algorithms - a number of unit tests have been implemented. These tests are included in the test.cpp file, which is also included in appendix 1 (section 6.1.1). These unit tests ensure that the algorithms produce the expected results for a number of predefined test cases.

5.2 Brute Force

5.2.1 Number of operations

The number of operations of the brute force median algorithm was calculated by running the algorithm on a random array of various sizes, and counting the number of times the comparison $A[j] < A[i]$ is performed.

This expected number of comparisons was calculated above to be $\frac{n^2}{2}$ for odd sized arrays, and $\frac{n^2+n}{2}$ for even sized arrays. This function has been graphed along the data that was collected from the tests - and clearly shows that the expected trend holds true.

It must be noted that the data sets used to test the number of operations is not tested for uniqueness - a single value may appear more than once. This causes some amount of discrepancy between the expected result and the actual result, but can mostly be corrected by increasing the maximum value a value can take. This can be seen in the source code in section 6.1.4, on line 97.

5.2.2 Computation time

The computation for the brute force algorithm was expected to run in Θn^2 time. In other words, the time for the algorithm to run is determined quadratically by the number of elements in the array. It can be seen in the included figure that the collected data follows a quadratic trend. Included on this graph also is the operation count again, which allows us to see that the average case running time is indeed correlated to the number of basic operations performed.

5.3 Selection

5.3.1 Number of operations

The selection algorithm was quoted to perform less than $4 \cdot n$ comparisons to determine the median value of array of size n on average. For comparison, the line $comparisons = 4 \cdot array_size$ has been graphed. It is clear that the recorded data fits the expected trend - as it almost always lies below the line. As the number of samples or averages is increased, the variance would become less, and the recorded data would trend closer and closer to a straight line which lies underneath the plotted function.

5.3.2 Computation time

An operating complexity of Θn was expected for the median selection. While the included graph of the data makes the linear trend clear - it is also apparent that the data is neither pretty nor smooth. This is due to the fact that the algorithm runs in a very short amount of time, and the number of averages was not high enough to accurately depict a trend for each array size.

5.4 Comparison

As the two algorithms perform the same operation, they are able to be quantitatively compared. In section 3, we have theoretically shown that the average computational complexity of the brute force and selection algorithms are Θn^2 and Θn respectively. In section 5, we have shown that these theoretical models fit the data that was collected quite well - which allows us to confidently compare the two algorithms.

It is clear that the selection algorithm is superior to the brute force method for an arbitrarily sized array - shown computationally and [lol] analytically. Included in the figures are graphs comparing the computation time and operation count for both algorithms, these allow us to see the drastic difference between these two (result wise) equivalent algorithms. It can be seen that the difference in linear and quadratic running time causes the brute force method to reduce the selection method to less than a percent of the graphs y axis.

6 Appendix

6.1 Source Code

6.1.1 Brute Force Median

```
1 #include "bruteForceMedian.h"
2
3 int bruteForceMedian(int *a, int num_elements, float *time_taken, int *num_operations){
4     int num_smaller, num_equal, operations_counter = 0, k;
5
6     clock_t start = clock();
7
8     k = ceil((float)num_elements/2.0);
9
10    for(int i = 0; i < num_elements; i++){
11        num_smaller = 0;
12        num_equal = 0;
13        for(int j = 0; j < num_elements; j++){
14            if (a[j] < a[i] & (operations_counter++ | 1)){
15                num_smaller++;
16            } else {
17                if(a[j] == a[i]){
18                    num_equal++;
19                };
20            };
21        };
22
23        if (num_smaller < k && k <= (num_smaller + num_equal)){
24            *num_operations = *num_operations + operations_counter;
25            *time_taken = (float)*time_taken + ((float)clock() - (float)start);
26
27            return a[i];
28        };
29    }
30 }
```

6.1.2 Selection Median

```
1 #include "selectionMedian.h"
2
3 int selectionMedian(int *a, int num_elements, float *time_taken, int *num_operations){
4     int answer, operations_counter = 0;
5
6     clock_t start = clock();
7
8     if (num_elements == 1){
9         answer = a[0];
10    } else {
11        answer = select(a, 0, floor(num_elements/2.0) , num_elements - 1, &←
            operations_counter);
12    }
13
14    *num_operations = *num_operations + operations_counter;
15    *time_taken = *time_taken + ((float)clock() - (float)start);
16
17    return answer;
18 }
19
20 int select(int *a, int lower, int middle, int upper, int *operations_counter){
21
22     int pos = partitionSort(a, lower, upper, operations_counter);
```

```

23
24     /*operations_counter = *operations_counter + 1;
25     if (pos == middle){
26         return a[pos];
27     }else{
28         /*operations_counter = *operations_counter + 1;
29         if (pos > middle){
30             return select(a, lower, middle , pos-1, operations_counter);
31         } else if (pos < middle){
32             return select(a, pos+1, middle , upper, operations_counter);
33         };
34     };
35 }
36
37 int partitionSort(int *a, int lower, int upper ,int *num_operations){
38     int pivot_val = a[lower];
39     int pivot_loc = lower;
40     int operations_counter = 0;
41
42     for (int i = lower+1; i <= upper; i++) {
43         if (a[i] < pivot_val & (operations_counter++ | 1)) {
44             pivot_loc++;
45             swapValues(&a[pivot_loc], &a[i]);
46         }
47     }
48     swapValues(&a[lower], &a[pivot_loc]);
49
50     *num_operations = *num_operations + operations_counter;
51     return pivot_loc;
52 }
53
54 void swapValues(int *val_1, int *val_2){
55     int *temp = val_1;
56     val_1 = val_2;
57     val_2 = temp;
58 }

```

6.1.3 Main

```

1 #include <iostream>
2 #include <fstream>
3 #include <ctime>
4 #include <cstdio>
5 #include <stdlib.h>
6 #include "bruteForceMedian.h"
7 #include "selectionMedian.h"
8 #include "tests.h"
9
10 #define MIN_INPUT_SIZE 990
11 #define MAX_INPUT_SIZE 1003
12 #define STEP_INPUT_SIZE 3
13 #define NUMAVERAGES 2000
14
15 #define MS_PER_SEC 1000
16
17 using namespace std;
18
19 // Creates an array of random numbers for a given size
20 int *createRandomArray(int num_elements);
21
22 // writes the results to a csv file
23 int writeResultsToFile(int num_elements, int *input_size, float *brute_time, int *↵
    brute_operations, float *select_time, int *select_operations);
24
25 int main()

```



```

26 {
27     int do_test;
28     int num_datapoints = (MAX_INPUT_SIZE - MIN_INPUT_SIZE)/STEP_INPUT_SIZE;
29
30     srand(clock());
31
32     cout << "Would you like to run the tests? (0 = No, 1 = Yes)"<< endl;
33     cin >> do_test;
34     if (do_test > 0){
35         run_tests();
36         return 0;
37     };
38
39     int *input_size = new int[num_datapoints];
40     float *brute_time = new float[num_datapoints];
41     float *select_time = new float[num_datapoints];
42     int *brute_operations = new int[num_datapoints];
43     int *select_operations = new int[num_datapoints];
44
45     // #pragma omp parallel for
46     for (int i=0; i < num_datapoints; i++) { // for each test
47         input_size[i] = MIN_INPUT_SIZE+(i*STEP_INPUT_SIZE);
48
49         cout << "Input Size = " << input_size[i] << endl;
50
51         brute_time[i] = 0;
52         select_time[i] = 0;
53         brute_operations[i] = 0;
54         select_operations[i] = 0;
55
56         for (int j = 0; j < NUM_AVERAGES; j++){
57             int *a = createRandomArray(input_size[i]);
58
59             bruteForceMedian(a, input_size[i], &brute_time[i], &brute_operations[i] );
60             selectionMedian(a, input_size[i], &select_time[i], &select_operations[i] );
61
62             delete [] a;
63         };
64
65         brute_time[i] = (brute_time[i]*(float)MS_PER_SEC)/((float)NUM_AVERAGES*(float)←
        CLOCKS_PER_SEC);
66         select_time[i] = (select_time[i]*(float)MS_PER_SEC)/((float)NUM_AVERAGES*(float)←
        CLOCKS_PER_SEC);
67
68         brute_operations[i] = (int)brute_operations[i]/NUM_AVERAGES;
69         select_operations[i] = (int)select_operations[i]/NUM_AVERAGES;
70
71         cout << "    Brute Force Algorithm" << endl;
72         cout << "        Time Taken(ms) = " << brute_time[i] << endl;
73         cout << "        Num Operations = " << brute_operations[i] << endl;
74
75         cout << "    Selection Algorithm" << endl;
76         cout << "        Time Taken(ms) = " << select_time[i] << endl;
77         cout << "        Num Operations = " << select_operations[i] << endl << endl;
78
79     };
80
81     writeResultsToFile(num_datapoints, input_size, brute_time, brute_operations, ←
        select_time, select_operations);
82
83     delete [] input_size;
84     delete [] brute_time;
85     delete [] select_time;
86     delete [] brute_operations;
87     delete [] select_operations;
88
89     return 0;
90 }

```

```

91 }
92
93 int *createRandomArray(int num_elements){
94     int *a = new int[num_elements];
95
96     for ( int i = 0 ; i < num_elements ; i++) {
97         a[i] = rand() % 100000000 + 1;
98     };
99
100     return a;
101 }
102
103 int writeResultsToFile(int num_elements, int *input_size, float *brute_time, int *brute_operations, float *select_time, int *select_operations){
104
105     ofstream myfile ("results.csv");
106     if (myfile.is_open())
107     {
108         myfile << "input size" << "," << "brute force - time taken" << "," << "brute force - number operations";
109         myfile << "," << "select - time taken" << "," << "select - number operations" << endl;
110         for ( int i = 0 ; i < num_elements ; i++) {
111             myfile << input_size[i] << "," << brute_time[i] << "," << brute_operations[i];
112             myfile << "," << select_time[i] << "," << select_operations[i] << endl;
113         };
114         myfile.close();
115     }
116     else cout << "Unable to open file";
117     return 0;
118 }

```

6.1.4 Unit Tests

```

1 #include "tests.h"
2
3 void run_tests(){
4     int test_size;
5
6     cout << "Test 1" << endl;
7     test_size = 5;
8     int test1 [5] = { 1, 2, 3, 4, 5 };
9     print_array(test1, test_size);
10    cout << "    Passed = " << get_results(test1, test_size, 3, 3) << endl;
11    delete [] test1;
12
13    cout << "Test 2" << endl;
14    test_size = 6;
15    int test2 [6] = { 1, 2, 3, 4, 5, 6 };
16    print_array(test2, test_size);
17    cout << "    Passed = " << get_results(test2, test_size, 3, 4) << endl;
18    delete [] test2;
19
20    cout << "Test 3" << endl;
21    test_size = 5;
22    int test3 [5] = { 1, 2, 3, 3, 5 };
23    print_array(test3, test_size);
24    cout << "    Passed = " << get_results(test3, test_size, 3, 3) << endl;
25    delete [] test3;
26
27 }
28
29 bool get_results(int* a, int num_elements, int brute_expected, int selection_expected){
30
31     bool brute_result = get_brute_results(a, num_elements, brute_expected);

```

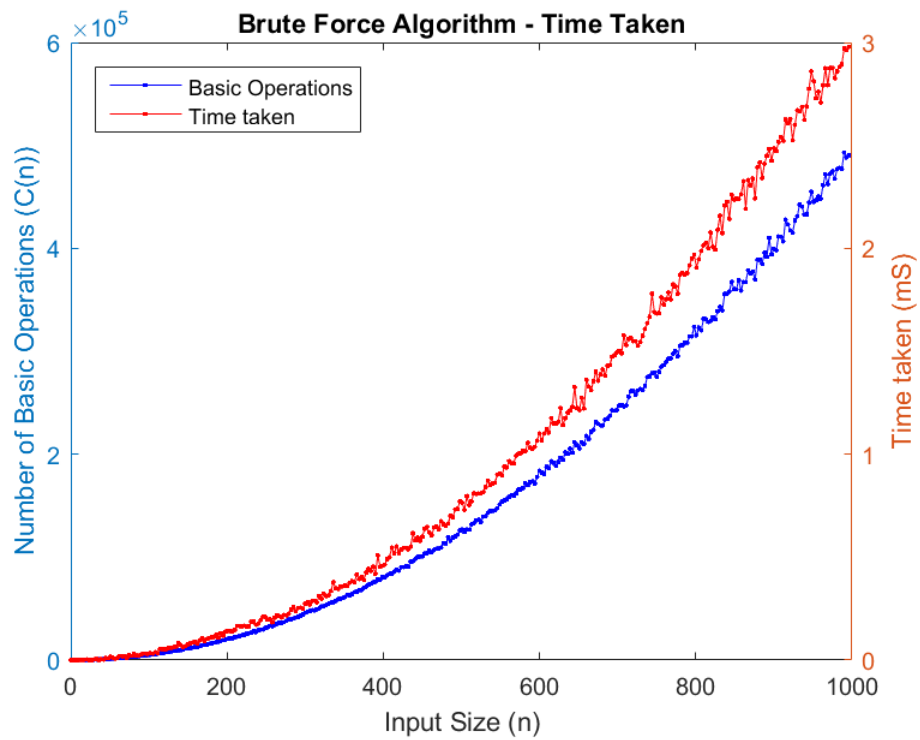
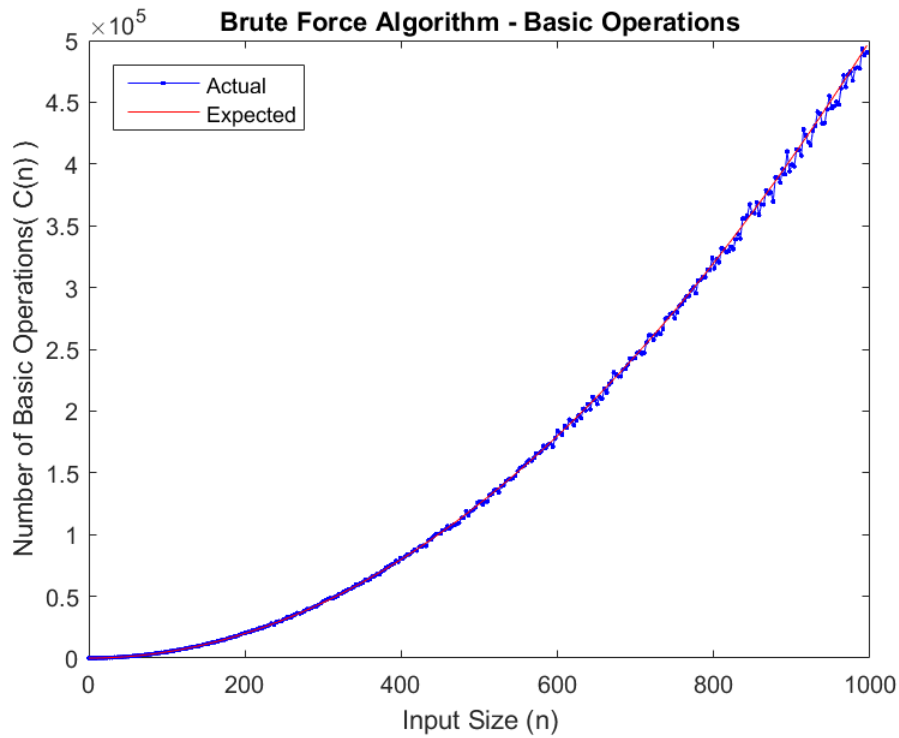
```

32
33     bool selection_result = get_selection_results(a,num_elements,selection_expected);
34
35     return (brute_result && selection_result);
36
37 }
38
39
40 bool get_brute_results(int* a, int num_elements, int expected){
41     int operations = 0;
42     float time = 0;
43     int ans;
44     bool result;
45
46     ans = bruteForceMedian(a, num_elements, &time, &operations);
47
48     result = (ans == expected);
49
50     cout << "          Brute Force Algorithm" << endl;
51     cout << "          Num Operations = " << operations << endl;
52     cout << "          Answer = " << ans << endl;
53     cout << "          Expected = " << expected << endl;
54     cout << "          Passed = " << result << endl;
55
56
57     return result;
58 }
59
60 // Returns true if the select median returns the correct expected result
61 bool get_selection_results(int* a, int num_elements, int expected){
62     int operations = 0;
63     float time = 0;
64     int ans;
65     bool result;
66
67     ans = selectionMedian(a, num_elements, &time, &operations);
68
69     result = (ans == expected);
70
71     cout << "          Selection Algorithm" << endl;
72     cout << "          Num Operations = " << operations << endl;
73     cout << "          Answer = " << ans << endl;
74     cout << "          Expected = " << expected << endl;
75     cout << "          Passed = " << result << endl;
76
77     return result;
78 }
79
80 void print_array(int* a, int length){
81     cout << "          Array = [ ";
82     for (int i = 0; i < length; i++){
83         cout << a[i] << " ";
84     };
85     cout << "]" << endl;
86 }

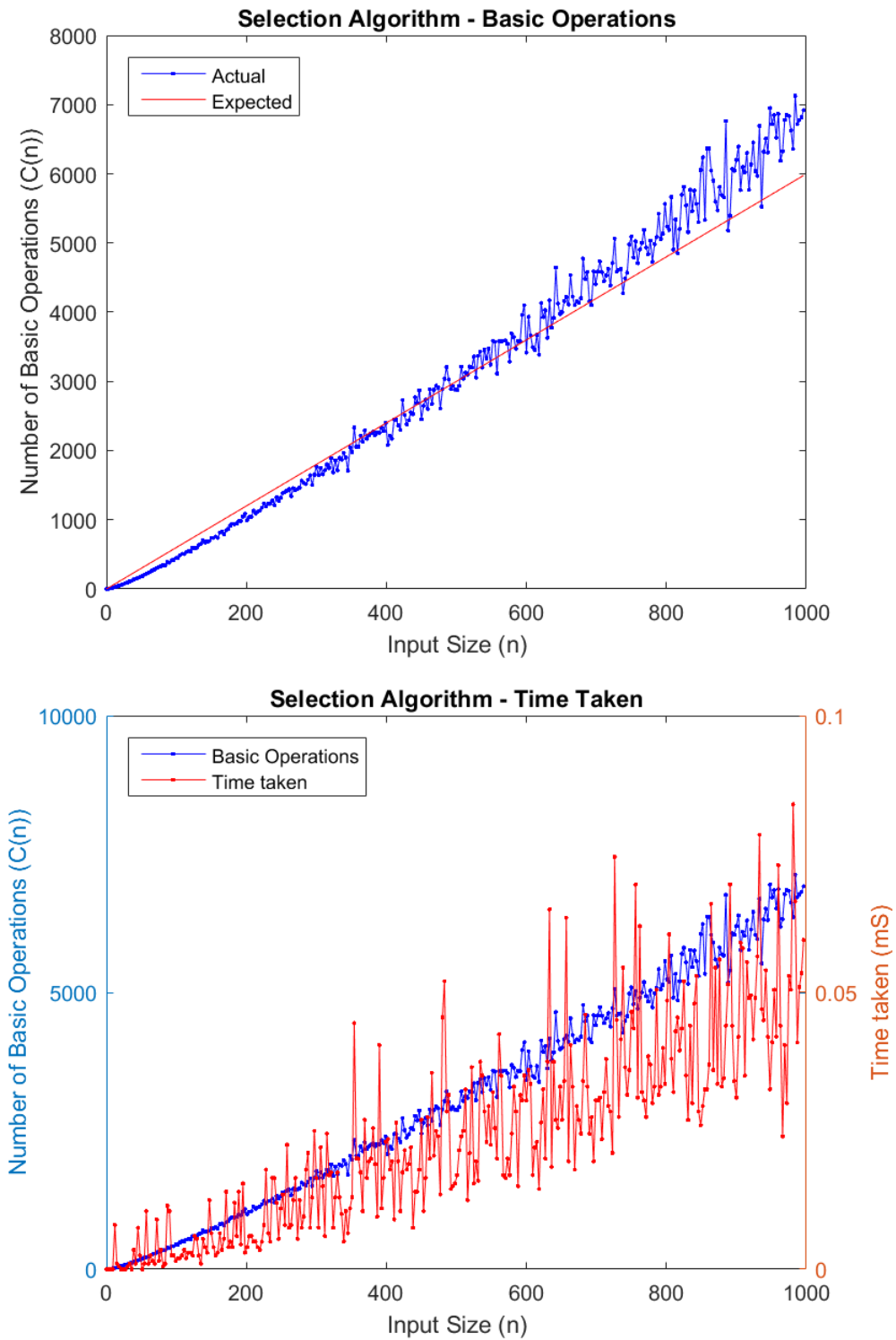
```

6.2 Figures

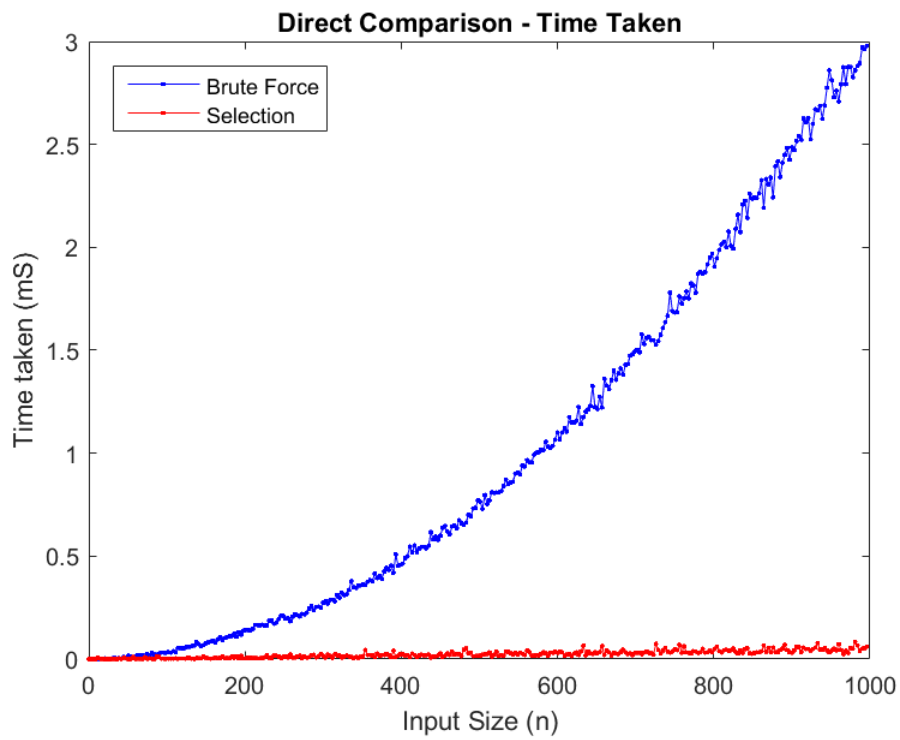
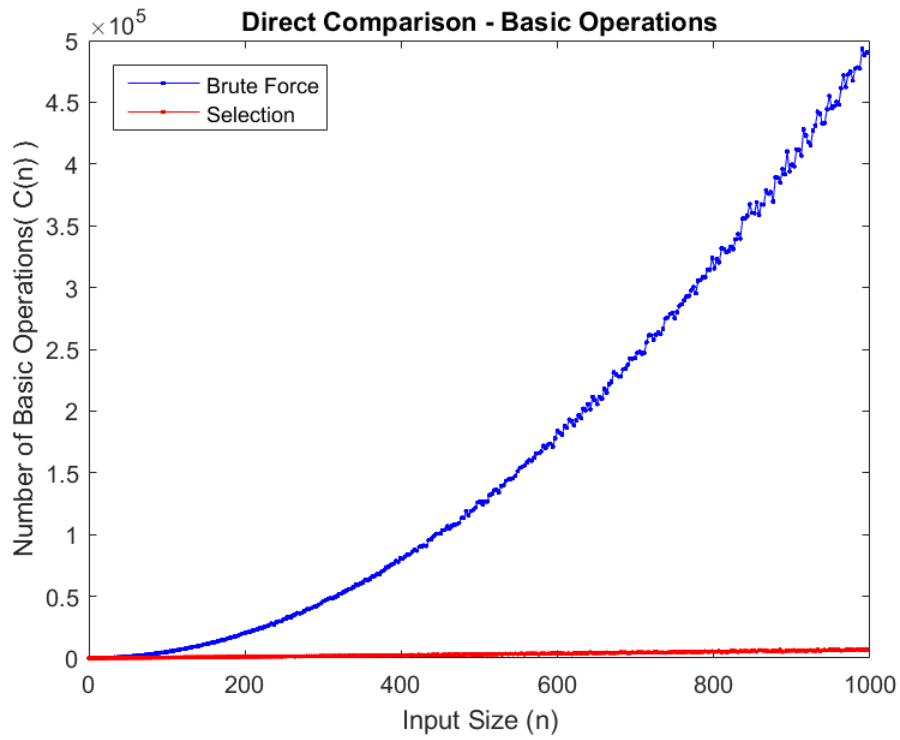
6.2.1 Brute Force



6.2.2 Selection



6.2.3 Comparison



7 Bibliography

References

- Balkcom, T. C. . D. (n.d.). *Analysis of quicksort*. <https://www.khanacademy.org/computing/computer-science/algorithms/a/analysis-of-quicksort>
(Online; accessed 15-May-2016)
- Paul, K. A. B. . J. L. (2004). *Algorithms: Sequential, parallel, and distributed* (1st ed., Vol. 1). Course Technology.