

IE531: Algorithms for Data Analytics
Spring, 2019
Programming Assignment 1: Median-of-Median Algorithms with
different Stopping Lengths
Due Date: February 8, 2019
©Prof. R.S. Sreenivas

1 Introduction

In Homework 1 you looked at the Median-of-Median Algorithm where the recursion was stopped as soon as the array-size was less-than-or-equal-to $m \in \{5, 7, 9, 11, 13\}$. We violated a bunch of “rules” by using the formulae for asymptotic running-time on arrays of small size (i.e. size m). If the mathematics of Homework 1 is to be believed, we expect the running-time to be the smallest for an optimal value of m in the range $\{5, 7, 9, 11, 13\}$. In this programming exercise you are going to experimentally verify/refute this claim.

Programming Assignment

You are going to modify the Python Code for the *Deterministic-Selection* (or *Median-of-Medians Algorithm*) written by me to create a version where the recursion stops whenever the array-size is less than m (cf. figure 1). Following this, you are going to implement the pseudo-code shown in figure 2 that verifies/refutes the theory developed in Homework 1.

Figure 3 shows the plots of mean-running-time versus m for different array-sizes (i.e. different values of n). It is not exactly clear that there is a unique value of m for which the mean-running-time is the smallest for all lengths. This is an artifact that comes out of using the asymptotic formulae (i.e. formulae that are correct only as n is very large) for cases where the array size (i.e. n) is small. It is hard to tell if repeating this experiment for very large values of n will confirm the theory discovered in Homework 1.

What I need from you

You can submit your Python code on Compass – make sure it runs without issues. In addition, I would like you to upload a PDF file that shows the plots (cf. figure 3) and a short explanation of why you were able to confirm/refute the theory developed in Homework 1.

Deterministic Select (vector **A** of size n , int $k \leq n$, int $m \leq n$)

```
1: if  $n \leq m$  then
2:   Sort A and return the element in  $k$ -th position.
3: end if
4: Partition A into vectors  $\{\mathbf{B}_i\}_{i=0}^{(n/5)-1}$ , where each vector  $\mathbf{B}_i$  has 5 elements.
5: for  $0 \leq i < n/5$  do
6:    $\mathbf{C}[i] = \text{Deterministic Select}(\mathbf{B}_i, 3)$ 
7: end for
   { /* C is a  $(n/5)$ -long vector, where the  $i$ -th entry is the median of  $\mathbf{B}_i$  */ }
8: (median-of-medians)  $p = \text{Deterministic Select}(\mathbf{C}, (n/10))$ 
9: Partition A into three sub-vectors  $\mathbf{A}_{<p}$ ,  $\mathbf{A}_{=p}$ , and  $\mathbf{A}_{>p}$ .
   { /*  $\mathbf{A}_{<p}$  has all elements of A that are less than  $p$ . */ }
   { /*  $\mathbf{A}_{=p}$  has all elements of A that are equal to  $p$ . */ }
   { /*  $\mathbf{A}_{>p}$  has all elements of A that are greater than  $p$ . */ }
10: if  $k \leq \text{length}(\mathbf{A}_{<p})$  then
11:   return  $\text{Deterministic Select}(\mathbf{A}_{<p}, k)$ 
12: else
13:   if  $k > \text{length}(\mathbf{A}_{<p}) + \text{length}(\mathbf{A}_{=p})$  then
14:     return  $\text{Deterministic Select}(\mathbf{A}_{>p}, k - \text{length}(\mathbf{A}_{<p}) - \text{length}(\mathbf{A}_{=p}))$ 
15:   end if
16: else
17:   return  $p$ 
18: end if
```

Figure 1: Pseudo-code for the *median-of-medians*, *deterministic-selection*, where m can be varied.

```

1: for  $1000 \leq n \leq 10,000$  in steps of 1000 do
2:   for  $m \in \{5, 7, 9, 11, 13\}$  do
3:     for  $1 \leq i \leq \#trials$  do
4:       Fill an array of size  $n$  with random values.
5:       Let  $k = \lceil \frac{n}{2} \rceil$  (i.e  $k$  is almost the median)
6:       Find the amount of time it took to find the  $k$ -th smallest element in the array
          (for the present value of  $m$ , array-size  $n$ , and trial  $i$ )
7:     end for
8:     Compute the mean-running-time of the  $\#trials$ -many experiments (for the
          present value of  $m$ , array-size  $n$ )
9:   end for
10:  Plot the mean-running-time as a function of  $m$  (for the present array-size  $n$ ).
11:  Check if there is a value of  $m$  for which the mean-running-time is the smallest
          (for the present array-size  $n$ ).
12: end for

```

Figure 2: Experimentally verifying/refuting the existence of an optimal value for m .

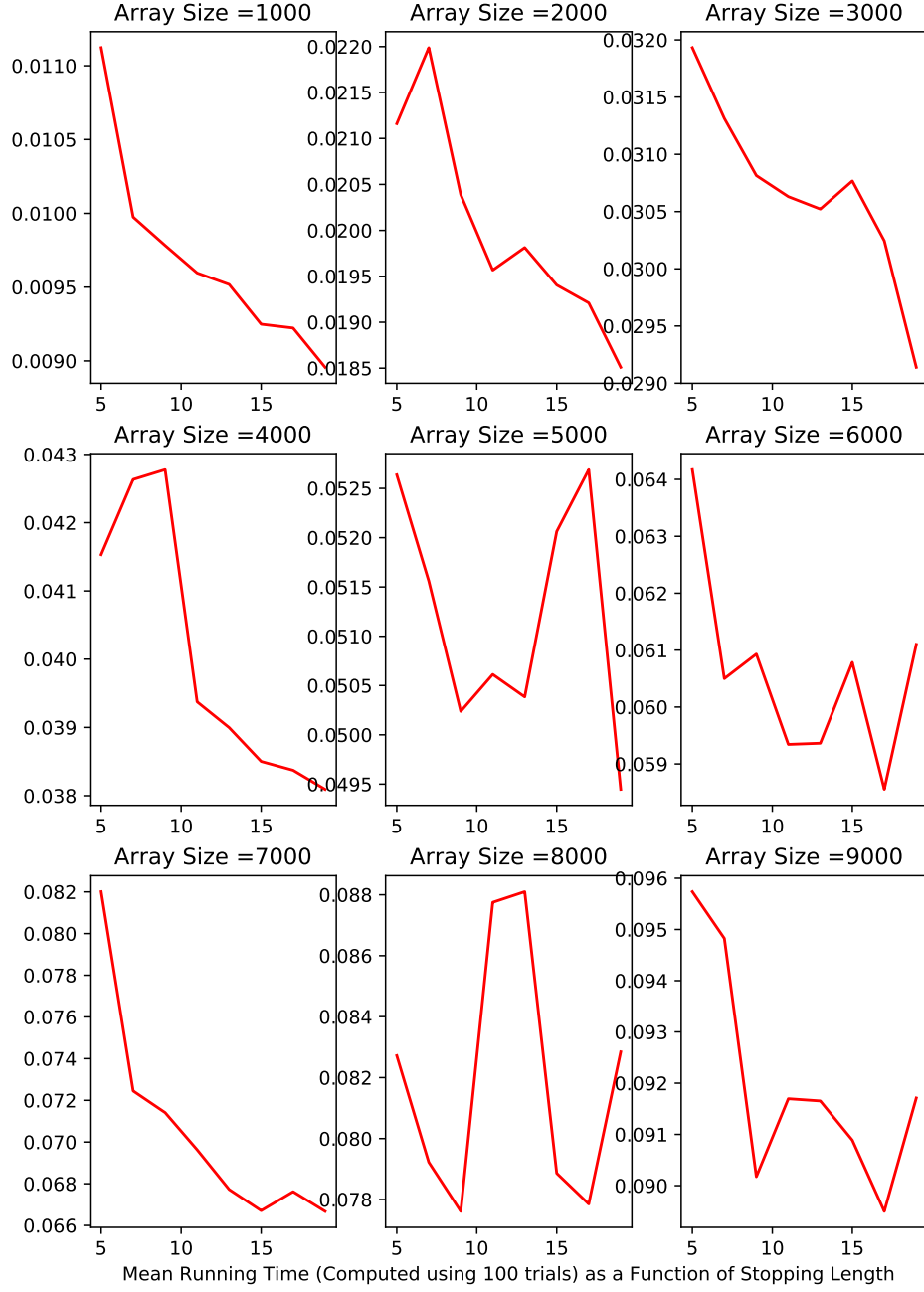


Figure 3: Mean-running-time (from 100 trials) vs. m plots for different values of n