

EE3980 Algorithms

Homework 7. Linear Sort

By 105061212 王家駿

1. Introduction

In the previous homework, the best time complexity of comparison-based sorting we got before is $O(n * \log n)$. In this homework, we try to use a non-comparison-based sorting algorithm to reach the time complexity of $O(n)$, which is linear. And then we compare the average time consumed for sorting strings by using linear sort and heap sort, respectively.

2. Implementation

For implementation, we first read in all the input data, and store them in a 2D array. Then, we use the radix sort, a non-comparison-based sorting algorithm, to sort the input strings in lexicographical order, and repeat the sorting for R times. Moreover, we also implement the heap sort from homework 2 for R times and get the average running time. Finally, we print out the array after sorted as well as the CPU time consumed for the two algorithms. The details are given below.

2.1. Method of storing data

In order to store the input strings, we use a 2D array of character. The size of the array is $N*14$, since there are N input strings and at most 14 characters in a string.

Now we examine the English characters in UTF-8 encoding rules. The lowercase letters are encoded from 97 for “a” to 122 for “z”, so we put the characters which are not belong to the string content, i.e. after ‘\0’, to be “^”, whose UTF-8 code is 96. After this step, we now have 14 characters for a string, and they are at the range from 96 to 122 in UTF-8 code.

The reason why we let the rest of the characters be “^” is that for the two strings with the same contents at the begin but one has extra letters, for example, “bit” and “bitwise”, the lexicographical order must be “bit” before “bitwise”. So if we put the “^” character at the end of “bit”, like “bit^^”, the two words have the same number of digits and it matches the alphabetical order in UTF-8. So, we could easily use it in the radix sort since the number of digits are same.

Also, we use a global variable to record the maximum string length. For the characters after the maximum string length, all of them are not the content of strings. That is, the characters after the maximum string length are all “^”. Thus, when we use the radix sort, we could neglect the digits after the maximum string length, since they are all the same kind and no need to be sorted.

2.2. Radix Sort

```
1. Algorithm RadixSort(A, d)
2. {
3.     for i := max_string_length to 1 step -1 do {
```

```

4.      Sort array A by digit i using CountingSort;
5.      }
6.  }

```

With the input data stored in the array, we use radix sort to sort it in lexicographical order. Radix sort applies sorting algorithm with respect to the certain digit from the LSB to MSB, and the sorting algorithm has to be stable. Thus, when we sort the higher digit, the lower digit we have sorted would not change its order if two of the higher digits are same. That is, the two string would arrange in the alphabetical order.

In our implementation, we use the counting sort for the stable sorting algorithm. And then we seem each character as a digit to be sorted, so we apply counting sort from the maximum string length (LSB) to the first character in the string (MSB). The result would be all the string arranged in lexicographical order.

2.3. Counting Sort

```

1. Algorithm CountingSort(A,B,n,k)
2. {
3.     // Initialize C to all 0
4.     for i := 1 to k do {
5.         C[i] := 0;
6.     }
7.
8.     // Count number elements in C[A[i]]
9.     for i := 1 to n do {
10.        C[A[i]] := C[A[i]] + 1;
11.    }

```

```

12.
13.    // C[i] is the accumulate number of elements
14.    for i := 1 to k do {
15.        C[i] := C[i] +C[i-1];
16.    }
17.
18.    // Store sorted order in array B
19.    for i := n to 1 step -1 do {
20.        B[C[A[i]]] := A[i];
21.        C[A[i]] := C[A[i]]-1;
22.    }
23. }

```

We use the counting sort as the stable sorting algorithm in the radix sort. First, we construct an array B with size N to store the sorting result, and an array C with size 27 to store the number of each element.

At the next step, we iterate through the array, and record the number of times each letter appears in C. Since we have only the lowercase letters and “\” we added, we need $26 + 1 = 27$ space to store the number of each letter in A. We use UTF-8 code to decide where we should record. Since we have known that the characters in the strings only have 96 to 122 in UTF-8 code. We stores the character with code 96 in C[0], and 91 in C[1], and so on. Thus, we could get all of the 27 indexes in the array C to record the appearing times of each letter respectively.

Then, we make the C array be the accumulation of the number of letters by adding the value from C[1] to C[i-1]. Thus, after the additions through the whole array, the content in the array C would be (not strictly) increasing.

At the final step, because we had known how many times each letter appears in the array, the only thing we have to do is place the string to array B directly. For a letter ltr, we would find C[ltr] to get the accumulation of index, and the index is where we should put the string in B. And then minus the accumulation by one, which indicates that the next string with the same letter would be place before the previous string. We keep executing this step until the iteration goes to the end, and it also means every component in array A was arranged well.

Finally, we copy the content in array B to A and finish the sorting of the digit.

2.4. Complexity analysis

For the time complexity, first in the radix sort function, the iteration at line 3 would run at most the time of maximum string length, which is no greater than 14. And at each iteration step, the counting would run for one time. Thus, the time complexity of radix sort is $14 * \text{time complexity of counting sort}$.

In the counting sort function, the iteration at line 4 and 14 would run for k times, where k is the number of possible values = 27. The iteration at line 9 and 19 would

run for n times. And there are only assignments and additions in the four loops. Thus, the time complexity is $O(4 * \max(n, k)) = O(4 * n) = O(n)$ since the data size n is often larger than k . Hence, the overall time complexity is $14 * O(n) = O(n)$, which is linear.

For the space complexity, we need two extra arrays B and C with size n and k for each counting sort, and the counting sort would execute for at most 14 times. Thus, the space complexity is $O(14 * (n + k)) = O(14 * 2n) = O(n)$.

Time complexity: $O(n)$

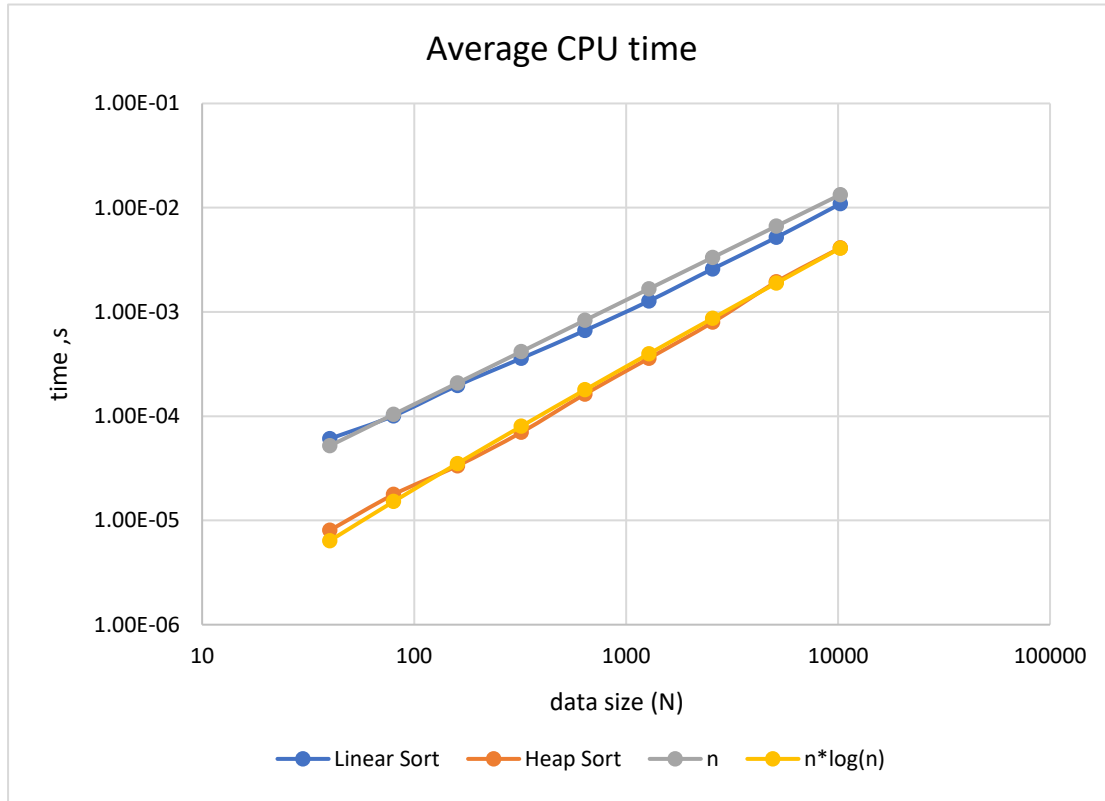
Space complexity: $O(n)$

3. Executing results

For repeating 500 times, we run the testing data from w11.dat to w19.dat with different input data size, and record the average CPU time used.

Data size (N)	Linear Sort	Heap Sort
40	60.74 μ s	8.036 μ s
80	100.8 μ s	17.72 μ s
160	196.7 μ s	33.41 μ s
320	358.5 μ s	69.90 μ s
640	663.8 μ s	163.0 μ s
1280	1.271ms	358.4 μ s

2560	2.588ms	795.4 μs
5120	5.187ms	1.944ms
10240	10.88ms	4.119ms



4. Result analysis and conclusion

From the graph, we could observe that the radix sort has linear trend, and the heap sort has the trend of $O(n * \log(n))$, which are same as our estimation.

However, for the case $N < 10000$, the radix sort is slower than the heap sort in spite of the fact that it has the lower time complexity. It may be the effect that at the each counting sort step in the radix sort, we have to construct two new arrays, and apply assignment through the whole array at least four times. These steps make the radix

sort slower than the heap sort in a small input data size. Yet, since the heap sort has a faster growing trend, it may become slower with regard to the radix sort if the input size keeps growing up.

We get the time complexity of $O(n)$ by neglecting the effect of string length. That is, assume the string length is far smaller than input data size. However, if the string length is close to, or greater than the data size, the time complexity would become $O(n * \text{string length}) = O(n^2)$. So, if the algorithm would become slower than the non-comparison-based sorting algorithm we used. Thus, the radix sort is only suitable when a small string length.

For the space complexity, the radix sort is $O(n)$, while the heap sort is $O(1)$ since it is in-place sorting, so the radix sort would take more space than the latter. In the end, we have to choose the algorithm carefully to be suitable for different problems.