

Efficient Data Structures for Range Selections Problem

Xiaodong Wang

Faculty of Mathematics & Computer Science
Fuzhou University, Fuzhou, China
Quanzhou Normal University

Jun Tian*

School of Public Health
Fujian Medical University
Fuzhou, China

Abstract—Building an efficient data structure for range selection problems is considered. While there are several theoretical solutions to the problem, only a few have been tried out, and there is little idea on how the others would perform. The computation model used in this paper is the RAM model with word-size $\Theta(\log n)$. Our data structure is a practical linear space data structure that supports range selection queries in $O(\log n)$ time with $O(n \log n)$ preprocessing time.

Keywords- range selection; RAM model; data structure

I. INTRODUCTION

In this work, we consider the problem of building an efficient data structure for range selection queries. The problem is to preprocess an input array A of n integers, such that given a query (i, j, k) , we can report the k' th smallest integer in the subarray $A[i], A[i+1], \dots, A[j]$ efficiently. In the rest of the paper, the subarray $A[i], A[i+1], \dots, A[j]$ is denoted as $A[i, j]$. A special case of the problem is known as range median query, which arises when k is fixed to $\lfloor (j-i+1)/2 \rfloor$. The prefix selection query is another special case of the problem, which arises when i is fixed to 0. These problems have many important applications in statistical analysis, and have been studied extensively in the last few years, see e.g. [1, 2, 3, 6].

We present a practical study on data structures for sequences supporting range selection queries. While there are several theoretical solutions to the problem, only a few have been tried out, and there is little idea on how the others would perform [5, 7, 8, 9]. The computation model used in this paper is the RAM model with word-size $\Theta(\log n)$. The data structure presented in this paper has the same basic approach as in [4]. We design a practical linear space data structure that supports range selection queries in $O(\log n)$ time.

The organization of the paper is as follows. In the following 3 sections we describe our general data structure design paradigm.

In section 2 we give an extremely simple data structure for answering range selection queries with $O(n \log n)$ time and space. Then the space consumption of the data structure is reduced to $O(n)$.

In section 3 we give a computational study of the presented data structure which demonstrates that the

achieved results are not only of theoretical interest, but also that the techniques developed may actually lead to a practical data structure for general range selection algorithms.

Some concluding remarks are in section 4.

II. THE ALGORITHM DESIGN

A. A Simple Data Structure

Let $A = (A[0], A[1], \dots, A[n-1])$ be the input array. Our data structure is a complete binary tree. We sort the array A and build a corresponding complete binary tree T that stores the n elements in the leaves in sorted order.

For a node v in the tree T , the subtree rooted at node v is denoted as T_v , and the number of leaves in T_v is denoted as $|T_v|$. In each node v of T , the $|T_v|$ elements in the leaves of subtree rooted at node v are stored in an array A_v of size $|T_v|$ sorted by their index in A .

The construction of the basic data structure can be described as follows.

Algorithm II.1: PREPROCESS(A, n)
comment: Build range select data structure
 $y \leftarrow$ INDEXSORT(A)
 $v \leftarrow$ BUILD($y, 0, n-1$)
return (v)

In the above algorithm, we first sort the input array A into an index array y by the algorithm *IndexsortA* such that the sequence $A[y[0]], A[y[1]], \dots, A[y[n-1]]$ is exactly the n elements $A[0], A[1], \dots, A[n-1]$ of the input array in sorted order. Then the actual construction of the tree T is performed by the recursive algorithm *Buildy, first, last* as follows.

Algorithm II.2: BUILD($y, first, last$)
comment: Construct the tree recursively
 $mid \leftarrow (first + last)/2$
if $last > first$
then
 $left \leftarrow$ BUILD($y, first, mid$)
 $right \leftarrow$ BUILD($y, mid + 1, last$)
 $v \leftarrow$ MAKE-TREE($left, right$)

```

Av ← MERGE(left.Av, right.Av)
return (v)
    
```

In the algorithm *Buildy, first, last*, the parameters *first* and *last* indicate the begin and end positions of array A_v in array y , which is the index array computed by *IndexsortA*. The tree T is built recursively in a bottom up manner. The sub-algorithm *Mergeleft.A_v, right.A_v* merges the two arrays in the subtrees of v into an array A_v of size $|T_v| = last - first + 1$ sorted by their index in A . This merge procedure is exactly the same as the merge procedure of merge sort algorithm.

Since our algorithm for constructing the basic data structure is essentially a sorting process, we can readily conclude that our algorithm uses $O(n \log n)$ time and $O(n \log n)$ words of space in the worst case.

If we visit the nodes of the complete binary tree T by an in-order traversal we can see that, at any node v of T , the elements in the array A_v are subdivided into two parts of (almost) equal size and stored in the left child and the right child nodes of v . The two parts in the subtrees are recursively subdivided further.

To answer a range selection query, we first visit the root of T and determine in which of its two subtrees the element of the desired rank lies. Once this is known, the search continues recursively in the appropriate subtree until a trivial problem of constant size is encountered.

The algorithm for answering the range selection queries (i, j, k) is described as follows.

Algorithm II.3: QUERY($v, first, last, i, j, k$)
comment: Answer range selection queries (i, j, k)
 $m \leftarrow (first + last)/2$
if $last = first$
then return ($A[y[first]]$)
 $l \leftarrow RANK(v.l, i - 1)$
 $r \leftarrow RANK(v.l, j)$
 $s \leftarrow r - l$
if $k \leq s$
then return (QUERY($v.l, first, m, l, r - 1, k$))
else return (QUERY($v.r, m + 1, last, i - l, j - r, k - s$))

In the above algorithm *Query_v, first, last, i, j, k*, we want to find the element of rank k in the subarray $A[i, j]$ from node v of T , where array A_v begins at position *first* and ends at position *last* of array y .

The number l is the rank of $i - 1$ in the array A_v of the left subtree $v.left$. It means that there are l elements in the left subtree $v.left$ are to the left of i in the subarray $A[i, j]$. The number r is the rank of j in the array A_v of the left subtree $v.left$. It means that there are r elements in the left subtree $v.left$ are to the left of and up to j in the subarray $A[i, j]$. Thus the number s is the number of elements in the array A_v of the left subtree $v.left$ contained in in the subarray $A[i, j]$.

If $k \leq s$ then the element of rank k in $A[i, j]$ is the element of rank k in the subarray $A_v[l, r - 1]$ of the left subtree $v.left$. Otherwise, the element of rank k is the element of rank $k - s$ in the subarray $A_v[i - l, j - r]$ of the right subtree $v.right$.

Thus the algorithm reduces the problem of finding an element of a given rank in the subarray $A[i, j]$ to the same problem, but on a smaller array. This reduction is applied recursively by the algorithm.

The number l and r can be found in $O(\log((last - first + 1)/2))$ time by a binary search. The query descends $\log n$ levels of recursion, so the algorithm *Query* would get a total execution time of up to $\sum_{i=1}^{\log n} O(\log n/2^i) = O(\log^2 n)$.

B. An Improved Data Structure

In the algorithm *Query_v, first, last, i, j, k* we can see that to find out in which subtrees the element of rank k lies, only the array stored in the left subtree is used. Therefore, we can lift it to the node v and thus save half of the total space costs. Furthermore, we note that to compute the numbers l and r , the information available in the arrays stored at the interior nodes of our data structure can be reduced further. We can use a bit-vector to store the information we need, where a 1-bit indicates whether an element of the original array is in A_v . Since we have n positions one very level, a total of $O(n \log n)$ bits are used in our data structure.

With these bit-vectors, the execution time of the algorithm *Query* can also be reduced.

In order to compute the number l and r in the query algorithm efficiently, we can store a table with ranks for indices that are a multiple of the size of machine word w . General ranks are then the sum of the next smaller table entry and the number of 1-bits in the bit-vector between this

rounded position and the query position. In this way, the number l and r can be computed in $O(1)$ time. Therefore the execution time of the algorithm *Query* can be reduced to $O(\log n)$.

Summing up, we have obtained a data structure for solving range selection problem with preprocessing time $O(n \log n)$ using $O(n)$ space and query time $O(\log n)$. This improves the space consumption compared to [16] by a factor $O(\log^2 n / \log \log n)$.

III. THE EXPERIMENTS

In this section we will describe the implementation of the data structure presented in last section. Based on the discussion above, we can design an node class for the nodes of the complete binary tree of our new data structure.

In the node class we store a bitvector *low* to indicate the elements in the left subtree of current node sorted by their index in A .

The vector t is a table with ranks for indices that are a multiple of the size of machine word w . With the vector t the rank of an index in the bitvector *low* can be computed in $O(1)$ time. The complete binary tree of our new data structure is built recursively in a bottom up manner. For the current node, the two arrays of its subtrees are merged into one array and then the bitvector *low* is formed. According to the information of *low*, the vector t can then be constructed readily. The merge procedure is exactly the same as the merge procedure of merge sort algorithm. In the algorithm *init*, the parameters *first* and *last* indicate the begin and end positions of array A_v in array y , which is the index array computed by the *Indexsort* algorithm. The size of bitvector *low* must be $last - first + 1$ bits.

A dequeue is used for merging two sorted arrays $y[first, mid]$ and $y[mid + 1, last]$. In the merge process, if the next element comes from the first array $y[first, mid]$, then the corresponding bit of bitvector *low* is marked true. When the two arrays are sorted, the bitvector *low* is built. The table t can then be constructed easily from the bitvector *low*.

With the class *node*, we can design a new class *rmedian* for our new data structure for general range selection query as follows. In the class *rmedian*, array *data* is used to store the input sequence. The arrays y and z are used to store the sorted index arrays of the input sequence. The main part of the class is a vector *bt* of element type *node*. This vector is used to store the

complete binary tree T that the n elements of the input sequence are stored in its leaves in sorted order. The vector *bt* is in fact an array indexed binary tree. The root of the tree is $bt[0]$. For a given index i of a node, the left and right child node of $bt[i]$ are $bt[2*i+1]$ and $bt[2*i+2]$ respectively. The parent node of $bt[i]$ is $bt[(i-1)/2]$.

The vector *bt* can be built recursively in a bottom up manner. Once an instance of *rmedian* is built, we can then answer any range selection query in $O(\log n)$ time.

For any range selection query *QUERY*(ind, first, last, left, right, rank), we want to find the element of rank *rank* in the subarray $A[left, right]$ from node $v = bt[ind]$ of *bt*, where array A_v begins at position *first* and ends at position *last* of array z .

We first find the number l and r , which are the ranks of $left - 1$ and $right$ in the array A_v of the left subtree of node v respectively.

Thus the number $length = r - l$ is computed, which is the number of elements in the array A_v of the left subtree of v contained in in the subarray $A[left, right]$.

If $rank \leq length$ then the element of rank *rank* in $A[left, right]$ is the element of rank *rank* in the subarray $A_v[l, r - 1]$ of the left subtree of v . Otherwise, the element of rank *rank* is the element of rank $rank - length$ in the subarray $A_v[left - l, right - r]$ of the right subtree of v .

The algorithm reduces the problem of finding an element of a given rank in the subarray $A[left, right]$ to the same problem, but on a smaller array. This reduction is applied recursively by the algorithm.

We implemented our data structure in C++ and tested them on a personal computer with Pentium(R) Dual Core CPU 2.10 GHz and 2.0 Gb RAM, using the Microsoft Visual C++ version 8.0 compilers. The word size of the processor is $w = 32$.

The experiment results show that our data structure is very practical for solving the range selection query problem.

We also performed some limited experiments on the relative performance of our data structure. The new data structure has similar or better speed than existing data structures but uses less space in the worst case.

IV. CONCLUDING REMARKS

We have presented new data structure for solving the range selection query problem. While there are several

theoretical solutions to the problem, only a few have been tried out, and there is little idea on how the others would perform. The computation model used in this paper is the RAM model with word-size $\Theta(\log n)$. Our data structure is a practical linear space data structure that supports range selection queries in $O(\log n)$ time with $O(n \log n)$ preprocessing time.

The computational experiments in Section 3 demonstrate that the achieved results are not only of theoretical interest, but also that the techniques developed may actually lead to considerably faster algorithms.

REFERENCES

- [1] M. J. Atallah and H. Yuan, Data structures for range minimum queries in multidimensional arrays, In Proceedings of the 20th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA), 150-160, 2010.
- [2] G. S. Brodal and A. G. Jorgensen, Data structures for range median queries, In Proceedings of the 20th International Symposium on Algorithms and Computation, 822-831, 2009.
- [3] M. Chan, Persistent predecessor search and orthogonal point location in the word RAM, In Proceedings of the 22nd ACM/SIAM Symposium on Discrete Algorithms (SODA), 1131-1145, 2011.
- [4] B. Gfeller and P. Sanders. Towards optimal range medians. In Proceedings of the 36th International Colloquium on Automata, Languages and Programming, 475-486, 2009.
- [5] D. Krizanc, P. Morin, and M. H. M. Smid, Range mode and range median queries on lists and trees. Nordic Journal of Computing, 12(1):1-17, 2005.
- [6] Kasper Green Larsen, The cell probe complexity of dynamic range counting, In Proceedings 44th ACM Symposium on Theory of Computing (STOC), 2012.
- [7] H. Petersen, Improved bounds for range mode and range median queries, In Proceedings of the 34th Conference on Current Trends in Theory and Practice of Computer Science , 418-423, 2008.
- [8] H. Petersen and S. Grabowski, Range mode and range median queries in constant time and sub-quadratic space, Information Processing Letters, 109(4):225-228, 2008.
- [9] D. E. Willard, Log-logarithmic worst-case range queries are possible in space $\Theta(n)$, Information Processing Letters, 17(2):81-84, 1983.