

A Research Project Report (Midsem)

On

**Reporting and counting maximal points in a query orthogonal rectangle**

By

Ananda Swarup Das, Prosenjit Gupta, Kishore Kothapalli, Kannan Srinathan

**PUBLISHER:** Journal of Discrete Algorithms

**DAS ANAMIKA ABHOYPAD**

**2014H313059H**

Under the supervision of

**TATHAGATA RAY**



**BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE**

**HYDERABAD CAMPUS**

## ABSTRACT

In this work we show that given a set  $S$  of  $n$  points with coordinates on an  $n \times n$  grid, we construct data structures for reporting and counting the maximal points in an axes-parallel query rectangle in sub-logarithmic time. We assume our model of computation to be the word RAM with size of each word being  $\log n$  bits.

## INTRODUCTION

Database system plays a major role in the industries as well as in our daily life. Hence, optimisation and efficiency are the important factors of the database system.

The given research paper focuses on database community problems. The domain relates to “skyline points”.

Lets walkthrough with an example,

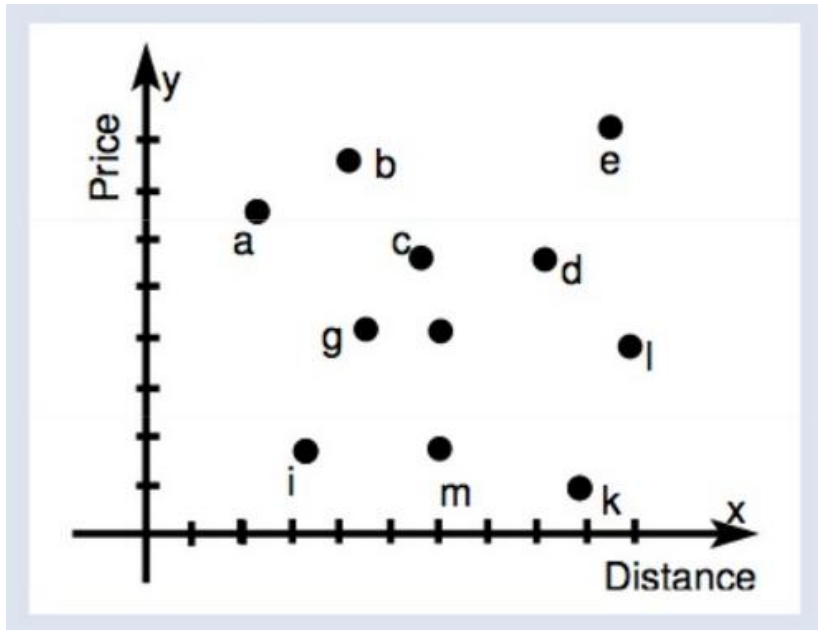
Suppose, you have a 10 days suspension from the mundane work and you have decided to go on a holiday. The favorite scenario that you prefer, is a beautiful sunset on a beach. You open some travel site and search for a hotel that is near to the beach! Plus the expense must be relative to your pocket size.

There are three possible search results:

1. Cheapest hotel but not close to the beach.
2. Closest to the beach but expensive hotel.
3. Neither close to the beach nor expensive hotel , relatively fits the criteria of your choice.

Lets take a look at the graph with the datasets that contains the details of the distance to the beach and the price of the hotel:

<http://www.slideshare.net/victormier/skyline-queries>



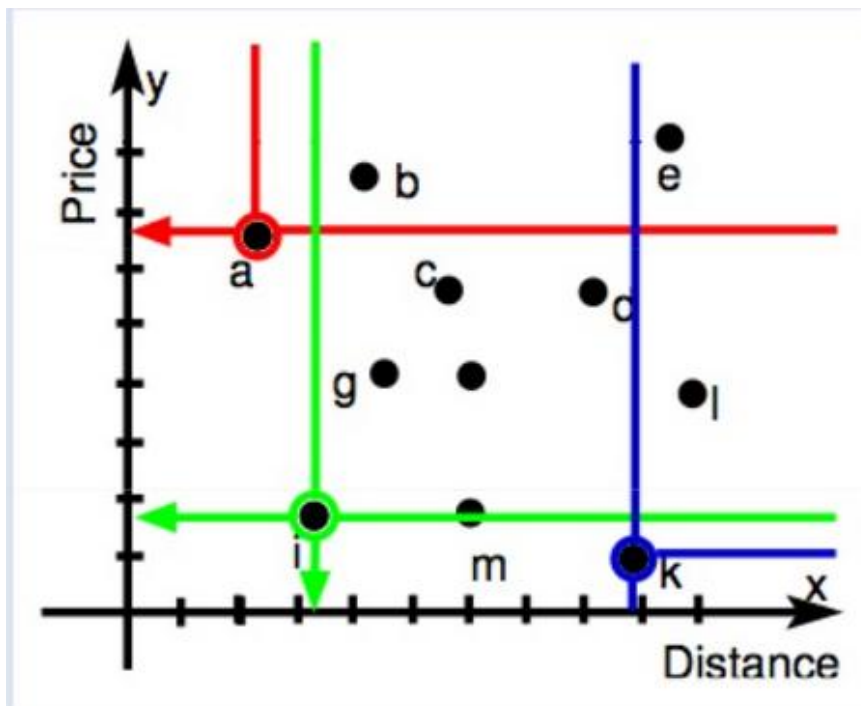
x-axis= Distance, y-axis= Price

The interesting points are a, i, k. Why??

We can see that point a has the least distance but high cost, point k has the lowest price with the not so close to the beach.

i neither has shortest distance nor lowest price. i has distance less than k and price more than a. i is not **dominated** by either of them

All the other points are dominated by the points a, i and k.



Hence the points that are not dominated by another points are called as skyline points. The skyline points for the above graph are a, i, k.

The dataset: (price , distance).

The skyline points are important to the large datasets that are stored in the database.

**DEFINITION:** Skyline is a set of tuples of information that are important or of special interests to us.

A point dominates another point if it is good in all dimensions or atleast in one dimension.

A point  $p(x_1, y_1)$  is dominated by point  $q(x_2, y_2)$  if  $x_1 < x_2$  and  $y_1 < y_2$  in terms of speciality.

The research paper has described **skyline points as maximal points**.

The reseach paper deals with 2 Dimensional skyline points i.e., the co-ordinates of x-axis and y-axis.

## PROBLEM DEFINITION

The user is not concerned with the entire dataset. A user enters his/her queries i.e. window size and accordingly the result is displayed. The simple approach goes as to check every dataset with the given queries and give the output. If the datasets are too large, the time taken for the search queries increases and hence expensive in terms of query time.

The solution uses a data structure of size  $O(n \log n)$  and can be queried to report maximal points in  $O(\log^2 n + k)$  time, where  $k$  is the size of the output.

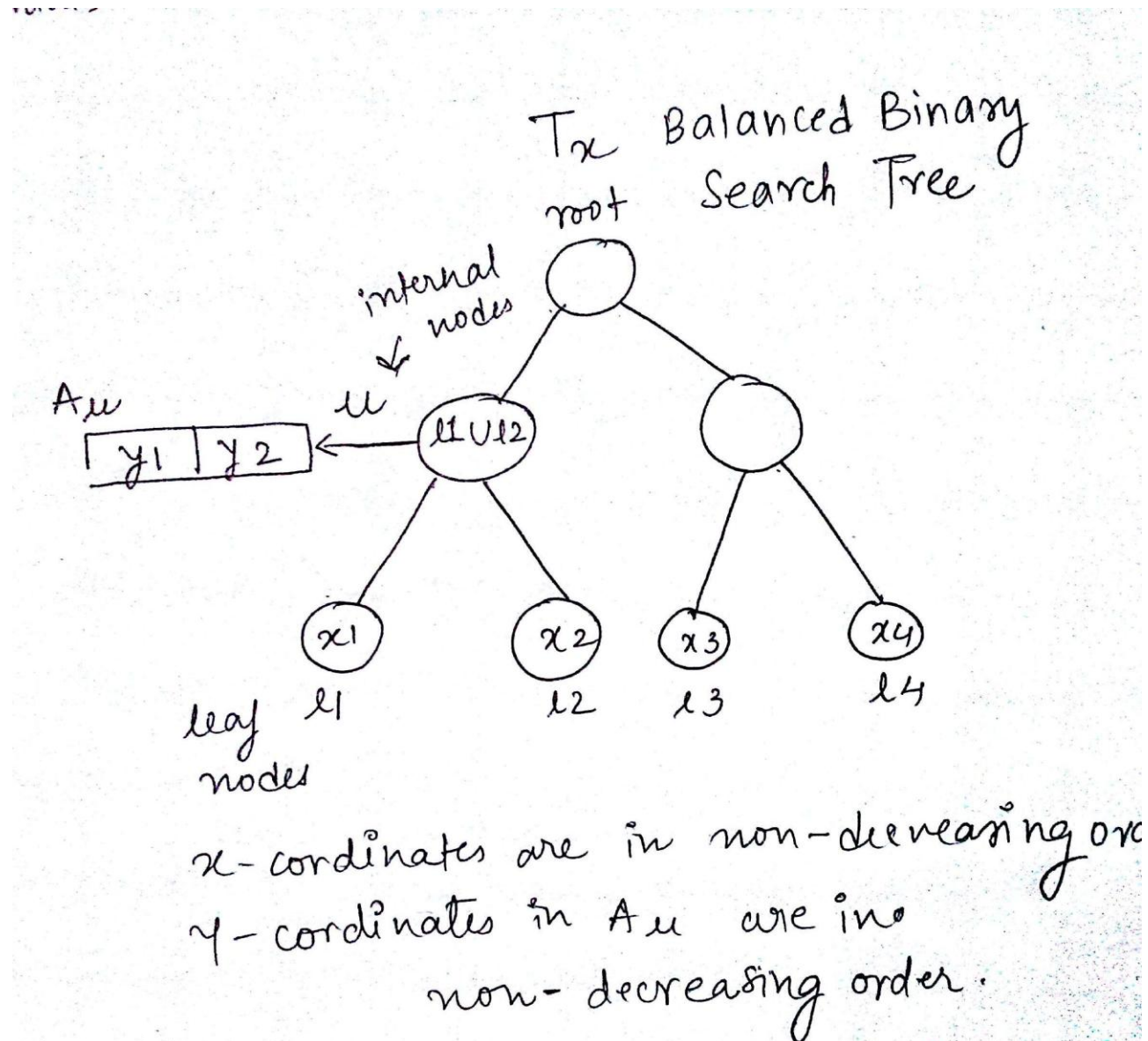
A proposed solution for the word RAM model of computation with size of each word being  $\theta(\log n)$  and propose a data structure of size  $O(n^*(\log n / \log \log n))$  that can be queried in  $O((\log n / \log \log n) + k)$  time.

The maximal point also reduces the size of the result. We have an algorithm that computes this size. It is useful in the context, there is a need to transfer the result using network bandwidth. We propose a data structure of size  $O((n^*(\log^3 n / \log \log n)))$  that can be queried in  $O(\log n / \log \log n)$  time.

## RANGE TREE WITH FRACTIONAL CASCADING

We present a data structure with  $O(n \log n)$  space and  $O(\log n + k)$  query time i.e., 2-dimensional range tree (with fractional cascading pointers) in order to support "2-dimensional range maximal points" reporting queries.

Given a Set  $S$  of points (dataset) with  $x$  and  $y$  coordinates and a query rectangle  $q$  (user queries). We can process the points that are in  $S \cap q$ ;



We have four points  $(x_1, y_1)$ ,  $(x_2, y_2)$ ,  $(x_3, y_3)$  and  $(x_4, y_4)$ .

We Build a tree  $T_x$  that is balanced binary tree that has leaf nodes  $l_1, l_2, l_3, l_4$ .

We store the  $x$  co-ordinates in non- decreasing order of their values at the leaf nodes i.e.,  $x_1 < x_2 < x_3 < x_4$

Let  $\mu$  be an internal node which is union of the discrete intervals associated with the leaf nodes of the subtree rooted at  $\mu$ .

Each internal node  $\mu$  also maintains an auxiliary array  $A_\mu$  which stores the  $y$ -coordinates of the leaf nodes present in the subtree rooted at  $\mu$ . These  $y$ -coordinates are sorted in non-decreasing order of their values i.e.,  $y_1 < y_2$

Leaf nodes of subtree rooted at  $\mu$  are  $l_1$  and  $l_2$ .

*For any node  $\phi$  belongs to  $T_x$ , it must be noted that  $int(\phi) = int(u) \cup int(w)$  where  $u, w$  are respectively the left and the right children for  $\phi$ . Also  $A_\phi = A_u \cup A_w$ . Thus, with each element  $y$  belongs to  $A_\phi$ , we store two pointers pointing to the smallest values greater than  $y$  in  $A_u$  and  $A_w$  respectively. The idea of maintaining such pointers is known as fractional cascading.*

The storage space needed by the range tree is  $O(n \log n)$ .