

A Research Project Report

Reporting and counting maximal points in a query orthogonal rectangle

By

Ananda Swarup Das: *IBM India Research Lab, New Delhi, India*

Prosenjit Gupta: *Heritage Institute of Technology, Kolkata, India*

Kishore Kothapali: *International Institute of Technology, Hyderabad, India*

Kannan Srinathan: *International Institute of Technology, Hyderabad, India*

PUBLISHER: Journal of Discrete Algorithms

ACCEPTED: 4th December 2014

Under the supervision of

TATHAGATA RAY



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

HYDERABAD CAMPUS

ANAMIKA DAS

2014H313059H

ACKNOWLEDGEMENT

I would like to express my gratitude to Prof Tathagarte Ray who gave me the opportunity to do this Research project on the topic of Counting the maximal points in an orthogonal rectangle, which helped me doing the Research and learn a new dimension from it.

I would express my gratitude to BITS Pilani Hyderabad Campus for providing me with the required opportunity and infrastructure.

ABSTRACT

In the given Research Paper, a proposal of a solution with sub-logarithmic query time for counting the number of maximal points in an axis parallel query rectangle.. The problem has been previously studied, to the authors' best knowledge, this is the first sub-logarithmic query time solution for the problem. The model of computation is the word RAM with word size of $\log(n)$ bits.

The main purpose of the project is to analyse, study and implement the algorithm to solve open problems.

This algorithm is useful in the context of database where enormous data are executed per second. There are customer inputs as to retrieve the results. Time complexity increases as the size of the data increases. To solve this problem, the author has given an efficient algorithm.

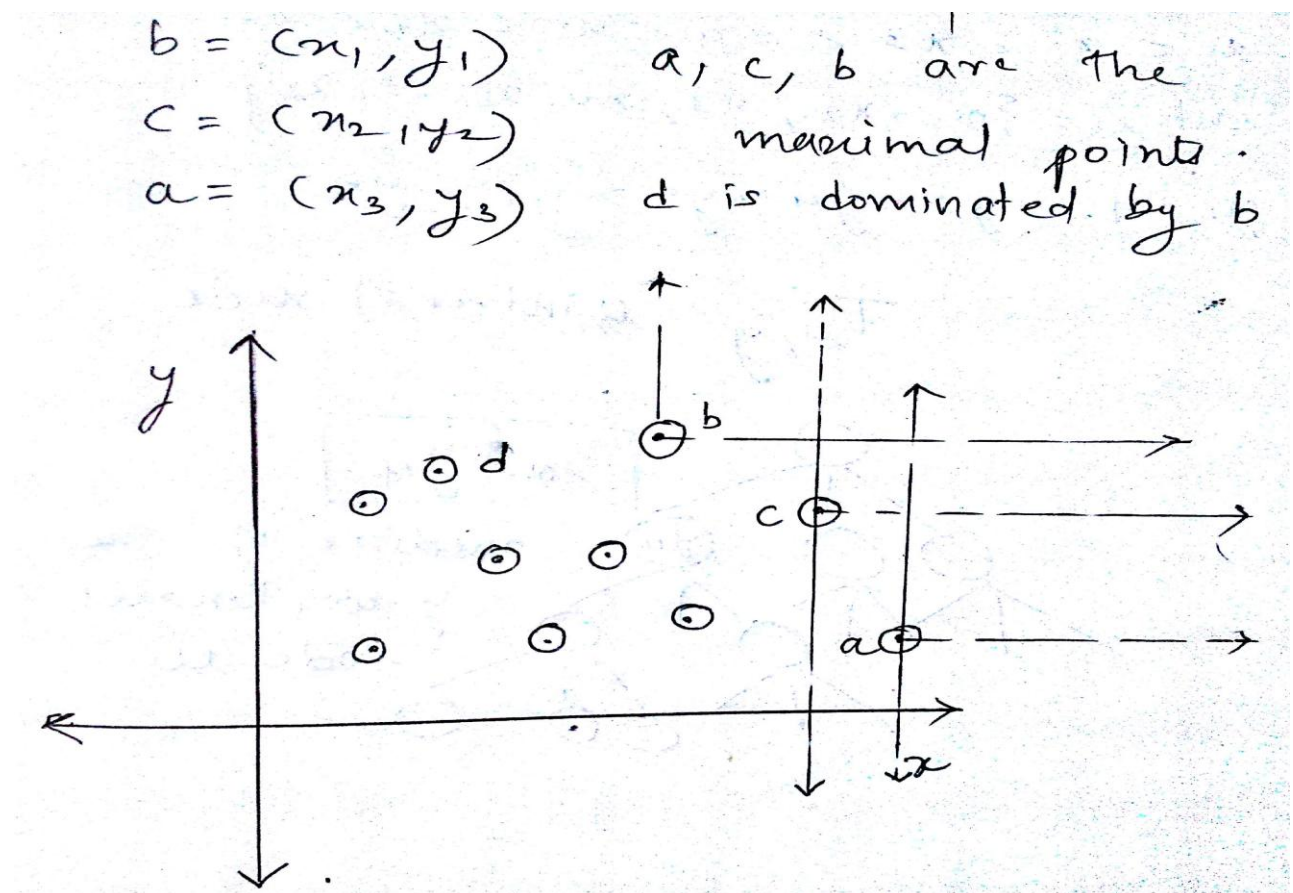
INTRODUCTION

The given research paper deals with the skyline points or maximal points.

According to the definition, we have two coordinates x and y .

Let us consider set of points in the 2 D space i.e., x and y coordinates.

According to the definition of maximal points, if point $q(x_1, y_1)$ is dominated by point $p(x_2, y_2)$ iff $x_1 \leq x_2$ and $y_1 \leq y_2$. This point p is called as maximal point.



From the above diagram, we can observe that point d is dominated point b . similarly, a, c are dominating points. a, c, b are the maximal points.

Our main context, we need to compute the count of maximal points from a given set of points S and a query rectangle R in $O(\log n / \log \log n)$.

PROBLEM DEFINITION

PROBLEM 1:

Given a set R of n points on an $n \times n$ grid, that is $R = \{(x_i, y_i) | x_i \in \{1, n\}, y_i \in \{1, n\}\}$, preprocess R into a data structure such that given an orthogonal query rectangle q , the maximal points in $q \cap R$ can be reported efficiently.

Assume all the points to have distinct x and y co-ordinates. This problem requires storage of size $O(n \log n / \log \log n)$ and query time $O((\log n / \log \log n) + k)$ time.

SOLUTION:

1. Consider two level binary search Tree T_x .
2. Given a query rectangle $q = [a, b] \times [c, d]$. the segment $[a, b]$ can be allocated to at most $O(\log n)$ canonical nodes of the tree T_x . Ordered position: right to left.
3. To achieve the desired complexity, increase the internal degree of each internal node of T_x from two to $O(\sqrt{\log n})$, thereby decreasing the height of the tree to $O(\log n / \log \log n)$
4. Given a query rectangle $[a, b] \times [c, d]$, allocate the segment $[a, b]$ to a node $\mu \in T_x$, if $\text{int}(\mu) \subset [a, b]$ but $\text{int}(\text{parent}(\mu)) \not\subset [a, b]$. The node μ is known as a *canonical node*.
5. The set of all such canonical nodes can be grouped into $O(\log n / \log \log n)$ groups with each group containing children v_1, \dots, v_k for some node v .
6. Node v is called as group leader for the group $g(i)$. Denoted by $GL(i)$.
7. Let $G = \{GL(1), \dots, GL(O(\log n / \log \log n))\}$ be the set of *group leaders* stored in order of their positions from right to left in the tree T_x .
8. If we visit each group leader node $GL(i)$ from right to left and the process count in $O(1)$ time the number of maximal points in q from the nodes of the group $g(i)$ where $GL(i)$ is the group leader, we can count the maximal points in q in $O(\log n / \log \log n)$.

In order to solve the problem 1, we need to solve the subproblems.

PROBLEM 2:

Given a set S of n points from a universe of $[1, \log^{(p)}n] \times [1, n]$ for $0 < p \leq 1/2$, preprocess the points into a data structure such that given an axis parallel rectangle

$q = [a, b] \times [c, d]$ for (a, b) belongs to $[1, \log^{(p)}n] \times [1, \log^{(p)}n]$ and (c, d) belongs to $[1, n] \times [1, n]$, we can efficiently count the number of maximal points in $S \cap q$.

SOLUTION:

There are 4 points: $(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4)$

1. Array $A = \{y_1, y_2, y_3, y_4\}$ where $y_1 > y_2 > y_3 > y_4$.

2. Construct a height-balanced binary search tree T_y .

3. Leaf nodes store the values in the array A , the level is same for every leaf node. At each internal node m belongs to T_y , store a key value y_m which is the median of the points stored in the subtree rooted at m .

For example: $y_m = (y_1 + y_2)/2$ where y_1 and y_2 are the leaf nodes of subtree m .

3. For each pair of possible points (a, b) belongs to $[1, \log^{(p)}n] \times [1, \log^{(p)}n]$, form an interval $[a, b]$.

4. For each value y_2 belongs to A and for each possible interval $[a, b]$, form 3-sided anchored rectangles $[a, b] \times [y_2, \infty)$ and $[a, b] \times (-\infty, y_2]$. Next, for the interval $[a, b]$ and the value y_2 .

(a) Visit the ancestors of y_2 in the tree T_y . At each ancestor m , find (y_m) .
Median of the root = $(y_2 + y_3)/2$

(b) If $y_m < y_2$, form an axis-parallel rectangle $R_1 = [a, b] \times [y_m, y_2]$.

i. For the points in $S \cap R_1$, compute the subset of points that are not dominated by any other point. This subset is called as maximal chain.

$p_{y_{max}}$ = the topmost point of the maximal chain

$p_{x_{max}}$ = the bottommost point of the maximal chain

ii. Count the number of maximal points in the chain $p_{x_{max}}$ to $p_{y_{max}}$. Denote the value as $count_m(y_2) = |p_{x_{max}}, p_{y_{max}}|$.

iii. Find the topmost point $pnodom$ in $[a, b] \times (-\infty, ym]$ such that $pnodom$ is not dominated by $pxmax$.

iv. Create a tuple $(countm(y2), pnodom)$ and store it with reference to the rectangle $R1$ in a lookup table.

Here the suffix m denotes the index for the node m belong to T_y that is an ancestor of the leaf node storing the value $y2$.

v. Special cases

A. If no points are present in the rectangle $R1$, store $(0, NULL)$.

B. If the point $pnodom$ does not exist, store $(countm(y2), NULL)$.

(c) If $ym > y2$

Rectangle $R2 = [a, b] \times [y2, ym]$

Find $(p'ymax, p'xmax)$

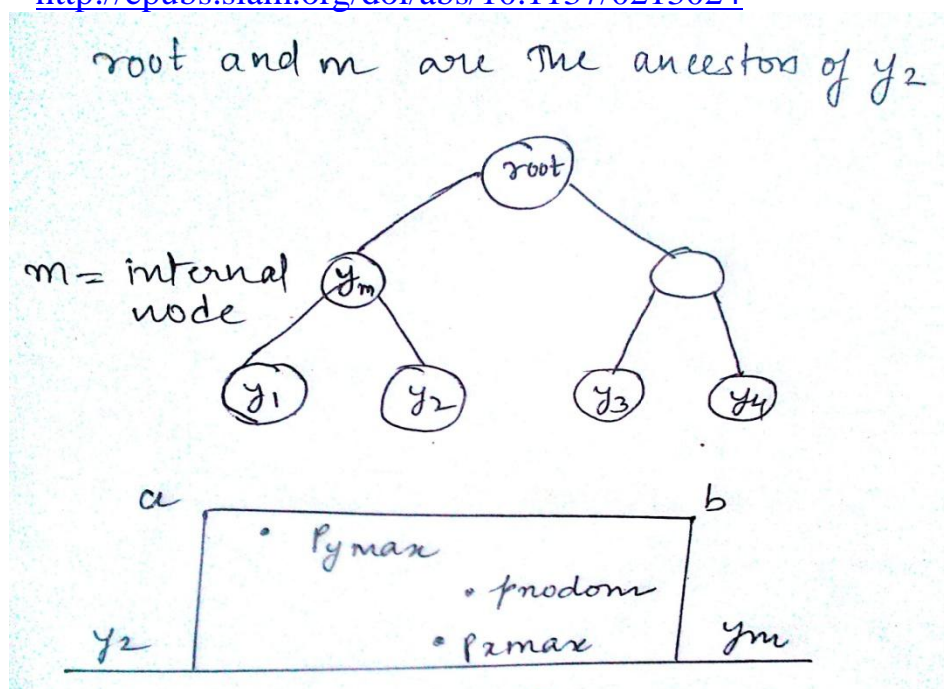
$count'm(y1) = |p'xmax, p'ymax|$.

Store it in a tuple $(count'm(y1), p'xmax)$.

(d) Special Case: If there are no points in $R2$, store $(0, NULL)$.

5. Maintain the data structure (*) in order to find the least common ancestor for two given leaf of the tree T_y in $O(1)$ time.

* <http://epubs.siam.org/doi/abs/10.1137/0213024>



LEMMA 1:

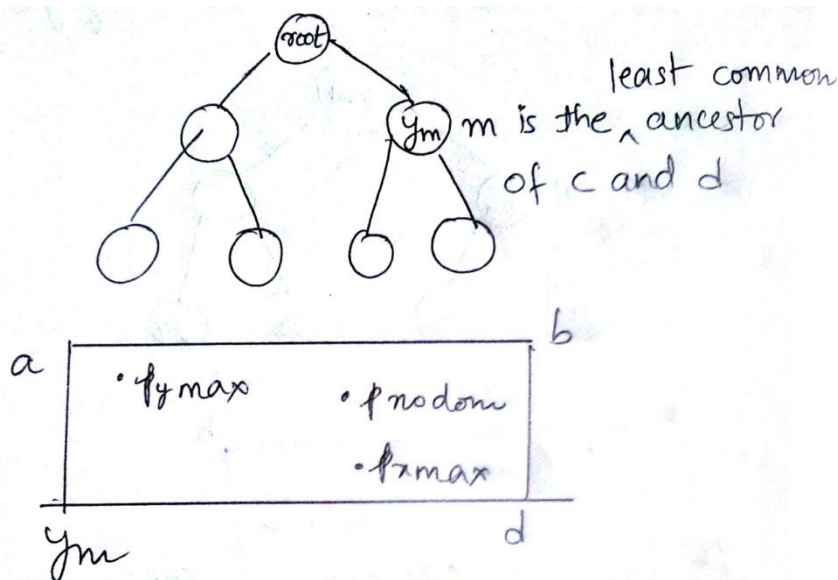
The total storage space needed by the data structure is $O(n \log(1+2\rho)n)$ words.

QUERY ALGORITHM:

1. Given a query rectangle $[a, b] \times [c, d]$ such that (a, b) belongs to $[1, \log n] \times [1, \log n]$ for $0 < \rho \leq 1/2$ and (c, d) belong to $[1, n] \times [1, n]$.

Find the least common ancestor m for the values d, c in the tree T_y . Let the key value stored at m be y_m .

2. Consider the rectangle $[a, b] \times [y_m, d]$ and the corresponding tuple $(count_m(d), p_{nodom})$. Where $count_m(d)$ has the same processor as $count_m(y_2)$.



3. Pseudo code

get_maximal_chain(Rectangle)

find pxmax and pymax and return the count of points between pxmax and pymax.

If ($pnodom \neq NULL$)

$pnodom(y) \Rightarrow$ y-coordinate of the point $pnodom$.

If ($c \leq pnodom(y)$)

Rectangle $\Rightarrow [a, b] \times [pnodom(y), ym]$.

Rectangle $\Rightarrow [a, b] \times [c, ym]$.

$count'm(pnodom(y)) \Rightarrow$ get_maximal_chain($[a, b] \times [pnodom(y), ym]$).

$count'm(c) \Rightarrow$ get_maximal_chain ($[a, b] \times [c, ym]$)

return $count'm(c) - count'm(pnodom(y)) + 1 + countm(d)$.

else if($pnodom(y) < c$)

return $countm(d)$.

else if ($pnodom=NULL$)

if ($countm(d) \neq 0$)

return $countm(d)$.

else

return $count'm(c)$.

// Special Cases

if($count'm(c)=0$) return $countm(d)$.

if ($countm(d)=0$) return $count'm(c)$.

if ($count'm(c) \neq 0$ and $pnodom= NULL$) return $countm(d)$.

PROBLEM 3:

Given an unsorted array A of n integers from the range of $[0, \sqrt{\log n}-1]$, preprocess A into a linear space data structure such that given two indices i, j and two values, a, b : (a, b) belongs to $[0, \sqrt{\log n}-1] \times [0, \sqrt{\log n}-1]$, we can efficiently report the smallest value t belongs to $A[i, j]$ for $a \leq t \leq b$.

This is clearly a variant of range successor problem and we can solve the problem as follows:

PREPROCESSING:

1. For $A[i]=x$, we create a 2D point (x, i) . Let S' be the set of such points.

2. We preprocess S into a data structure denoted by RS .

3. RS is an instance of the data structure of Theorem 4.

4. From $[*]$, RS queries of the form $[x_1, \infty) \times [y_1, y_2]$:

x_1 belongs to $[1, \log n / \log \log n]$ and (y_1, y_2) belongs to $[1, n] \times [1, n]$ queried in $O(1)$ time using a linear space data structure for a set of n 2D points with coordinates from a rectangle of $[1, \log n / \log \log n] \times [1, n]$.

In the given query, for any point (x, i) belongs to S' , x belongs to $[0, \sqrt{\log n}-1]$ and i belongs to $[1, n]$.

Theorem 4:

If the coordinates of the points are on an integer grid of $[1, \log n / \log \log n] \times [1, n]$, then there exists a linear space data structure that can be queried in $O(1)$ time to report the point with smallest x -coordinate in $[a, \infty) \times [c, d]$ where a belongs to

$[1, \log n / \log \log n] \times [1, \log n / \log \log n]$ and (c, d) belongs to $[1, n] \times [1, n]$.

RS: Range Successor

[*]<http://www.sciencedirect.com/science/article/pii/S0925772110000696>

QUERY ALGORITHM:

1. Given the two indices $i, j : (i, j)$ belongs to $[1, n] \times [1, n]$ and two values $a, b : (a, b)$ belongs to $[0, \sqrt{\log n} - 1] \times [0, \sqrt{\log n} - 1]$, we construct a three sided query rectangle of the form $[a, \infty) \times [i, j]$.

2. We run the range successor query on the data structure RS and find the smallest value t in $[a, \infty) \times [i, j]$.

If $(t \leq b)$ return t

else return *null*.

PROBLEM 1 SOLUTION

$R = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$

1. Construct a tree T_x , the leaf nodes are at the same height.
2. The leaf nodes store the x-coordinates of the points in the set R in non-decreasing order.
3. Each internal node μ belong to T_x has $O(\sqrt{\log n})$ children. Left most child numbered as $\sqrt{\log n} - 1$ and right most child numbered as 0.
4. μ is an internal node. $\text{int}(\mu) = \text{union of } x_1, x_2, x_3, \dots$. The values at the leaf node subrooted at μ .
5. For the internal node except root:

(a) create an auxillary array $A(\mu)$ which stores the y coordinates of the x coordinates subrooted at μ in decreasing order.

$A(\mu) = \{y_1, y_2, y_3, \dots\}$ where $y_1 > y_2 > y_3 > \dots$ and $A(\mu) = U(A(i))$: $i=0, 1, \dots, \sqrt{\log n} - 1$ where U is the union and v_i is a child of node μ .

(b) Each element of $A(i)$, $i=0, 1, \dots, \sqrt{\log n} - 1$ will point to its corresponding position in the array $A(\mu)$.

(c) A data structure $D(\mu)$ for $A(\mu)$: $A(\mu)[j] = y_j$.

Create a 2D point (v_i, y_j) where y_j belongs to $A(i)$, $A(i)$ = an auxillary array of the child v_i of the node μ and v_i belongs to $[0, \sqrt{\log n} - 1]$. These points belong to $[0, \sqrt{\log n} - 1] \times [1, n]$.

(d) Each child v_i of μ maintains a binary string lookup of size $|A(\mu)|$ where $A(\mu)[z] = 1$ if belongs to $A(i)$.

(e) lookup supports $\text{rank}()$ and $\text{select}()$ queries

Rank query: A rank query denoted by $\text{rank}_c(S, i)$ reports the number of c in S from position 1 to position i . Let $S = 222111331$. The query $\text{rank}_2(S, 4) = 3$

Select query: A select query $\text{select}_c(S, j)$ returns the position of the j th occurrence of the character c in S . Let $S = 222111331$ The query $\text{select}_2(S, 2) = 2$.

(f) $\text{RMA}(\mu)$, a range maxima data structure: return the maximum x coordinates whose y coordinates are between $A(\mu)[i]$ to $A(\mu)[j]$.

(g) For the values of $A(\mu)$:

$VT(\mu) \Rightarrow$ van Emde Boas tree and $T(\mu, y) \Rightarrow$ height balanced binary tree for any node n of $T(\mu, y)$ stores the median of the values stored in the leaf nodes.

4. For the root :

(a) Each index i of the array at the root node $A(root)$ will have $2 \sqrt{\log n}$ pointers of which $\sqrt{\log n}$ pointers will be pointing to the smallest elements greater than $A(root[i])$ in each of the arrays $A[j]$ for j being a child of the root node.

(b) The other $\sqrt{\log n}$ pointers will be pointing to the largest elements greater than $A(root[i])$ in the arrays $A[j]$ for j being a child of the root node.

(c) Construct a range maximum data structure $RMA(root)$:return the maximum x coordinates whose y coordinates are between $A(root[i])$ to $A(root[j])$.

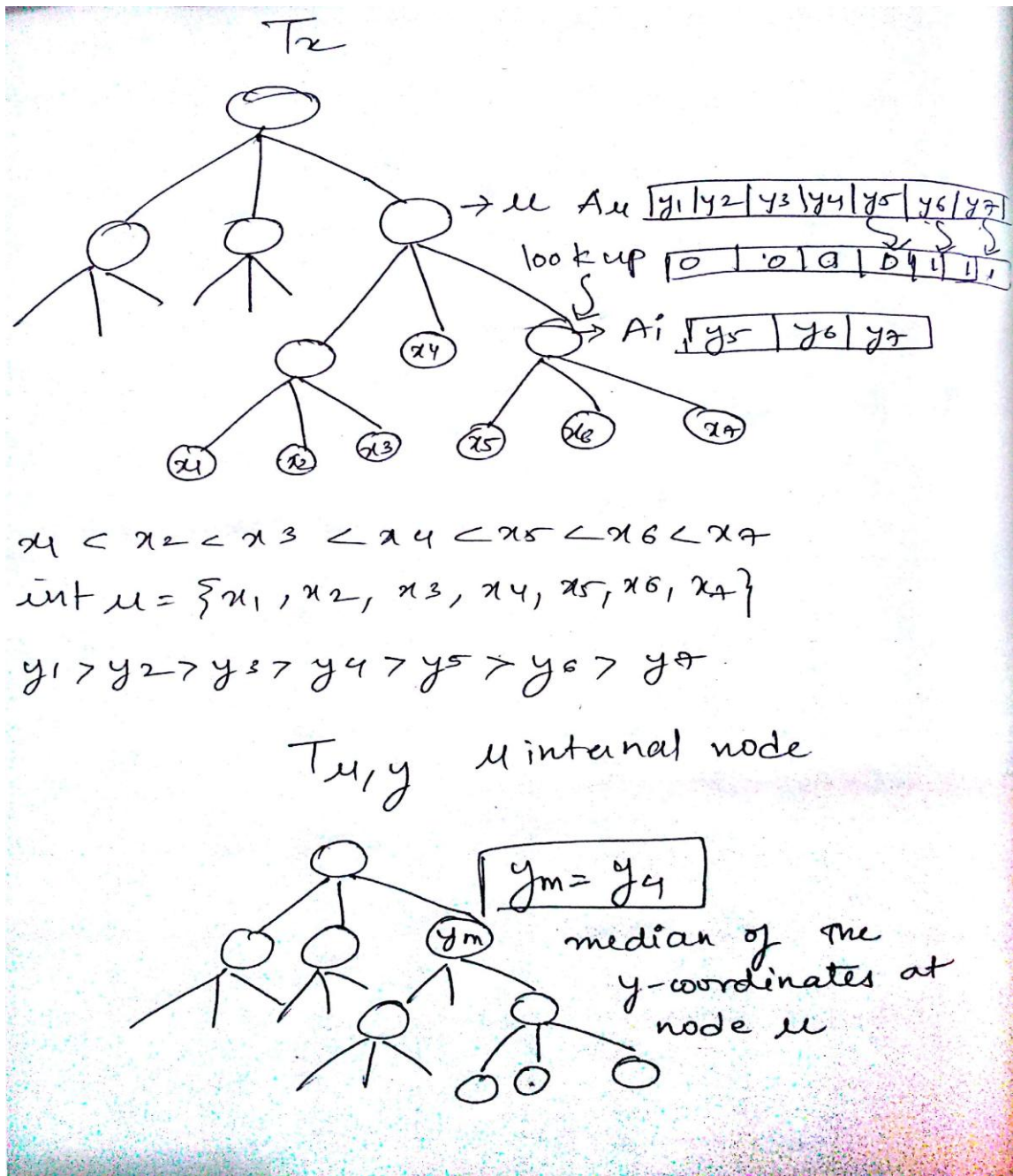
5. At each internal node μ belongs to Tx :

S' be the set of children of the node μ where $children = v_i, \dots, v_j$.

(a) $R1 = [i, j] \times [y, y_m]$, $y_m > y$ or $R1 = [i, j] \times [y_m, y]$, $y_m < y$ where y_m is stored in tree $T(\mu, y)$ for each element y belong to $A(\mu)$.

(b) $Find_maximal_chain(Rectangle\ R1, set\ S')$ denote it $ptop$ and $pbottom$ and count the maximal points.

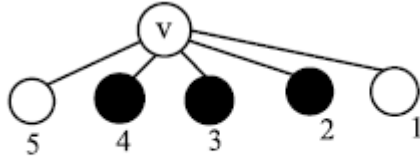
(c) $Tuple(R1) = (ptop, pbottom, count)$ for Rectangle $R1$



Lemma2 : The total storage space needed by the data structure for the counting problem will be $O(n \log^3 n / \log \log n)$.

DECOMPOSITION OF RECTANGLE

1. Given a rectangle $[a,b] \times [c,d]$, the segment $[a, b]$ is allocated to a node μ in T_x . If $\text{int}(\mu)$ subset of $[a,b]$ but $\text{int}(\text{parent}(\mu))$ doesnot belong to $[a, b]$. There will be $O(\log^{3/2}/\log\log n)$ such points. The set of such nodes be V .



. The segment $[a, b]$ is allocated to the nodes colored black.

2. These are canonical points and group it as $GL(1)$, .. $GL(\log n / \log \log n)$ and GL is group leader.

3. To count the maximal point, we decompose the rectangle into smaller rectangle. The number of rectangles will be $O(\log n / \log \log n)$.
Rightmost Rectangle: $R1$ and Leftmost Rectangle: R_z .
All the rectangles have same height.

For width:

Let v_i be the rightmost point in T_x among all the nodes in V then

(a) Consider $GL(1)$ for which the node v_i is a child and $g(1)$ consists of the children of $GL(1)$.

$R1$ width interval = $U(\text{int}(v_j))$ where U is the union. Iff v_j is a child of $GL(1)$ and $\text{int}(v_j)$ is a subset of $[a, b]$ or v_j belongs to $g(1)$.

(b) $V' = V - g(1)$ for this v_k be the rightmost node. Consider $GL(2)$ and set $g(2)$ for the node v_k . Width of $R2 = U(\text{intervals of the nodes in } g(2))$ where U is the union.

(c) Thus we can calculate the width of all the rectangle.

COUNTING ALGORITHM

Tx tree

$lca = \text{Find_Least_Common_Ancestor}(a, b)$

$\text{Get}(A(lca))$

// array of lca

$A(lca)[i]$ (smallest value) $\geq c$ and $A(lca)[j]$ (largest value) $\leq d$

//use van Emde Boas tree at the node lca.

$GL(1) \Rightarrow$ leader node subtree rooted at the child node vm for node lca .

$t = \text{number_of_ones_lookup_till_}(i-1)\text{th msb.}$

$\text{select}(1)(A(vm), (t+1)) = \text{position } z$

//use select operation

The largest value $\leq d = \text{position } z$

find $y'2 \leq d$ and $y'1 \geq c$ in $A(GL(1))$.

The children of $GL(1)$ are in the set $g(1)$ // v_i to $v_{\sqrt{\log n} - 1}$.

$\text{Search}(m \text{ in } D(GL(1)))$ m is the smallest index such that the node vm has the rightmost maximal point in the rect $[a, b] \times [y'1, y'2]$. run RS query to find m .

$\text{Total1} = \text{Count_maximal_chain}$ in the rectangle $[m, \sqrt{\log n} - 1] \times [y'1, y'2]$

$\text{pymax} = (\text{pymax}(x), \text{pymax}(y))$

$GL(2) \Rightarrow$ leader node and $g(2) = vw, \dots, vq$

$\text{Total2} = \text{Count_maximal_chain}$ in the rectangle $[w, q] \times [\text{pymax}(y), d]$

$GL(2)$ satisfy one of the following:

(a) $GL(2)$ is in the path from $GL(1)$ to lca .

$A(GL(1))$ has a pointer to the parent of $GL(1)$ array.

$y'2 \leq d$ in $A(GL(2))$ is easy to find.

(b) $GL(2)$ is in the path from lca to $GL(1)$.

First move to the node lca by following pointers from $GL(1)$

Then proceed with (a)

(c) $GL(2)$ is the lca

Move to the node lca by following pointers from $GL(1)$.

Repeat for all the group leader nodes

Return $\text{Total} = \sum i \text{ Total}_i: i=1, \dots, O(\log n / \log \log n)$.

QUERY TIME ANALYSIS AND CORRECTION PROOF

The query algorithm correctly counts the number of maximal points inside the query rectangle.

Divide the problem of reporting maximal points for a query rectangle in a general setting to $O(\log n / \log \log n)$ instances of the problem of counting maximal points for a query rectangle (as per problem 2).

By Lemma 1, the correct count of maximal points can be calculated.

Thus, the query algorithm takes $O(\log n / \log \log n)$ time to count the number of maximal points in the query rectangle.

The storage space can be improved without sacrificing the query time efficiency by reducing the degree of each internal node of the tree to $\log^{(p)} n$ for $0 < p < 12$.

CONCLUSION

According to the theorem and analysis of the algorithm, we can implement this algorithm in many practical areas of database management system for speeding the query time. We have also learnt that we can reduce the storage capacity without compromising the query time.

For future work, we can increase the dimension i.e experiment with 3D coordinates and analyze the time complexity of the algorithm.

Here is the link of the ongoing code:

<https://drive.google.com/file/d/0B6O5O3fDqH3Ld1JWX21YNmF1SzQ/view?usp=sharing>

REFERENCES

- <http://epubs.siam.org/doi/abs/10.1137/0213024>
- https://en.wikipedia.org/wiki/Van_Emde_Boas_tree
- [Books](#)
- <http://www.sciencedirect.com/science/article/pii/S0925772110000696>
- [http://refhub.elsevier.com/S1570-8667\(14\)00095-1/bib7975s1](http://refhub.elsevier.com/S1570-8667(14)00095-1/bib7975s1)