

# Database Implementation System

## Assignment 3

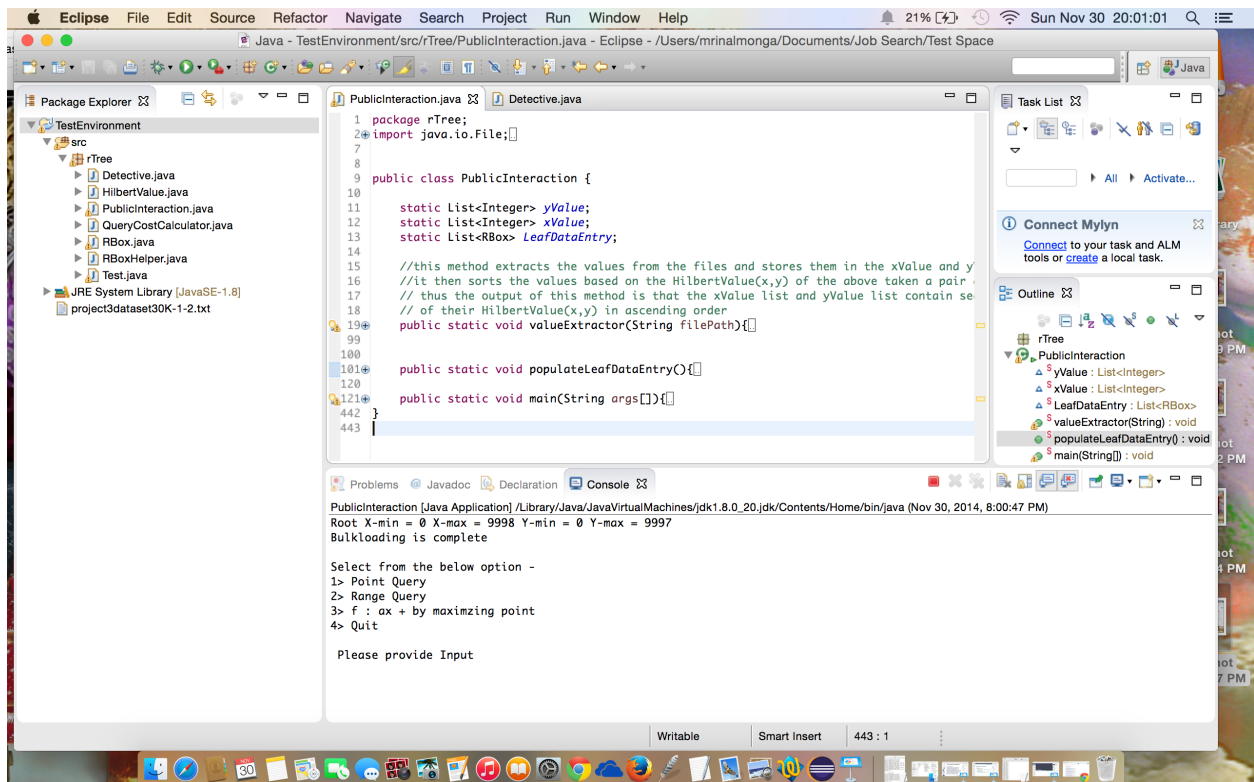
### Mrinal Monga

#### Minimum System Requirements

1. Free RAM space of at least 60 MB
2. Developed and Tested using **Java 8 SDK**

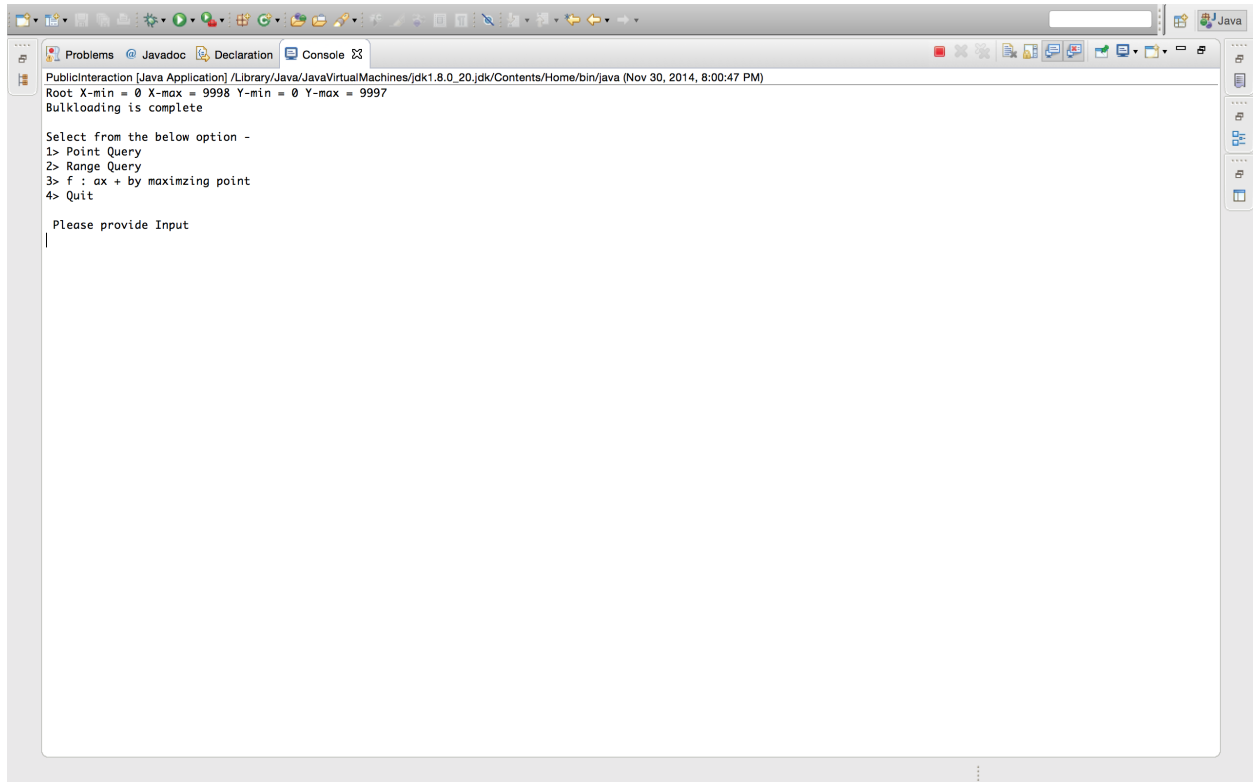
#### Compilation Instructions -

1. The zip folder is a java project. Use the java project to run the program. The class which must be run is the **PublicInteraction.java** file. This class is responsible for interacting with the user. It is also responsible for bulkloading R-Tree. Please take great care to note the location of the file project3dataset30K-1-2.txt which contains the dataset. **Please carefully note the location of project3dataset30K1-2.txt with respect to the package rTree.**
2. Once you have a Java Project in the Eclipse as shown in the image below run the file PublicInteraction.java



# User Interaction

1. The following screen is seen



2. Note carefully the options it provides. 1 for Point Query 2 for Range Query and so on.
3. Some sample inputs outputs are given.

## POINT QUERY

The screenshot shows the Eclipse IDE interface. At the top, there are tabs for "Problems", "Javadoc", "Declaration", and "Console". The "Console" tab is active, displaying the output of a Java application. The application prompts the user to "Please provide Input". The user enters "2335" for X and "500" for Y. The program then searches for rectangles containing the point (2335, 500). It lists several rectangles found, such as those with corners at (0, 0), (998, 4093), (2096, 4094), etc. Finally, it outputs the result of the point query, stating that the point was found within one of the rectangles. The console window has a scroll bar on the right side. The background of the IDE is dark gray, and the text in the console is white. The title bar of the window indicates it's running on a Mac OS. The overall layout is typical of a professional development environment.

Note only if a data point is found then FOUND key word along with the Data Entry details is shown. The DESCRIPTION is also shown. If data point is not found then FOUND Key word followed by the Data Entry details will not be observed.

## RANGE QUERY

```
PublicInteraction [Java Application] /Library/Java/JavaVirtualMachines/dk1.8.0_20.jdk/Contents/Home/bin/java (Nov 30, 2014, 8:16:04 PM)
Root X-min = 0 X-max = 9998 Y-min = 0 Y-max = 9997
Bulkloading is complete

Select from the below option -
1> Point Query
2> Range Query
3> f : ax + by maximizing point
4> Quit

Please provide Input
2
Enter x-min. Please ensure x-min <= x-max
2300
Enter x-max. Please ensure x-max >= x-min
2400
Enter y-min. Please ensure y-min <= y-max
500
Enter y-max. Please ensure y-max >= y-min
700
Now Searching for - X-min = 2300 X-max = 2400 Y-min = 500 Y-max = 700
HIT ENTER TO SEARCH ...

-----
Now search will be conducted. If the Data Points are FOUND
then the FOUND data points will be displayed with the details
The DESCRIPTION will also be displayed
Number of pages is dispalyed.

-----

FOUND the Data Entry X-min = 2335 X-max = 2335 Y-min = 500 Y-max = 500
DESCRIPTION of the Data Entry aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa

-----

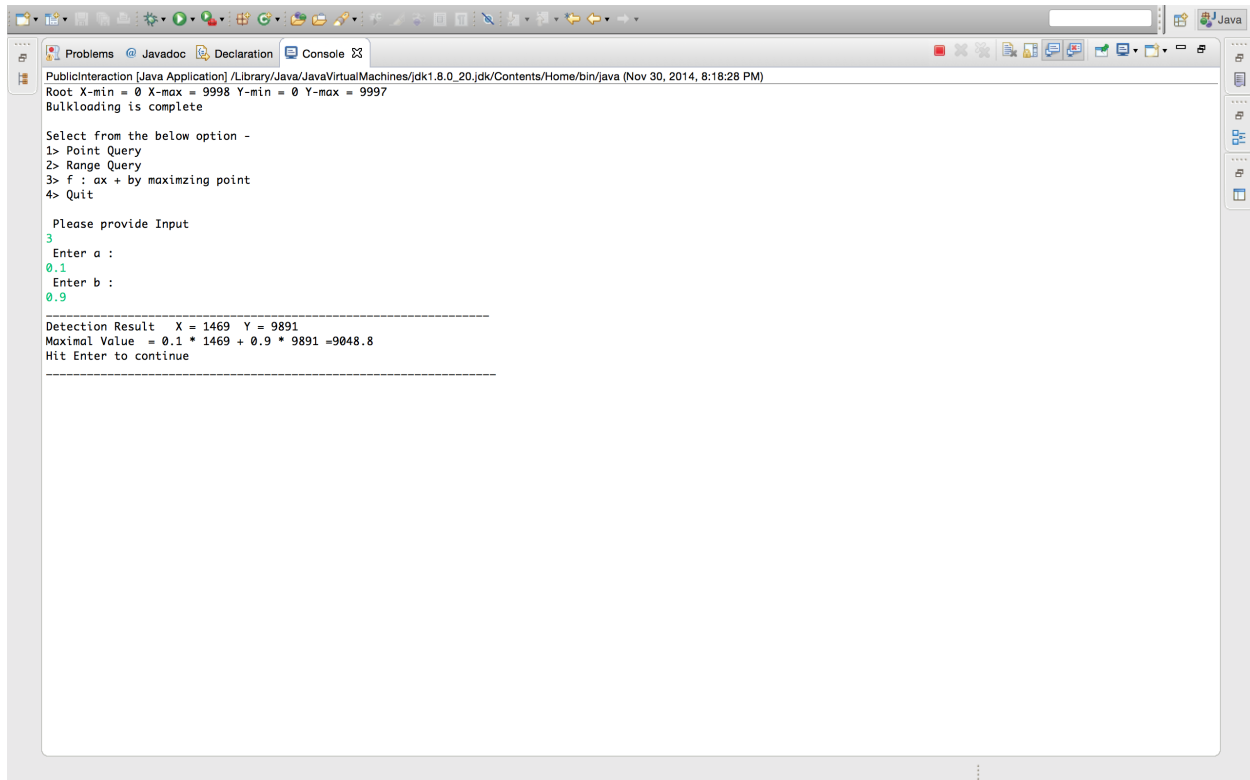
FOUND the Data Entry X-min = 2305 X-max = 2305 Y-min = 550 Y-max = 550
DESCRIPTION of the Data Entry aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa

-----

OUTPUT of Range Query
The found points will have FOUND key word at the beginning. If FOUND key word is NOT in the above messages than no point could be located
Search Complete
Number of Pages Accessed = 12
Hit Enter to continue
```

Note - Here only those points are shown which exist in the given range. If no point is found then the FOUND key word followed by Data Entry details will not be observed.

f:ax + by QUERY



```
PublicInteraction [Java Application] /Library/Java/JavaVirtualMachines/jdk1.8.0_20.jdk/Contents/Home/bin/java (Nov 30, 2014, 8:18:28 PM)
Root X-min = 0 X-max = 9998 Y-min = 0 Y-max = 9997
Bulkloading is complete

Select from the below option -
1> Point Query
2> Range Query
3> f : ax + by maximizing point
4> Quit

Please provide Input
3
Enter a :
0.1
Enter b :
0.9

-----
Detection Result  X = 1469  Y = 9891
Maximal Value   = 0.1 * 1469 + 0.9 * 9891 =9048.8
Hit Enter to continue
-----
```

## RTree Details

1. RBox is the class which represents a rectangle.
2. All rectangles in this program are RBox objects.
3. Data Entry Rectangles contain description of 500 char. This variable is null for all rectangles other than the Data Entry Rectangles. This implies that the root, index nodes and leaf nodes will have description as null. The Data Entry Rectangles are given the description while they are bulkloaded.
4. The RBox has List<RBox> Entry which points to child rectangles or RBox objects of a given RBox object.
5. For Leaf Rectangles the Entry points to data entries.
6. For Index Rectangles the Entry points to child RBox objects.
7. For Data Entry Rectangles Entry is null;
8. All data is held in the main memory.
9. Size of data entry =  $4 * 4 \text{ bytes} + 500 * 2 \text{ bytes} = 1016 \text{ bytes}$   
 $= 4 * \text{size of an int} + 500 * \text{size of char}$   
Therefore 1 page of size 4096 byte has  $[4096/1016] = 4$  data entries  
Therefore there are 7500 Leaf Rectangles
10. Size of an Index Entry =  $4 * \text{int} + 4 \text{ byte for object reference or pointer}$   
Similarly it is seen that 1 Index Node has  $[4096 / 20]$  or 204 entries

There are 37 Index Rectangle

11. Thus in this implementation every thing is in the main memory.

## Class Description

1. The name of the classes used for the implementation of R-Tree

(1) **RBox** - The RBox class represents a class of the Rectangle Box. These are used to create rectangles. Note that due to the nature of R-Tree one rectangle can contain as its children many other rectangles. Also the data entries are also rectangles.

(2) The data members of the RBox class are -

```
int xmax;  
int xmin;  
int ymax;  
int ymin;
```

```
public char[] description; // This is used only by a Data Entry.  
                          Not used by any other type of RBox.
```

```
List<RBox> Entry = null; // for data entry its value is null
```

```
List<RBox> listOfPotentialRBoxes = null;
```

(3) Note that as the name suggests xmin, xmax, ymin, ymax contain the the minimum and maximum x and y values respectively corresponding to the rectangle. The char[] description is a data member which will ONLY be assigned value for the data entries of the leaf nodes. char[] description is null for all other types of rectangles.

(4) List<RBox> Entry is the array list which holds the list of children rectangles pointed to by a rectangle . For DATA ENTRY rectangles Entry would be null as they don't point to anything. The LEAF NODES would point to the DATA ENTRY rectangles. They INDEX NODES would point to other INDEX NODES, LEAF NODES.

(5) The List<RBox> listOfPotentialRBoxes is an array list which is used by the **pSearch(RBox b, QueryCostCalculator QueryCost)** and the **public rSearch(RBox b, QueryCostCalculator QueryCost)** methods and used for the Point and Range Queries respectively.

(6) The **public void dimensionCalculation()** method is used to assign the values to the xmax, xmin, ymax, ymin variables. This method takes into consideration the xmax, xmin, ymax, ymin values of the rectangles in the Entry array list.

(7) The **public void pSearch(RBox b, QueryCostCalculator QueryCost)** method is used for the point queries. When the user inputs a point, it is converted into an RBox object - that is a rectangle. This rectangle RBox b is sent to the method along with a QueryCostCalculator object which helps to keep track of the pages accessed during the search.

This method first check whether it is inside a rectangle which is a Leaf Data Entry or not. If it is not a Data Entry rectangle then it prepares a list of potentialRectangles which is basically those rectangles which are children of the current rectangle and which overlap the RBox -

It displays as output all the nodes that it visits and if it is able to reach inside a Data Entry Rectangle which is the input being searched for, then it displays it with FOUND key word. If for a query no FOUND key word is in the output then that point couldn't be found. It also mentions the **Number of Pages Accessed** at the bottom.

```
Java - TestEnvironment/src/rTree/PublicInteraction.java - Eclipse - /Users/mrinalmanga/Documents/Job Search/Test Space  
PublicInteraction [Java Application] /Library/Java/JavaVirtualMachines/jdk1.8.0_20.jdk/Contents/Home/bin/java (Nov 30, 2014, 4:09:05 PM)  
Please provide Input  
1  
Enter X =  
2335  
Enter Y =  
500  
  
Now Searching for - X-min = 2335 X-max = 2335 Y-min = 500 Y-max = 500  
HIT ENTER TO SEARCH ...  
  
-----  
Now search will be conducted. Note that all the Rectangles being  
traversed will be displayed. If the Entered Data Point is FOUND  
then the FOUND will be displayed with the details of the rectangle  
corresponding to the Data Point. The DESCRIPTION will also be displayed  
Number of pages is dispalyed.  
  
-----  
Now Searching   inside RBoxX-min = 0 X-max = 9998 Y-min = 0 Y-max = 9997  
Now Searching   inside RBoxX-min = 2096 X-max = 4094 Y-min = 4 Y-max = 4093  
Now Searching   inside RBoxX-min = 2110 X-max = 4085 Y-min = 46 Y-max = 2141  
Now Searching   inside RBoxX-min = 2048 X-max = 3192 Y-min = 7 Y-max = 2047  
Now Searching   inside RBoxX-min = 2312 X-max = 2343 Y-min = 39 Y-max = 500  
  
FOUND the Data Entry : X-min = 2335 X-max = 2335 Y-min = 500 Y-max = 500  
DESCRIPTION of the Data Entry = aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa  
  
-----  
Now Searching   inside RBoxX-min = 2068 X-max = 2379 Y-min = 175 Y-max = 627  
Now Searching   inside RBoxX-min = 2069 X-max = 2998 Y-min = 27 Y-max = 1054  
Now Searching   inside RBoxX-min = 0 X-max = 4080 Y-min = 42 Y-max = 5999  
Now Searching   inside RBoxX-min = 2 X-max = 4080 Y-min = 42 Y-max = 4193  
Now Searching   inside RBoxX-min = 0 X-max = 8190 Y-min = 0 Y-max = 9997  
Now Searching   inside RBoxX-min = 30 X-max = 8176 Y-min = 80 Y-max = 8204  
  
OUTPUT result of Point Query  
The found points will have FOUND key word at the beginning. If FOUND key word is NOT in the above messages than no point could be located  
Search Complete  
Number of Pages Accessed = 12  
Hit Enter to continue
```

(8) The `public void rSearch(RBox b, QueryCostCalculator QueryCost)` method answers the Range Queries. It is slightly different in its logic then the method for answering point queries. In this method once the list of potential rectangles which could contain the required RBoxes-Rectangles is formed, then the entries of the leaf nodes of the subtree of the potential rectangles are checked those entries are only displayed which actually are found to be CONTAINED by the rectangle corresponding to the input.

### A Sample output -





Explanation - Every node in the R-Tree is a rectangle - here implemented by the R-Box. Each R-Box or rectangle can have children rectangles which may contain the points which could maximize the  $f:ax+by$ . So a method is there in the R-Box class and thus in each object of the class called the **public** `Detective maximizingEntryDetector(Detective detective)` method.

This method is both the point of insertion and point of extraction of the of the Detective probe into the R-Tree rooted at that R-Box or the rectangle under consideration. Once the Detective probe is initialized with  $a$  and  $b$  from the user, the R-Tree's root RBox or rectangle is asked to consume/ingest this Detective probe. The root then runs an algorithm explained below and passes this Detective probe to its such children rectangles or RBoxes which most likely contain the maximizing point. Each of the children rectangles then subsequently follow the same algorithm to identify their most likely descendant children and the process goes on till Leaf Rectangle is reached.

The Detective probe on reaching a Leaf Rectangle will scan the entire list of entries to identify such a Data Entry which maximizes the  $f:ax+by$ . It records the  $x$ ,  $y$  and `maximalValue`. Subsequently when the Detective probe visits other leaf nodes it behaves similarly, **but** it would now record the  $x,y$  and maximal value  $ax+by$  **if and only if** the new maximal value found is **bigger** then the previous maximal value. So the Detective node in effect traverses the likely leaf rectangles based upon the algorithm, however at any point in time it would only have the  $x,y$  and maximal value corresponding to such a Data Entry which had maximized  $f:ax + by$  till that point in execution. Also once it finds that for a particular sub-tree the value  $f:ax+by$  is not greater than the `maximalValue` it has found till that point in execution, it would no longer go down that sub-tree. This is taken care of by the algorithm explained below.

Subsequently the Detective probe is returned by the same method **public** `Detective maximizingEntryDetector(Detective detective)` method at which it was inserted to into the root of the R-Tree, **where it was ingested by the R-Tree by for the very first time**. After extraction the Detective probe is examined to see what it found inside the R-Tree.

**Algorithm which runs inside the **public** `Detective maximizingEntryDetector(Detective detective)` method**

- 1. Check if the existing RBox is a Leaf Node**
- 2. If it is not a Leaf Node then -**

Extract the  $a$ ,  $b$  and `maximalValue` (found till this point in execution) from the Detective probe.

For the existing RBox find the value of  $ax+by$  by analyzing the centroids of all the children rectangles. Every RBox has a **public float** `centroidX()` and **public float** `centroidY()` method. These methods return the  $X$  and  $Y$  coordinates of the centroid of the rectangle.

For a child RBox - rectangle its centroid's coordinates are put in the  $f:ax+by$  to find the value.

Scan all the children rectangles **to observe and find out if at all** any of the child rectangles or RBoxes gives a value of  $f: ax+by$  greater than that found till this point in execution which is stored in Detective probe's variable `maximalValue`.

This is implemented by substituting the centroid of a RBox or rectangle in the  $f:ax+by$  and comparing it with the “maximalValue” found till now, which was extracted from the Detective probe. The detective probe on its way from other Rectangles to the current rectangle stores the  $x,y$ , and the maximalValue found till this point in execution.

The algorithm is designed in such a manner so that if it was found that the new maximalValue found was greater than that stored inside the Detective probe then again scan all the children RBoxes - rectangles to find those rectangle for which the value obtained by putting their centroid's coordinates for  $x, y$  in  $f:ax+by$  gives a value equal to the new greater maximalValue found now.

The list of all such children rectangles the centroids of which when put in  $f:ax+by$  give the maximalValue found till now is prepared.

Thereafter again scan the entire list of children rectangles of the current Rectangle to find all such rectangles which **overlap** the list of rectangles the centroids of which give the maximalValue when substituted for  $x,y$  in the  $f:ax + by$ .

**The above step is done to ensure that the rectangle the centroid of which maximizes the  $f:ax+by$  may not contain that data entry in its sub-tree which actually maximizes  $f:ax+by$ . So all those rectangles which overlap the rectangle are also added to the search space.**

Prepare a list of all of the above rectangles, remove any redundancy which may be there on account of counting a rectangle more than once because of several scans.

Then send the Detective probe to every such rectangle from the list prepared.

**Note - We send the Detective probe down every possible sub-tree but the probe need not actually reach the Leaf Rectangles of each of the sub-tree as the algorithm takes care to let the probe go down further if and only if the value  $f:ax+by$  obtained by scanning the children rectangles or Rboxes of a rectangle is greater than the maximalValue found till that point in execution.**

**This means even if we send the Detective probe down every potential sub-tree it would actually go down only if the condition new-maximal-Value-found > maximalValue-found-till-now is met. Other wise it would NO LONGER CONTINUE down that subtree.**

### **3. If a Leaf Rectangle is reached**

Scan all the Data Entries of the leaf node to check whether the new value found for  $f:ax+by$  by substituting the  $x,y$  of the Data Entry exceeds the maximalValue found till this point in execution inside Detective probe. If it does exceed record the new  $x, y$ , and maximalValue found else don't store anything.

### **4. Running Time of the Algorithm -**

**In worst case - Height of tree levels would need to be traversed -  $\log(n)$  base - branching factor. maximum data entries which could ever be scanned in a tree =  $n$**

**Thus Running Time -  $O(n * \log(n))$  (ASSUME base = branching factor of R-Tree)**

**(11) Class RBoxHelper**

This class has two methods

One is `static boolean doTheseOverlap(RBox b1, RBox b2)` which tells whether two rectangles overlap. The other is `static boolean contains(RBox b1, RBox b2)` which checks whether one rectangle contains the other.

**(12) QueryCostCalculator Class**

The object of this class is used to keep track of the number of pages traversed .

**(13) PublicInteractionClass**

This interacts with the user. This is responsible for bulk loading of the tree and interacting with the user.