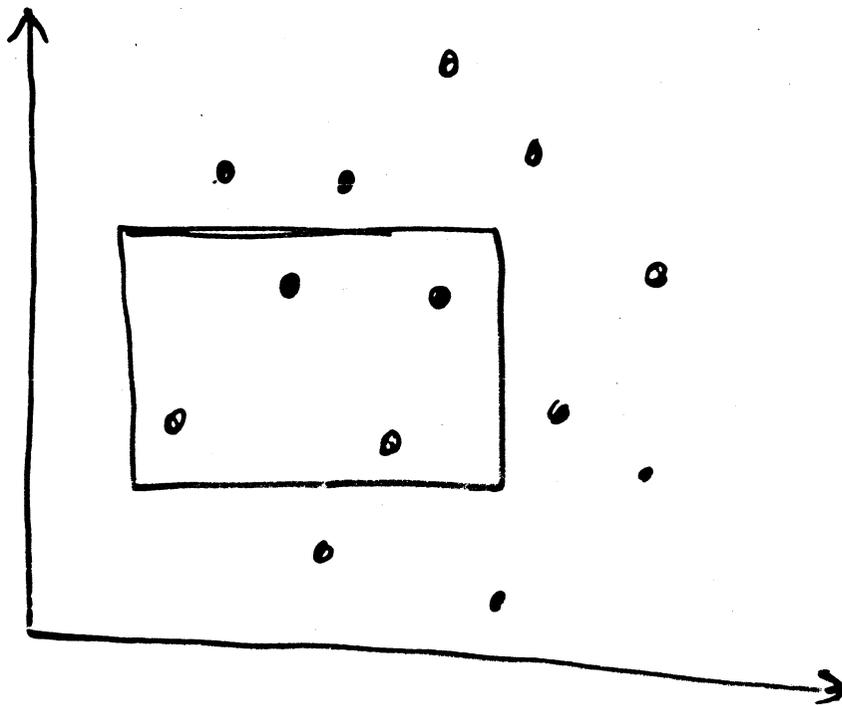


# Orthogonal Range Searching (Querying a database)

Week 4

- Data structure for a set of objects (points, rectangles, polygons) for efficient range queries.
  - Given a set of  $n$  points, build a data structure that for any query rectangle  $R$ , reports all points in  $R$ .



- Efficiency : depends on type of objects & queries.
- Time-space tradeoff : the more we preprocess the data & store, the faster we can solve a query
- Orthogonal range searching deals with points sets and axis-aligned rectangle queries.

## Kd-trees (Bentley)

- Not the most efficient solution in theory.
- Everyone uses it in practice.
- Algorithm:

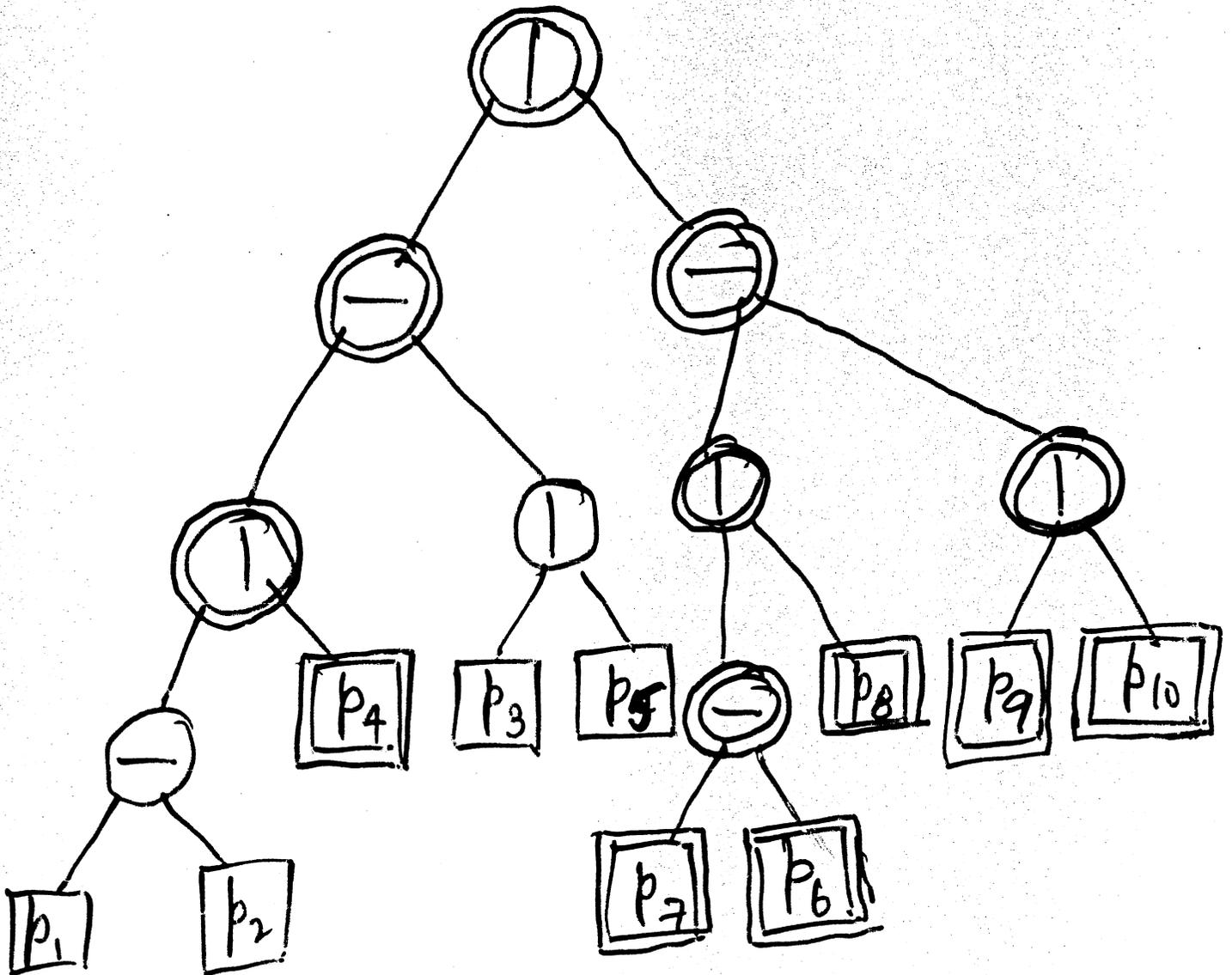
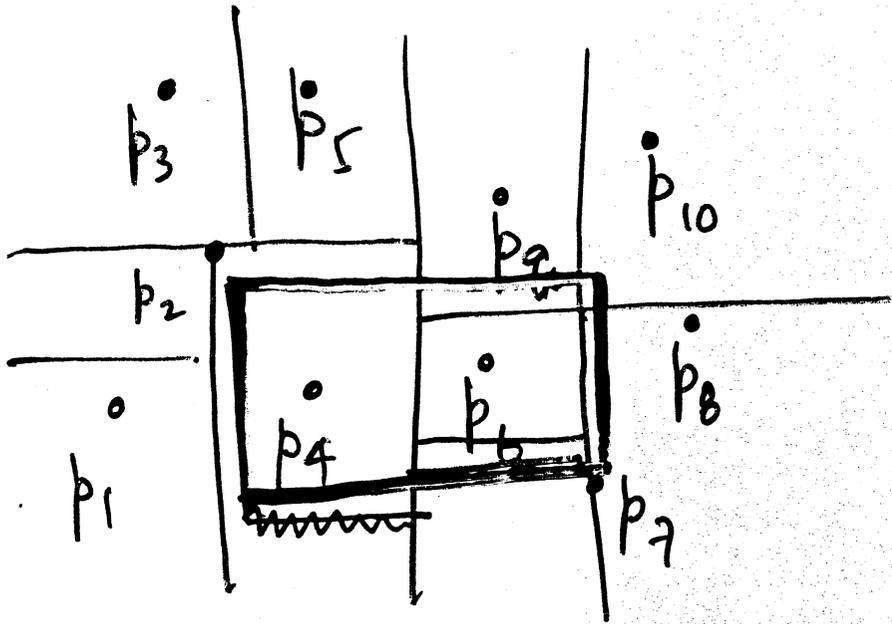
- Choose  $x$  or  $y$  coordinate  
(alternate)

- Choose the median of the coordinate.

(The point set gets split into two equal halves.)

This defines a horizontal or a vertical line.

- Recurse on both sides.



## Kd-tree : Range Queries

1. Recursive procedure, starting from  $v = \text{root}$ .

2. Search  $(v, R)$ :

a) If  $v$  is a leaf, report the pt. stored in  $v$  if it lies in  $R$ .

b) Otherwise, if  $\text{Region}(v)$  is contained in  $R$ , report all the points in the subtree of  $v$ .

c) Otherwise :

i) If  $\text{Region}(\text{left}(v))$  intersects  $R$ ,  
 $\text{search}(\text{left}(v), R)$

ii) If  $\text{Region}(\text{right}(v))$  intersects  $R$ ,  
 $\text{Search}(\text{right}(v), R)$



- When  $\text{Region}(v)$  is contained in  $R$ , the complexity is linear in output size.
- We just need to ~~know~~ bound the number of nodes  $v$  st.  $\text{Region}(v)$  intersects  $R$  but is not contained in  $R$ . In other words, the boundary of  $R$  intersects the boundary of  $\text{Region}(v)$ .
- We will make a gross overestimation. We will bound the number of  $\text{Region}(v)$  which are crossed by any one of the 4 horizontal/vertical lines of  $R$ .

- Let  $l$  be the line defining one side of  $R$ .
- How many ~~set~~ ~~Region~~  $\text{Region}(v)$  can  $l$  intersect?
- ~~Consider~~ Consider two levels at a time
- Suppose the first cut is vertical, & second horizontal. We have 4 regions, each with  $\frac{n}{4}$  points
- A line intersects exactly two <sup>regions</sup> ~~sets~~ (line is either vertical or horizontal)  
The other regions will be either outside or entirely inside  $R$ .

The recurrence is

$$Q(n) = \begin{cases} 1 & \text{if } n=1 \\ 2Q\left(\frac{n}{4}\right) + 2 & \text{otherwise} \end{cases}$$

$$\therefore Q(n) \in O(\sqrt{n}).$$

- ~~A~~ Kd-tree is an  $O(n)$  space data structure that solves the orthogonal range query in worst-case time  $O(\sqrt{n+k})$ , where  $k$  is the output size.

## Extending Kd-tree to higher dimensions?

- Try 3-dimensional case, & then generalize.

- The recurrence is

$$Q(n) = 2^{d-1} Q\left(\frac{n}{2^d}\right) + 2^{d-1}$$

which solves to  $Q(n) \in O(n^{1-\frac{1}{d}})$

( $d$  is fixed)

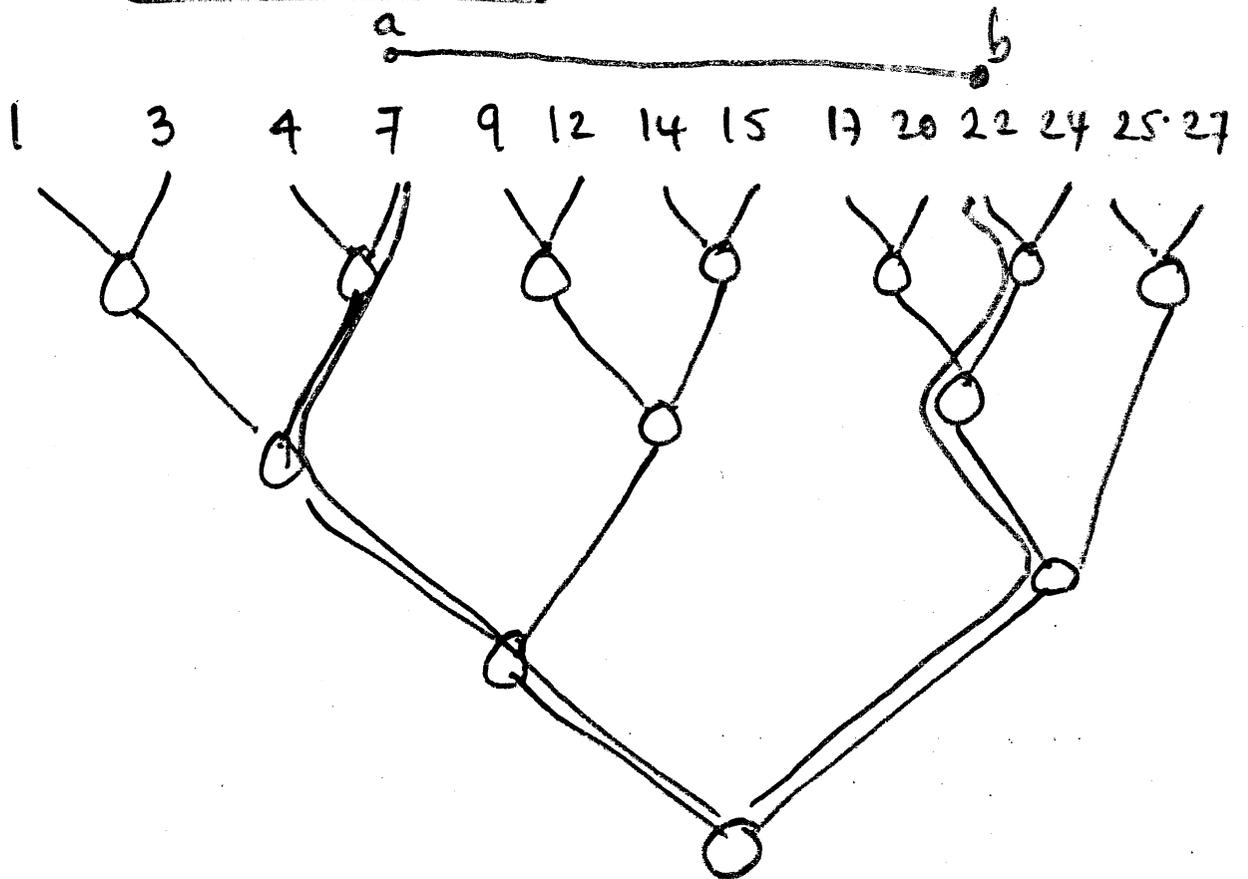
- Kd-tree is an  $O(dn)$  space data structure that solves  $d$ -dim range query in worst-case time  $O(n^{1-\frac{1}{d}} + k)$ , where  $k$  is the output size.

## 1-dimensional range search.

- Points on a line
- Queries are intervals.
- Preprocessing: Sort the points
- ~~Each~~ Each query can be answered in  $O(\log n + k)$  time where the interval contains  $k$  points.

(Does not extend to higher dimension.)

# Tree Search.



- build a balanced binary tree on the sorted list of points.
- given an interval  $[a, b]$ , search down the tree for  $a$  &  $b$ .
- all leaves between the two form the answer.

## Analysis

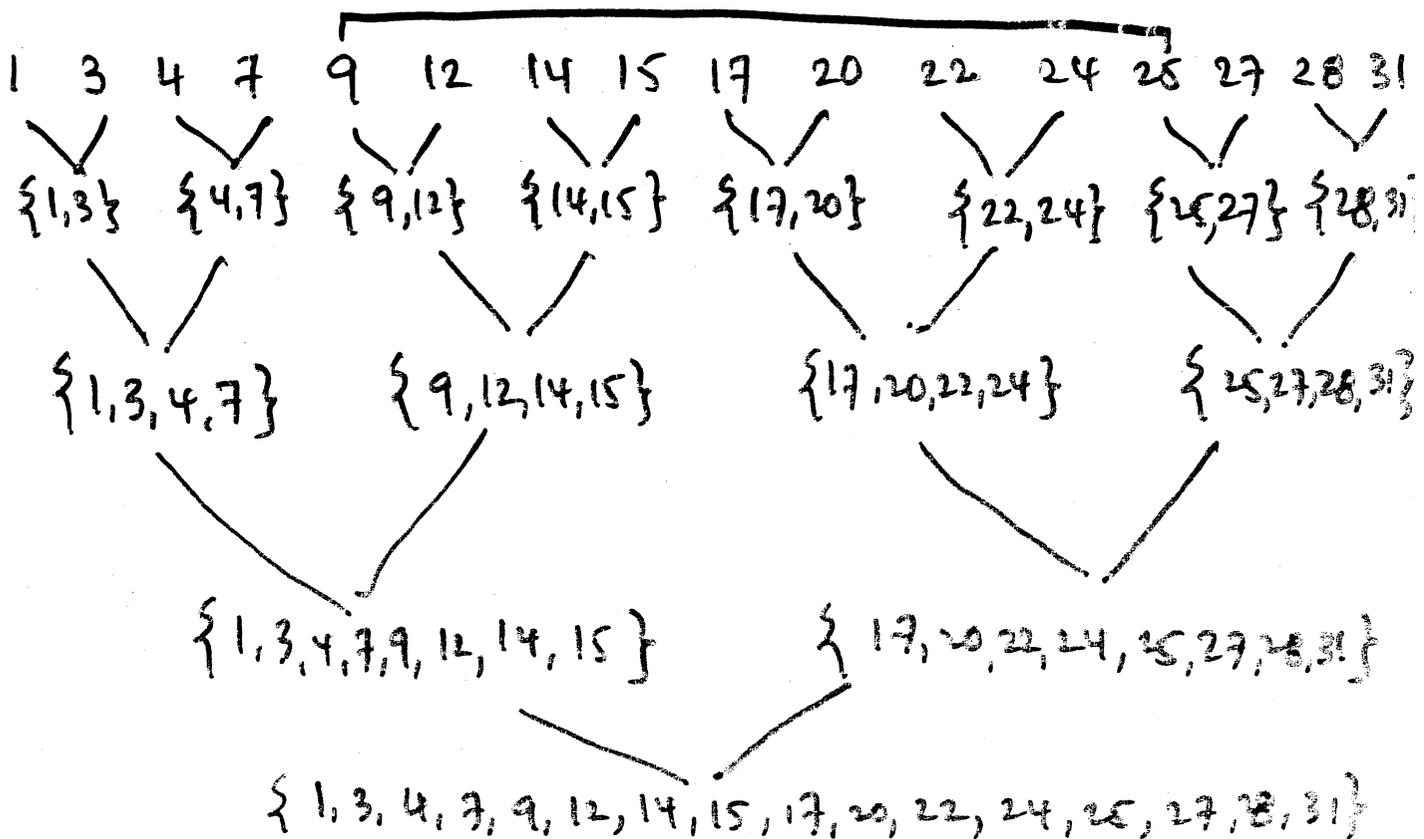
- Search time :  $O(\log n + k)$
- Preprocessing time :  $O(n \log n)$
- Storage space :  $O(n \log n)$   
(canonical sets)

## 1-d Range query

- Given query  $[a, b]$ , search down the tree for the leftmost leaf  $u \geq a$  and the ~~right~~ leftmost leaf  $v \geq b$ .
- All leaves between  $u$  &  $v$  are in the range.
- If  $u = a$  or  $v = b$ , include the canonical set stored at  $u$  &  $v$ .
- The remainder nodes are the union of at most  $2 \log n$  canonical sets.

# Canonical Subsets

- $S_1, S_2, \dots, S_E$  are Canonical Subsets  $S_i \subseteq P$ , if the answer to any range query can be written as the disjoint union of some  $S_i$ 's.
- The Canonical subsets may overlap.
- Only  $O(\log n)$  subsets are reported for any range query.



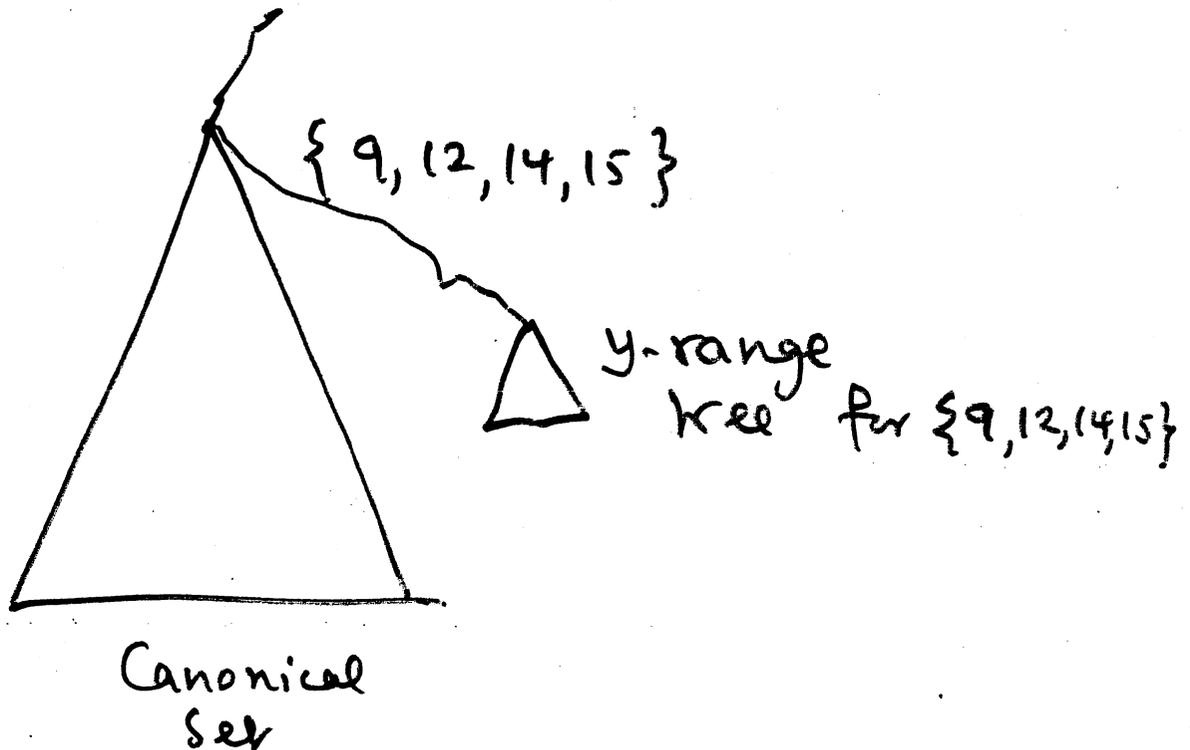
## Orthogonal Range Trees

- Generalizes 1-D search trees to dimension  $d$ .
- Each search recursively decomposes into multiple lower dimensional searches.
- Search complexity is  $O((\log n)^d + k)$ , where  $k$  is the answer size.
- Space and time complexity is  ~~$O((\log n)^d)$~~   $O(n(\log n)^{d-1})$ .
- Fractional Cascading eliminates one  $\log n$  factor from search time.
- We focus on 2-d. The generalization step is straightforward.



## Level 2 trees

- Collect points of each canonical set & build a y-range tree.



- Search each of the  $O(\log n)$  canonical sets that include points for x-range  $[x_l, x_r]$  using their y-range trees for range  $[y_l, y_r]$
- The y-range searches list out the points of  $P$  in the query rectangle (No duplicates)

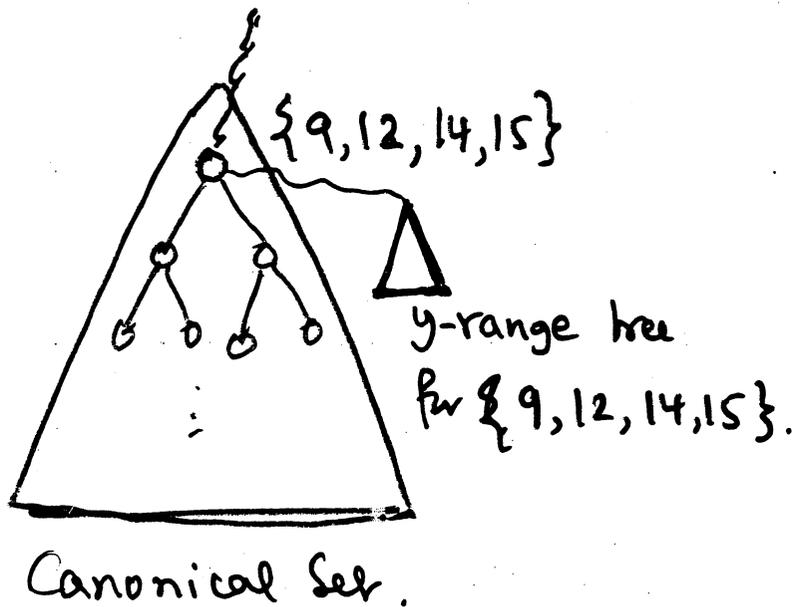
# Analysis

- Query time:  $O(\log^2 n)$ 
  - $O(\log n)$  x-range query
  - each canonical set y-range query takes  $O(\log n)$  time.
- Space complexity:  $O(n \log n)$ 
  - each point appears in  $O(\log n)$  canonical sets.
- Preprocessing complexity:  $O(n \log n)$



## Level 2 trees

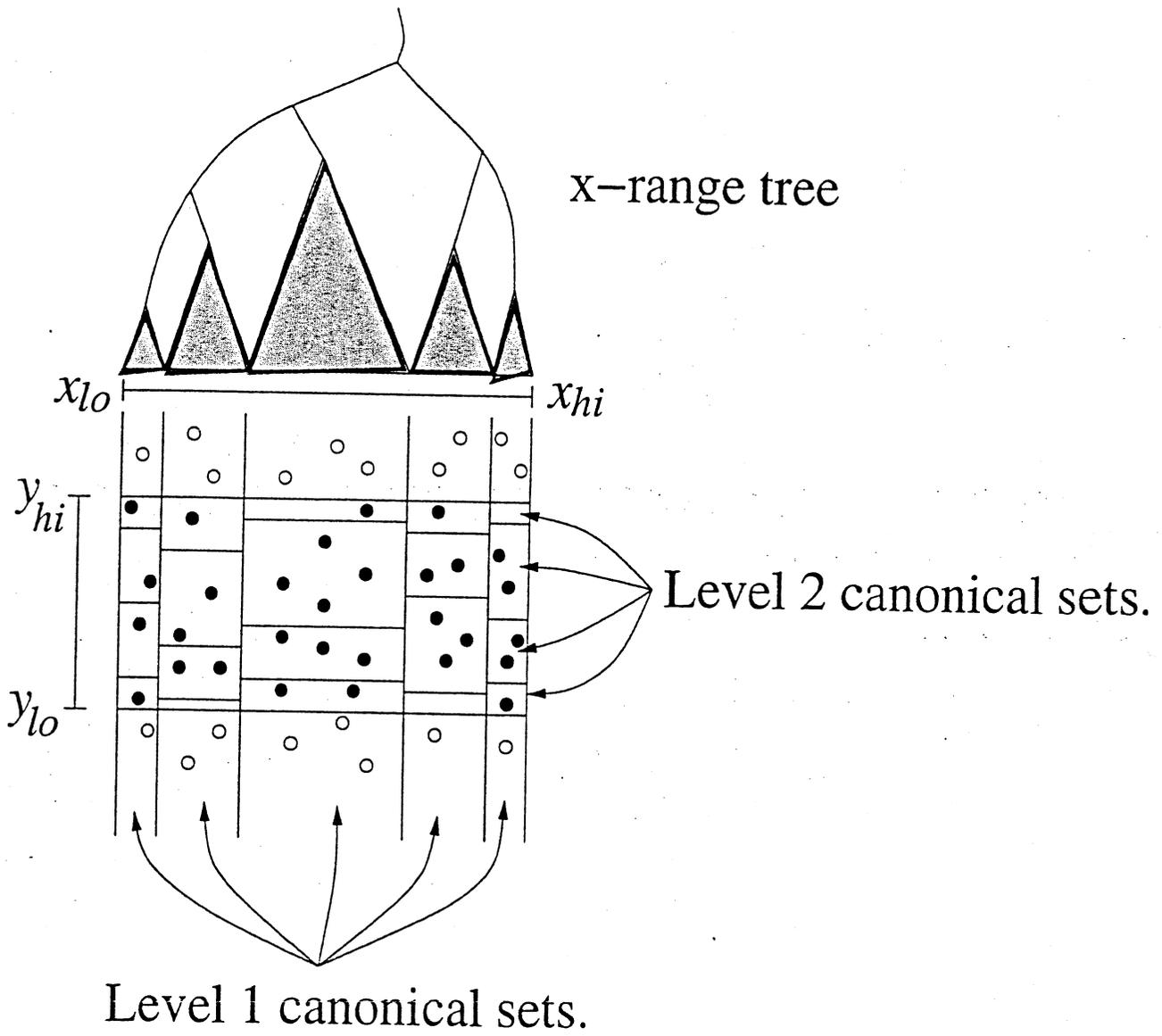
- Collect points of each canonical set, & build a y-range tree



- We search each of the  $O(\log n)$  canonical sets that include points for x-range  $[x_{l_0}, x_{h_1}]$  using their y-range trees for range  $[y_{l_0}, y_{h_1}]$ .
- The y-range searches list out the points of  $P$  in the query rectangle.  
(No duplicates.)

# Canonical Sets

---



## Analysis

- Query time is  $O(\log^2 n)$ 
  - $O(\log n)$  x-range query
  - Each canonical set y-range query takes  $O(\log n)$  time.
- Space complexity is  $O(n \log n)$ .
  - Total canonical set size is  $O(n \log n)$   
{ Each point of  $S$  appears in  $O(\log n)$  canonical sets. }
  - Each canonical set of size  $m$  requires  $O(m)$  space for the y-range tree.

Preprocessing time complexity.  $O(n \log n)$ .

- $x$ -range tree can be built in  $O(n \log n)$  time.

(Building them bottom up)

- $y$ -range trees of all canonical sets can be constructed in  $O(n \log n)$  time.

## d-dimensional range trees

- The multilevel range tree extends naturally to any dimension  $d$ .
- Build the  $x$ -range tree on the first coordinate.
- For each canonical set build the  $(d-1)$ -dimensional range tree on the remaining  $(d-1)$ -dimensions;  $x_2, x_3, \dots, x_d$ .  
 $T_d(n) = O(n \log n) + O(\log n) \cdot T_{d-1}(n)$
- Complexity grows by one  $\log n$  factor for each dimension.
- Search cost is  $O((\log n)^d)$  (query cost)
- Space complexity:  $O(n (\log n)^{d-1})$   
Preprocessing Time:  $O(n (\log n)^{d-1})$ .

# Fractional Cascading

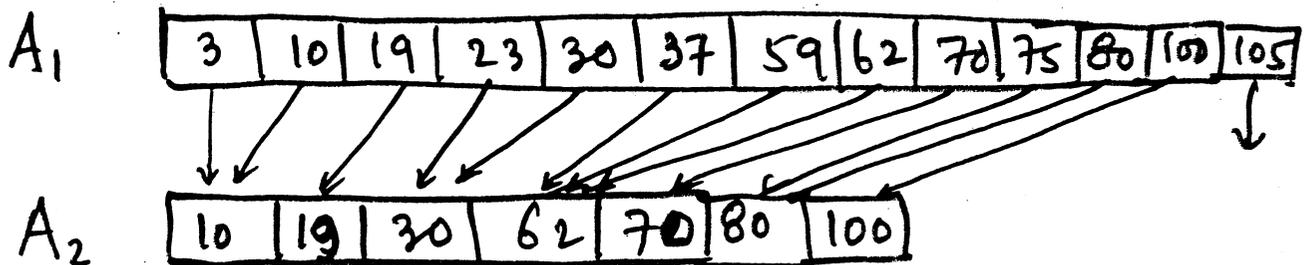
- Allows the 2-dimensional problem to be solved in
  - a)  $O(\log n + k)$  query time
  - b)  $O(n \log n)$  storage space & preprocessing time.
- This allows the  $d$ -dimensional orthogonal range tree problem to be solved in
  - a)  $O(\log^{d-1} n + k)$  query time
  - b)  $O(n \log n)^{d-1}$  storage space & preprocessing time.

## Basic idea.

- $x$ -range tree first finds the set of canonical sets lying in  $[x_{lo}, x_{hi}]$ .
- Each canonical set is searched using the ~~tree~~  $y$ -range tree for the same range  $[y_{lo}, y_{hi}]$ .
- Since each set is searched for the same key, we can improve the <sup>amortized</sup> search cost to  $O(1)$  time per set.
- Key is to attach pointers ~~to~~ linking the search structures for the canonical sets.

# Simple Example

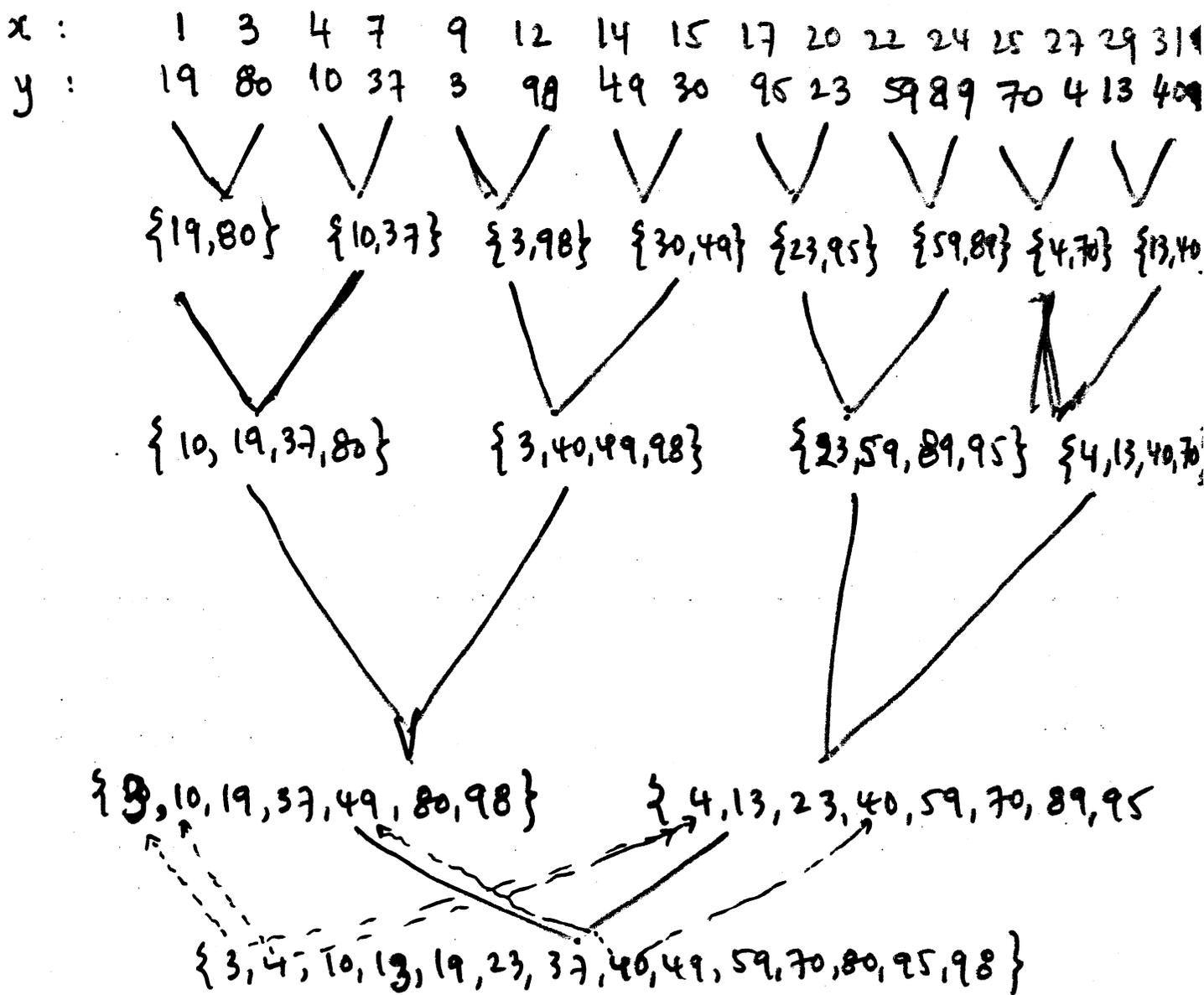
- We have two sets of numbers  $A_1, A_2$ ,  $A_1 \supseteq A_2$ ; both ~~set~~ sorted.
- Given a range  $[x, x']$ , we want to report all keys in  $A_1, A_2$  that lie in the range  $[x, x']$ .
- $2 \log n + k$  solution is easy.



- Point  $A_1[i]$  to  $A_2[j]$  ~~for all i where~~  
~~for all i where~~ for all  $i$  where  $A_1[i]$   
is the smallest key  $\geq A_1[i]$ .

- Suppose we want keys in range  $[y, y']$ .
- Search  $A_1$  for  $y$  and walk until past  $y'$ .  $O(\log n + k_1)$
- $A_1$  search <sup>for  $y$</sup>  ended ~~at~~ at  $A_1[i]$ .  
Use the pointer at  $A_1[i]$  to start search in  $A_2$ . This takes  $O(1 + k_2)$  time.

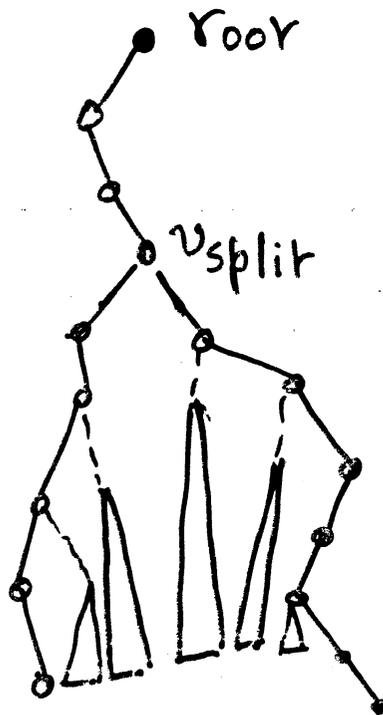
# Fractional Cascading in Range Trees



- Canonical subsets  $S(\ell(v))$  &  $S(r(v))$  are subsets of canonical set  $S(v)$ .
- Each entry in  $A(v)$  stores two pointers into arrays  $A(\ell(v))$  &  $A(r(v))$ .  $A(v)[i]$  points to the left most entry of  $A(\ell(v))[i]$  with smallest value  $\geq A(v)[i]$ . Same for  $A(r(v))$ .

# Fractional Cascading Search.

- Consider range  $R = [x, x'] \times [y, y']$ .
- Search  $x, x'$  in the main  $x$ -range tree.
- Let  $v_{\text{split}}$  be the node where the two search paths diverge.



- At  $v_{\text{split}}$ , do binary search to locate  $y$  in  $A(v_{\text{split}})$ .
- We can then search (for  $y$  in each) canonical set attached to the paths in  $O(1)$  time per canonical set.
- Report all points in rectangle query in  $O(\log n + k)$  time.

◦◦  $d$ -dimensional orthogonal range search problem can be solved in

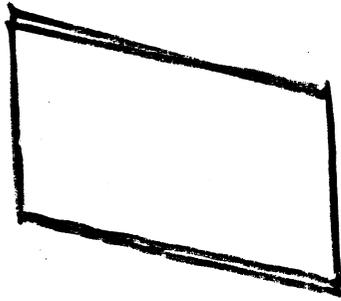
a)  $O((\log n)^{d-1} + k)$  query time

b)  $O(n(\log n)^{d-1})$  storage space

c)  $O(n(\log n)^{d-1})$  preprocessing time

Other query object ?

1) Parallelogram of the type.

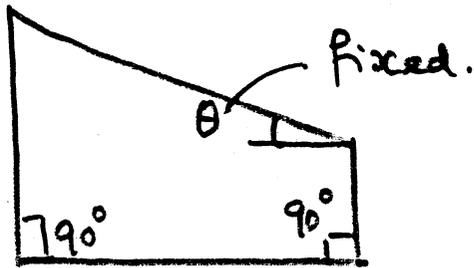


In 2-d

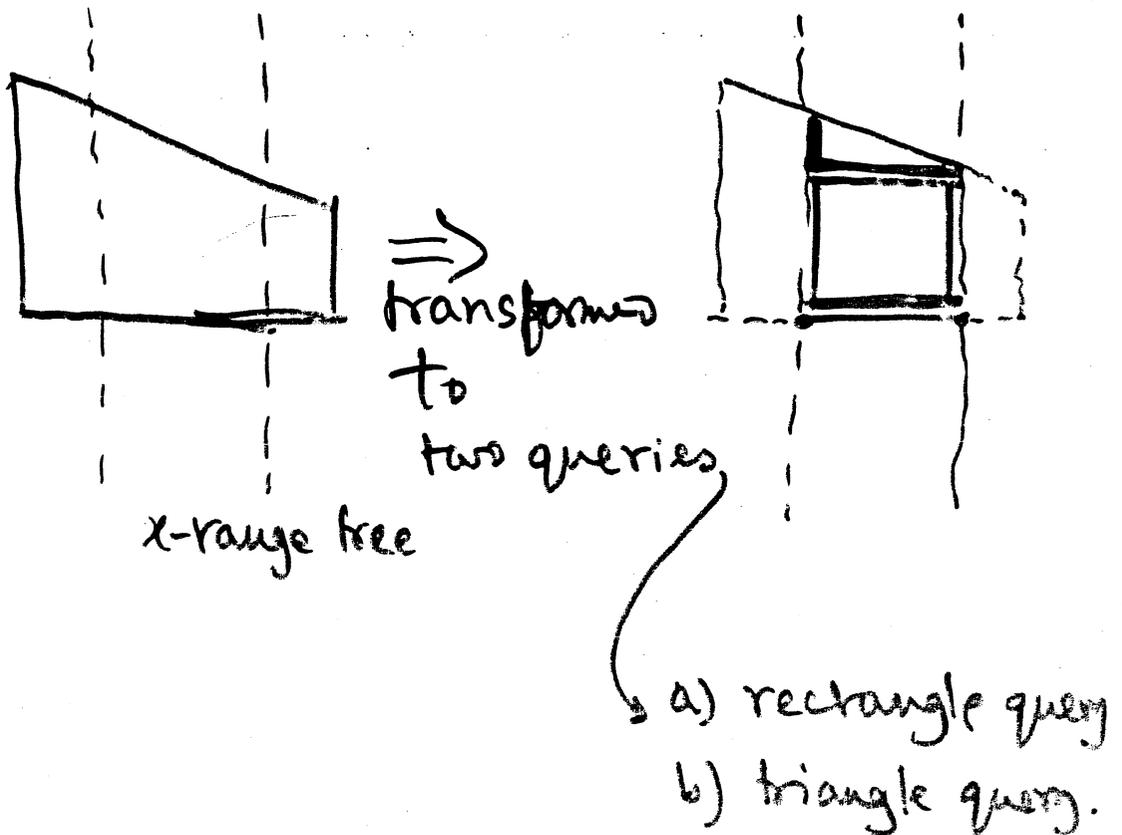
a)  $O(\log n + k)$  query time

b)  $O(n \log n)$  storage space  
+ preprocessing time.

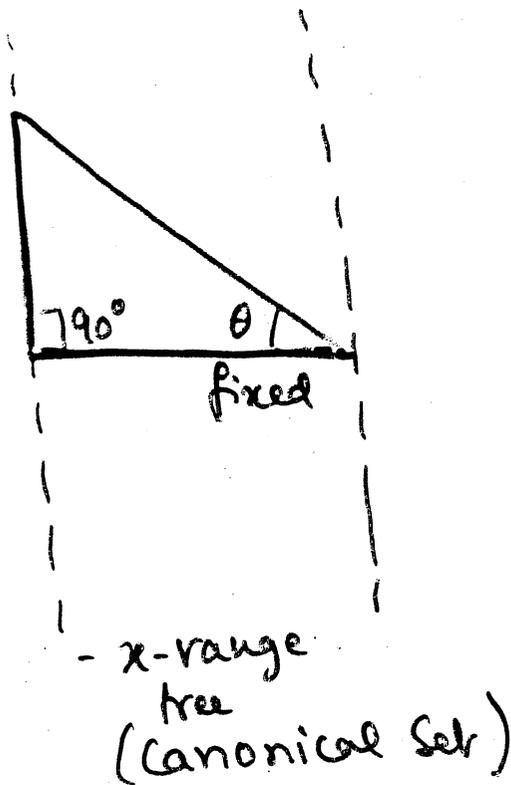
## 2. Trapezoid query object ?



Searching within a canonical set.



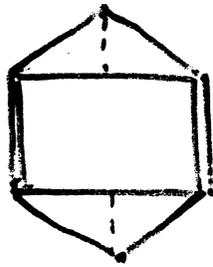
## ⊗ (20) Right-angled triangle.



Using the priority-search tree data structure, it is possible to answer the query (i.e. reporting the points inside the triangle) in  $O(\log n + k)$  time.

∴ The trapezoid <sup>query</sup> ~~set~~ can be answered in  $O(\log^2 n + k)$  time using  $O(n \log n)$  storage space + preprocessing time.

(4) Regular Hexagon, query object  
(Octagon)



Can be answered by querying

a) a rectangle

b) 4 right angled triangle.

with one fixed acute angle.

Query time :  $O(\log^2 n + k)$

Storage space :  $O(n \log n)$

Preprocessing time :  $O(n \log n)$

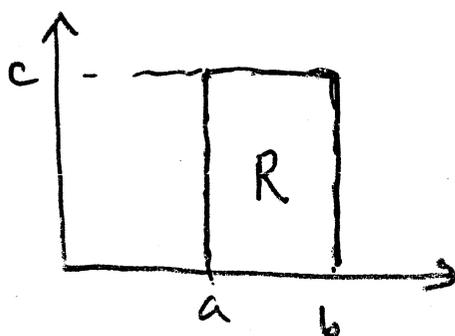
# McCreight's Priority Search Trees

Priority Search Trees, SIAM J. on Computing

Vol. 14, 257-276, 1985.

Problem Given a set of points  $(x, y)$ , find all points inside a testing rectangle.

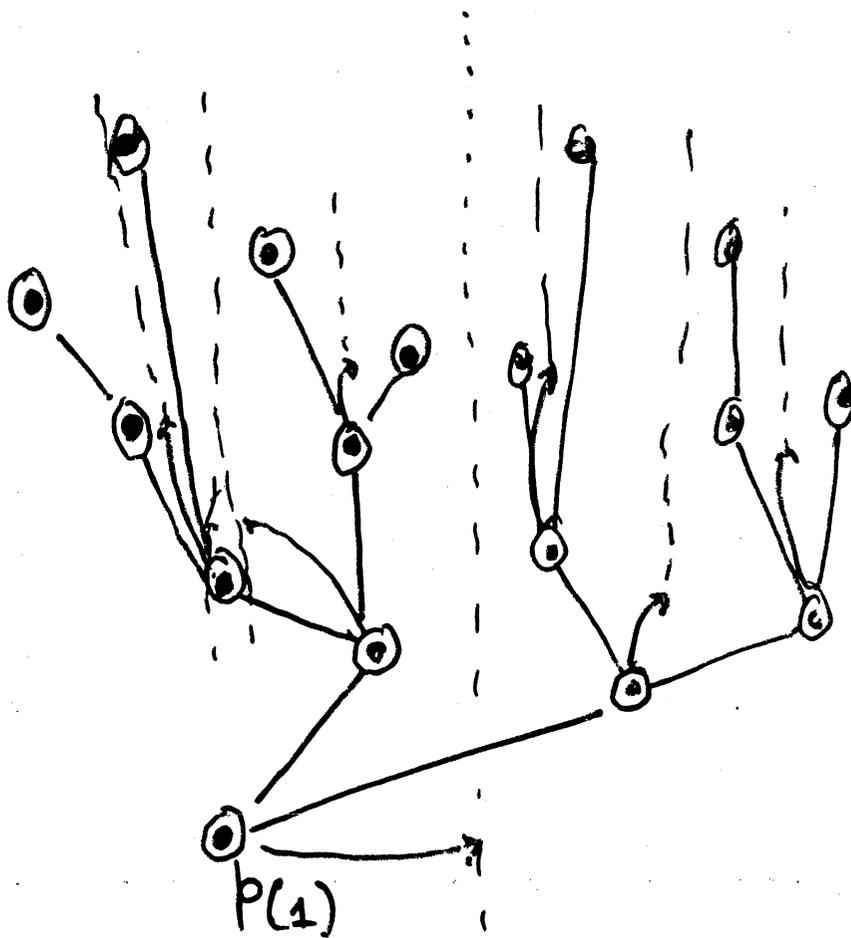
- allows insertion & deletions of points in  $O(\log n)$  time.
- $O(\log n + k)$  query time for rectangle  $R = \{x \mid a \leq x \leq b \text{ and } 0 \leq y \leq c\}$



- Search tree allows in  $O(\log n)$  time the following operations (Given a rectangle query)
  - ~~minimum~~ point with minimum x-coordinate
  - largest x-coordinate point in R
  - ~~largest~~ smallest y-coordinate point in R (No largest y-coordinate point)

Preprocessing Stage ( $S = \{p_1, \dots, p_n\}$  is the set of points)

- ~~Step~~ Select the point with the smallest y-coordinate. (say  $p_{(1)}$ ).  $p_{(1)}$  is the root of the search tree.
- The remaining points  ~~$S - \{p_{(1)}\}$~~   $S - \{p_{(1)}\}$  are assigned to the left or to the right subtree of  $p_{(1)}$  according to whether they are to the left or to the right of a vertical line determined by the x-median of  $S - \{p_{(1)}\}$  points. The root also contains the x-range of  $S$ , along with the vertical line.
- Recursively construct the left priority search tree of  $p_{(1)}$
- Recursively construct the right priority search tree of  $p_{(1)}$ .



Height of the tree is  $O(\log n)$  (Balanced)

Search for a point  $p = (x_p, y_p)$

- Check if  $y_p = y_{p(1)}$

- If yes, stop.

- Else if  $x_p < \text{divider}(p_i)$ ,

search left <sup>sub</sup> tree

else search right <sup>sub</sup> tree.

Insert a point  $p = (x_p, y_p)$

- If  $y_p < y_{p(i)}$ , (a) exchange  $p$  with  $p(i)$   
(b) Insert  $p(i)$  in  $S - \{p(i)\}$ .
- If  $y_p \geq y_{p(i)}$ , insert  $p$  in  $S - \{p(i)\}$ .

The search will terminate at a null ~~link~~ link which gets replaced by the point being inserted.

Cost:  $O(\log n)$  time.

## Delete a point p

- Locate  $p$  in the tree.
- Delete  $p$  from the tree. (say at node  $v$ )  
(Leaves a hole in the structure)
- Fill it up by promoting the child of  $v$  with the smallest  $y$ -coordinate.  
(The hole moves down to the position of that child)
- Continue until the hole moves out of the bottom of the tree.

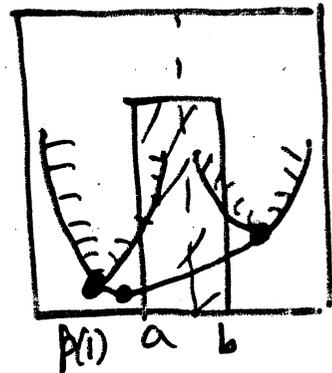
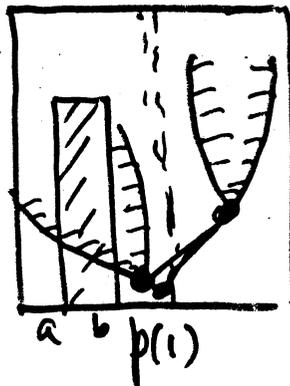
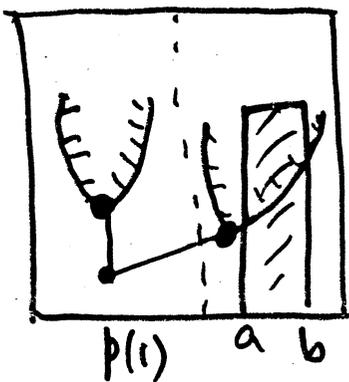
Cost  $O(\log n)$

(Balancing the tree is required after insertion & deletion operations are performed.)

Balancing can be avoided if all the points are known before hand.

Find the point of  $S$  with smallest  $y$ -coordinate inside the query rectangle  $R = [a, b] \times [c, e]$ .

- The root is in  $R$ : Stop,  $p_{(1)}$  is the point
- Otherwise we have three possibilities.

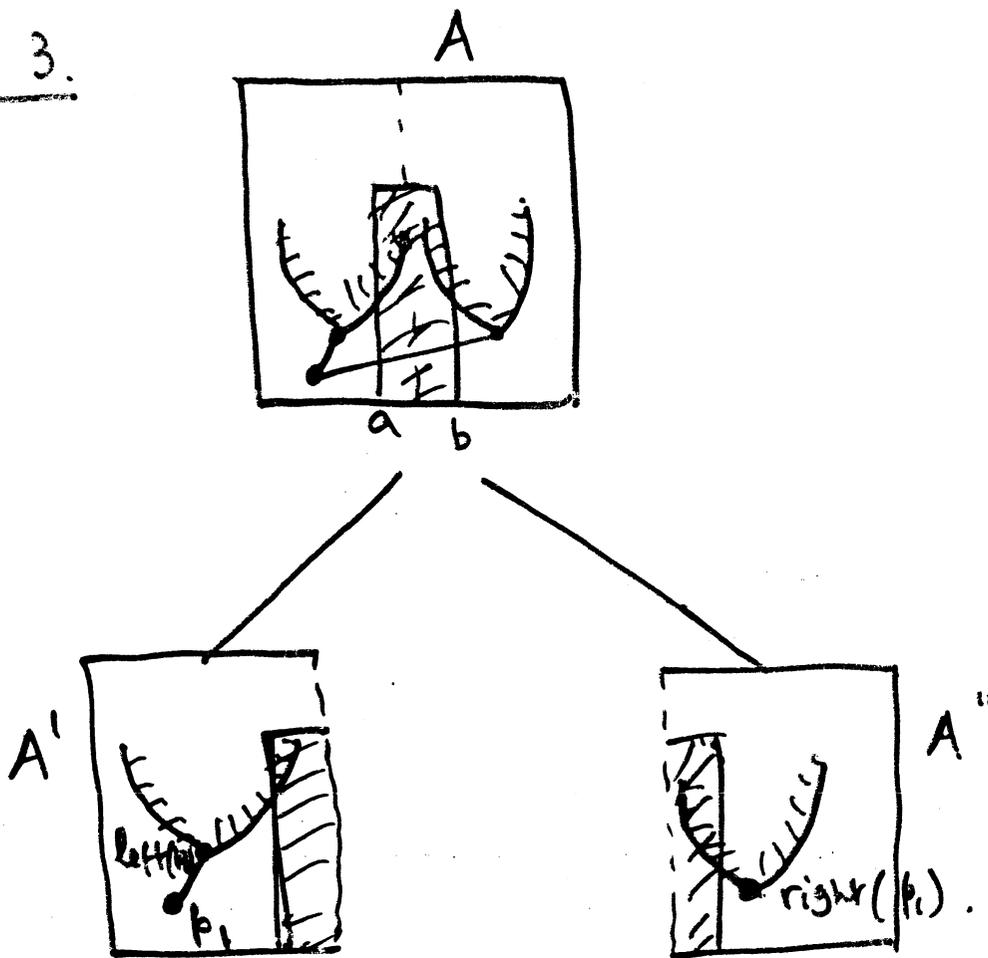


First two cases: ignore either the left or the right.

Third case: (a) Find the lowest point in each subtree.

(b) Pick the smallest of the two.

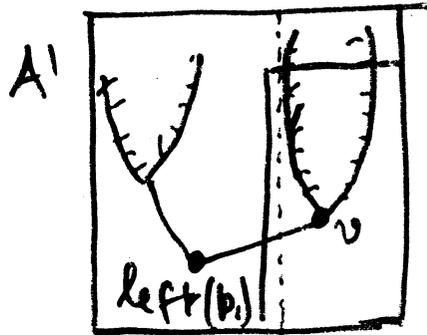
Case 3.



~~R~~ The problem reduces to two smaller similar subproblems.

In  $A'$  if the vertical line for the left subtree ~~R~~  $\text{left}(p_i)$  does not intersect the shaded rectangle, the left subtree can be rejected. If the line intersects the rectangle, the ~~right~~ shaded region contains the  $x$ -range of the right subtree of ~~R~~  $\text{left}(p_i)$ .

### Case 3 (contd.)



We then test  $v$ , the right child of  $\text{left}(p_i)$   
↓ See if it is inside  $R$ . If yes, ~~search in left subtree~~ <sup>search in left subtree</sup>.

Otherwise we ignore the search tree rooted at  $v$ . Search in left subtree.

In any case we are left with one subtree of  $\text{left}(p_i)$ . Total time spent at  $\text{left}(v_i)$  is  $O(1)$ .

Similarly, ~~after~~ after  $O(1)$  time, we are left with one subtree of  $\text{right}(p_i)$ .

oo

Lemma: The point with smallest  $y$  coordinate in a query rectangle  $R$  can be found in  $O(\log n)$  time.