

5800

data structures

apr8/apr11 2022
shelat

Dictionary

data structure

$\text{insert}(\underline{\text{key}}, \underline{\text{value}})$: add an item to the data structure

$\text{lookup}(\text{key})$: returns $(\text{key}, \text{value})$ if it was previously inserted into the data structure

$\text{Delete}(\text{key})$:

$\text{FINDNEXT}(\text{key})$: finds the $(\text{key}', \text{value})$ entry where key' is the smallest value inserted for which $\text{key}' > \text{key}$

MIN, MAX

DICTIONARY

insert(key, value)

delete(key)

lookup(key)

findnext(key)

DICTIONARY

standard solution: hashtable

insert(key, value)

delete(key)

lookup(key)

findnext(key)

We expect $\Theta(1)$ performance.

But usually this only happens in expectation. worst case performance could be $O(n)$.

→ might require scanning all keys.
 $O(n)$

Hashtables are tricky

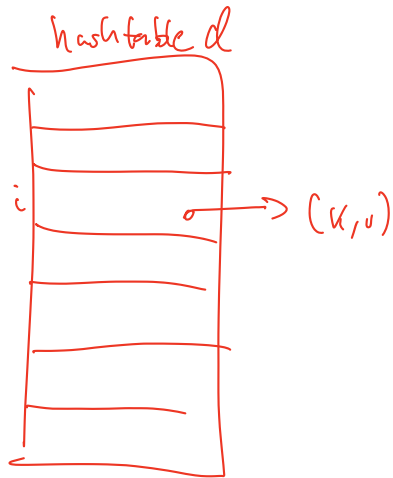
```
1
2 import time
3 import sys
4 import d
5
6 dd = {}
7
8 # make a dictionary with elements from the list
9 for l in d.list:
10     dd[l] = l
11
12 def lookup(v):
13     start = time.time()
14     t = 0
15     for j in range(10000):
16         if v in dd:
17             t = t + 1
18     end = time.time()
19     print(end - start)
20     return t
21
```

$d[k] = v$
 $i = \text{hash}(k)$.

if i is free
add (k, v)
at i

else

= keep looking for
a new spot
"open addressing"



Hashtables are tricky

```
1
2 import time
3 import sys
4 import d
5
6 dd = {}
7
8 # make a dictionary with elements from the list
9 for l in d.list:
10     dd[l] = l
11
12 def lookup(v):
13     start = time.time()
14     t = 0
15     for j in range(10000):
16         if v in dd:
17             t = t + 1
18     end = time.time()
19     print(end - start)
20     return t
21
```

This is a trivial lookup experiment.
Looking up 1 key takes 2000x longer.

```
MacBook-Pro-2:hashing abhi$ python3 bad.py
size of dictionary: 43689
Starting experiment to lookup 1000:
0.0005161762237548828
Starting experiment to lookup 100000:
1.0303189754486084
MacBook-Pro-2:hashing abhi$
```

Hashtables are tricky

```
1
2 import time
3 import sys
4 import d
5
6 dd = {}
7
8 # make a dictionary with elements from the list
9 for l in d.list:
10     dd[l] = l
11
12 def lookup(v):
13     start = time.time()
14     t = 0
15     for j in range(10000):
16         if v in dd:
17             t = t + 1
18     end = time.time()
19     print(end - start)
20     return t
21
```

This is a trivial lookup experiment.
Looking up 1 key takes 2000x longer.

```
MacBook-Pro-2:hashing abhi$ python3 bad.py
size of dictionary: 43689
Starting experiment to lookup 1000:
0.0005161762237548828
Starting experiment to lookup 100000:
1.0303189754486084
MacBook-Pro-2:hashing abhi$ █
```

Worst case performance: $O(n)$

DICTIONARY

new constraint: keys belong to limited range:

$$\{ \underline{1}, \dots, \underline{u} \}$$

$$u \sim \underline{\underline{2^{32}}}$$

insert(key, value)

delete(key)

lookup(key)

findnext(key)

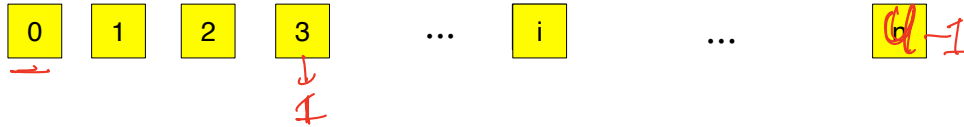
$\Theta(\log \log u)$ performance.

$$\log \log 4 = 511$$

Note: $\Theta(\log n)$ or $\Theta(\log u)$ performance
can be achieved with
balanced binary trees.

A simple solution: bit vector

Maintain an array of bits



insert(key, value)
delete(key)
lookup(key)
findnext(key)

$\Theta(1)$
 $\Theta(n)$ in worst case

• We can still use this DS as a base case. for small u .

CAN WE DO BETTER THAN $O(N)$ FINDNEXT?

van emde Boas Q VEB Q

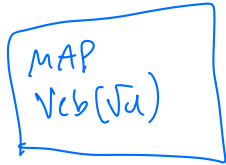
THE BIG IDEA: a datastructure for universe $1..U$ that is comprised of small versions of the same datastructure for universe $1..√U$.

VEB for U

sz min max

points to smaller VEB($√U$) queues.

Keeps track
of which
ptr are
not
empty



van emde Boas Q

THE BIG IDEA:

Use recursion for a data structure.

A data structure that handles $1..n$ can be designing using several smaller versions of the same structure.

VEB queue

$\text{VEB}_{(N)}$

VEB queue

$\text{VEB}_{(N)}$

SZ, MIN, MAX

VEB queue

$\text{VEB}_{(N)}$

SZ, MIN, MAX

BASE CASE: 1 BIT VECTOR.

VEB queue

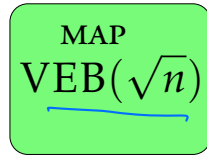
VEB_(N)

SZ, MIN, MAX

BASE CASE: 4 BIT VECTOR.

for $u \leq 4$, bit vector.

NORMAL CASE:



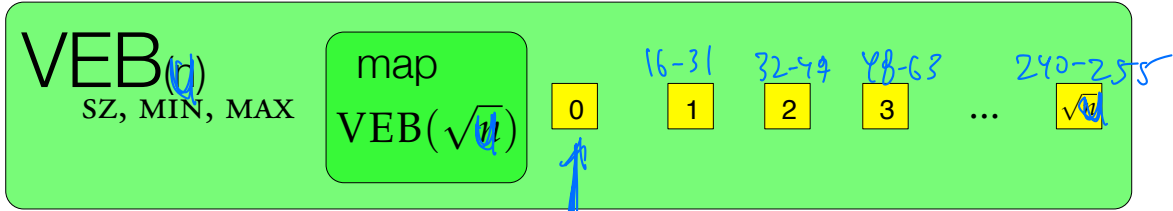
Pointers to recursive, smaller instances of VEB.



Keeps track of which ptrs
are not empty.

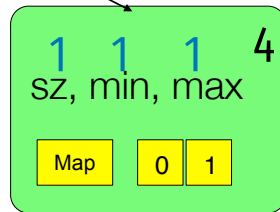
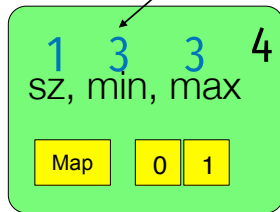
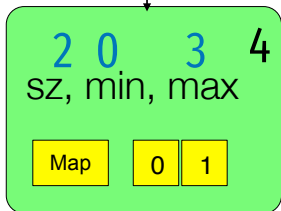
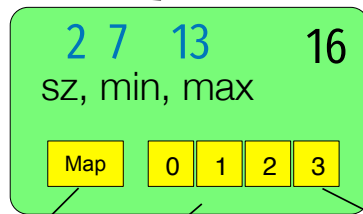
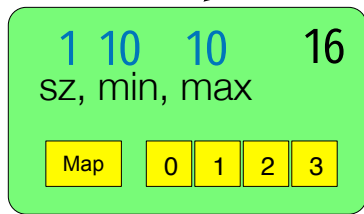
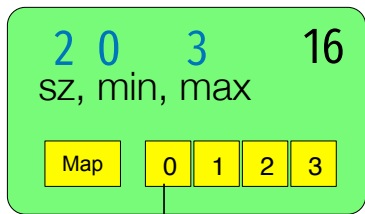
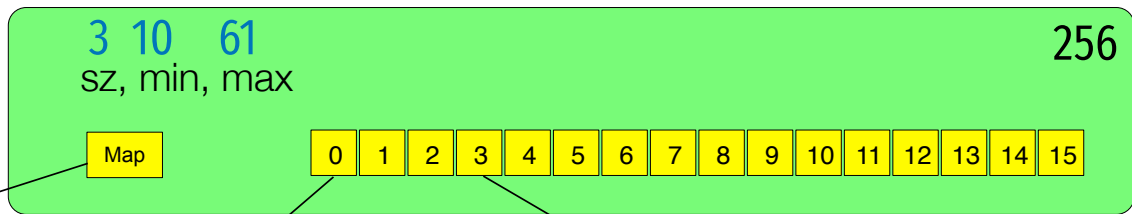
EXAMPLE:

$u = 256$

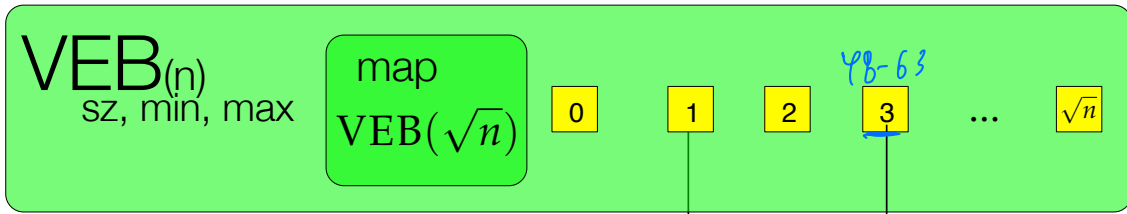


this pointer would point to
a VEB queue that would store
any entries between 0...15

Example $n=256$, keys={10, 55, 61}



(Roughly correct)

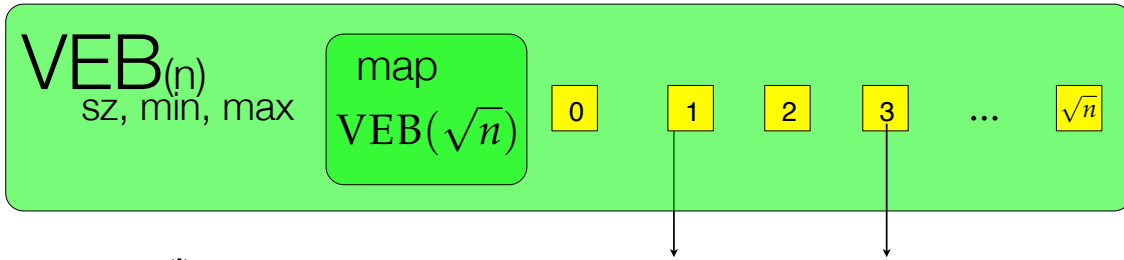


LOOKUP(i)

$\Theta(u)$ $\left\{ \begin{array}{l} \text{write } i = a \cdot \sqrt{u} + b \text{ for } a, b \in [0 \dots \sqrt{u}] \text{ eg } 55 = 3 \cdot \sqrt{256} + 7 \\ \langle \text{base case} \rangle \rightarrow \text{for } u \leq 4, \text{ bit vector} \\ \text{If } a = \text{null return false} \end{array} \right.$

else return $a \cdot \text{lookup}(b)$

$$T(u) = T(\sqrt{u}) + \Theta(1) = \Theta(\log \log u)$$



LOOKUP(i)

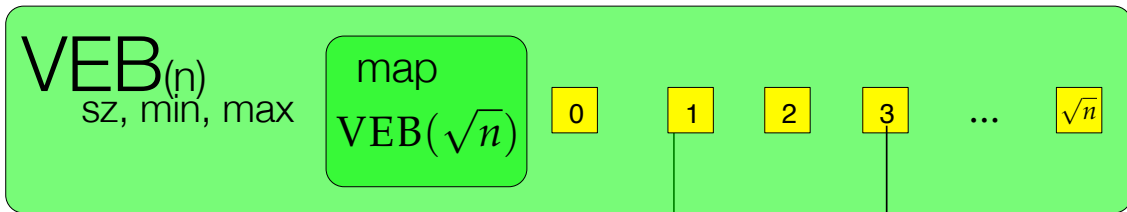
WRITE $i = a\sqrt{n} + b$

IF <BASE CASE>: CHECK BIT VECTOR

IF SIZE = 0 OR **a**.SIZE = 0 THEN RETURN FALSE

ELSE RETURN **a**.LOOKUP(b)

(Almost right, we will have to slightly change this later.)



LOOKUP(*i*)

WRITE $i = a\sqrt{n} + b$

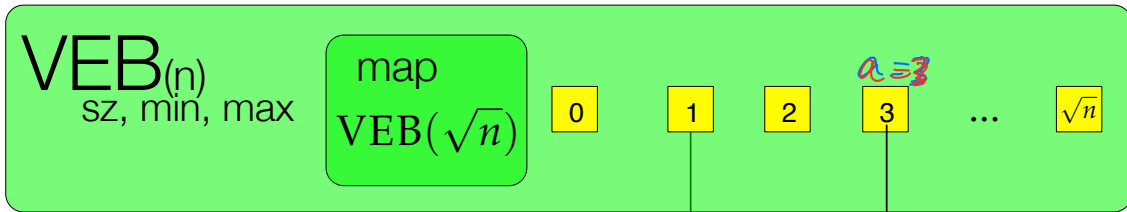
IF <BASE CASE>: CHECK BIT VECTOR

IF SIZE = 0 OR **a**.SIZE = 0 THEN RETURN FALSE

ELSE RETURN **a**.LOOKUP(*b*)

Running time: $T(n) = T(\sqrt{n}) + \Theta(1) = \Theta(\log \log n)$

(Almost right, we will have to slightly change this later.)



FINDNEXT(i)

EXAMPLE $\text{findnext}(55) = 3 \cdot 16 + 7$

MAX
55

IDEA:

write $i = a\sqrt{n} + b$

CASE 1

if the next element occurs in bucket a

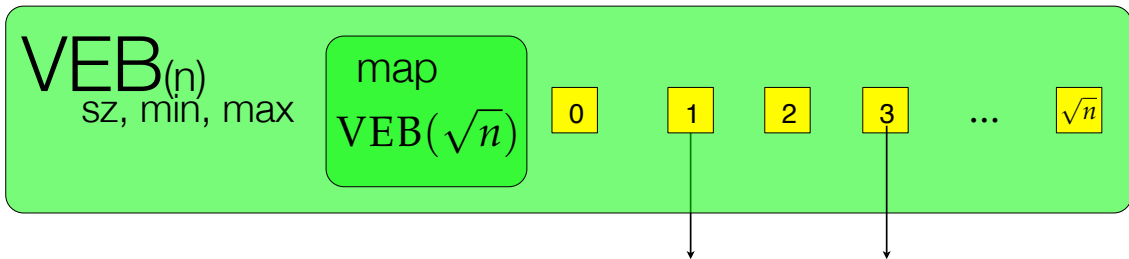
easy: we can use $a.\text{findnext}(b)$

use
 $a.\text{max}$
to precu

CASE 2

the next element occurs in the next non-empty bucket

use $\text{map}.\text{findnext}(a) \cdot \text{map}$



FINDNEXT(*i*)

IDEA:

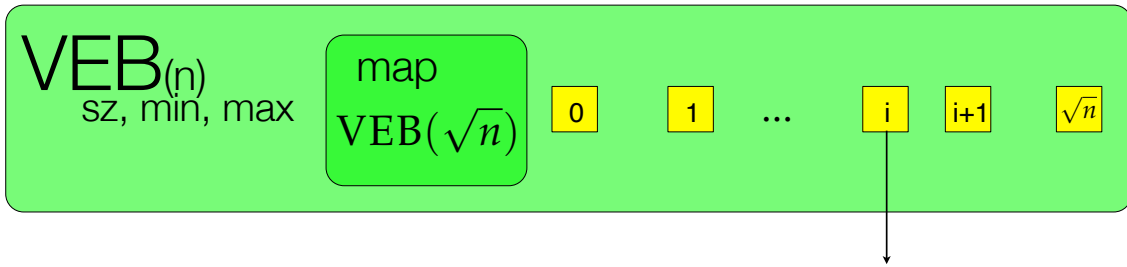
Write $i = a\sqrt{n} + b$ as usual.

Case 1: Bucket *a* has the next value.

Recursively use findnext_{*a*}(*b*)

Case 2: Bucket *a* does not have the next value.

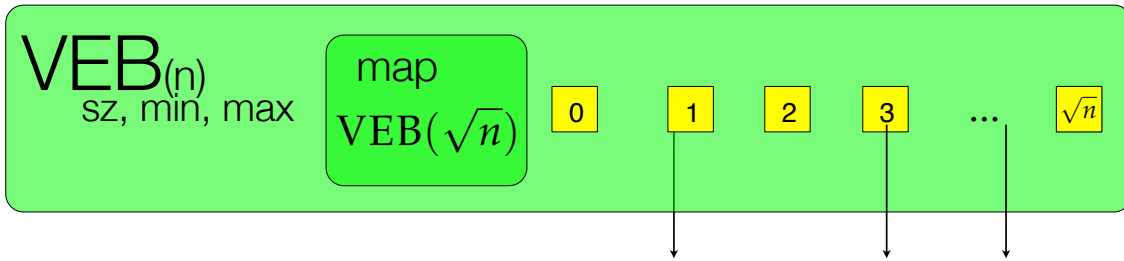
Use $x = \text{findnext}_{\text{map}}(a)$, return $x.\text{min}$.



FINDNEXT(i)

$\Theta(u)$ { write $i = a \cdot \sqrt{u} + b$
 base case, use the bitvector if $u \leq 4$
 IF $a \cdot \max > b$ then
 return $a \cdot \text{findnext}(b)$ ← only do one
 else
 return $\text{map} \cdot \text{findnext}(a) \cdot \min$ ←

$T(u) = T(\sqrt{u}) + \Theta(u)$
 $= \Theta(\log \log u)$



`FINDNEXT(i)`

WRITE $i = a\sqrt{n} + b$

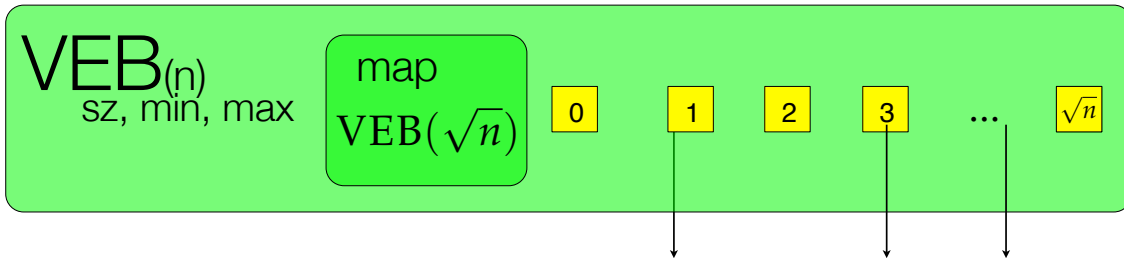
<BASE CASE IF SIZE IS ZERO>

IF `a`.MAX > b THEN

RETURN `a`.FINDNEXT(b)

ELSE

RETURN `MAP`.FINDNEXT(`a`).MIN



FINDNEXT(i)

WRITE $i = a\sqrt{n} + b$

<BASE CASE IF SIZE IS ZERO>

IF a .MAX $> b$ THEN

RETURN a .FINDNEXT(b)

ELSE

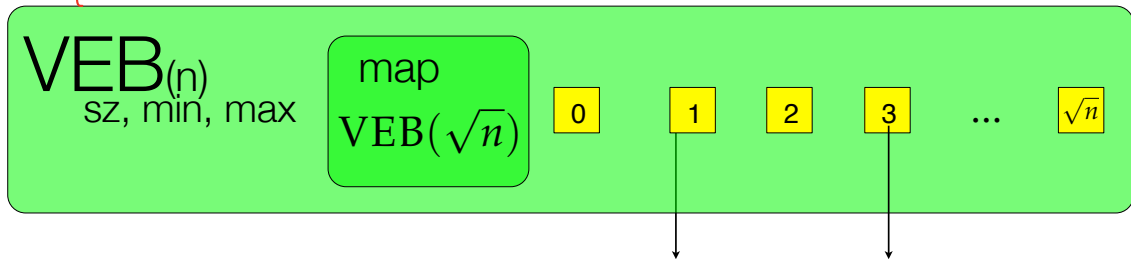
RETURN MAP .FINDNEXT(a).MIN

Running time:

$$T(n) = T(\sqrt{n}) + \Theta(1)$$

$$\Theta(\log \log n)$$

First and last



INSERT(i)

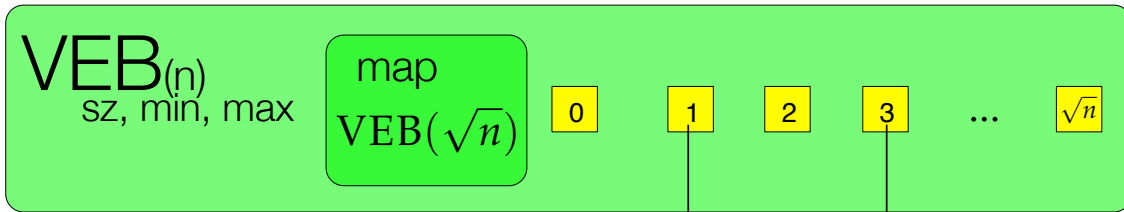
WRITE $i = a\sqrt{n} + b$

a. insert(b)

map.insert(a)

$$T(u) = 2T(\sqrt{u}) + \Theta(1)$$
$$= O(\log u)$$

$$S(n) = 2S(\frac{n}{2}) + \Theta(1)$$
$$O(n)$$

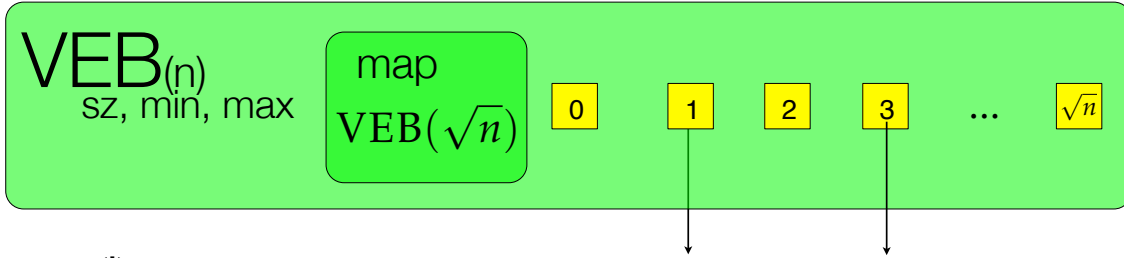


INSERT(*i*)

WRITE $i = a\sqrt{n} + b$

A.INSERT(*B*)

MAP.INSERT(*A*)



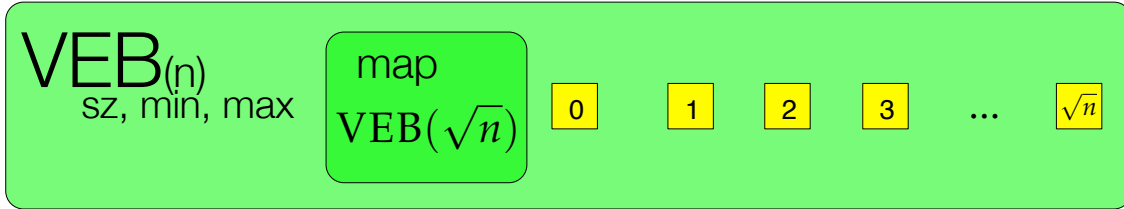
INSERT(*i*)

WRITE $i = a\sqrt{n} + b$

A.INSERT(*B*)

MAP.INSERT(*A*)

WHAT IS THE PROBLEM WITH THIS?



INSERT(i)

WHAT IS THE PROBLEM WITH THIS?

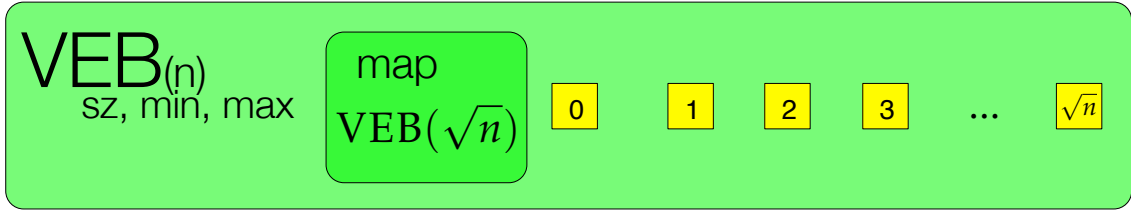
WRITE $i = a\sqrt{n} + b$

A.INSERT(B)

MAP.INSERT(A)

HOW CAN WE GET AROUND THE PROBLEM OF
INSERTING TWICE?

ANSWER: LAZY INSERTS. HOW MANY TIMES DO WE NEED
TO INSERT INTO MAP?



INSERT(i)

~~WRITE $i = a\sqrt{n} + b$~~

se case w/ BIT vector)
IF SZ==0 THEN

→ SZ=(MIN = MAX = i

$\Theta(1)$ for insert into empty datastructure.

ELSE

IF MIN > i SWAP(i, MIN)

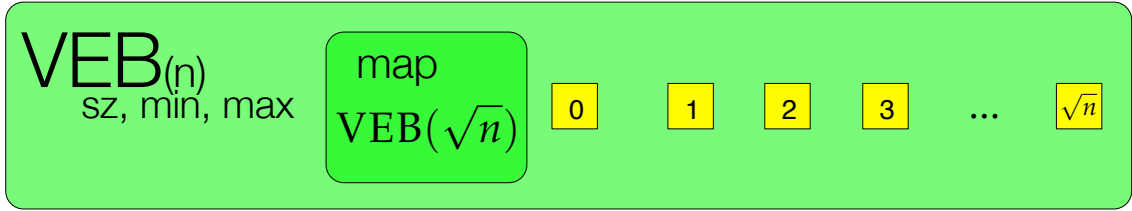
write $i = a\sqrt{n} + b$

IF a.size == 0 {map.insert(a)}

a.insert(b)

update SZ AND MAX

$$T(n) = T(\sqrt{n}) + \Theta(1)$$



INSERT(i)

IF $SZ == 0$ THEN UPDATE $SZ = 1, MIN = MAX = i$

ELSE

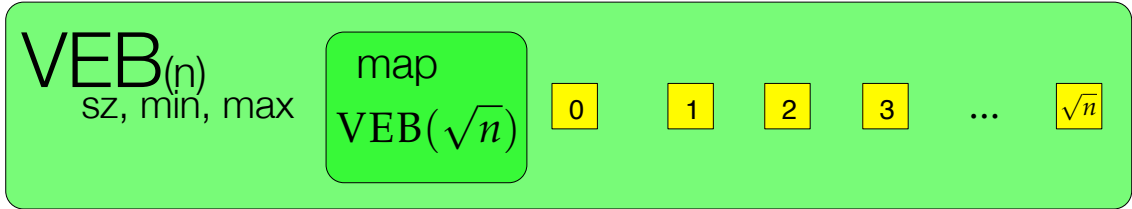
IF $MIN > i$ SWAP(i, MIN)

WRITE $i = a\sqrt{n} + b$

IF a .SZ == 0 THEN MAP[MAP].INSERT(a).

a .INSERT(b)

UPDATE SZ, MAX



INSERT(i)

IF SZ==0 THEN UPDATE SZ=1, MIN=MAX=i

ELSE

IF MIN>i SWAP(i,MIN)

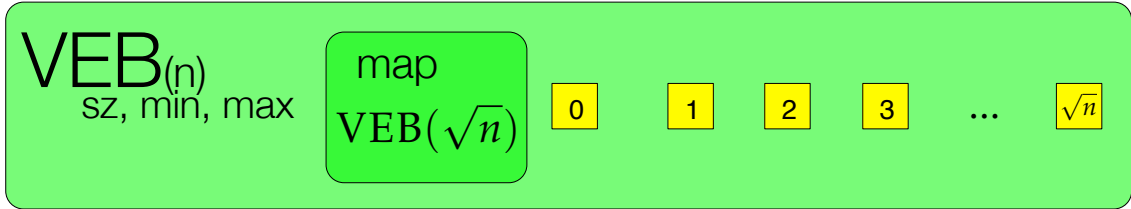
WRITE $i = a\sqrt{n} + b$

IF **a**.SZ==0 THEN **MAP**.INSERT(a).

a.INSERT(b)

UPDATE SZ, MAX

If a is empty:
then 1 full recursive call + 1 base case



INSERT(i)

IF SZ==0 THEN UPDATE SZ=1, MIN=MAX=i

ELSE

IF MIN>i SWAP(i, MIN)

WRITE $i = a\sqrt{n} + b$

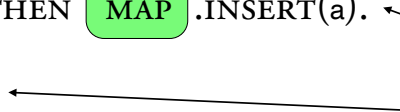
IF **a**.SZ==0 THEN **MAP**.INSERT(a).

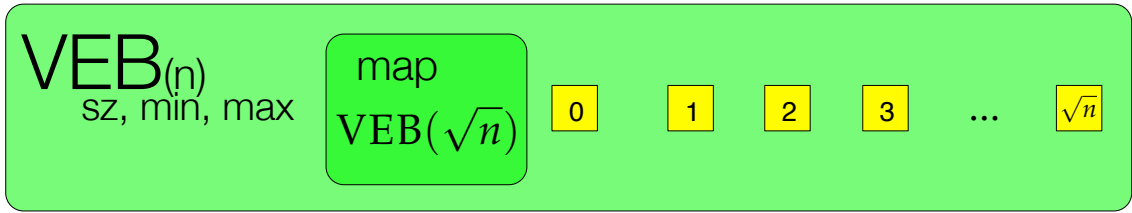
a.INSERT(b)

UPDATE SZ, MAX

If a is empty:
then 1 full recursive call + 1 base case

If a is not empty:
Then this line does not run
but 1 full recursive call is made





INSERT(i)

IF SZ==0 THEN UPDATE SZ=1, MIN=MAX=i

ELSE

IF MIN>i SWAP(i,MIN)

WRITE $i = a\sqrt{n} + b$

IF **a**.SZ==0 THEN **MAP**.INSERT(a).

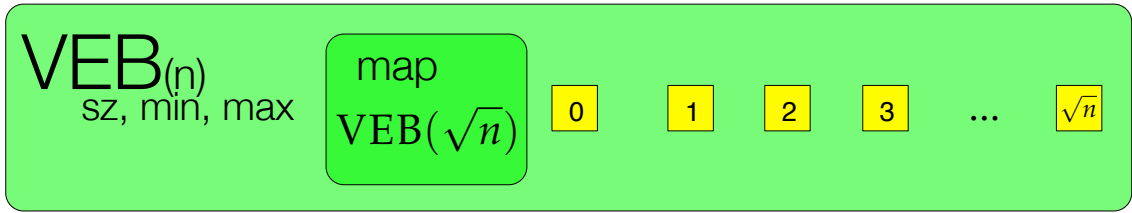
a.INSERT(b)

UPDATE SZ, MAX

If a is empty:
then 1 full recursive call + 1 base case

If a is not empty:
Then this line does not run
but 1 full recursive call is made

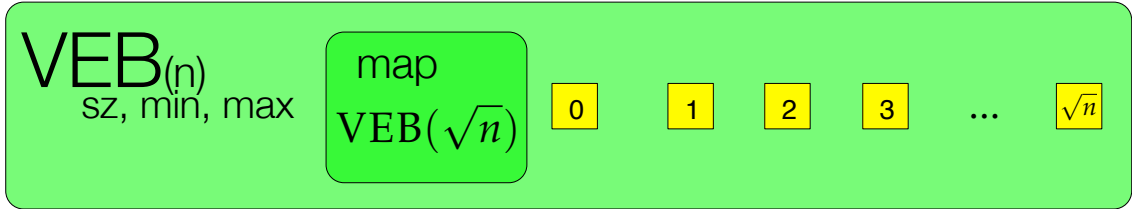
$$T(n) = T(\sqrt{n}) + \Theta(1)$$



LOOKUP(*i*)

WRITE $i = a\sqrt{n} + b$

We need to fix the Lookup to work with Lazy inserts.



LOOKUP(i)

WRITE $i = a\sqrt{n} + b$

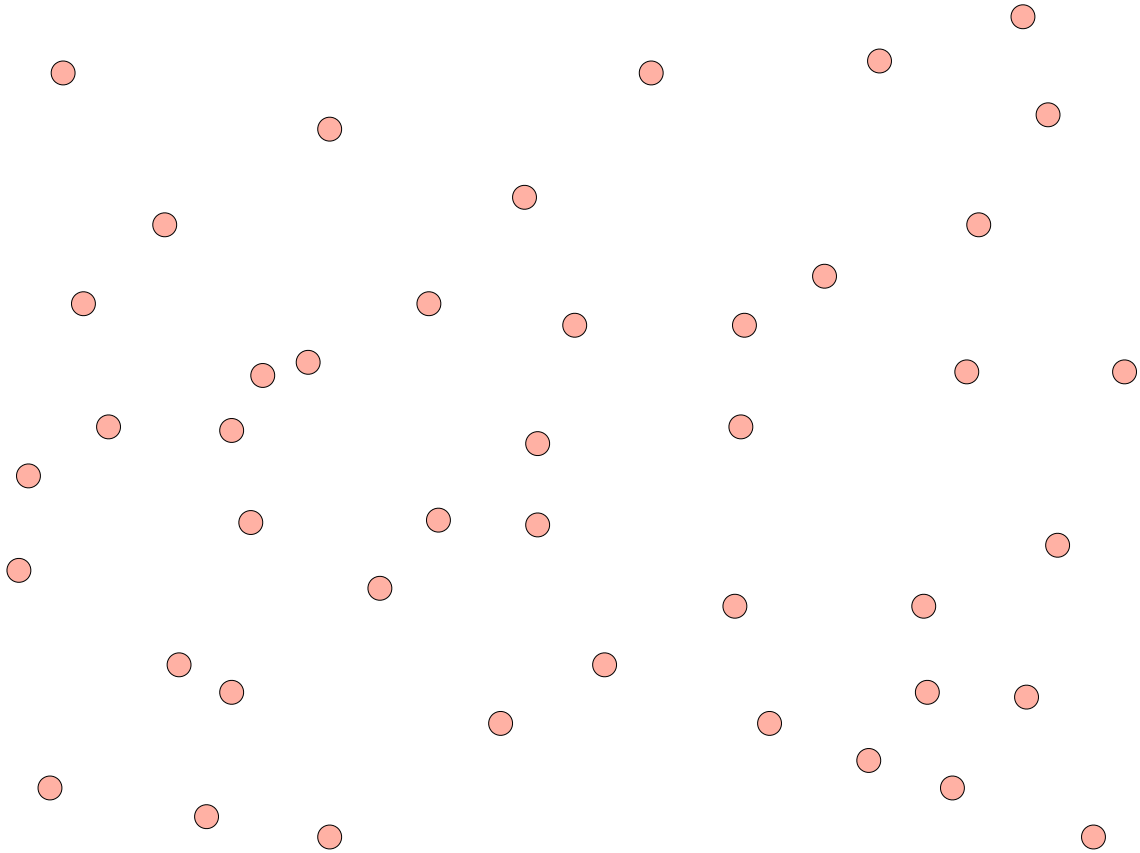
IF SIZE==0 RETURN FALSE

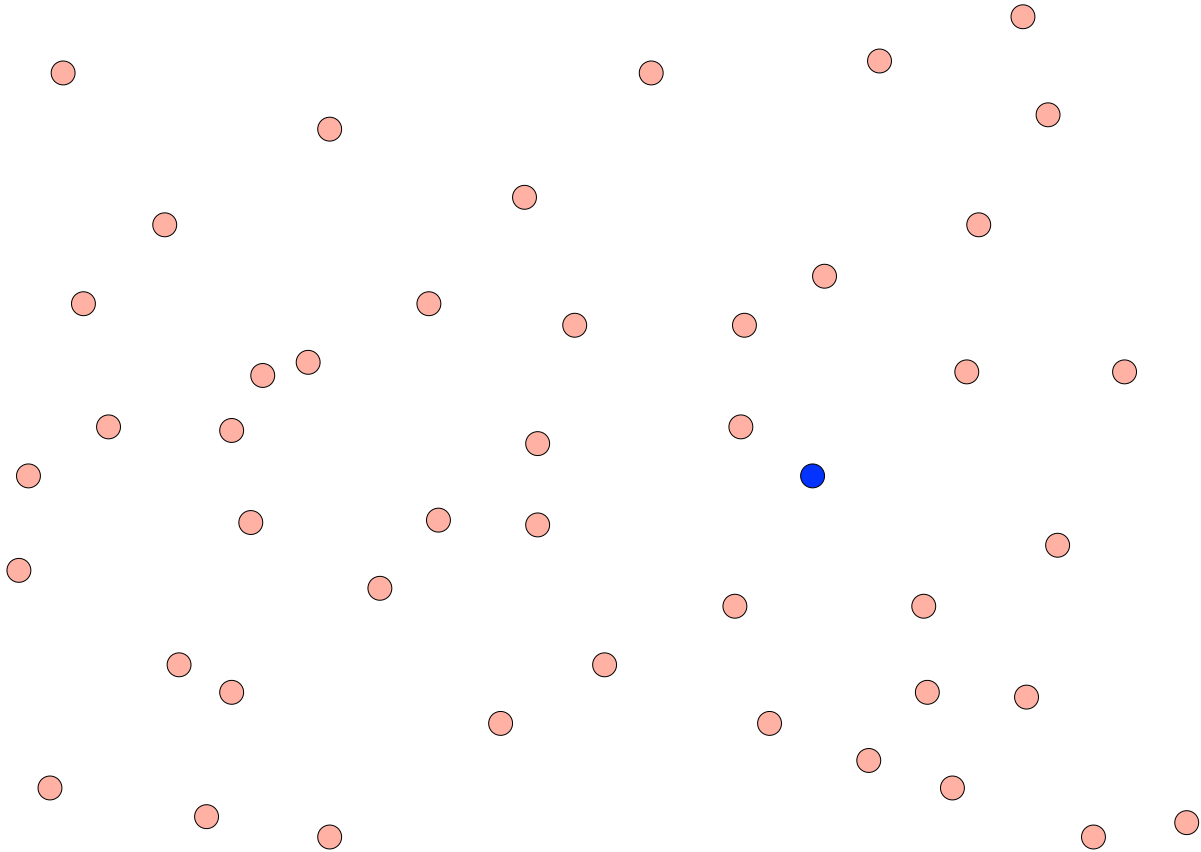
IF i ==MIN RETURN TRUE

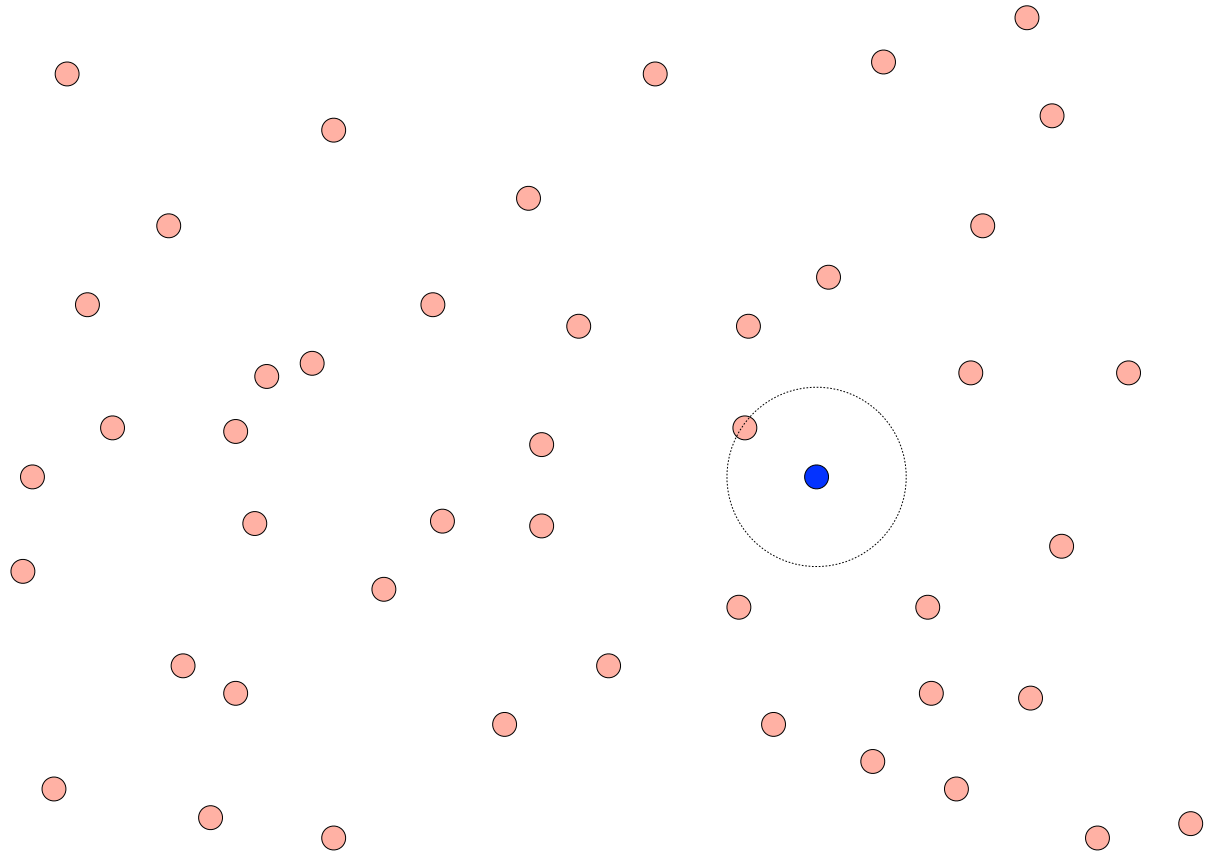
ELSE RETURN a .LOOKUP(b)

We need to fix the Lookup to work with Lazy inserts.

Nearest
neighbor
queries



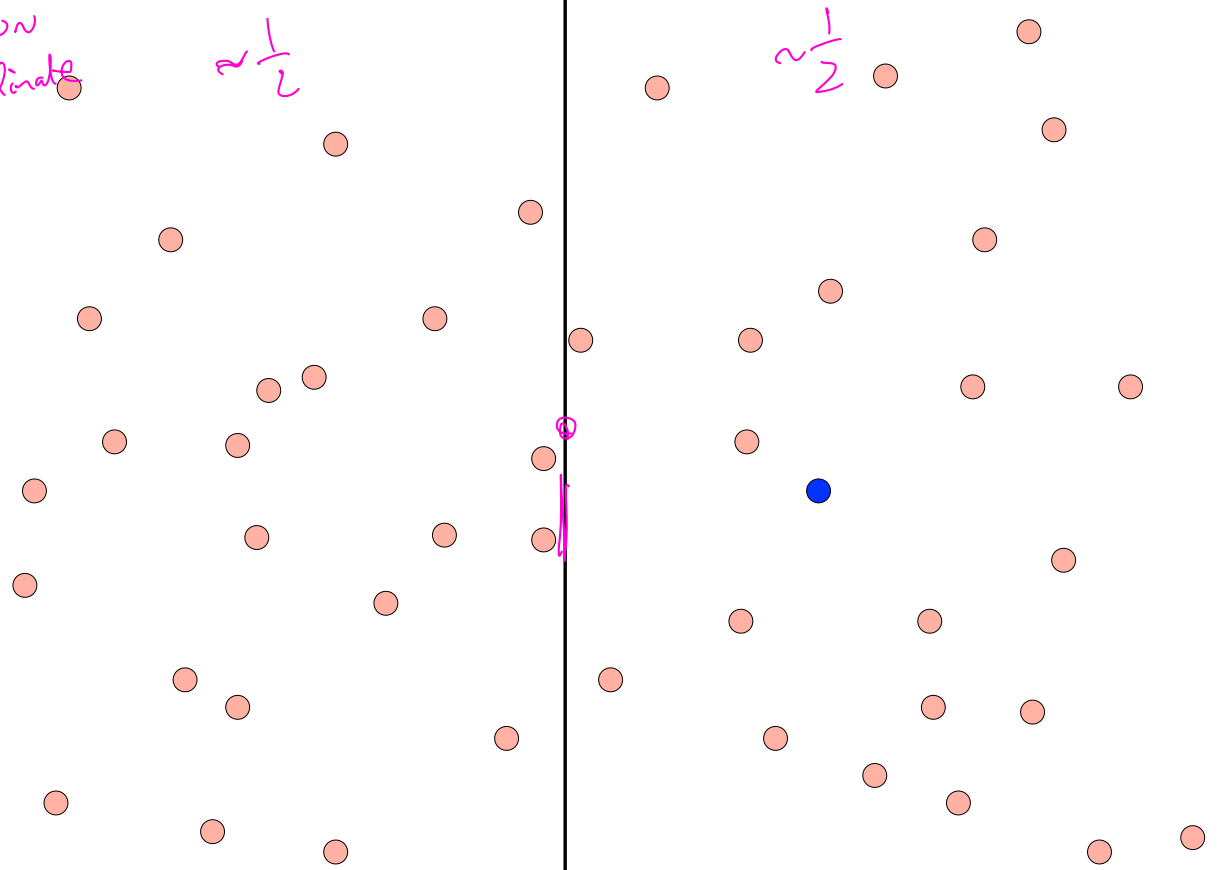


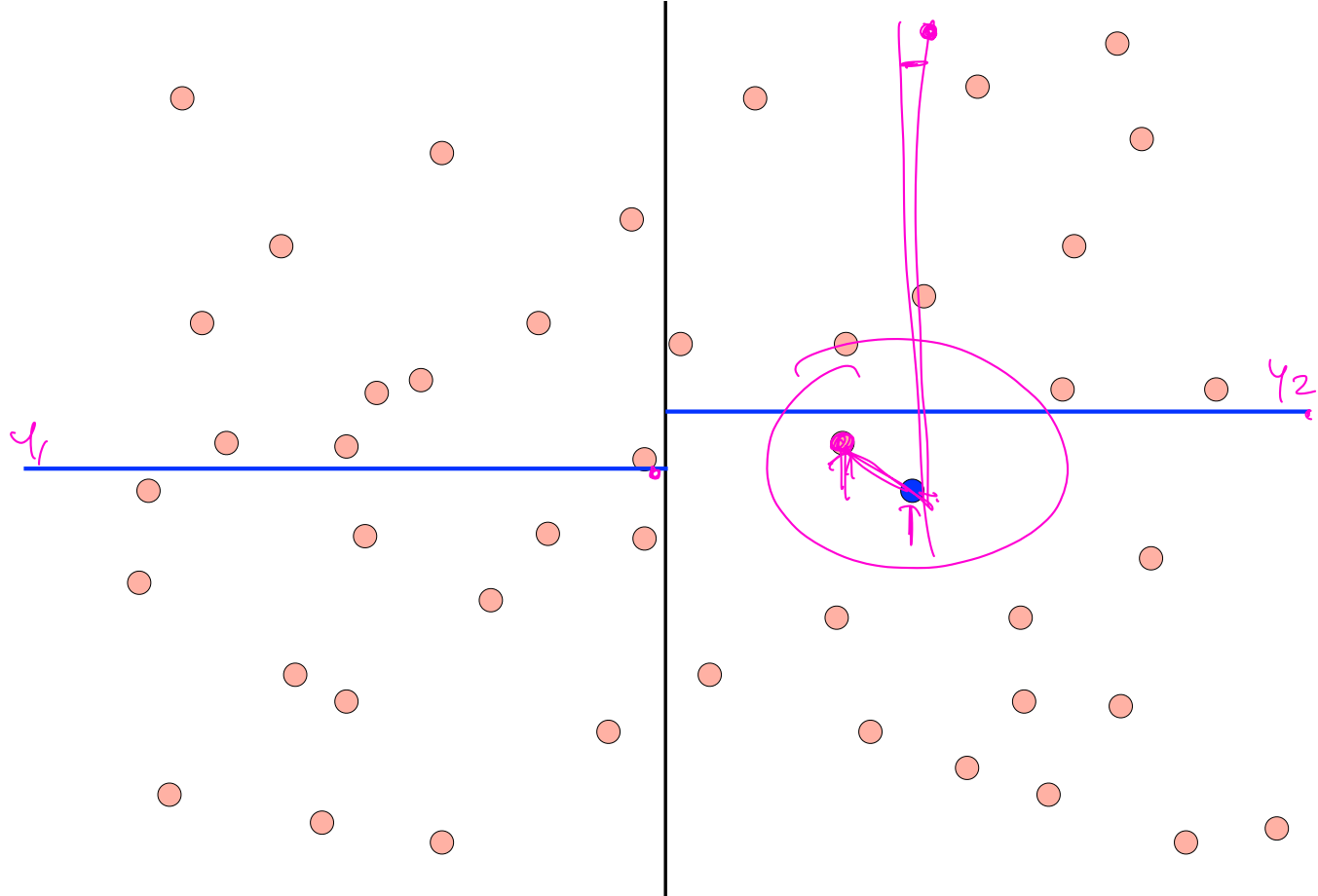


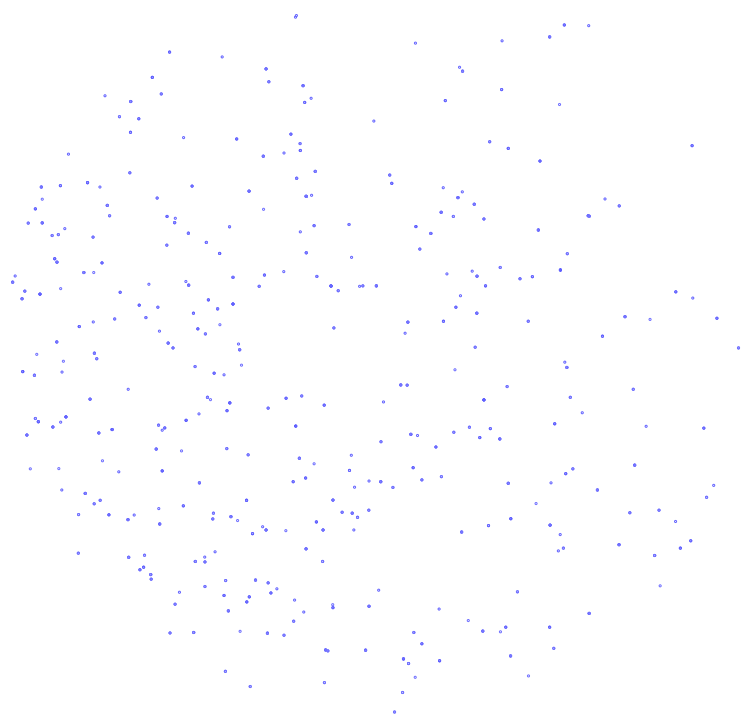
SORT ON
X coordinate

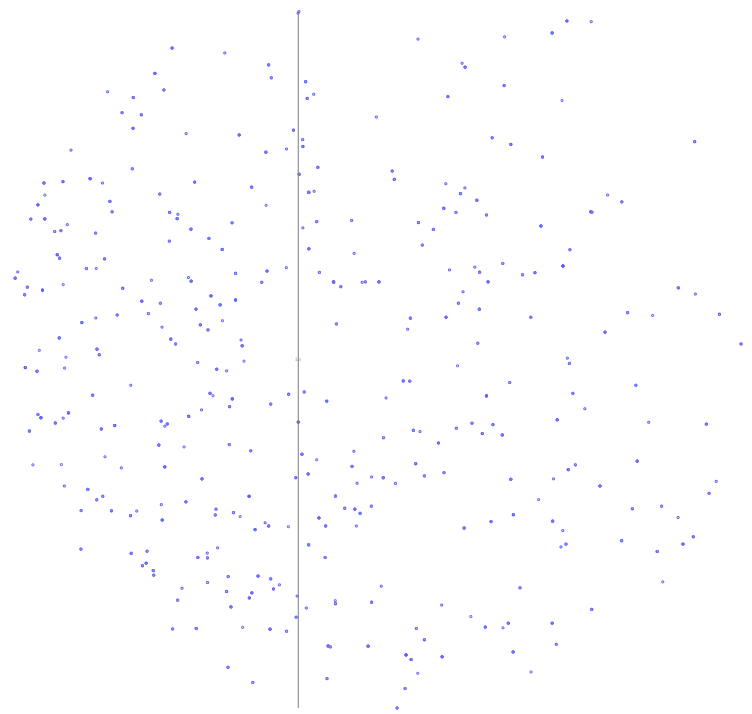
$\approx \frac{1}{2}$

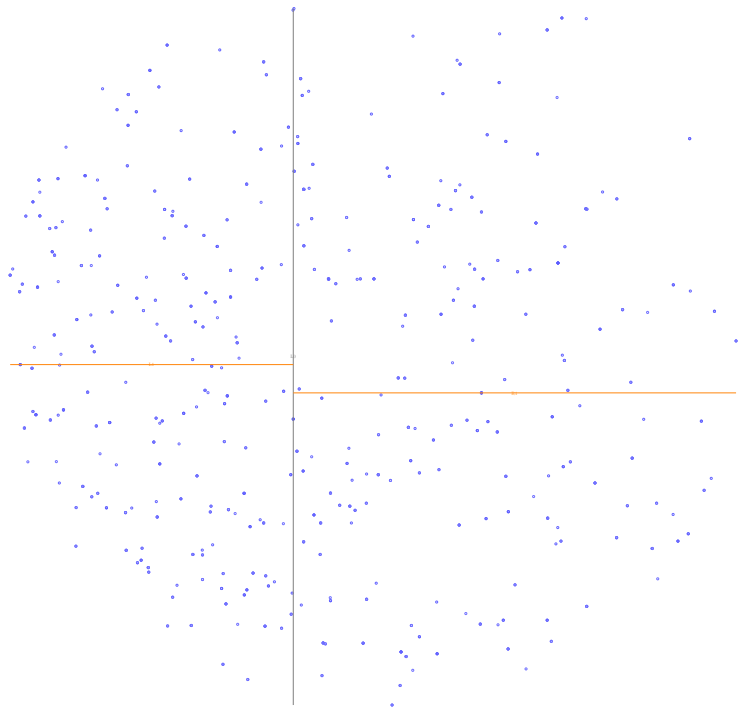
$\approx \frac{1}{2}$

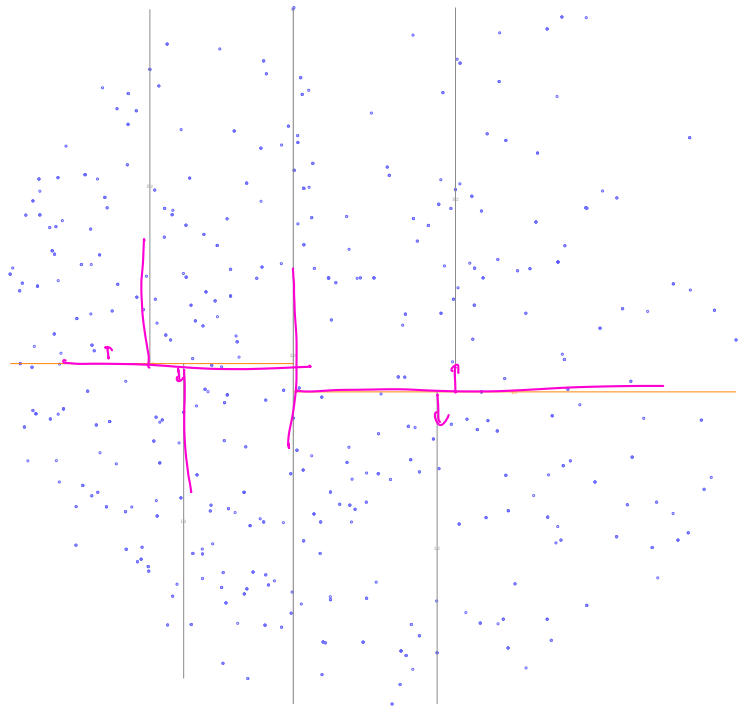


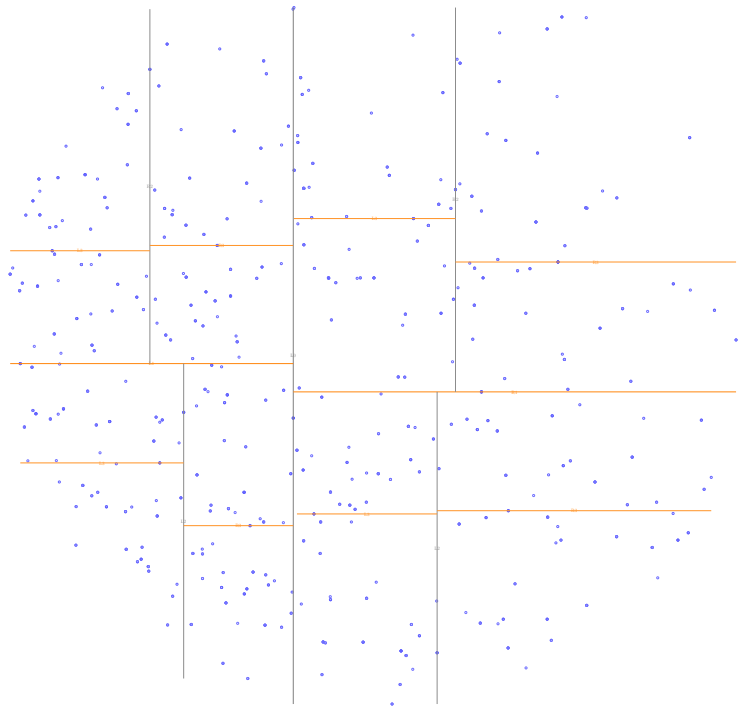


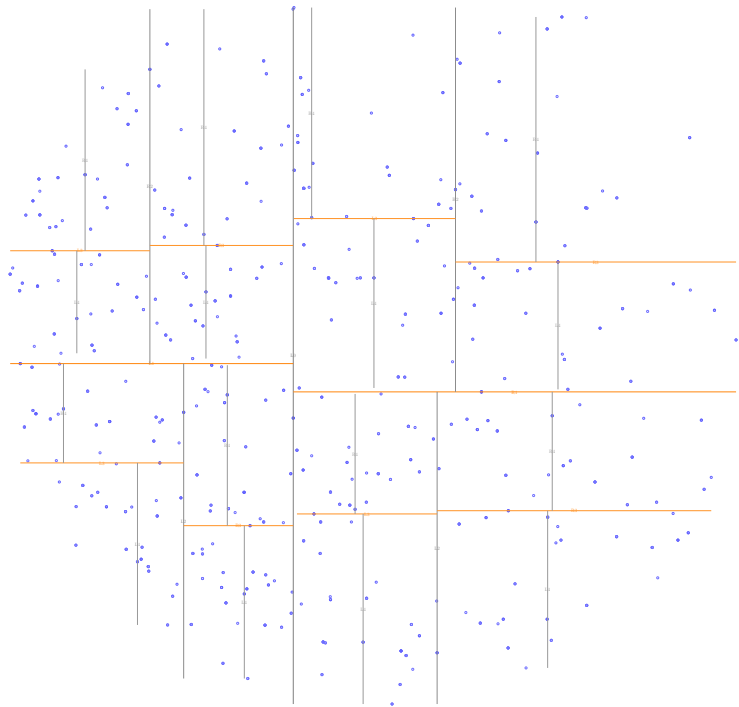


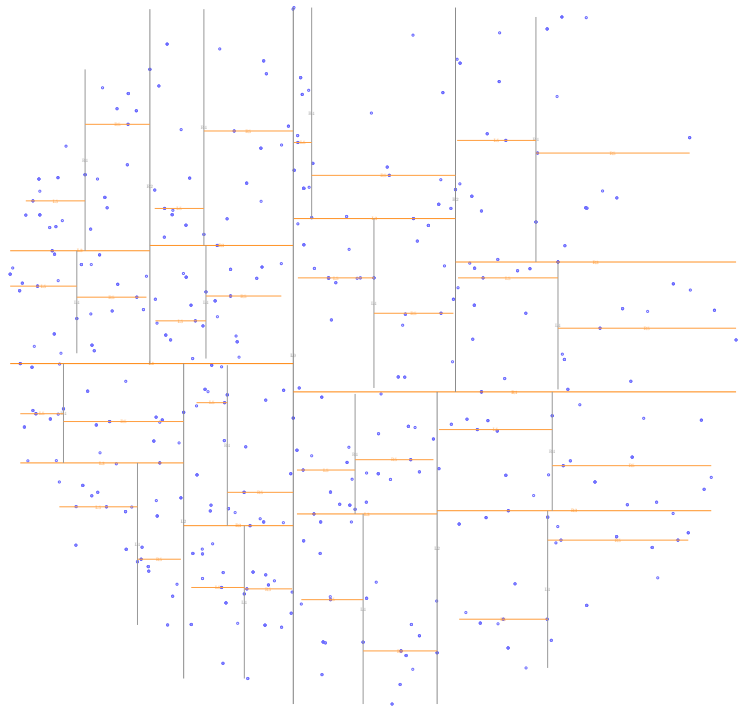








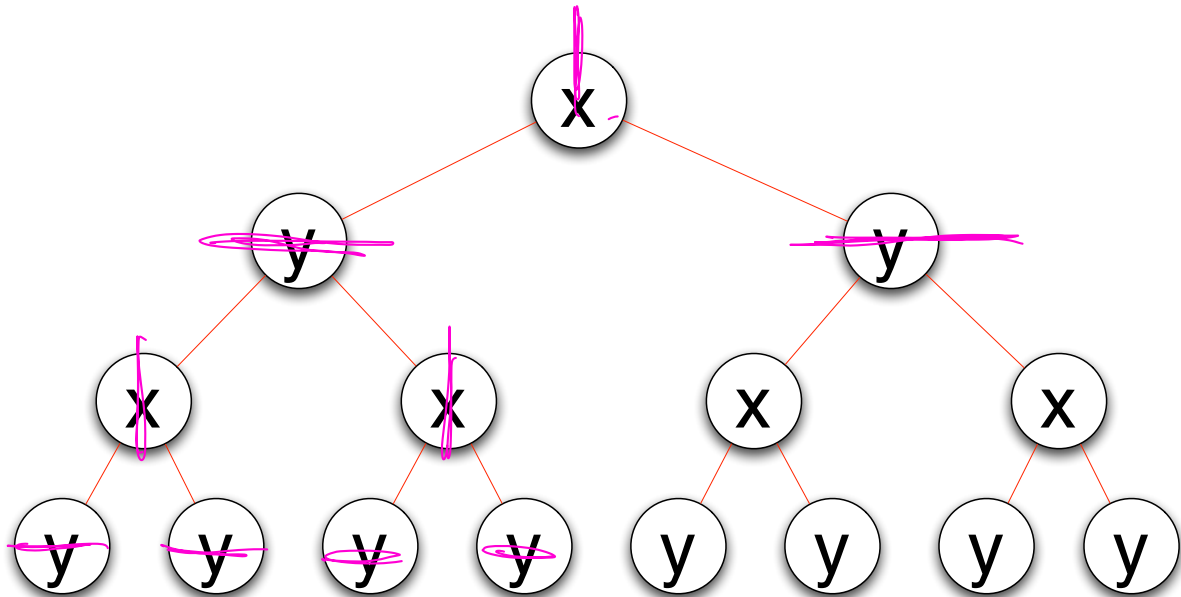


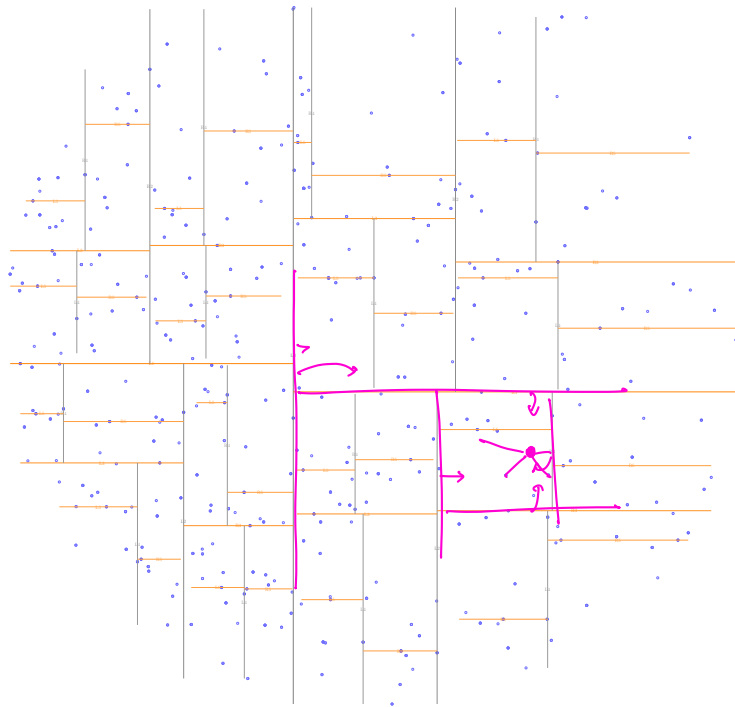


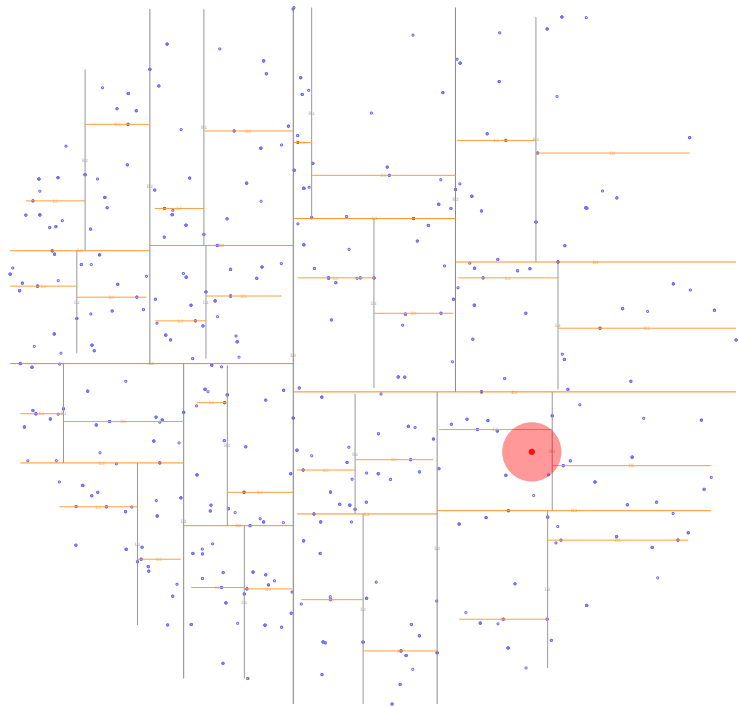
KD-Tree

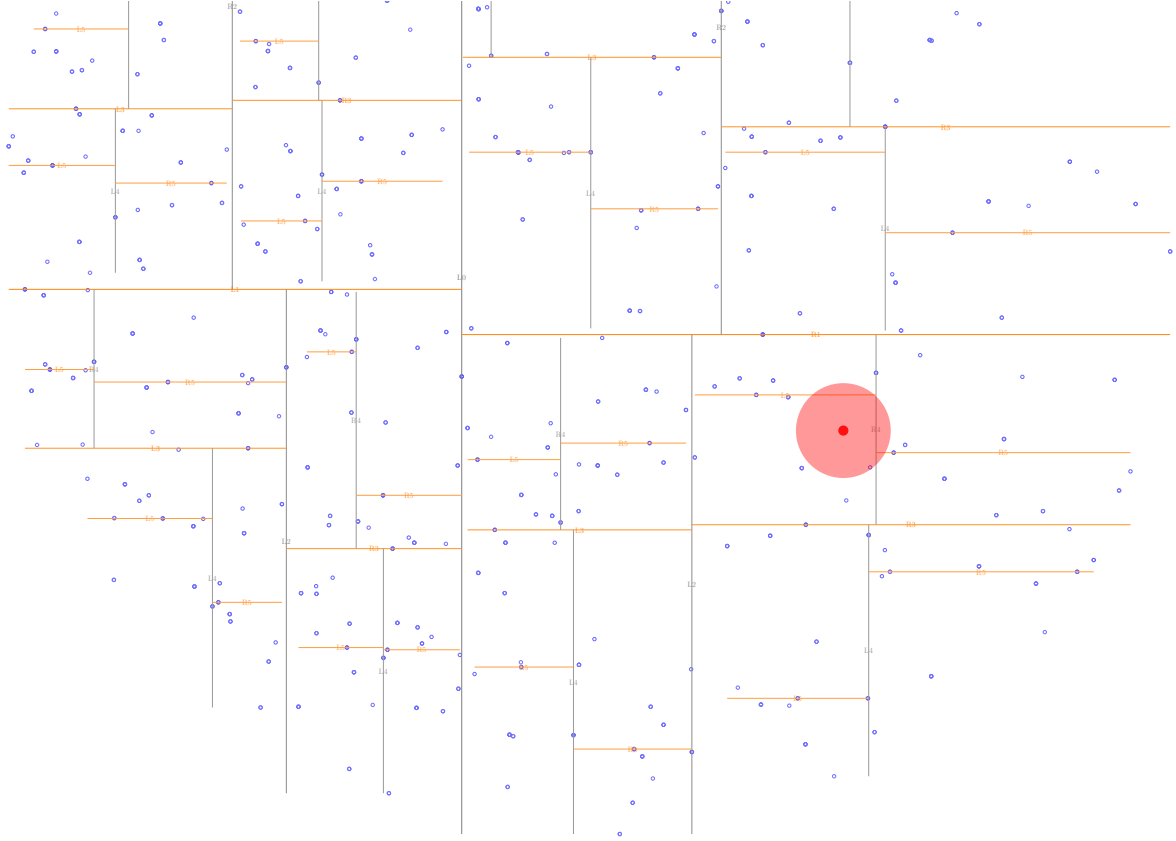
Each node in tree maintains variable “box”

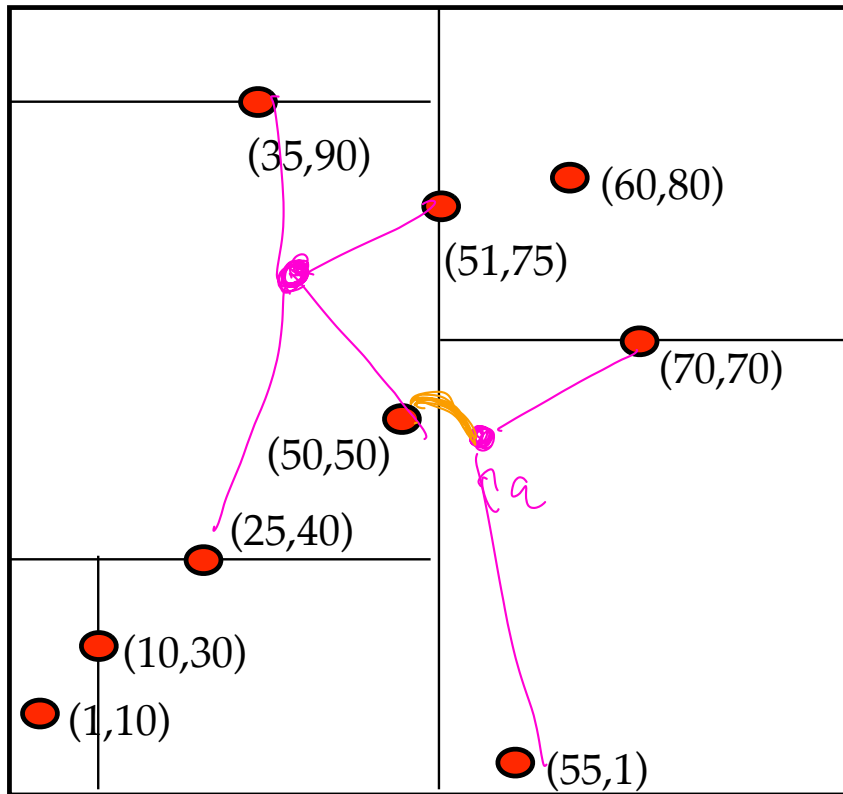
```
node {  
  rect box  
  point split  
  node* left  
  node* right  
}
```











NN(q, tree, dir, closest-so-far)

if empty(tree) or dist(q, tree.box) > closest return

→ if dist(q, tree.root) < closest { update closest} *heuristically*

if q.dir < tree.dir {

NN(q, tree.left, nextdir, closest)

NN(q, tree.right, nextdir, closest)

} else {

NN(q, tree.left, nextdir, closest)

NN(q, tree.right, nextdir, closest)

}

$$T(n) = \frac{2}{n} T\left(\frac{n}{2}\right) + \Theta(1)$$

$$\log(n)$$