

CSE520 Computational Geometry
Lecture 23
Quadtrees and Approximate Nearest Neighbor
Searching

Antoine Vigneron

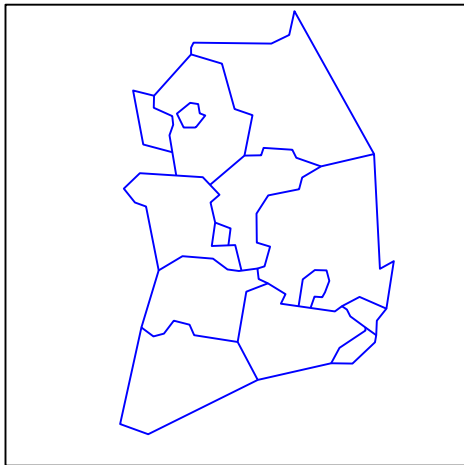
Ulsan National Institute of Science and Technology

June 16, 2020

Course Organization

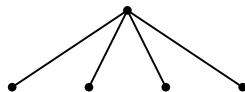
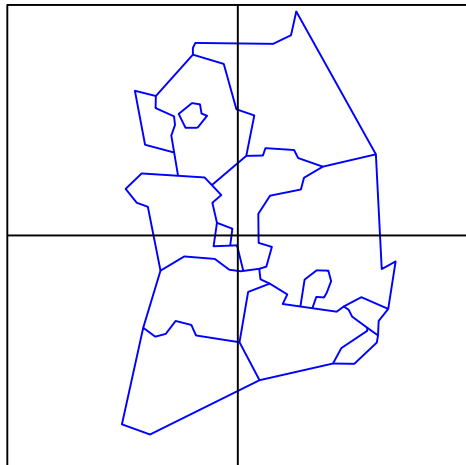
- Today, I will will present *quadtrees*.
- A quadtree is a tree data structure that is suitable for geometric data in low dimension.
- I will also present an application to near-neighbor searching.
- References:
- Sariel Har Peled's [book](#), chapters 2 and 17.

Example



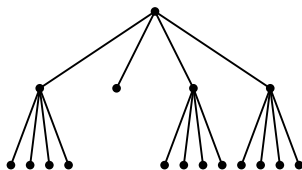
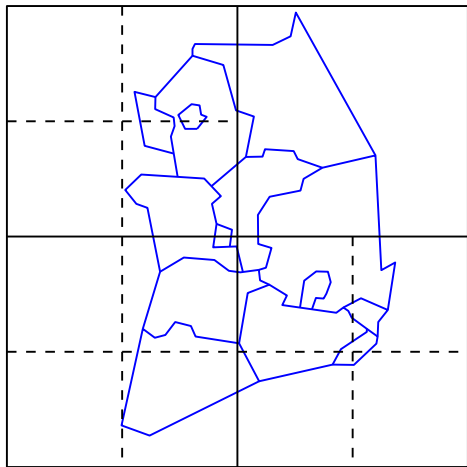
- Find a suitable square box containing the input subdivision.

Example



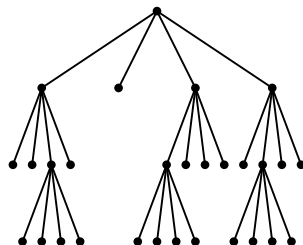
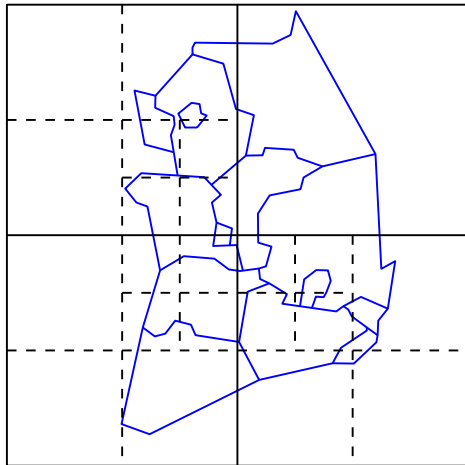
- Partition the outer box into 4 boxes.

Example



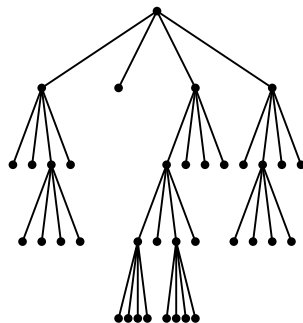
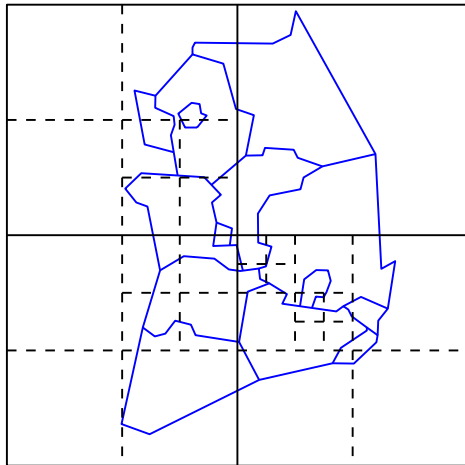
- Keep subdividing any box that overlaps with 5 regions or more.
- Remark: We could have used any other constant instead of 5.

Example



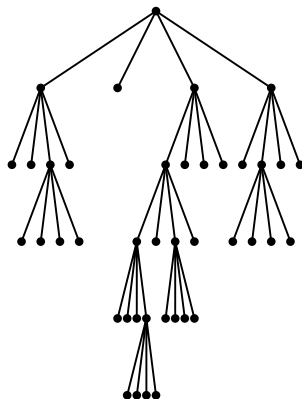
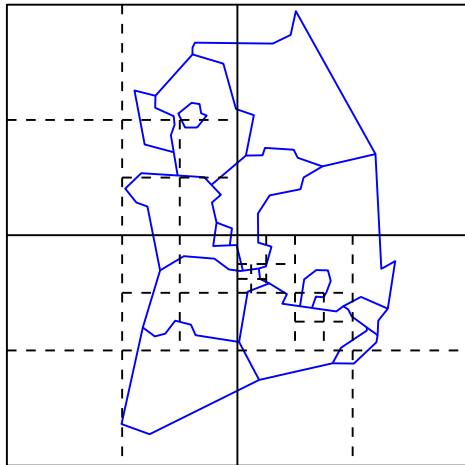
- 3 regions are subdivided.

Example



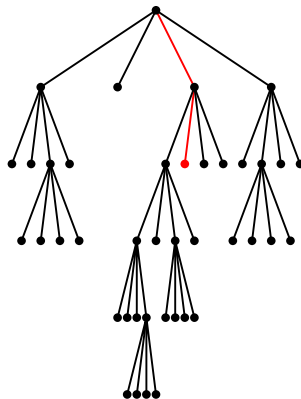
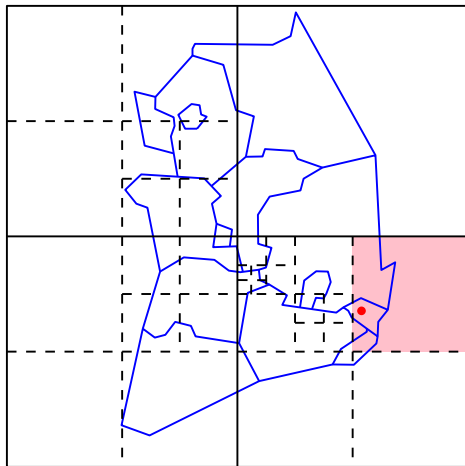
- 2 regions are subdivided.

Example



- Last step of the construction.

Example

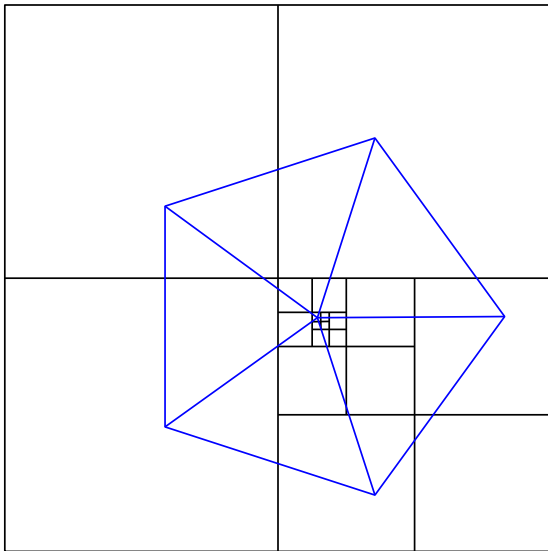


- In order to locate the red point, we go down to the leaf containing it, and check each of the region overlapping with this leaf box.

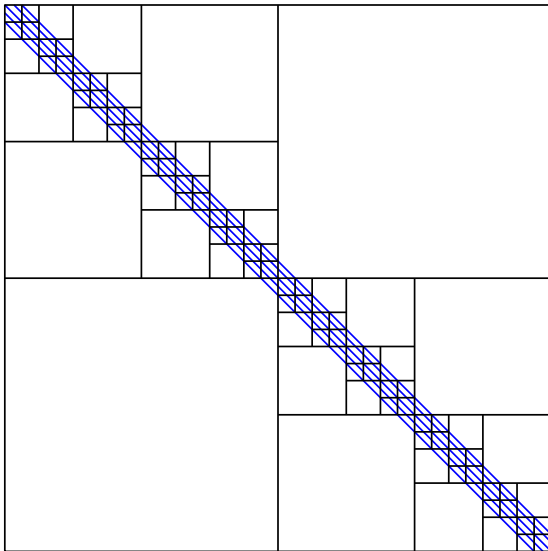
Quadtrees

- Each node of the graph is associated with a conflict list.
- The conflict list of a node u contains all the input regions that overlap with the box corresponding to u .
- Each box that overlaps with more than 4 input regions is subdivided into 4 boxes of half the size.
- We only need to record the conflict list of the leaves.
- This data structure is a *quadtree*.
- In practice, this data structure can be good if the regions have a nice shape.
- But for some input subdivision it may have a very large size. (See next two slides.)

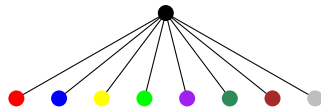
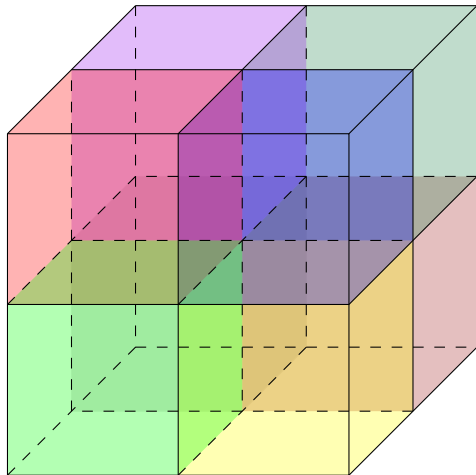
Quadtrees



Quadtrees



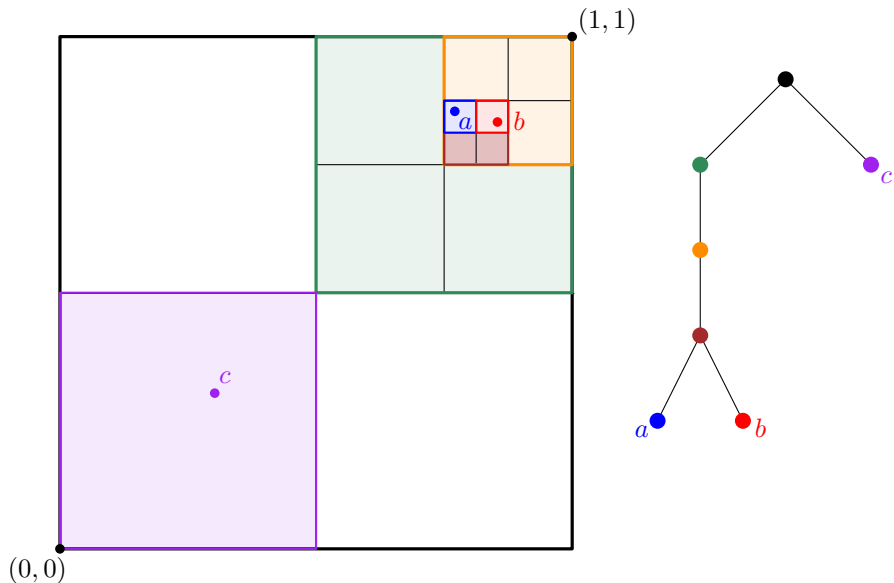
Quadtrees



Quadtrees

- Quadtrees generalize to arbitrary dimension.
- A box is split according to each dimension.
- A node can have up to 2^d children.
- So it is suitable for *fixed* dimension, i.e. $d = O(1)$.
- From now on, we will assume $d = O(1)$ and we will focus on quadtrees for set of points.
- We will subdivide a box whenever it contains more than one input point.

Recording a Point-Set in a Quadtree



Recording a Point-Set in a Quadtree

- Let P be a set of n points in \mathbb{R}^d .
- We assume that P is contained in the box $[0, 1]^d$.
- A *quadtree box* is a box obtained by recursively partitioning the box $[0, 1]^d$ into 2^d boxes. In other words, a quadtree box is of the form

$$\left[\frac{k_1}{2^i}, \frac{k_1 + 1}{2^i} \right] \times \cdots \times \left[\frac{k_d}{2^i}, \frac{k_d + 1}{2^i} \right]$$

where $k_1, \dots, k_d \in \{0, \dots, 2^i\}$ and $i \in \mathbb{N}$.

- The *level* of this box is i .

Recording a Point-Set in a Quadtree

- Each node v of the quadtree records a quadtree box denoted $\text{box}(v)$.
- The root node is associated with the box $[0, 1]^d$.
- Each leaf corresponds to exactly one point of P .
- Whenever $\text{box}(v)$ contains more than one point of P , we create 2^d children of v corresponding to the 2^d boxes obtained by partitioning $\text{box}(v)$.
- In order to save space, we delete the nodes v such that $\text{box}(v) \cap P = \emptyset$. (See Slide 15.)
- For instance, we replace them with a null pointer.

Recording a Point-Set in a Quadtree

- The *spread* of P is

$$\Phi(P) = \frac{\text{diam}(P)}{\text{dmin}(P)}$$

where

$$\text{diam}(P) = \max_{p,q \in P} d(p, q)$$

is the diameter of P and

$$\text{dmin}(P) = \min_{\substack{p,q \in P \\ p \neq q}} d(p, q).$$

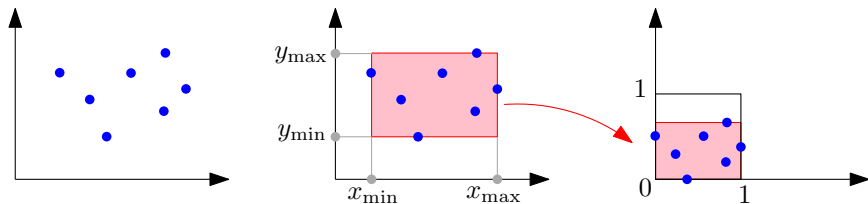
- When there is no ambiguity, we write Φ instead of $\Phi(P)$.

Recording a Point-Set in a Quadtree

Theorem

Let P be a set of n points contained in $[0, 1]^d$, and such that $\text{diam}(P) \geq 1$. Then the quadtree recording P has height $O(\log \Phi)$, size $O(n \log \Phi)$, and can be constructed in $O(n \log \Phi)$ time.

- We can always enforce the conditions $P \subset [0, 1]^d$ and $\text{diam}(P) \geq 1$ by scaling and translating P appropriately. It can be done in $O(n)$ time:



Proof

- We first prove that the height of the quadtree is $O(\log \Phi)$.
- Let u be an internal node at level i .
- Then $\text{box}(u)$ contains at least 2 points of P , so

$$\begin{aligned} \text{dmin}(P) &\leq \text{diam}(\text{box}(u)) \\ &= \frac{\sqrt{d}}{2^i} \leq \text{diam}(P) \frac{\sqrt{d}}{2^i} \end{aligned}$$

so $2^i \leq \Phi \sqrt{d}$, and thus $i \leq \log(\Phi) + \log(d)/2$.

- As $d = O(1)$, it means that $i = O(\log \Phi)$.
- So all internal nodes are at level $O(\log \Phi)$.
- It implies that the height of the quadtree is $O(\log \Phi)$.

Proof

- As the quadtree has height $O(\log \Phi)$ and has n leaves, it has at most n nodes per level, and thus it has $O(n \log \Phi)$ nodes.
- As each node records at most one point, this data structure uses $O(n \log \Phi)$ space.
- It can be constructed in $O(n \log \Phi)$ time by calling the procedure below with $Q = P$, $b = [0, 1]^d$ and $v = \text{root}$.

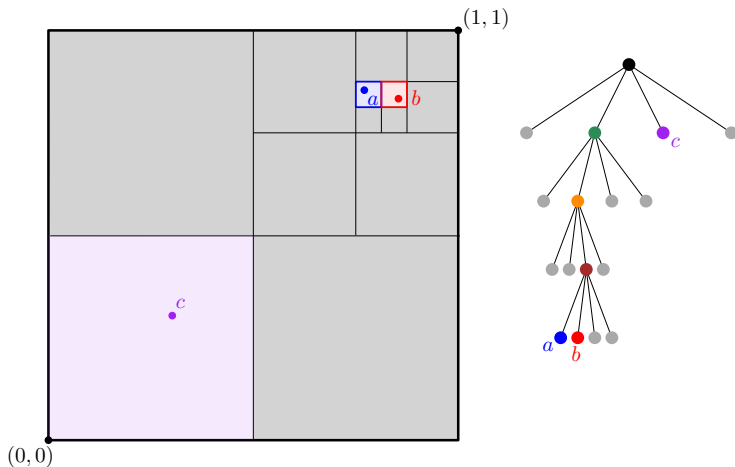
Proof

Pseudocode

```
1: procedure CONSTRUCTQUADTREE( $Q, b, v$ )
2:   if  $|Q| = 1$  then
3:     record the point  $q \in Q$  at node  $v$ 
4:     return
5:   let  $b_1, \dots, b_m$  be the  $m = 2^d$  sub-boxes of  $b$ 
6:   for  $i \leftarrow 1, m$  do
7:     if  $Q \cap b_i \neq \emptyset$  then
8:       create a child  $v_i$  of  $v$ 
9:       CONSTRUCTQUADTREE( $Q \cap b_i, b_i, v_i$ )
```

- As $d = O(1)$, each level of the quadtree contains at most n nodes, and the height of the quadtree is $O(\log \Phi)$, the total construction time is $O(n \log \Phi)$.

Point Location in a Quadtree



- In this section, we include in the quadtree the empty leaf nodes (grey).
- It increases the number of nodes by a factor less than 4.

Point Location in a Quadtree

Theorem

Given a quadtree of height h , we can preprocess it in linear time so that given a query point q , we can find the leaf containing q in $O(\log h)$ time.

- Remark: We will need to be able to perform hashing in $O(1)$ time per query, which is a common assumption in algorithms design.
- Remember that a quadtree box is of the form

$$\left[\frac{k_1}{2^i}, \frac{k_1 + 1}{2^i} \right] \times \cdots \times \left[\frac{k_d}{2^i}, \frac{k_d + 1}{2^i} \right]$$

where $k_1, \dots, k_d \in \{0, \dots, 2^i\}$ and $i \in \mathbb{N}$.

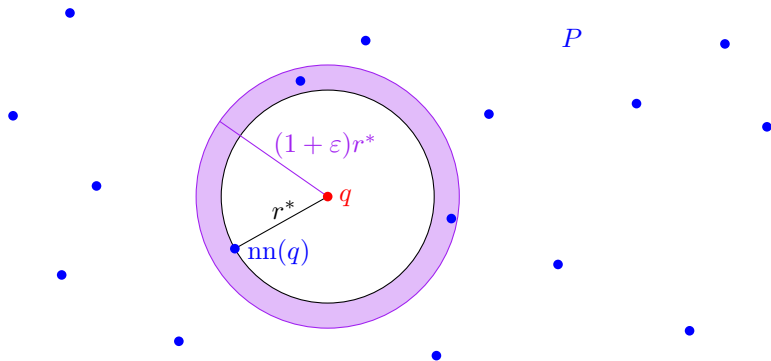
Point Location in a Quadtree

- We record $\text{box}(u)$ for all u in a hash table, using as keys the integers k_1, \dots, k_d described above.
- As $d = O(1)$, we can construct this data structure in linear time, and each query takes $O(1)$ time.
- If $q = (q_1, \dots, q_d)$, the quadtree box containing q at level i satisfies $k_i = \lfloor q_i/2^i \rfloor$ for all i .
- If there is a node u at level i containing q , we can find it in $O(1)$ time by hashing.
- So we can find the leaf node containing u in $O(h)$ time by trying each level separately.

Point Location in a Quadtree

- Better approach: Do binary search on the level of the leaf node containing q .
- So we check whether there is a node u at level $\lfloor h/2 \rfloor$ such that $q \in \text{box}(u)$.
- If it is not the case, recurse on the first $\lfloor h/2 \rfloor$ levels.
- If it is the case, and if this node is a leaf, return it.
- If it is an internal node, recurse on levels deeper than $\lfloor h/2 \rfloor$.
- This approach takes $O(\log h)$ time.

Approximate Nearest Neighbor Searching



- A $(1 + \varepsilon)$ -ANN query for q may return any of the points in the purple annulus.

Approximate Nearest Neighbor Searching

- P is still a set of n points in \mathbb{R}^d .
- Given a query point q , the distance from q to a closest point in P is denoted

$$r^* = \min_{p \in P} d(q, p).$$

- We denote by $\text{nn}(q)$ a closest point $p \in P$ such that $d(q, \text{nn}(q)) = r^*$. This point p is also called a *nearest neighbor*.
- Let $0 < \varepsilon \leq 1$ be a relative error ratio.
- A $(1 + \varepsilon)$ -*approximate nearest neighbor* of q is a point $p \in P$ such that $d(q, p) \leq (1 + \varepsilon)r^*$.
- We also write it $(1 + \varepsilon)$ -ANN.

Approximate Nearest Neighbor Searching

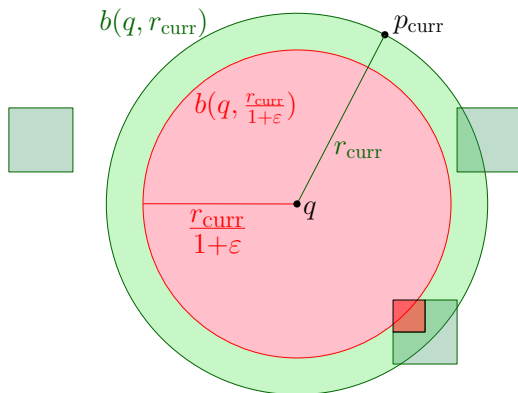
- Why should we use ANN searching?
- Alternative: exact near neighbor using Voronoi diagrams.
- Problem: in dimension d , the Voronoi diagram has combinatorial complexity $\Theta(n^{\lfloor d/2 \rfloor})$.
- In practice, we only compute it in dimension 2 or 3.
- But ANN queries can be answered efficiently using quadtrees in any fixed dimension:

Theorem

The quadtree described above allows us to answer $(1 + \varepsilon)$ -ANN queries in time $O(1/\varepsilon^d + \log \Phi)$.

- Remember we can compute this quadtree in $O(n \log \Phi)$ time.
- We will first describe the algorithm, and then analyze it.

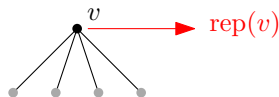
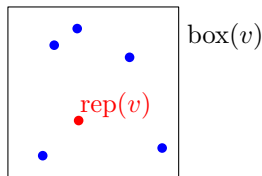
Algorithm



- Maintain a collection A_i of quadtree nodes at level i . (Green boxes.)
- Maintain the closest point found so far. (Point p_{curr} .)
- Recurse to the children of nodes in A_i whose box overlaps with the red ball. (Here, only the red box.)

Algorithm

- We denote by \mathcal{T} the quadtree recording P .
- An *empty* node v of \mathcal{T} is a node such that $\text{box}(v) \cap P = \emptyset$, i.e. the corresponding box contains no input point.



- For each non-empty node v of \mathcal{T} , record a *representative point* $\text{rep}(v) \in \text{box}(v) \cap P$.
- We can compute these representative points in time $O(n \log \Phi)$ for the whole quadtree by a bottom-up construction.

Algorithm

Pseudocode

```
1: procedure ANN(quadtrees  $\mathcal{T}$ , query point  $q$ ,  $\varepsilon$ )
2:    $A_0 \leftarrow \{\text{root of } \mathcal{T}\}$ ,  $p_{\text{curr}} \leftarrow \text{rep}(\text{root of } \mathcal{T})$ ,  $i \leftarrow 0$ 
3:   while  $A_i \neq \emptyset$  do
4:      $A_{i+1} \leftarrow \emptyset$ 
5:     for each node  $u \in A_i$  do
6:       for each non-empty child  $v$  of  $u$  do
7:         if  $d(q, \text{rep}(v)) \leq r_{\text{curr}}$  then
8:            $p_{\text{curr}} \leftarrow \text{rep}(v)$ ,  $r_{\text{curr}} \leftarrow d(q, \text{rep}(v))$ 
9:       for each node  $u \in A_i$  do
10:        for each non-empty child  $v$  of  $u$  do
11:          if  $\text{box}(v) \cap b(q, r_{\text{curr}}/(1 + \varepsilon)) \neq \emptyset$  then
12:            insert  $v$  into  $A_{i+1}$ 
13:      $i \leftarrow i + 1$ 
14:   return  $p_{\text{curr}}$ 
```


Proof of Correctness

- During our traversal of \mathcal{T} , we discard a non-empty node v , and the subtree rooted at it, if $\text{box}(v)$ does not overlap with the ball $b(q, r_{\text{curr}}/(1 + \varepsilon))$.
- Let w be the node discarded by the algorithm such that $\text{nn}(q) \in \text{box}(w)$.
- Let c be the point of $\text{box}(w)$ that is closest to q .

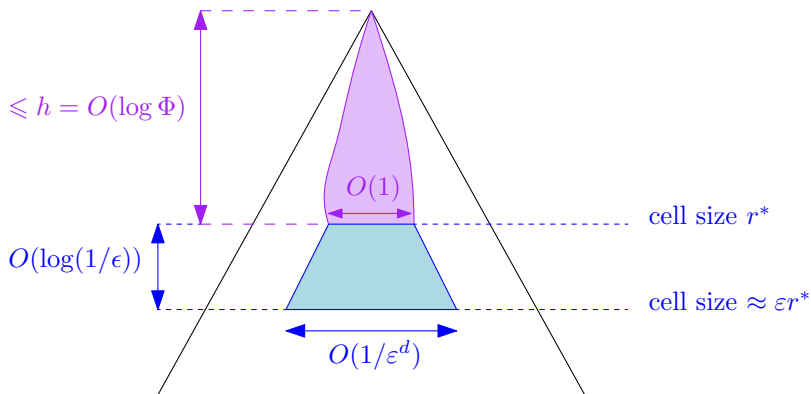
$$\begin{aligned} r^* = d(q, \text{nn}(q)) &\geq d(q, c) \\ &\geq \frac{r_{\text{curr}}}{1 + \varepsilon} \end{aligned}$$

so p_{curr} is a $(1 + \varepsilon)$ -ANN.

- After that, p_{curr} may be updated, but it only gets closer to q , so it is still a $(1 + \varepsilon)$ -ANN.

Analysis

- Idea: At lower levels i , the set A_i contains $O(1)$ nodes.
- At later stages, it contains $O(1/\varepsilon^d)$ nodes in total.



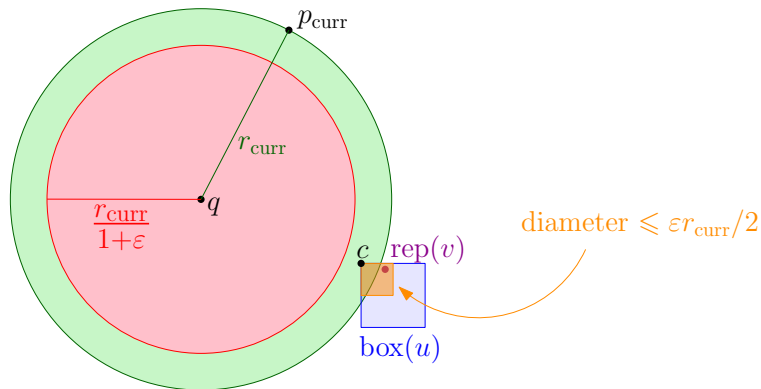
Analysis

Lemma

During the course of the algorithm, the size of the quadtree boxes of the nodes we visit is $\Omega(\varepsilon)$. More precisely, we have $\frac{1}{2^i} \geq \frac{\varepsilon r^}{4\sqrt{d}}$.*

- Suppose that, at line 9 and 12 of the algorithm, we visit a child v of a node $u \in A_i$.
- Suppose that $\text{box}(v)$ has side length $1/2^{i+1} < \varepsilon r^*/(4\sqrt{d})$.
- Then $\text{box}(u)$ has side length $1/2^i < \varepsilon r^*/(2\sqrt{d})$.
- As $r^* \leq r_{\text{curr}}$, it implies $1/2^i \leq \varepsilon r_{\text{curr}}/(2\sqrt{d})$.

Analysis



- Therefore the diameter $\sqrt{d}/2^i$ of $\text{box}(u)$ satisfies $\frac{\sqrt{d}}{2^i} \leq \frac{\epsilon r_{\text{curr}}}{2}$.
- As v is visited at lines 5–8, we must have $d(q, \text{rep}(v)) \geq r_{\text{curr}}$.

Analysis

- Let c be the point of $\text{box}(v)$ that is closest to q .

$$\begin{aligned} d(q, c) &\geq d(q, \text{rep}(v)) - \text{diam}(\text{box}(v)) && \text{by triangle inequality} \\ &\geq r_{\text{curr}} - \frac{\varepsilon r_{\text{curr}}}{2} \\ &= \left(1 - \frac{\varepsilon}{2}\right) r_{\text{curr}} \\ &\geq \frac{r_{\text{curr}}}{1 + \varepsilon} && \text{because } 0 < \varepsilon \leq 1 \end{aligned}$$

- It means that v will not be inserted in A_{i+1} .
- So the algorithm does not reach stage $i + 1$, which completes the proof of the lemma.

Analysis

Lemma

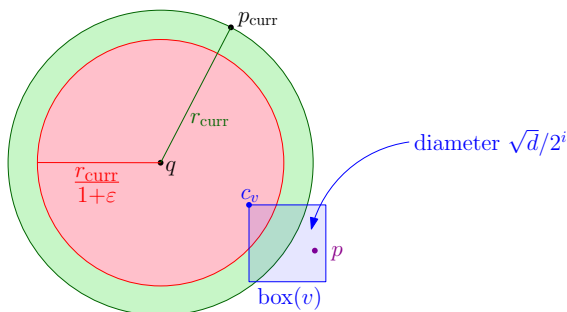
For any node $v \in A_i$, we have $\text{box}(v) \subset b(q, \ell_i)$ where

$$\ell_i = \frac{2\sqrt{d}}{2^i} + (1 + \varepsilon)r^*.$$

- We now prove this lemma.
- There are two cases.
- **Case 1:** There exists $w \in A_i$ such that $\text{nn}(q) \in \text{box}(w)$.
- Since $\text{diam}(\text{box}(w)) = \sqrt{d}/2^i$, we have $d(\text{nn}(q), \text{rep}(w)) \leq \sqrt{d}/2^i$.
- It follows that

$$d(q, \text{rep}(w)) \leq d(q, \text{nn}(q)) + d(\text{nn}(q), \text{rep}(w)) \leq r^* + \frac{\sqrt{d}}{2^i}$$

Analysis



- So at Line 5 of the $i - 1$ th step, just before we construct A_i , we have

$$r_{\text{curr}} \leq d(q, \text{rep}(w)) \leq r^* + \frac{\sqrt{d}}{2^i}$$

- For any node $v \in A_i$, the ball $b(q, r_{\text{curr}}/(1 + \epsilon))$ must intersect $\text{box}(v)$, so there is a point $c_v \in \text{box}(v)$ such that $d(q, c_v) \leq r_{\text{curr}}$.

Analysis

- Then any point $p \in \text{box}(v)$ satisfies

$$\begin{aligned} d(q, p) &\leq d(q, c_v) + d(c_v, p) && \text{by the triangle inequality} \\ &\leq r_{\text{curr}} + \sqrt{d}/2^i && \text{see previous slide} \\ &\leq r^* + 2\sqrt{d}/2^i \\ &\leq \ell_i \end{aligned}$$

- In other words, $\text{box}(v) \subset b(q, \ell_i)$, which completes the proof of Case 1.

- **Case 2:** For all $v \in A_i$, we have $\text{nn}(q) \notin \text{box}(v)$.
- So we have discarded the box containing $\text{nn}(q)$, and thus by the argument on Slide 33, we have $r_{\text{curr}} \leq (1 + \varepsilon)r^*$.
- It follows that $\text{box}(v) \subset b(q, (1 + \varepsilon)r^* + \sqrt{d}/2^i) \subset b(q, \ell_i)$.
- This completes the proof of the lemma on Slide 38.

Analysis

Lemma

There are two constants α_d, β_d depending on d only such that

- (a) If $1/2^i > r^*$, then $|A_i| \leq \alpha_d$.*
- (b) If $1/2^i \leq r^*$, then $|A_i| \leq \beta_d(2^i r^*)^d$.*

- **Proof of (a):** By the lemma on Slide 38, for any $v \in A_i$, we have $\text{box}(v) \subset b(q, \ell_i)$ where $\ell_i = \frac{2\sqrt{d}}{2^i} + (1 + \varepsilon)r^*$.
- As $\varepsilon < 1$ and $1/2^i > r^*$, we have $\ell_i < 2(1 + \sqrt{d})/2^i$.
- The volume of $b(q, \ell_i)$ is $C_d \ell_i^d$ for some constant C_d depending on d only. The volume of $\text{box}(v)$ for any $v \in A_i$ is $(1/2^i)^d$.
- As these boxes are disjoint, their number is at most

$$\frac{C_d \ell_i^d}{(1/2^i)^d} < C_d(2(1 + \sqrt{d}))^d.$$

- So we can just choose $\alpha_d = C_d(2(1 + \sqrt{d}))^d$.

Analysis

- **Proof of (b):** By the lemma on Slide 38, for any $v \in A_i$, we have $\text{box}(v) \subset b(q, \ell_i)$ where $\ell_i = \frac{2\sqrt{d}}{2^i} + (1 + \varepsilon)r^*$.
- Since $1/2^i \leq r^*$ and $\varepsilon < 1$, we have $\ell_i < 2(\sqrt{d} + 1)r^*$.
- So the volume of each $\text{box}(v)$ where $v \in A_i$ is at most $(1/2^i)^d$.
- So the number of these boxes is at most

$$\frac{C_d(2(\sqrt{d} + 1)r^*)^d}{(1/2^i)^d} = \beta_d(2^i r^*)^d$$

where β_d is a constant depending only on d .

Analysis

Lemma

Our algorithm returns a $(1 + \varepsilon)$ -ANN in time $O(1/\varepsilon^d + \log(1/r^))$*

- In the first part of the traversal, when $1/2^i > r^*$, we have $i < \log(1/r^*)$, so we traverse at most $\lceil \log(1/r^*) \rceil$ levels. By previous lemma (a), we visit $O(1)$ nodes per level, so it takes time $O(\log(1/r^*))$.
- In the second part of the traversal, by the lemma on Slide 35, we have

$$\frac{\varepsilon r^*}{4\sqrt{d}} \leq \frac{1}{2^i} \leq r^*$$

and at each level we visit at most $\beta_d(2^i r^*)^d$ nodes.

Analysis

- So the number of nodes we visit is at most $\sum_{i=0}^m \beta_d(r^*)^d (2^d)^i$ where $m = \lfloor \log(4\sqrt{d}/(\varepsilon r^*)) \rfloor$.
- This is the sum of a geometric series with common ratio 2^d , so it is dominated by its last term. Therefore, this sum is

$$O\left((r^*)^d (2^m)^d\right) = O\left((r^*)^d (4\sqrt{d}/(\varepsilon r^*))^d\right) = O(1/\varepsilon^d).$$

- So overall, the algorithm visits $O(1/\varepsilon^d + \log(1/r^*))$ nodes, which completes the proof of the lemma.
- We prove Theorem 29 by using the upper bound $i \leq h = O(\log \Phi)$ for the first part of the traversal.