

Nearest Neighbor Search

Input: $S \subseteq \mathbb{R}^d$ $|S| = n$

Output: Search data structure DS

$\forall x \in \mathbb{R}^d$ DS(x) $\in S$ closest point

Cost:

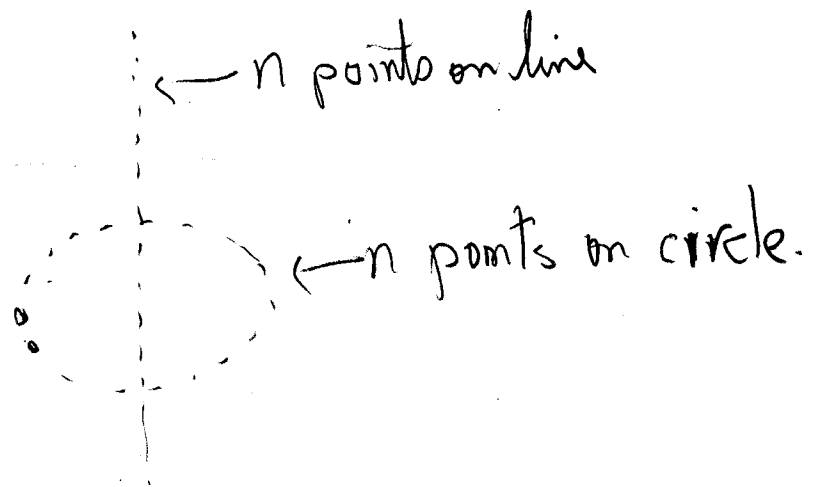
- 1) Time to compute DS
- 2) Size of DS
- 3) Query time

Equivalent to Finding the Voronoi cell containing x .

Prob: # of Voronoi points on ∂ maybe of size $O(n^{1+d/2})$

of Delaunay Simplices

Preparata's Ex.



Claim: neig pair on circle + neig pair on line
form a Del-Tet

$$\therefore n^2 \text{ Tets} \quad d=3 \quad n^{\lceil d/2 \rceil} = n^2$$

Naive DS may have size $\Omega(n^2)$ in 3D.

Work around: Approximate NN

Def $p \in S$ is a $(1+\epsilon)$ -approx NN of q if

$$\frac{\text{dist}(p, q)}{\text{dist}(p', q)} \leq (1+\epsilon) \quad \forall p' \neq p \in S$$

Goal for d fixed

$O(n \log n)$ preprocessing

$O(n)$ space DS

$O(\log n)$ query time $(1+\epsilon)$ -approx

4

Assume $P \subseteq \mathbb{R}^d$ $P \subseteq \text{Unit box}$

Build Quadtree (P)

Unit $B \equiv \text{Unit Box}$ & $P \subseteq P$

Root = (B, P) ; $\text{rep}(B) = P$

While a box $|B' \cap P| \geq 2$

Split B' into B_1, B_2, B_3, B_4

for each i s.t. $|B_i \cap P| \geq 1$ pick $P_i \in B_i \cap P$

add box (B_i, P_i) ; $\text{rep}(B_i) = P_i$

ANN-Search(g, Q) $g \equiv \text{point}$
 $Q \equiv \text{Quad tree DS}$

Init $\text{Current} = \text{rep}(B)$
 $r_{\text{cur}} = d(g, \text{current})$
 $\text{Queue}_1 = \{ \text{children}(B) \} \quad i=1$

While $\text{Queue}_i \neq \emptyset$
 Search(Queue_i); $i = i+1$ (**)

Return current

Search(Queue_i)

While $\text{Queue}_i \neq \emptyset$ do $B = \text{Dequeue}(\text{Queue}_i)$

1) If $d(g, \text{rep}(B)) < r_{\text{cur}}$ then

$\text{current} = \text{rep}(B)$; $r_{\text{cur}} = d(g, \text{rep}(B))$

2) If $d(g, \text{rep}(B)) - \text{dia}(B) < (1 - \epsilon/2)r_{\text{cur}}$ (*)
 then add $\text{children}(B)$ to Queue_{i+1}

Claim If $d(B) < (\frac{\epsilon}{2})\delta$ then B is Discarded.

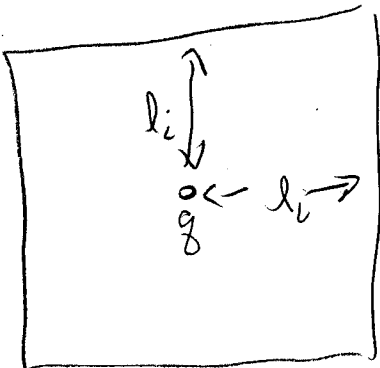
$$\begin{aligned} d(q, \text{rep}(B)) - d(B) &\geq d(q, \text{rep}) - (\frac{\epsilon}{2})\delta \\ &\geq r_{\text{cur}} - (\frac{\epsilon}{2})r_{\text{cur}} \\ &= (1 - \frac{\epsilon}{2})r_{\text{cur}} \end{aligned}$$

Claim After Search (Q_i) $l_i = \delta + d(B_i) = \delta + \sqrt{d}2^{-i} \geq r_{\text{cur}}$.

Consider box B s.t. $n(q) \in B$

$$r_{\text{cur}} \leq d(q, \text{rep}(B)) \leq \delta + d(B)$$

Claim Live box at time i contained in $2l_i \times 2l_i$ box \bar{B} centered at q_i .



Correctness

To show It is OK to discard a box

ie Let B be box containing $NN(B) \in \mathcal{P}$.

If B is discarded then $\text{rep}(B)$ is an $(1+\epsilon)$ -ANN.

Suppose B is discarded

$$(1-\epsilon/2)r_{\text{cur}} \leq d(g, \text{rep}(B)) - \text{dia}(B)$$

$$\leq d(g, nn(B)) = \delta$$

$$\frac{\delta}{(1-\epsilon/2)} \geq r_{\text{cur}}$$

$$\text{but } \frac{1}{(1-\epsilon/2)} \leq 1+\epsilon \quad \text{for } 0 < \epsilon \leq 1$$

$$r_{\text{cur}} \leq (1+\epsilon)\delta$$

$$\text{Vol}(\bar{B}) = (2l_i)^d$$

$$\text{Vol}(B_i) = \left(\frac{l_i}{2}\right)^d \quad n_i \leq \left(\frac{2l_i}{l_i}\right)^d$$

$$= \left\lfloor \frac{2\left(\delta + \sqrt{d} \frac{\delta}{2^i}\right)}{\frac{l_i}{2^i}} \right\rfloor^d = O\left(\left(1 + \frac{\delta}{2^i}\right)^d\right) = O\left(1 + 2^i \delta\right)^d$$

$$= O\left(1 + (2^i \delta)^d\right)$$

Total # Box

$$\sum_{i=0}^h n_i = O\left(\sum_0^{\lfloor \log(\frac{\epsilon}{2}\delta) \rfloor} \left[1 + (2^i \delta)^d\right]\right)$$

$$= O\left(\log\left(\frac{1}{\epsilon\delta}\right) + \left(\frac{\delta}{\frac{\epsilon}{2}\delta}\right)^d\right)$$

$$= O\left(\log\left(\frac{1}{\epsilon\delta}\right) + \frac{1}{\epsilon^d}\right)$$

Lemma $P \subseteq \mathbb{R}^d$ is n points \wedge $Q \in \Sigma$ of P
 $\text{dia}(P) \approx 1$, then $(1+\epsilon)$ -ANN queries in

$$O\left(\frac{1}{\epsilon^d} + \log(\bar{\Phi})\right) \quad \bar{\Phi} \equiv \text{spread of } P.$$