

# Approximate Nearest Neighbors via Point Location in Balls

- **Nearest Neighbor Problem:** Given a set of  $n$  points,  $P$  in a metric space  $M$ ; preprocess it to be able to determine the closest point in  $P$ , to a query point  $q \in M$ .

Requires prohibitive pre-processing time. Need to compute Voronoi diagrams and process for point location queries.  $O(n^{\lceil d/2 \rceil})$  time.

So we ..... Approximate!

Let  $m(q, P)$  be said nearest neighbor, and  $d(q, P)$  be the distance between  $q$  and  $m(q, P)$ .

A point  $s \in M$  is a  $(1+\epsilon)$ -approximate nearest neighbor if  $d_M(q, s) \leq (1+\epsilon) \cdot d(q, P)$ .

- Ball is defined by its center  $p$  & radius  $r$ .

$$b(p, r) = \{q \in M \mid d_M(p, q) \leq r\}$$

- Target Ball of  $q$  in a given set  $B$ ,  $\odot_B(q)$ , is the smallest ball in  $B$  that contains the query  $q$ .

- **Objective:** Show that  $(1+\epsilon)$ -ANN queries can be reduced to target ball queries on a "small" # of balls.

- **Notation:** Union of balls of radius  $r$ , centered at points of  $P$ .  $\bigcup_{\text{balls}}(P, r) := \bigcup_{p \in P} b(p, r)$ .

⇒ Quadratic # of balls

Lemma 841: Given  $P, q$  in  $M$ . Let  $B = \bigcup_{i=-\infty}^{\infty} \mathcal{U}_{\text{balls}}(P, (1+\epsilon)^i)$ .  
 Let  $b = \mathcal{O}_B(q)$  be the target ball and  $p \in P$  its center.  
 Then  $p$  is  $(1+\epsilon)$ -ANN in  $P$  to  $q$ .

Pf: Let  $s = \text{nn}(q, P)$  be the nearest neighbor, and  $i$  be the index such that

$$r_1 := (1+\epsilon)^i < d_M(s, q) \leq (1+\epsilon)^{i+1} =: r_2$$

No ball of radius  $r_1$  can contain the query, but  $\exists$  a ball  $b(s, (1+\epsilon)^{i+1})$  that does contain it.  
 So, the target ball must have radius  $\leq (1+\epsilon)^{i+1}$ .

$$\alpha = \frac{r(\text{target ball})}{r(\text{OPT ball})} \leq \frac{(1+\epsilon)^{i+1}}{(1+\epsilon)^i} \leq 1+\epsilon$$



So, is the objective achieved? Nope as balls is not "small number" by any measure.

Consider two points  $u, v$ .

Case 1:  $d_M(q, u) \geq \frac{2d_M(u, v)}{\epsilon}$



Case Close(d):  $d_M(q, u) \leq d_M(u, v)/4$

Intuitively, only  $d(q, \{u, v\}) \in \underbrace{[1/4, 2/\epsilon]}_{\mathcal{Z}(u, u)} d_M(u, v)$  hard to decide.

→ Construction.

$$B(u, v) = \{ b(u, r), b(v, r) \mid r = (1+\epsilon)^i \in Z(u, v) \text{ and } i \in \mathbb{Z} \}$$

For every pair  $u, v \in P$ , add balls  $B(u, v)$  to the set  $B$  being constructed.

→ Size  $\in O(n^2 \epsilon^{-1} \log \epsilon^{-1})$ .

Ratio between high and low values

$$R = \frac{2d_M(u, v)}{\epsilon} \Bigg/ \frac{d_M(u, v)}{4} \in O(1/\epsilon)$$

Total no. of balls is,

$$O(n^2 \log_{1+\epsilon} R) = O(n^2 \epsilon^{-1} \log R) = O(n^2 \epsilon^{-1} \log \epsilon^{-1})$$

→ Correctness.

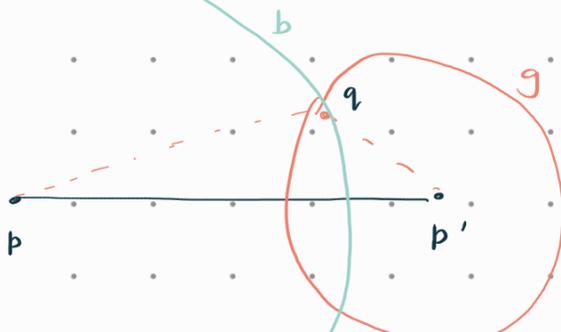
Let  $b = O_B(q)$  be the target ball, and  $p$  its center.

Let  $p' = nn(q, P)$  be the nearest neighbor.

• Far Case:  $d_M(q, p') \geq 2d_M(p, p')/\epsilon$

$$\frac{d(q, p)}{d(q, p')} \leq \frac{d(q, p') + d(p', p)}{d(q, p')} \leq 1 + \frac{\epsilon}{2}$$

• Close Case:  $d_M(q, p') \leq d_M(p, p')/4$



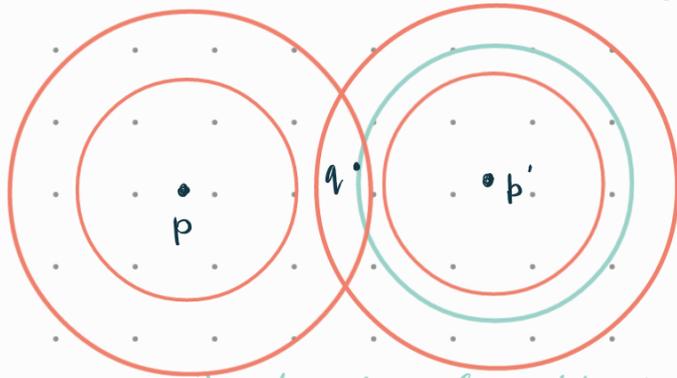
$$\text{rad}(g) = d_M(p', q) \leq d_M(p, p')/4$$

$$\text{but } \text{rad}(b) \geq 3/4 d_M(p, p')$$

$b$  can never be returned  $\Rightarrow \epsilon$ .

• Bad Case:  $d_m(q, p') \in [1/4, 2/\epsilon]$   $d_m(p, p')$

then by construction, there is a good ball.



Are we done now? Kind of. Maybe we can still improve on the quadratic size?

⇒ Auxilliaris

→ Proximity Queries (a.k.a. Near Neighbor Data Structure)

• Given point set  $P$ , query pt.  $q$ , and  $\epsilon > 0 \in \mathbb{R}$ .

$D_{\text{near}}(P, \epsilon)$  is a data structure, which decides whether  $d(q, P) \leq \epsilon$  or  $d(q, P) > \epsilon$ .

In the former case, it returns a witness point.

• Trivial realization: balls of radius  $\epsilon$ , with centers at  $P$ .

Target ball query for checking.

→ Interval Query (a.k.a. Interval Near Neighbor) to check whether query lies in dist.  $[a, b]$  from  $P$ .

$\hat{I}(P, a, b, \epsilon) := \{N_0, N_1, \dots, N_M\}$

where  $N_i = D_{\text{near}}(P, \epsilon_i = \min((1/2)^i a, b))$ , &  $M = \lceil \log_{1/2}(b/a) \rceil$ .

• Constructed with  $O(\epsilon^{-1} \log(b/a))$  PQ data structures.

- Contains  $O(\epsilon^{-1} n \log(b/a))$  balls overall.
- Requires 2 queries to report not in range, and requires  $O(\log(\epsilon^{-1} \log(b/a)))$  PQ to report  $(1+\epsilon)$ -ANN in range.

→ Far Case.

Let  $Q$  be a set of  $m$  balls. And  $\epsilon > 0 \in \mathbb{R}$  such that

$\bigcup_{\text{balls}}(Q, \epsilon)$  is a connected set. Then,

- i) any points  $p, q \in Q$  are within  $2\epsilon(m-1)$  distance.
- ii) for query  $q \in M$  s.t.  $d(q, Q) \geq \text{diam}(Q)/\delta$ , we have that any point of  $Q$  is a  $(1+\delta)$ -ANN to  $q$  in  $Q$ .

→ Binary Hierarchically well Separated Tree (BHST)

- Given a set of points  $P$ . Consider a <sup>binary</sup> tree  $H$ , having elements of  $P$  as leaves.  $H$  defines an **HST** over points of  $P$  if

- $\forall u \in H$ , we associate a label  $\Delta(u) \geq 0$ .
- $\Delta(u) = 0 \iff u$  is a leaf.
- $u$  is child of  $v \implies \Delta(u) \leq \Delta(v)$
- distance between leaves  $u, v \in H$  defined as  $\Delta(\text{lca}(u, v))$ .

- distances satisfy triangle inequality, hence HST defines metric. useful, since can be manipulated algorithmically.

- Quadtree with 'diameter of point within cell' as label defines an HST.

- In linear time HST can be converted to BHST, while

retaining distances.

• Metric  $N$   $t$ -approximates a metric  $M$ , if

$$d_M(u, v) \leq d_N(u, v) \leq t \cdot d_M(u, v).$$

• Any  $n$  point metric can be  $(n-1)$  approximated by some HST.

$\Rightarrow$  Linear # of balls.

• We use intuition of close/far/bad cases from earlier to refine our construction. We construct a tree to search for the ANN, such that traversal is based on that case work.

• Given set of points  $P$ , and a  $t$ -approximate BHST  $H$  of  $P$  for a leaf  $u$  of  $H$ ,  $\text{rep}_u$  is the point stored at  $u$ , and for an internal node,  $\text{rep}_u$  is the representative of one of its children.

For a subtree  $H' \subseteq H$ ,  $P(H') := \{\text{rep}_u \mid u \in H'\}$

not just rooted subtrees!

• A separator is a node  $u$  of tree  $H$ , such that its removal breaks  $H$  into connected components, each of size at most  $|V(H)|/2$ .

It can be found in linear time.

$\rightarrow$  Construction

- Initialize subtree  $S = H$ .

- Set root of tree  $T$  as node  $v = v(S)$ . Notation:  $n_v$  is # of nodes of  $S$ , and  $P^v = P(S) = \bigcup_{u \in S} \text{rep}_u$ .

- Find  $u^v \in V(S)$  as a separator of  $S$ . Let  $S_L$  and  $S_R$  be the subtrees rooted in the left and right child of  $u^v$  resp. And let  $S_{out} := S \setminus \{S_L \cup S_R\}$ .

define  $P_L^v = P(S_L)$ ,  $P_R^v$ ,  $P_{out}^v$  accordingly.

- Build Interval Query data structure  $\hat{I}(P^v, r_v, R_v, \epsilon/4)$  and store it in node  $v$ . Here

$$r_v = \Delta(u^v)/4t \quad R_v = \mu \Delta(u^v)$$

global parameter  $\mu \in O(\epsilon^{-1} \log n)$  to be set later.

$\Delta(u^v)$  is (upper bound)  $t$ -approx diameter of subtree at  $u^v$ .

- Recursively build tree on  $S_L, S_R, S_{out}$  and attach as children  $v_L, v_R, v_{out}$  of  $v$ .

→ Query

Perform query at  $\hat{I}_{root}$ :

i).  $d(q, P^v) \leq r_v$ : then  $q \in \mathcal{U}_{balls}(P^v, r_v)$ , and DS returns a point  $p \in P^v$  s.t.  $d_u(q, p) \leq r_v$ . If  $p \in P_L^v$  or  $P_R^v$  recurse

there, otherwise recurse in  $v_{out}$ .

ii).  $d(q, P^v) \in (r_v, R_v]$ : DS already returns an ANN.

iii).  $d(q, P^v) > R_v$ : recurse in  $v_{out}$ .

→ Correctness

• Observations:

$$- P_R^v \cap P_L^v = \emptyset$$

$$- P_R^v \cap (P_{out}^v \setminus \{rep_u\}) = \emptyset$$

$$- P_L^v \cap (P_{out}^v \setminus \{rep_u\}) = \emptyset$$

$$- d_u(P_L^v \cup P_R^v, P_{out}^v \setminus \{rep_u\}) \geq \Delta(u^v)/t.$$

$$- d_u(P_L^v, P_R^v) \geq \Delta(u^v)/t.$$

$$- \text{diam}(P_L^v \cup P_R^v) \leq \Delta(u^v).$$

For some node  $v$  on the search path.

- Close Case  $d(q, P^v) \leq r_v$ : we continue search in the subtree that contains req. ANN, & hence incur no loss.

Say we recurse on  $v_L$ , then  $d(q, P_L^v) \leq r_v$   
but  $d_m(P_L^v, P \setminus P_L^v) \geq \Delta(u^v)/t$

Let  $q_L = \text{nn}(q, P_L^v)$ . So

$$\begin{aligned} d(q, P \setminus P_L^v) &\geq d(q_L, P \setminus P_L^v) - d_m(q, q_L) \\ &\geq 4r_v - r_v > r_v. \end{aligned}$$

- Far Case  $d(q, P^v) > R_v$ : If  $s = \text{nn}(q, P^v)$  is in  $\gamma_{\text{out}}$ , we are fine (no loss). If not, if  $s \in P_L^v \cup P_R^v$ ?

we know  $\text{diam}(P_L^v \cup P_R^v) \leq \Delta(u^v) = R_v/\mu$

$$\Rightarrow \text{diam} \frac{(P_L^v \cup P_R^v)}{\mu} \leq R_v$$

$\Rightarrow$  Any point of  $P_L^v \cup P_R^v$  is  $(1 + 1/\mu)$ -ANN.

and  $\text{rep}_v$  (still in search space) is in  $P_L^v \cup P_R^v$ .

So we introduce  $(1 + 1/\mu)$  error.

- Bad Case: actually pretty good. We output an answer introducing  $(1 + \epsilon/4)$  error.

Now, depth of tree  $N = 2 \ln 2n \in O(\log n)$ . Let  $\mu = 4N/\epsilon$ .

Quality of ANN,

$$\left(1 + \frac{1}{\mu}\right)^N \left(1 + \frac{\epsilon}{4}\right) \leq \exp\left(\frac{N}{\mu} + \frac{\epsilon}{4}\right) = \exp\left(\frac{N}{4N/\epsilon} + \frac{\epsilon}{4}\right) = e^{\epsilon/2} \leq 1.1\epsilon.$$

→ # of balls.  $O(n/\epsilon \cdot \log^2 n)$ .

per structure.

$$U(n_v) = n/\epsilon \cdot \log(b/a) = n_v/\epsilon \cdot \log\left(\frac{R_v}{r_v}\right) = n_v/\epsilon \cdot \log\left(\frac{t \log n}{\epsilon}\right)$$

recurrence

$$B(n) = U(n) + \underbrace{B(n_{in}) + B(n_{out}) + B(n_{out})}_{\text{each } n_i \leq n/2 + 1}$$

each  $n_i \leq n/2 + 1$ .

set  $t = n^{O(1)}$  and  $1/\epsilon \in O(n)$ , to get

$$B(2n) = O(n/\epsilon \cdot \log^2 n)$$

→ # of Proximity Queries.  $O(\log(n/\epsilon))$ .

every internal node on the search path takes at most 2 queries. The final query takes,

$$O(\log(\epsilon^{-1} \log(b/a))).$$

$$= \log\left(\frac{\log(R_v/r_v)}{\epsilon}\right) = \log 1/\epsilon + \log \log\left(\frac{t \log n}{\epsilon}\right)$$

$$= \log 1/\epsilon + \log \log n.$$

at most  $O(\log n)$  internal nodes in the query, so overall

$$O(\log(n/\epsilon)).$$