

L9: Approximate Nearest Neighbors

Feb 5, 2025



Jeff M. Phillips

Large Set Data $P = \{p_1, p_2, \dots, p_n\} \in \mathcal{X}$

Given a distance $D: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$

Nearest Neighbor for a query $g \in \mathcal{X}$

$$\underset{P}{\Phi}(g) = \arg \min_{P \in P} D(p, g)$$

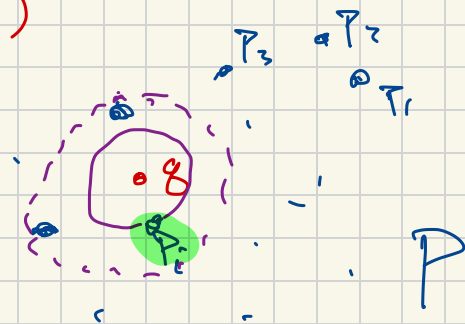
k-Nearest Neighbors $\Phi_{P,k}(g)$

query g , $k \geq 1$ $P_i = \Phi_P(g)$

Set $S \subseteq P$ s.t. $|S| = k$

no point $p' \in P \setminus S$ has

$D(g, p') < D(g, s)$ so $s \in S$



Algo for NN

$|P| = n$

1. Baseline checks all $O(|P|)$, return smallest

Pre-process : Spend $O(n)$ time
 $O(n \log n)$ time
near-linear

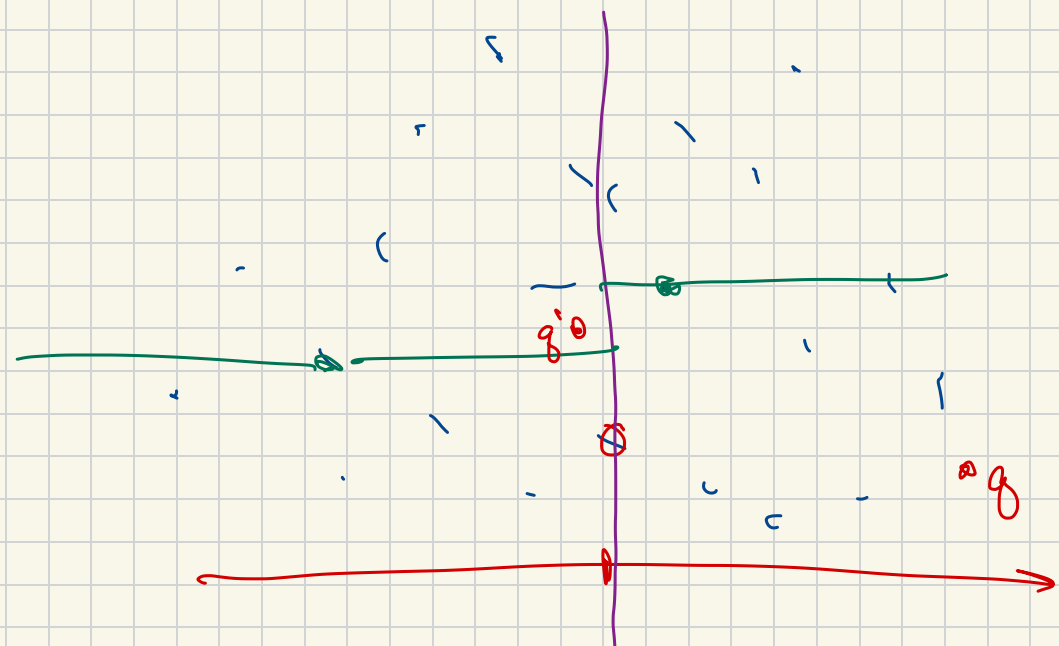
Build Data Structure $A(P)$

Then can ask NN queries of $A(P)$
quickly.

Ex: if PCR: Build binary tree, queries $O(\log n)$

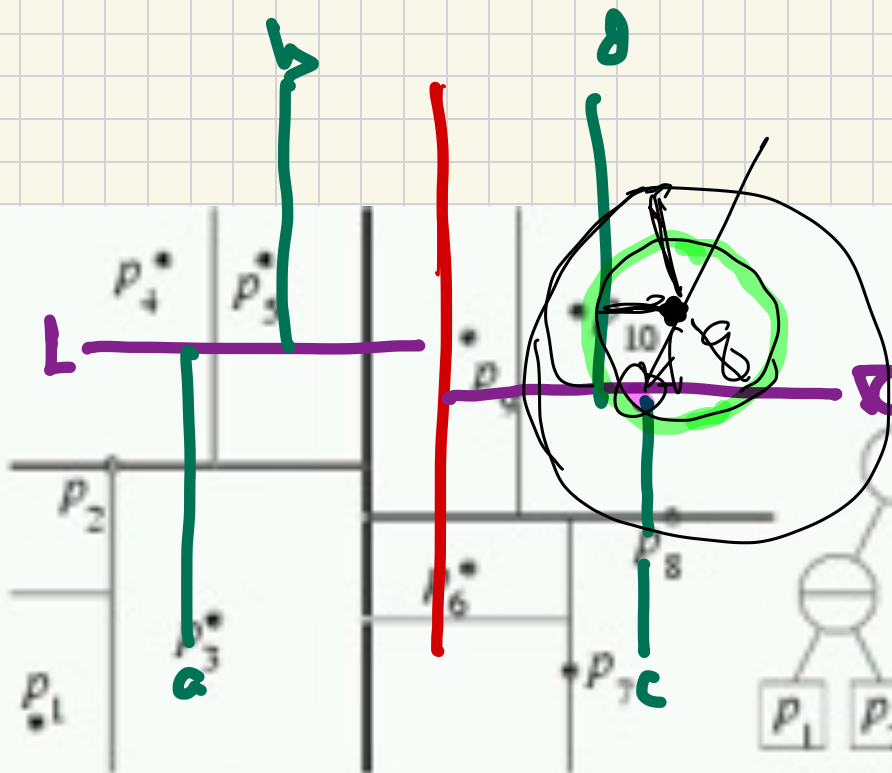
Classic Idea Applies $\mathbb{P}(\mathbb{R}^d)$ d , small

\checkmark $d=2,3$
or $d=4, \dots, (2)$ \rightarrow Hierarchy on Data
extending binary tree idea

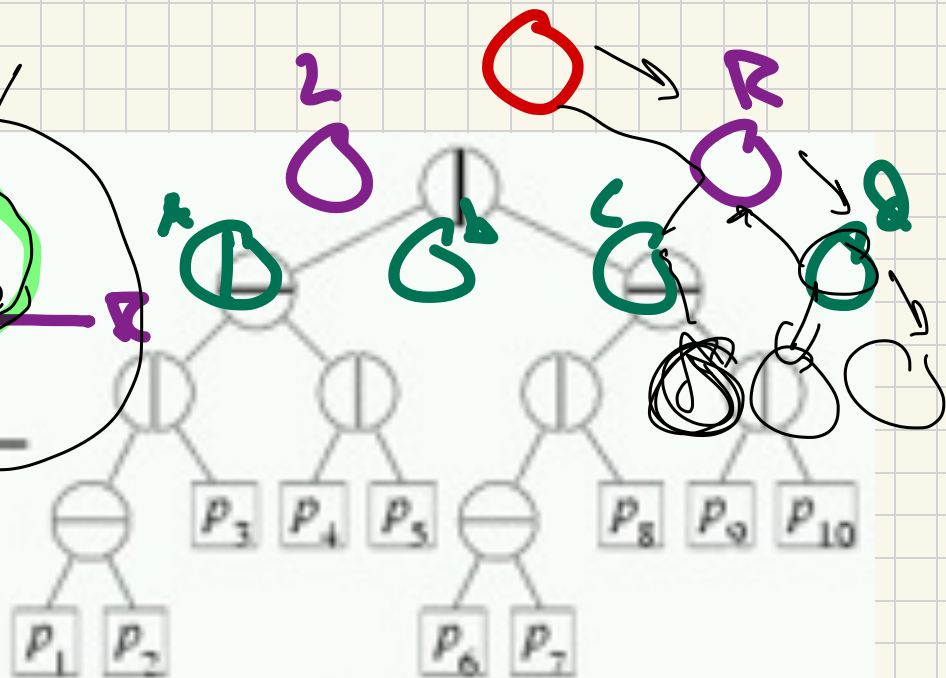


kd-tree

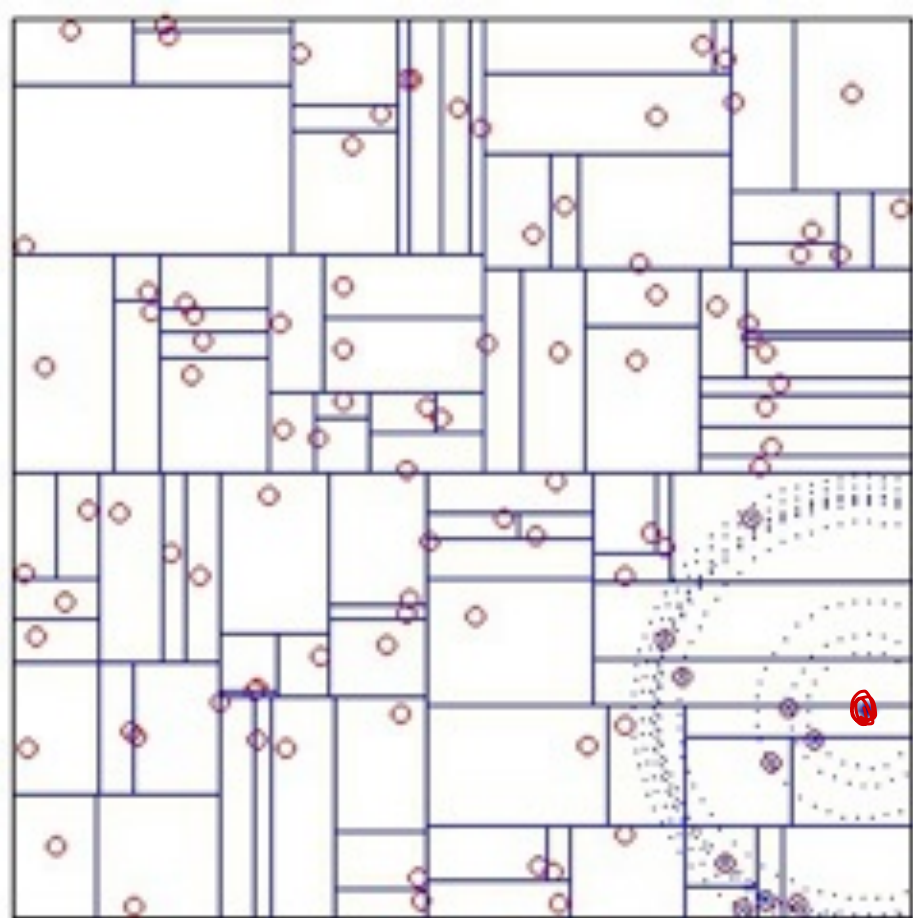
kd-tree



Subdivision



Tree structure



Approximate Nearest Neighbor

ϵ error tolerance $\epsilon = (0.05, 1)$

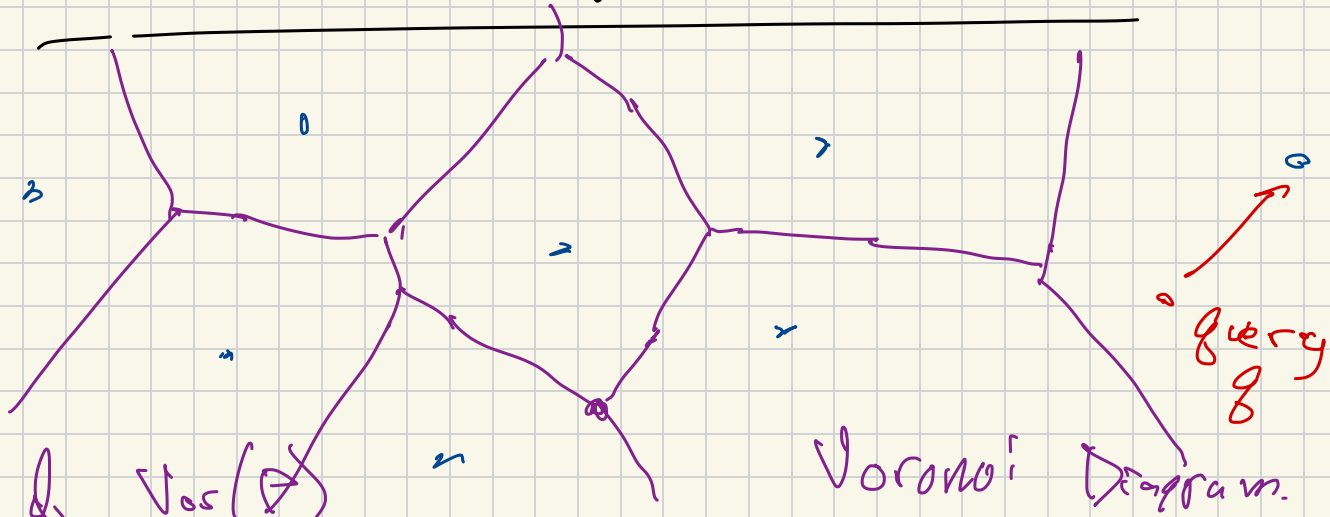
$$D(\Phi_P(q), q) = v$$

or with some $p' \in P$ so

$$D(q, p') \leq (1 + \epsilon) v$$

then p' ϵ -approx NN

Nearest Neighbor Problem

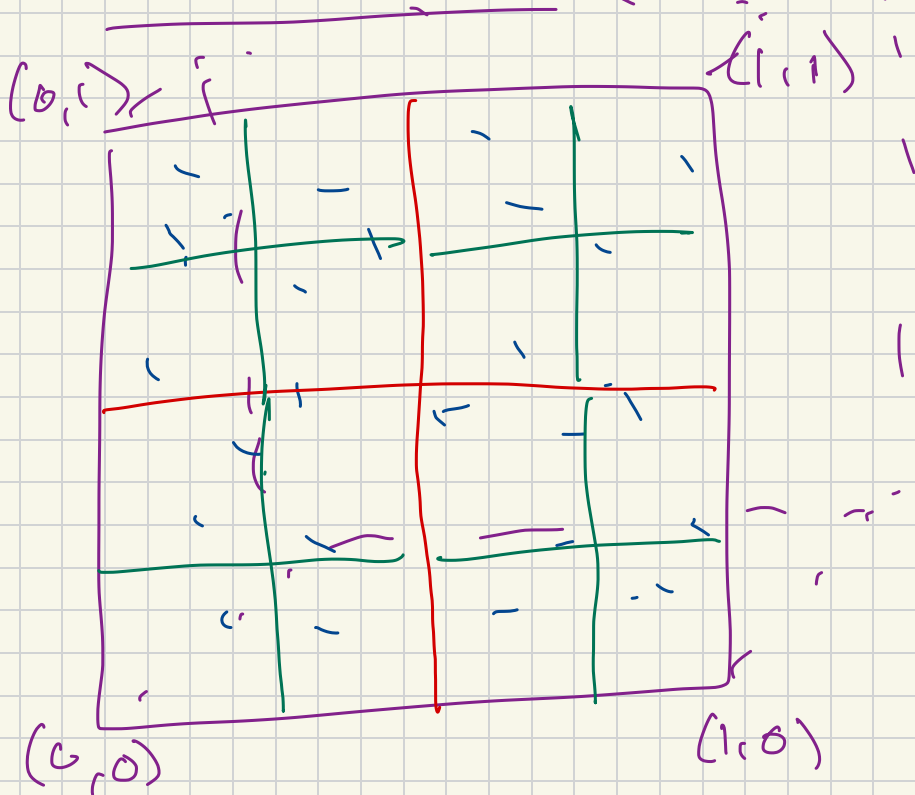


Size of $Voc(P)$
 $\Theta(|P|^{d/2})$

Voronoi Diagram.

query q



Quad Tree - Assumes \mathbb{R}^d

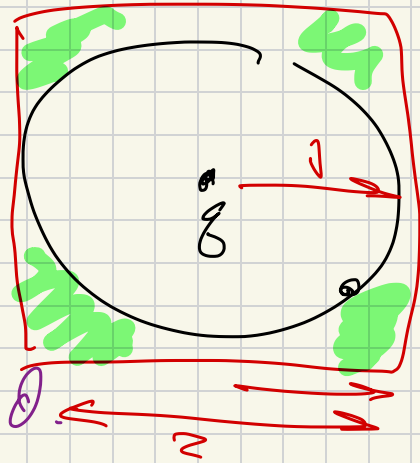


Library
ANN

Challenges in High Dimension

Using Euclidean Distance

Approximating   Rectangle
poss. in each dimension d



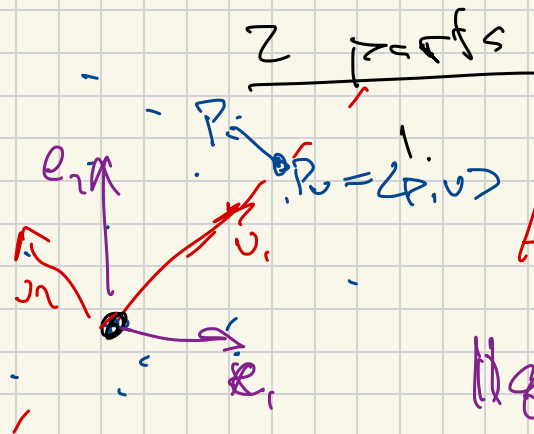
$$\text{Vol}(\text{Ball, rad}=1) = \frac{\pi^{d/2}}{\Gamma(d/2+1)} \approx \frac{\pi^{d/2}}{(d/2)!} \rightarrow \text{shrinks as } d \rightarrow \infty$$

$$\text{Vol}(\text{[0,1]}^d \text{ box}) = 2^d \rightarrow \text{grows exponentially in } d$$

Locality Sensitive Hashing

hash function

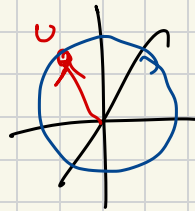
$$D_{\text{Eucl}}(p, g) = \|p - g\|$$



2 parts

Choose

$$u \sim \text{Unif}(\mathbb{S}^{d-1})$$



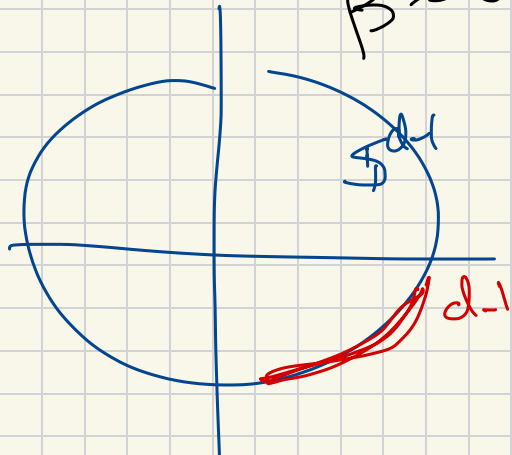
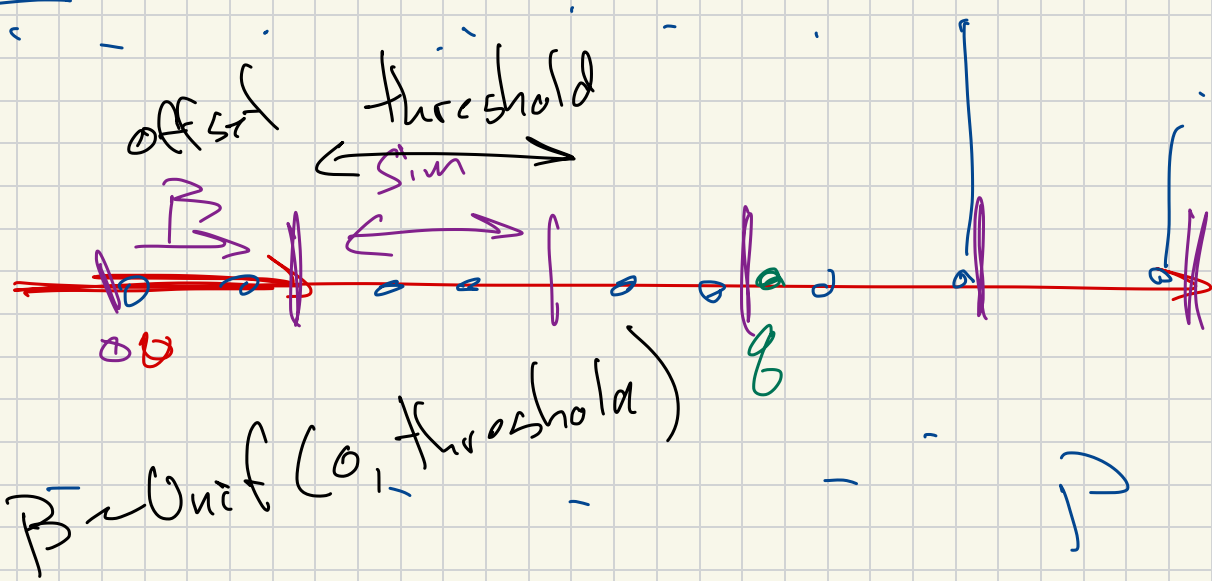
All of P project $P_u = \{(P_1, u), (P_2, u), \dots\}$

$$\|g - P\|^2 = \sum_{j=1}^d (P_j - g_j)^2$$

$$\|g - P\|^2 = \sum_{j=1}^d \langle (P - g), u_j \rangle^2$$

$$E_{u \sim \text{Unif}}[\langle P - g, u \rangle^2] = \frac{1}{d} \|P - g\|^2$$

FALCONN



$$\text{Sim}(p, g) = \max\{0, 1 - (p - g)\}$$

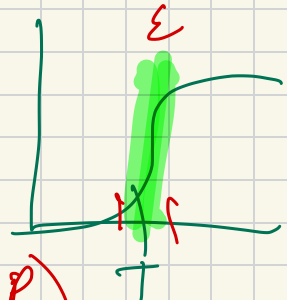
LSH

- For NN queries, need right threshold

Binary search on threshold until
get $O(r) > 0$ points.

What preprocessing
+ query time?

$$\tilde{O}(n \cdot r)$$
$$\tilde{O}(r)$$



t = hash function,

b in each bin superhash

r bins/super hash
query each

if $\epsilon = 1$

$$r = n^{1/\epsilon} = \Theta(n^{\epsilon})$$

$$n = (0, \text{billion}) \Rightarrow n^{1/\epsilon} = 10$$

Graph Descent NN

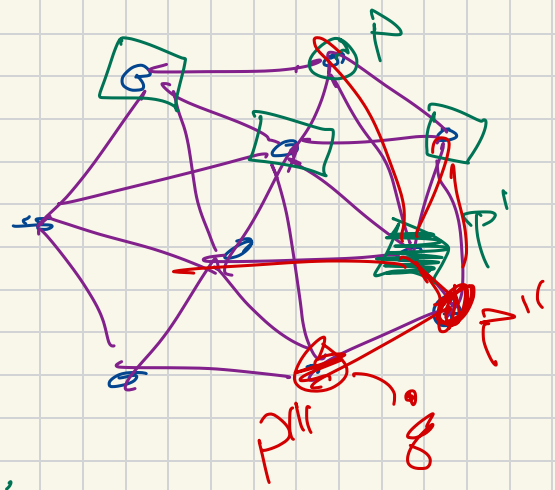
Data P , build graph on P

k -NN graph

edge p_i to p_j if $p_j \in \Phi_{p_i, k}(p_i)$

On going g
start at any $p \in P$

1. check if any neighbor d_p is closer to g .
2. Move to closest \rightarrow go to 1 or (if not) stop return p .



• Maintain k nearest neighbors
always check all neighborhoods
until can't improve.

• HNSW Build hierarchy over graph

DistANN

FAISS

Pinecone

