

Approximate Nearest Neighbor Search

Given points $p_1 \dots p_n \in \mathbb{R}^d$

Preprocess data set to answer nearest neighbor queries:

Given query point $q \in \mathbb{R}^d$, find $\arg \min_{p_i} d(p_i, q)$

	Brute Force	BF + dim red ⁿ
Storage	nd	$n \log(n)/\epsilon^2$
Preprocessing	—	$n d \log(n)/\epsilon^2$
Query	nd	$n \log(n)/\epsilon^2$

Can we answer queries in sublinear time

NNS in high dimⁿ: many data structures have exponential dependence on dimension

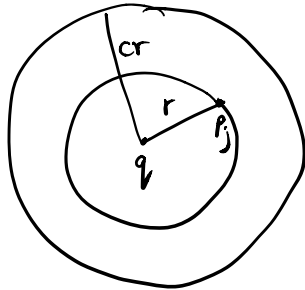
— Curse of dimensionality

JL dimension reduction not good enough

Approximate NNS

Relaxed Goal: Return p_c so that

$$d(q, p_c) \leq c \cdot \min_{p_j} d(q, p_j) \quad c \geq 1$$



Intrinsic dimensionality of data:

Complexity of data could be low even if points lie in high dimⁿ
eg. low dimensional manifold.

Doubling dimension $\dim(X)$ of metric space (X, d)
in min value of p st. every ball in X can be covered by
 2^p balls of half the diameter.

$X = \mathbb{R}^k$ with any norm $\dim(X) = \theta(k)$

Result: [Krauthgamer, Lee '02]

Navigating Nets

size: $2^{O(\dim(S))} \cdot n$

$(1+\epsilon)$ -approx nearest neighbor in time
 $2^{O(\dim(S))} \cdot \log \Delta + \left(\frac{1}{\epsilon}\right)^{O(\dim(S))}$

Aspect ratio $\Delta = \frac{d_{\max}}{d_{\min}}$

Claim: Exact nearest neighbor in time
 $2^{O(\dim(S))} \cdot \log n$

[KOR '98] [IM '98]

$(1+\epsilon)$ approx nearest neighbor in $\text{polylog}(n)$ query time
 $\text{poly}(n)$ preprocessing & storage.

$\hookrightarrow n^{O(1/\epsilon^2)}$

[IM '98] Locality Sensitive Hashing:

Defⁿ: Family of Hash fn is (r, cr, p_1, p_2) -LSH

with $p_1 \geq p_2, c > 1$ if

(a) $\Pr[h(x) = h(y)] \geq p_1$ when $d(x, y) \leq r$ (close points)

(b) $\Pr[h(x) = h(y)] \leq p_2$ when $d(x, y) > cr$ (distant points)

Can be used to design algo for approximate NNS
 Focus on following problem:

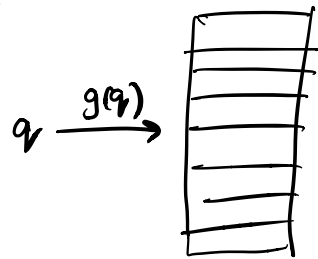
If \exists point within distance r of q , return a point within cr

To solve NNS, multiple copies for different values of r

Sample k hash h_i from (r, cr, p_1, p_2) -LSH family
 concatenate to get hash value

$$g(x) = h_1(x) h_2(x) \dots h_k(x)$$

$$g(x) = g(y) \quad h_i(x) = h_i(y) \quad \forall i$$



$$d(x, y) \leq r \quad \Pr[g(x) = g(y)] \geq p_1^k$$

$$d(x, y) \geq cr \quad \Pr[g(x) = g(y)] \leq p_2^k$$

Choose k so that $p_2^k = \frac{1}{n}$

$$k = \frac{\log n}{\log(\frac{1}{p_2})}$$

$$\mathbb{E}[\text{distant point collisions}] \leq 1$$

$$p_1 = p_2^p \quad p = \frac{\log(V/p_1)}{\log(V/p_2)}$$

$$p_2^k = \frac{1}{n} \Rightarrow p_1^k = \frac{1}{n^p}$$

$$\Pr[\text{good point collision}] \geq \frac{1}{n^p}$$

Repeat n^p times

n^p hash tables

n^p query time

n^{1+p} storage

(r, cr, p_1, p_2) -LSH

p : best parameter over LSH schemes for all r

P : fn of c and metric d

Hamming Metric:

$$p_i \in \{0, 1\}^d$$

hash fn: pick random coord

$$d(x, y) \leq r \quad \Pr[h(x) = h(y)] \geq 1 - \frac{r}{d} \approx e^{-r/d}$$

$$d(x, y) \geq cr \quad \Pr[h(x) = h(y)] \leq 1 - \frac{cr}{d} \approx e^{-cr/d}$$

$$p_1 = p_2^{1/c} \quad p = \frac{1}{c}$$

Best possible for Hamming

l_1 norm: same result

What is the hash fn? Pick random coord i , threshold t

$$\mathbb{I}_{\{x_i \geq t\}}$$

$$p = \frac{1}{c}$$

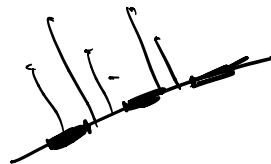
Euclidean space: l_2 norm

l_2 embeds into l_1 isometrically, so we can get $p = \frac{1}{c}$

Can do better

[DIIM'04] map to random line, split into buckets of width w

$$h_L(x) = \left\lfloor \frac{x \cdot L + \alpha}{w} \right\rfloor \quad \alpha \in \mathbb{R}([0, w])$$



[AI '06] map $\pi: X \rightarrow \mathbb{R}^d$
 random seq $s_1, s_2, \dots \in \mathbb{R}^d$, fix radius S

$$h_s(x) = \arg \min_{i > 0} \pi(x) \in B_S(s_i)$$

index of first ball containing $\pi(x)$

$$\rho \rightarrow \frac{1}{c^2} \text{ as } d \rightarrow \infty$$

$$\rho = \frac{1}{c^2} \text{ optimal}$$

Data Dependent Hashing

[AR '15] Can achieve $\rho = \frac{1}{2c-1}$ for Hamming
 $= \frac{1}{2c^2-1}$ for Euclidean

decompose dataset into pseudo-random sets.

Other Hash families:

Collection of sets

$$\text{sim}(A, B) = \frac{|A \cap B|}{|A \cup B|} \quad \text{Jaccard coeff.}$$

Min-hash scheme [Broder '97]

$$f: U \rightarrow 2^{64} \quad (\text{assume no collisions})$$

$$h(A) = \min_{a \in A} f(a)$$

$$\Pr \left[\underbrace{\min_{a \in A} f(a)}_{h(A)} = \underbrace{\min_{b \in B} f(b)}_{h(B)} \right] = \frac{|A \cap B|}{|A \cup B|}$$

"distance" for $1 - \frac{|A \cap B|}{|A \cup B|}$

Sketches of docs for estimating similarity

Sim-Hash

unit vectors with angular distance



$$h(u) = \text{sign}(\langle r, u \rangle) \quad r: \text{random vector}$$

$$\Pr[h(u) \neq h(v)] = \frac{\theta}{\pi} \quad \theta = \angle(u, v)$$

used at Google

Cross-Polytope Hash [Andoni et al '15]

unit vectors.

pick random rotation & return index of largest magnitude coordinate

$$\text{Achieves } \rho \approx \frac{1}{2c^2 - 1}$$