

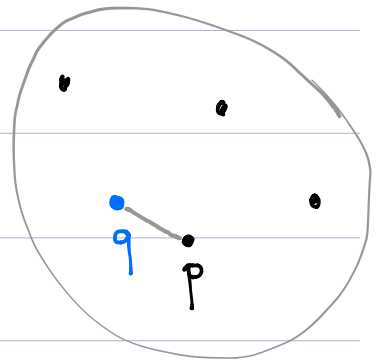
Locality Sensitive Hashing

Motivation: nearest neighbor search (NNS)

Preprocess: set S of n pts in metric space

Query: given pt q , find $p \in S$ s.t.

$$d(p, q) \text{ min}$$



Docs, images, etc often represented as pts in very high dimensional space

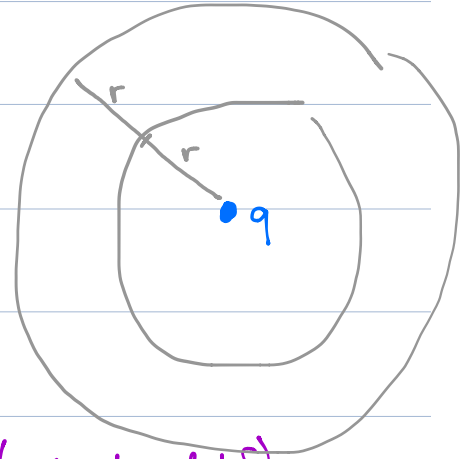
Similarity is fundamental problem in applications ranging from e-commerce to medical imaging, bioinformatics, astrophysics, finance web search,

2D - Voronoi diagram	$O(n)$ space	query time $O(\log n)$
In high dimensions	query time $O(\log n \cdot d)$	space $n^{O(d)}$
with preprocessing	$O(\log n \cdot d)$	$n^{O(d)}$
without	$O(nd)$	$O(nd)$

← curse of dimensionality

Approximate NNS

C-approximate r-neighbor search



Given query pt q , return

- all pts p s.t. $d(p, q) \leq r$ (each with prob $1-\delta$)
- may return some pts s.t. $d(p, q) \leq cr$

Do this with LSH [Indyk, Motwani]

This work won 2012 Kanellakis Theory and Practice Award.

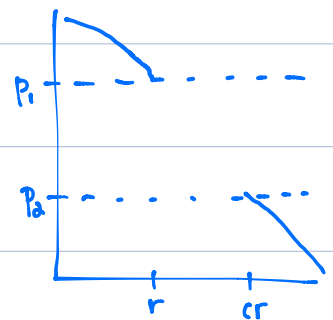
Locality sensitive hashing LSH

\mathcal{H} : family of hash fns mapping pts in metric space $\rightarrow \mathbb{R}$
for now \mathbb{R}^d

\mathcal{H} is (r, cr, p_1, p_2) sensitive $\forall p, q \in \mathbb{R}^d$

If $d(p, q) \leq r \Rightarrow \Pr(h(p)=h(q)) \geq p_1$
random $h \in \mathcal{H}$

If $d(p, q) \geq cr \Rightarrow \Pr(h(p)=h(q)) \leq p_2$



Example:

Suppose pts $e \in \{0,1\}^d$

$$\begin{aligned}d(p,q) &= \text{Hamming distance} \\ &= \# \text{ bits where } p \text{ \& } q \text{ differ} \\ &= \sum_{1 \leq i \leq d} \mathbb{1}_{p_i \neq q_i}\end{aligned}$$

$$h(p) = \{h_i(p) = p_i \mid 1 \leq i \leq d\}$$

$$d(p,q) \leq r$$

$$\begin{aligned}p_1 &= \Pr(h(p) = h(q)) \geq 1 - \frac{r}{d} \\ &\approx e^{-\frac{r}{d}}\end{aligned}$$

$$d(p,q) \geq cr$$

$$\begin{aligned}p_2 &= \Pr(h(p) = h(q)) \leq 1 - \frac{cr}{d} \\ &\approx e^{-\frac{cr}{d}}\end{aligned}$$

How do we use (r, cr, p_1, p_2) class of hash fns?

- amplify diff between p_1 & p_2 - concatenate hash fns in 3 steps

1. reduce prob that distant pts hash together

AND

2. increase prob that close pts hash together

OR

3. iterate to get failure prob below target

AND

Alg for approx NNS given (r, cr, p_1, p_2) sensitive family \mathcal{H}

each of these
selected uniformly at random from \mathcal{H}

Define $g_i(p) = [h_1(p), \dots, h_k(p)]$

e.g. Hamming dist $g: \{0,1\}^d \rightarrow \{0,1\}^k$

$\forall p \in S$ hash $g_i(p) \rightarrow$ table T_i of size n
to reduce
space usage.

Let p, q be 2 pts st. $d(p, q) \geq 2r$
far

$\Pr(g_i(p) = g_i(q)) \leq p_2^k$ choose k so this is $\leq \frac{1}{n}$

Why? so expected # of far pts
reported on query to q is ≤ 1

Defn: \mathcal{P} satisfies $p_1 = p_2^{\mathcal{P}}$

Suppose $d(\tilde{p}, q) \leq r$

Then $\Pr(g_i(\tilde{p}) = g_i(q)) \geq p_1^k = (p_2^{\mathcal{P}})^k = \frac{1}{n^{\mathcal{P}}}$

To catch close pairs, create $L = O(n^{\mathcal{P}})$ tables

and return \tilde{p} if $g_i(q) = g_i(\tilde{p}) \quad \forall 1 \leq i \leq L$

"signature"
 $g_1(p) = [h_{11}(p), \dots, h_{1k}(p)] \Rightarrow \text{hash to } T_1 \text{ of size } n$
 $g_2(p) = [h_{21}(p), \dots, h_{2k}(p)] \Rightarrow \text{hash to } T_2 \text{ of size } n$
 \vdots
 $g_L(p) = [h_{L1}(p), \dots, h_{Lk}(p)] \Rightarrow \text{hash to } T_L \text{ of size } n$

each used fresh random ness

$L = n^{\mathcal{J}}$ hash tables

On query q , compute $T_1(g_1(q)), \dots, T_L(g_L(q))$

If p in any of those buckets, compute $d(p, q)$

output all close points

Alg fails when an r -near neighbor p^* not in any of these buckets

$$\text{Prob of failure} = (1 - p_i)^k = (1 - p_2^{\mathcal{J}k})^{n^{\mathcal{J}}} \leq \frac{1}{n^{\mathcal{J}}}$$

recall $L = n^{\mathcal{J}}$ $p_i = p_2^{\mathcal{J}}$

Runtime: • hash fn evaluation $O(Lk)$

• distance computations to pts in buckets

Distance computations:

care only about far pts ($> cr$)

$$\Pr(\text{far pt collides}) \leq p_2^k = \frac{1}{n}$$

$$E(\# \text{ far pts in a bucket}) \leq n \cdot \frac{1}{n} = 1$$

$$\Rightarrow E(\# \text{ far pts is } L)$$

$$\text{Total: } O(Lk + Ld) = O(n^{\beta} (\log n + d))$$

$$p_2^k = \frac{1}{n}$$

$$k \log(p) = -\log n$$

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

Space: $O(nL) = O(n^{4\beta})$

plus space to store pts

Application: Hamming distance

Suppose pts $\in \{0,1\}^d$

$d(p,q)$ = Hamming distance
= # bits where p & q differ
= $\sum_{1 \leq i \leq d} \mathbb{1}_{p_i \neq q_i}$

$\mathcal{H} = \{h_i(p) = p_i \mid 1 \leq i \leq d\}$

$$d(p,q) \leq r$$

$$p_1 = \Pr(h(p) = h(q)) \geq 1 - \frac{r}{d} \\ \approx e^{-\frac{r}{d}}$$

$$d(p,q) \geq 2r$$

$$p_2 = \Pr(h(p) = h(q)) \leq 1 - \frac{2r}{d} \\ \approx e^{-\frac{2r}{d}}$$

Recall \mathcal{F} defined by $p_i = p_i^{\mathcal{F}}$

$$\Rightarrow \mathcal{F} = \frac{1}{2}$$

For example with $c=2$: query time $O(\sqrt{n}(\log n + d))$
space: $O(n^{\frac{3}{2}}) + \text{pts}$

LSH families

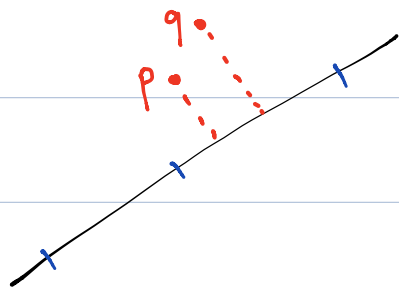
Euclidean distance in \mathbb{R}^d

Density of $N(0,1)$

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

Let $\vec{g} = (g_1, \dots, g_d)$ iid $N(0,1)$ random vars

Fix $w \gg r$



$$h(p) = \left\lfloor \frac{p \cdot g + b}{w} \right\rfloor = \left\lfloor \frac{\sum p_i g_i + b}{w} \right\rfloor$$



$p \cdot g$ projects p onto random line

$$(p - q) \cdot g = \sum_i \underbrace{(p_i - q_i) g_i}_{N(0, (p_i - q_i)^2)} \quad \text{weighted sum of } N(0,1) \text{ r.v.'s}$$

$$X_1 + \dots + X_k \quad X_i \sim N(\mu_i, \sigma_i^2) \quad \Rightarrow \quad X_1 + \dots + X_k \sim N(\mu_1 + \dots + \mu_k, \sigma_1^2 + \sigma_2^2 + \dots + \sigma_k^2)$$

Projected distance

$$\Rightarrow \sum_i (p_i - q_i) g_i \text{ is } N(0, \overbrace{\sum (p_i - q_i)^2}^{\text{variance}})$$

\Rightarrow exp distance² between projections

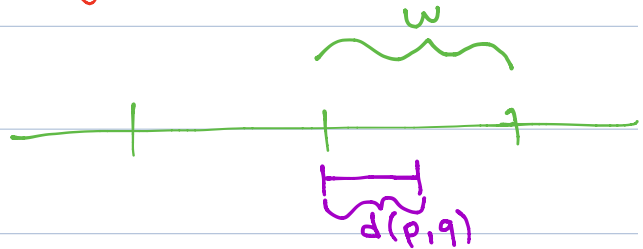
$$= E[(p \cdot g - q \cdot g)^2] = \text{Var}(p \cdot g - q \cdot g) = \underbrace{\sum (p_i - q_i)^2}_{\text{Euclidean distance}^2}$$

Concentrated around expectation

projection approximately preserves distance

random shift ensures that likely to go to same bucket

$\beta < \frac{1}{2}$ for w carefully chosen



Jaccard

A, B docs \equiv elts of $\{0,1\}^{|U|}$

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|} \quad d(A,B) = 1 - J(A,B)$$

Idea: permute rows of matrix at random

i.e. minhash $h_{\pi} = \min \{\pi(a) \mid a \in A\}$

$$\Pr(h_{\pi}(A) = h_{\pi}(B)) = \frac{|A \cap B|}{|A \cup B|}$$

$$d(A,B) \leq r$$

$$P_1 = \Pr(h_{\pi}(A) = h_{\pi}(B)) \geq 1 - r$$

$$d(A,B) \geq cr$$

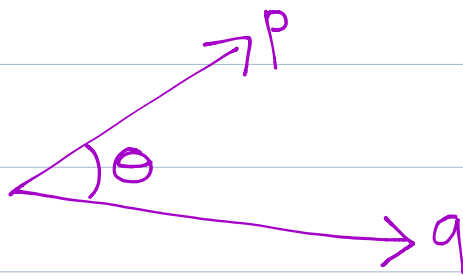
$$P_2 = \Pr(h_{\pi}(A) = h_{\pi}(B)) \leq 1 - cr$$

Cosine distance

used in information retrieval / data mining

points $\in \mathbb{R}^d$

$$d(p, q) = \arccos \left(\frac{p \cdot q}{\|p\| \|q\|} \right)$$



pick random hyperplane thru origin

$$\vec{g} = (r_1, \dots, r_n) \quad r_i \sim N(0, 1)$$

$$h(p) = \text{sign}(\vec{g} \cdot \vec{p})$$

(rotationally symmetric)

$$\Pr(h(p) = h(q)) = 1 - \frac{2\theta}{2\pi} = 1 - \frac{d(p, q)}{\pi}$$

Other distances:

- l_1

- edit distance