

Approximate Nearest Neighbor in High Dimensions

what if d is not constant?

[lots of apps -]
["curse of dimensionality"
all methods floundering
have exp. dep.
on d]

Method 1: Dimension Reduction (Indyk-Motwani '98)

Def Given n pts in \mathbb{R}^d ,

an ϵ -embedding into k dims is a map $f: P \rightarrow \mathbb{R}^k$

st. $\forall p, q \in P$,

$$\epsilon_0(1-\epsilon)\|p-q\| \leq \|f(p)-f(q)\| \leq \epsilon_0(1+\epsilon)\|p-q\|$$

Johnson-Lindenstrauss Lemma⁽¹⁹⁸⁴⁾ For any n pts in \mathbb{R}^d ,

\exists ϵ -embedding into $O\left(\frac{1}{\epsilon^2} \log n\right)$ dims.

Pf: Very Simple Randomized Algm:

just choose a random projection f !

more precisely, let $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$ be

$$f(x_1, \dots, x_d) = \left(\sum_{j=1}^d X_{1j} x_j, \dots, \sum_{j=1}^d X_{kj} x_j \right)$$

where coefficients X_{ij} are indep rand vars
of normal distribution $\sim N(0, 1)$

↑ mean
↑ variance

Analysis: (sketch)

Fix $p, q \in P$.

$$\text{Then } \|f(p)-f(q)\|^2 = \left\| \sum_{i=1}^k y_i \right\|^2$$

where $y_i =$ i th coord of $f(p)-f(q)$

$$y_i = \sum_{j=1}^d X_{ij} (p_j - q_j)$$

is normally distributed $\sim N\left(0, \sum_{j=1}^d (p_j - q_j)^2\right)$

$$= N\left(0, \|p-q\|^2\right)$$

(linear comb.
of normal
vars
is normal)
(Var[\sum
of vars])

$$\text{Let } \mu = E\left[\sum_{i=1}^k y_i^2\right] = \sum_{i=1}^k \text{Var}[y_i] \\ = k \|p - q\|^2$$

By Chernoff bound (for squared normal distribution), $-k\epsilon^2/4$

$$\Pr\left\{\sum_{i=1}^k y_i^2 \notin [(1-\epsilon)\mu, (1+\epsilon)\mu]\right\} \leq e^{-k\epsilon^2/4}$$

$$\text{Set } k = \frac{12}{\epsilon^2} \log n \leq e^{-3 \log n} = \frac{1}{n^3}$$

$$\Rightarrow \Pr\left\{\|f(p) - f(q)\|^2 \notin [(1-\epsilon)k\|p-q\|^2, (1+\epsilon)k\|p-q\|^2] \text{ for some } p, q \in P\right\} \leq \frac{1}{n} \quad \square$$

To solve approx decis problem for fixed r :

1. apply embedding
2. then use grid method with

query time $\tilde{O}(1)$ (\tilde{O} hides poly(d) factors)

$$\text{Space } O\left(\left(\frac{1}{\epsilon}\right)^k n\right) \\ = O\left(\frac{1}{\epsilon}\right)^{O\left(\frac{1}{\epsilon^2} \log n\right)} n \\ = \boxed{n^{O\left(\frac{1}{\epsilon^2} \log \frac{1}{\epsilon}\right)}}$$

polynomial space for any const ϵ regardless of d !

Remark - extends to ANN queries

- many other embedding results for diff. metrics (extends to $L_p, 1 \leq p \leq 2$)

- Alon-Chazelle '06: faster JL transform ($O(kd)$ time $\rightarrow \tilde{O}(d)$)

Method 2: Locality Sensitive Hashing (Indyk-Motwani)

Consider special case of Hamming space $\{0, 1\}^d$:

given $p, q \in \{0, 1\}^d$, $d(p, q) = \#$ bits that are different

Consider approx decs. problem for fixed r , with approx factor c

Preproc. Algm:

Fix k to be set later

define hash function $h: \{0, 1\}^d \rightarrow \{0, 1\}^k$ by

projecting to k random dims

build hash table

\Rightarrow space $O(n)$

Query(q):

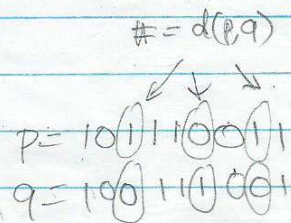
for each $p \in P$ with $h(p) = h(q)$

if $d(p, q) \leq cr$ return yes

return no

Facts:

$$0. \Pr[h(p) = h(q)] = \left(1 - \frac{d(p, q)}{d}\right)^k$$



all k chosen dims
are not among
circled ones

1. If $d(p, q) > cr$, then

$$\Pr[h(p) = h(q)] \leq \left(1 - \frac{cr}{d}\right)^k$$

$$\leq e^{-\frac{cr}{d}k}$$

by picking $k = \frac{d \ln n}{cr}$

$$= \frac{1}{n}$$

2. If $d(p, q) \leq r$, then

$$\begin{aligned} \Pr [h(p) = h(q)] &\geq \left(1 - \frac{r}{d}\right)^k \\ &\approx e^{-\frac{r^k}{d}} \\ &= \frac{1}{n^{1/c}} \end{aligned}$$

→ expected query time

$$\begin{aligned} &= \tilde{O} \left(\mathbb{E} \left[\# \text{ pts } p \in P \text{ with } h(p) = h(q) \right. \right. \\ &\quad \left. \left. \& d(p, q) > cr \right] \right) \\ &= \tilde{O} \left(n \cdot \frac{1}{n} \right) = \tilde{O}(1) \end{aligned}$$

Error Analysis:

Case 1. $d(p, q) > cr$ for all $p \in P$

algm is always correct

Case 2. $d(p, q) \leq r$ for some $p \in P$

\Pr [algm is correct]

$$\geq \Pr [h(p) = h(q)] = \frac{1}{n^{1/c}} \quad \text{tiny!}$$

final idea → repeat $t = 100n^{1/c}$ times

to lower error prob

$$\Rightarrow \text{err prob} \leq \left(1 - \frac{1}{n^{1/c}}\right)^t \leq e^{-100}$$

$$\Rightarrow \text{space } \tilde{O}(n^{1+1/c})$$

$$\text{query time } \tilde{O}(n^{1/c})$$

for approx factor c

(e.g. $c=2$!
= space $\tilde{O}(n^{3/2})$
query $\tilde{O}(\sqrt{n})$)

Remarks - immediately extends to L_1 metric in $[U]^d$

$$d(p, q) = \sum_{i=1}^d |p_i - q_i|$$

by first mapping $[U]^d \rightarrow \{0, 1\}^{dU}$

$p = (2, 3) \rightarrow (11000, 1100)$
 $q = (5, 2) \rightarrow (11111, 1100)$

L_1 dist $3+1=4$

ATI

3/18/16

- extends to Euclidean metric by modifying hash fn (Datar-Indyk et al. '04)

take rand. projection $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$

$$f(u_1 \rightarrow u_d) = \left(\sum_{j=1}^d X_{1j} u_j, \dots, \sum_{j=1}^d X_{kj} u_j \right)$$

where $X_{ij} \sim N(0, 1)$

take randomly shifted grid in \mathbb{R}^k of suff. large
side length m

define $h(u) =$ grid cell containing $f(u)$.

$$\text{Then } \Pr \left[\begin{array}{l} f(p), f(q) \text{ in diff} \\ \text{grid cell in } i^{\text{th}} \text{ coord} \end{array} \right] \approx \frac{E[|Y_i|]}{m} = \frac{\sqrt{2/\pi} \|p-q\|}{m}$$

$$\begin{aligned} \text{where } Y_i &= \text{ } i^{\text{th}} \text{ coord of } f(p) - f(q) \\ &= \sum_{j=1}^d X_{ij} (p_j - q_j) \\ &\sim N(0, \|p-q\|^2) \end{aligned}$$

$$\Rightarrow \Pr[h(p) = h(q)] \approx \left(1 - \frac{\sqrt{2/\pi} d(p,q)}{m} \right)^k$$

Similar to Fact 0

rest is as before

$$\begin{aligned} \Rightarrow \text{space } &O^*(n^{1+1/c}) \\ \text{query time } &O^*(n^{1/c}) \quad \text{dyn } O^*(n^{1/c}) \end{aligned}$$

Andoni-Indyk '06 improve to space $O^*(n^{1+1/c^2})$
time $O^*(n^{1/c^2})$

Andoni-Razenshteyn '15
(data-dependent LSH)

space $O^*(n^{1+1/(2c^2-1)})$
time $O^*(n^{1/(2c^2-1)})$

(Hamming $n^{1/(2c-1)}$)

Pagh '16 / Wei '19: Las Vegas

Andoni et al. '17: space/time trade-offs