

Near-Neighbour Search

&
locality-sensitive Hashing

TSS 20/09/23

- Deepanshu
Kush

Outline

- Problem Definition & Background
- Approximate Version
- LSH Families (Simple)
- Data Structure Construction
- LSH Families (Complex)
- Recent Developments

Nearest Neighbour Search

Nearest Neighbour Search

Given:

1. Database $X = \{x_1, \dots, x_n\} \subset U \subseteq \mathbb{R}^d$ (eg. $U = \{0, 1\}^d$)
2. A distance metric $d(a, b)$ for points $a, b \in U$

Nearest Neighbour Search

Given:

1. Database $X = \{x_1, \dots, x_n\} \subset U \subseteq \mathbb{R}^d$ (eg. $U = \{0, 1\}^d$)
2. A distance metric $d(a, b)$ for points $a, b \in U$

Goal:

1. Construct a data structure for X (pre-processing, extra space allowed)

Nearest Neighbour Search

Given:

1. Database $X = \{x_1, \dots, x_n\} \subset U \subseteq \mathbb{R}^d$ (eg. $U = \{0, 1\}^d$)
2. A distance metric $d(a, b)$ for points $a, b \in U$

Goal:

1. Construct a data structure for X (pre-processing, extra space allowed)
2. Given any query pt. $q \in U$, find a nearest pt. in X
(i.e. $\operatorname{argmin}_{x \in X} d(q, x)$)

What parameters do we care about?

What parameters do we care about?

> space complexity & query time

What parameters do we care about?

- > space complexity & query time
- > some trade off to be expected

What parameters do we care about?

- > space complexity & query time
- > some trade off to be expected

Naively: Just store the n pts using
 $O(dn)$ space & do linear search ie. $O(n)$ query time.

What parameters do we care about?

- > space complexity & query time
- > some trade off to be expected

Naively: Just store the n pts using
 $O(dn)$ space & do linear search ie. $O(n)$ query time.

So, goal: try for $O(n)$ time, more space

What parameters do we care about?

- > space complexity & query time
- > some trade off to be expected

Naively: Just store the n pts using
 $O(dn)$ space & do linear search ie. $O(n)$ query time.

So, goal: try for $O(n)$ time, more space
& minimize pre-processing time for practical apps

d = 1

$$\underline{d=1}$$

> Pre-processing: Sort the n pts

$d=1$

- > Pre-processing: Sort the n pts
- > Given a query q , binary search to find a nearest pt

$d=1$

- > Pre-processing: Sort the n pts
- > Given a query q , binary search to find a nearest pt
- > query-time: $O(\log n)$, space $O(n)$,
pre-processing $O(n \log n)$

Higher d?

Higher d?

k-d trees:

Higher d?

k-d trees: Partition \mathbb{R}^d using coord-aligned
planes, chosen suitably.

Higher d?

k-d trees: Partition \mathbb{R}^d using coord-aligned planes, chosen suitably.

> Unfortunately, fail to beat naive as soon as $d = \Omega(\log n)$

Higher d?

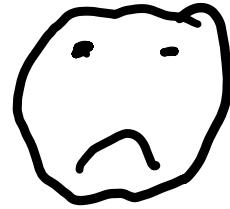
k-d trees: Partition \mathbb{R}^d using coord-aligned planes, chosen suitably.

- > Unfortunately, fail to beat naive as soon as $d = \Omega(\log n)$
- > Moreover, in all known approaches, size of DS / q-time grows like $\exp(d)$

Curse of Dimensionality

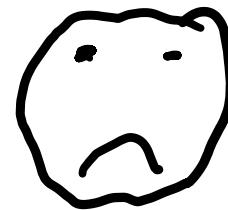


Curse of Dimensionality

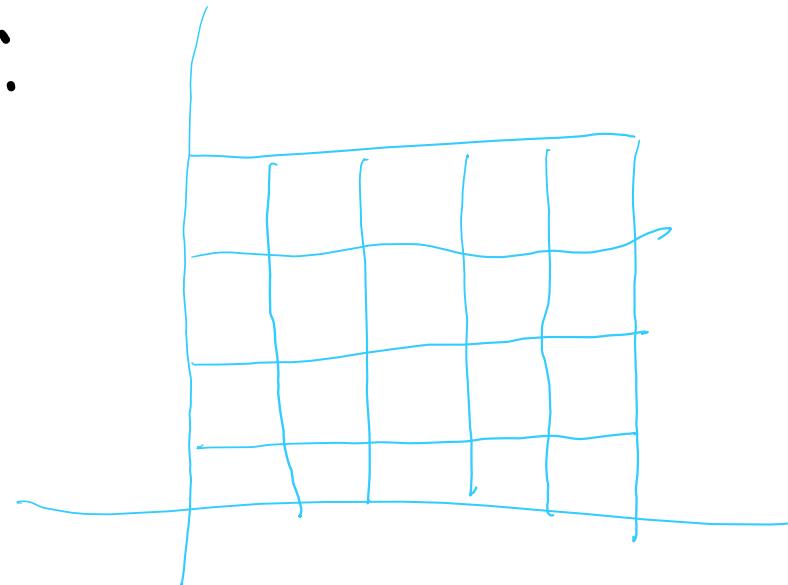


- Well-known phenomenon observed in many geometric problems, including NNS.

Curse of Dimensionality

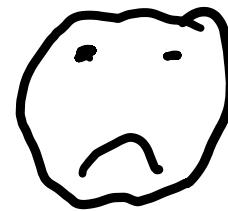


- Well-known phenomenon observed in many geometric problems, including NNS.
- Intuition:



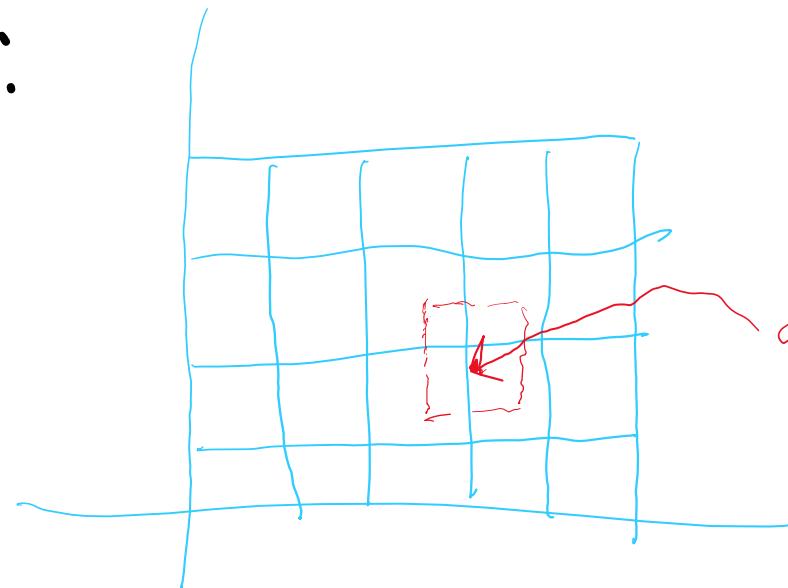
grid of side a .

Curse of Dimensionality



- Well-known phenomenon observed in many geometric problems, including NNS.

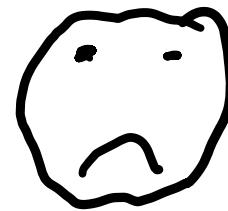
Intuition:



grid of side a .

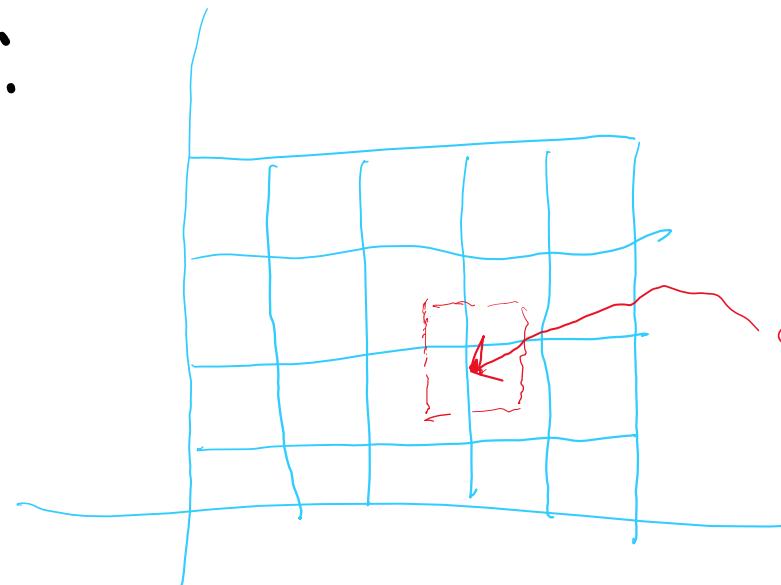
q: every pt will intersect
 2^d grid cells

Curse of Dimensionality



- > Well-known phenomenon observed in many geometric problems, including NNS.

Intuition:



grid of side a .

q: every pt will intersect
 2^d grid cells

- > NNS based on k-d trees will need $\mathcal{O}(2^d)$ lookup time

Example Applications

Example Applications

> Shazam!

Example Applications

- > Shazam!
- > Task: Identify a song from an audio recording
of a short clip from the song

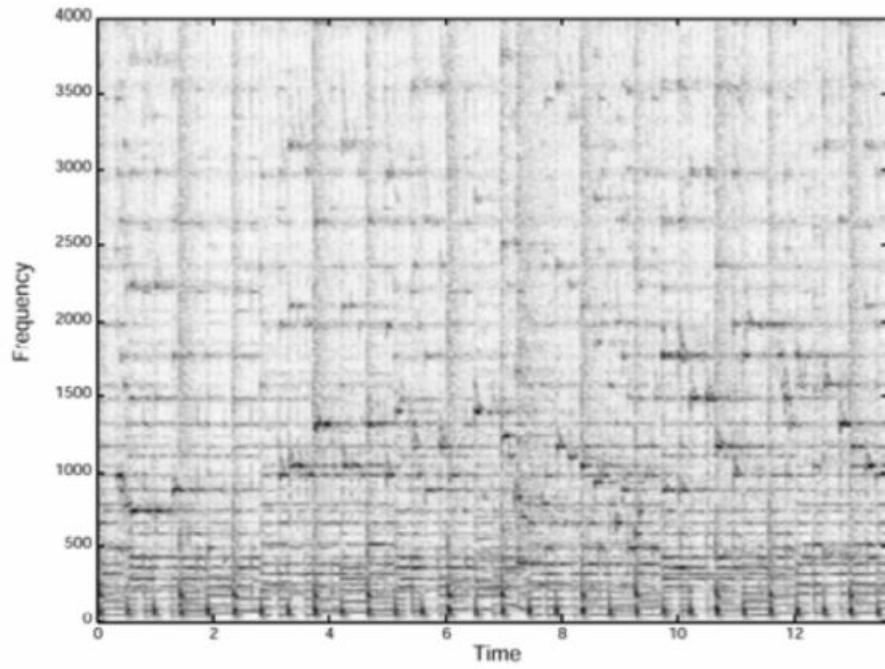
Example Applications

- > Shazam!
- > Task: Identify a song from an audio recording
of a short clip from the song
- > Idea (from Ansg hi-Chun Wang '03, ICMIR):

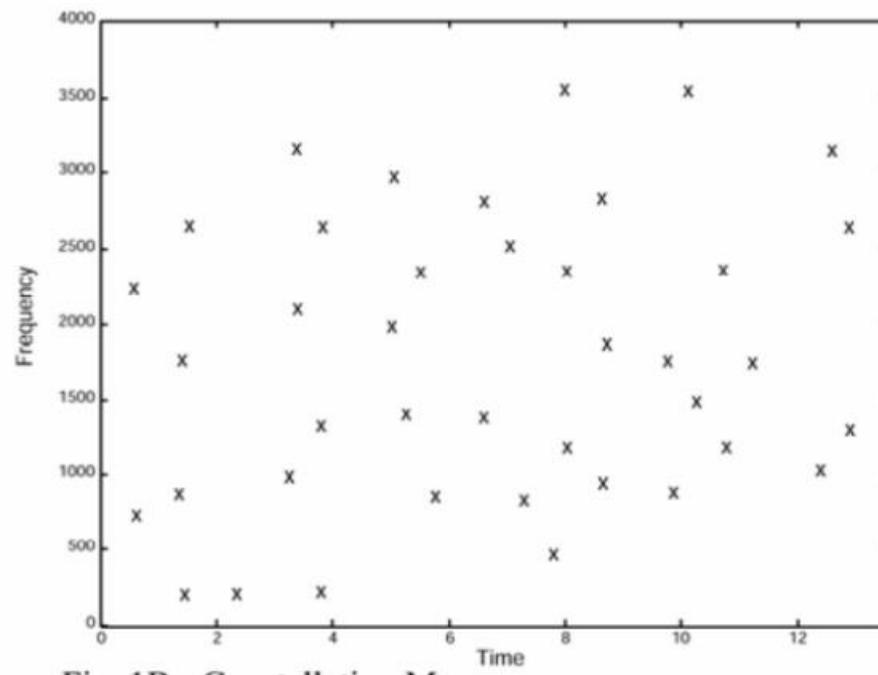
Example Applications

- > Shazam!
- > Task: Identify a song from an audio recording of a short clip from the song
 - Idea (from Ansgy Li-Chun Wang '03, ICMIR): first construct a "spectrogram" from the clip. (some windowed FFT). Then identify "peaks" in the spectrogram to construct $q \in \mathbb{A}^d$. Then apply NNS for a huge database of already processed clips.

Picture:



(a) Spectrogram extracted from clip of an audio file.



(b) Post-processed spectrogram which can be used to construct an audio “fingerprint” $q \in \mathbb{N}^d$ for identifying a song clip.

Figure 1: Images from the paper outlining Shazam’s approach to audio retrieval [2]. Once features are extracted from an audio clip, a hashing based algorithm is used to perform approximate nearest neighbor search in a database of known songs.

More Examples:

More Examples:

- > Used in detecting early warning signs of earthquakes and other seismic events

More Examples:

- > Used in detecting early warning signs of earthquakes and other seismic events
- > Again, spectrogram of a seismogram is computed in real-time & compared (quickly) against a database of surrounding major seismic events

More Examples:

- > Used in detecting early warning signs of earthquakes and other seismic events
- > Again, spectrogram of a seismogram is computed in real-time & compared (quickly) against a database of surrounding major seismic events
- > Also a very high-dim problem

Approximate NNS

→ Fortunately, break the curse if we're willing to tolerate some level of approximation.

Approximate NNS

- Fortunately, break the curse if we're willing to tolerate some level of approximation.
- (ϵ, r) -NNS: given $X \subset \mathbb{R}^d$ of n pts,
construct a DS s.t. for any query pt q ;
if $\exists p \in X$ s.t. $d(p, q) \leq r$, output a pt p s.t.
 $d(p, q) \leq \epsilon r$.

Approximate NNS

- > Fortunately, break the curse if we're willing to tolerate some level of approximation.
- > (ϵ, r) -NNS: given $X \subset \mathbb{R}^d$ of n pts,
construct a DS s.t. for any query pt q ;
if $\exists p \in X$ s.t. $d(p, q) \leq r$, output a pt p s.t.
 $d(p, q) \leq \epsilon r$.
- > if there is no such p , no requirements.

Picture:

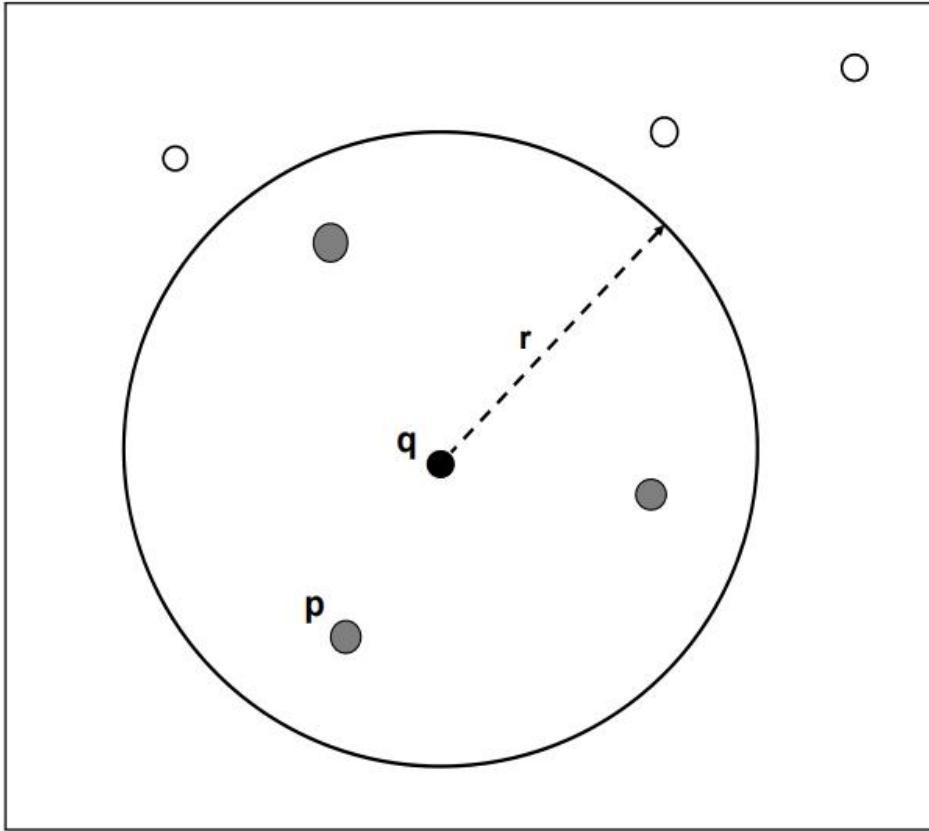


Figure 1: Near Neighbor problem, given q find p s.t. $D(q, p) \leq r$

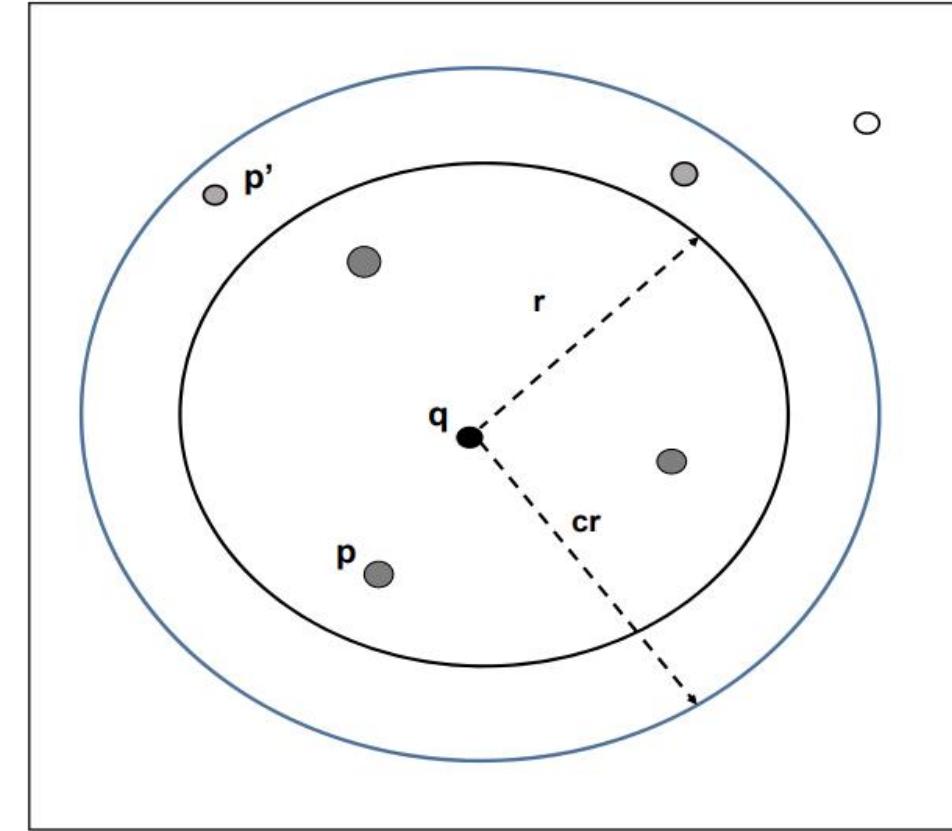


Figure 2: Approx Near Neighbor problem, given q find a p' s.t. $D(q, p') \leq cr$

(Approximately) Solve NNS using (c, r) -NNS

(Approximately) Solve NNS using (c, r) -NNS

→ guess $\min_{p \in X} d(p, \mathcal{Q})$ up to a mult. factor $1 \pm \epsilon$.

(Approximately) Solve NNS using (c, r) -NNS

- > guess $\min_{p \in X} d(p, \mathcal{Q})$ up to a mult. factor $1 \pm \epsilon$.
- > Assume $\text{diam}(P) = \max_{p, p' \in X} d(p, p') \leq 1$

(Approximately) Solve NNS using (c, r) -NNS

- > guess $\min_{p \in X} d(p, \mathcal{Q})$ up to a mult. factor $1 \pm \epsilon$.
- > Assume $\text{diam}(P) = \max_{p, p' \in X} d(p, p') \leq 1$
- > let $s = \min \text{dist. in } X$

(Approximately) Solve NNS using (c, r) -NNS

- > guess $\min_{p \in X} d(p, \mathcal{Q})$ up to a mult. factor $1 \pm \epsilon$.
- > Assume $\text{diam}(P) = \max_{p, p' \in X} d(p, p') \leq 1$
- > let $\delta = \min \text{dist. in } X$, so $1/\delta = \begin{matrix} \text{"bit-precision"} \\ \text{of dataset } X \end{matrix}$

(Approximately) Solve NNS using (c, r) -NNS

- > guess $\min_{p \in X} d(p, \mathcal{Q})$ up to a mult. factor $1 \pm \varepsilon$.
- > Assume $\text{diam}(P) = \max_{p, p' \in X} d(p, p') \leq 1$
- > let $\delta = \min \text{dist. in } X$, so $1/\delta = \begin{matrix} \text{"bit-precision"} \\ \text{of dataset } X \end{matrix}$
- > Next, solve $(c(1-\varepsilon), r)$ -NNS for
 $\gamma = \delta, (1+\varepsilon)\delta, (1+\varepsilon)^2\delta, \dots, 1$

(Approximately) Solve NNS using (c, r) -NNS

- > guess $\min_{p \in X} d(p, q)$ up to a mult. factor $1 \pm \varepsilon$.
- > Assume $\text{diam}(P) = \max_{p, p' \in X} d(p, p') \leq 1$
- > let $\delta = \min \text{dist. in } X$, so $1/\delta = \text{"bit-precision of dataset } X$
- > Next, solve $(c(1-\varepsilon), r)$ -NNS for
 $r = \delta, (1+\varepsilon)\delta, (1+\varepsilon)^2\delta, \dots, 1$
- > Report the minimal value of r for which we find a pt at dist $c(1-\varepsilon)$ of q .

Redⁿ to (C, r)-NNS (contd)

> This redⁿ imposes an additional $O(\log \frac{1}{\epsilon})$ overhead to the query time & memory

Redⁿ to (C, r)-NNS (contd)

> This redⁿ imposes an additional $O(\log \frac{1}{\epsilon})$ overhead to the query time & memory
(As we maintain a separate DS for each value of r above)

Redⁿ to (C, r) -NNS (contd)

- > This redⁿ imposes an additional $O(\log \frac{1}{\delta})$ overhead to the query time & memory
(As we maintain a separate DS for each value of r above)
- > So, for the rest of the talk, we focus on (C, r) -NNS !

LSH Families

LSH Families

> Kind of hash family that is sensitive to distance

LSH Families

- > Kind of hash family that is sensitive to distance
- > Ideally, want a hash fn. that maps "close points" to the same value w.h.p and maps "far points" to different values

LSH Families

- > Kind of hash family that is sensitive to distance
 - > Ideally, want a hash fn. that maps "close points" to the same value w.h.p and maps "far points" to different values. Formally,
- $\mathcal{H} = \{h: X \rightarrow \mathbb{Z}\}$ is $(c, c.r, p_1, p_2)$ -LSH if
 $\forall P, P' \in X:$

LSH Families

> Kind of hash family that is sensitive to distance
> Ideally, want a hash fn. that maps "close points" to the same value whp and maps "far points" to different values. Formally,

$\mathcal{H} = \{h: X \rightarrow \mathbb{Z}\}$ is (c, cr, p_1, p_2) -LSH if
 $\forall p, p' \in X:$

① $d(p, p') \leq r \Rightarrow \Pr[h(p) = h(p')] \geq p_1$

② $d(p, p') > cr \Rightarrow \Pr[h(p) = h(p')] \leq p_2$

LSH Families

- > Kind of hash family that is sensitive to distance
- > Ideally, want a hash fn. that maps "close points" to the same value whp and maps "far points" to different values. Formally,

$\mathcal{H} = \{h: X \rightarrow \mathbb{Z}\}$ is (c, cr, p_1, p_2) -LSH if
 $\forall p, p' \in X:$

$$\textcircled{1} \quad d(p, p') \leq r \Rightarrow \Pr[h(p) = h(p')] \geq p_1$$

$$\textcircled{2} \quad d(p, p') > cr \Rightarrow \Pr[h(p) = h(p')] \leq p_2$$

where the probabilities are over $h \sim \mathcal{H}$.

LSH Families (Contd)

$\mathcal{H} = \{h: X \rightarrow \mathbb{Z}\}$ is (c, cr, p_1, p_2) -LSH if
 $\forall p, p' \in X:$

① $d(p, p') \leq r \Rightarrow \Pr[h(p) = h(p')] \geq p_1$

② $d(p, p') \geq cr \Rightarrow \Pr[h(p) = h(p')] \leq p_2$

where the probabilities are over $n-h$.

→ Ideally, we want $p_1 \gg p_2$, but as we'll see,
this highly depends on c .

LSH Families (Contd)

$\mathcal{H} = \{h: X \rightarrow \mathbb{Z}\}$ is (c, cr, p_1, p_2) -LSH if
 $\forall p, p' \in X:$

① $d(p, p') \leq r \Rightarrow \Pr[h(p) = h(p')] \geq p_1$

② $d(p, p') \geq cr \Rightarrow \Pr[h(p) = h(p')] \leq p_2$

where the probabilities are over $n - H$.

- > Ideally, we want $p_1 \gg p_2$, but as we'll see, this highly depends on c .
- > Main idea of IM98: even if $p_1 > p_2$ slightly, can boost $p_1 \approx 1$ & $p_2 \approx 1/n$.

LSH family for binary vectors

LSH Family for binary vectors

> Say $U = \{0, 1\}^d$, $\text{dist}(p, p') := \|p - p'\|_1$, ie.
of coordinates in which p, p' differ

LSH Family for binary vectors

- > Say $U = \{0, 1\}^d$, $\text{dist}(p, p') := \|p - p'\|_1$, ie.
of coordinates in which p, p' differ
- > Consider $f = \{\phi_i\}_{i=1}^d$, where
 $\phi_i(p) = p_i$

LSH Family for binary vectors

- > Say $U = \{0, 1\}^d$, $\text{dist}(p, p') := \|p - p'\|_1$, ie.
of coordinates in which p, p' differ
- > Consider $h = \{h_i\}_{i=1}^d$, where
$$h_i(p) = p_i$$
- > Observe: for $p, p' \in U$,
$$\Pr[h(p) = h(p')] = \frac{\# \text{ agreeing words}}{\text{total # words}} = \frac{d - \|p - p'\|_1}{d} = \frac{1 - \|p - p'\|_1}{d}$$

LSH Family for binary vectors

- > Say $U = \{0, 1\}^d$, $\text{dist}(p, p') := \|p - p'\|_1$, ie. # of coordinates in which p, p' differ
- > Consider $h = \{h_i\}_{i=1}^d$, where $h_i(p) = p_i$
- > Observe: for $p, p' \in U$,
- $$\Pr[h(p) = h(p')] = \frac{\# \text{ agreeing words}}{\text{total \# words}} = \frac{d - \|p - p'\|_1}{d} = \frac{d - \|p - p'\|_1}{d}$$
- > So, $\Pr[h(p) = h(p')] = \begin{cases} \geq 1 - \frac{r}{d} \approx e^{-r/d}, & \text{if } d(p, p') \leq r \\ \leq 1 - \frac{cr}{d} \approx e^{-c r / d}, & \text{if } d(p, p') \geq c r \end{cases}$

LSH Family for binary vectors

- > Say $U = \{0, 1\}^d$, $\text{dist}(p, p') := \|p - p'\|_1$, ie. # of coordinates in which p, p' differ
- > Consider $\mathcal{H} = \{h_i\}_{i=1}^d$, where $h_i(p) = p_i$

Conclusion: \mathcal{H} is $(c, cr, e^{-r/d}, e^{-cr/d})$ -LSH
- > Observe: for $p, p' \in U$,
 $\Pr[h(p) = h(p')] = \frac{\# \text{ agreeing words}}{\text{total \# words}} = \frac{d - \|p - p'\|_1}{d} = \frac{1 - \frac{\|p - p'\|_1}{d}}{d}$
- > So, $\Pr[h(p) = h(p')] = \begin{cases} \geq 1 - \frac{r}{d} \approx e^{-r/d}, & \text{if } d(p, p') \leq r \\ \leq 1 - \frac{cr}{d} \approx e^{-cr/d}, & \text{if } d(p, p') \geq cr \end{cases}$

LSH Family \Rightarrow NNS - data structures!

LSH Family \Rightarrow NNS - data structures!

> Suppose we had a (r, cr, p_1, p_2) -LSH fam. st.
 $p_1 \times 1$, $p_2 \approx 0$. How to construct DS?

LSH Family \Rightarrow NNS - data structures!

- > Suppose we had a (r, cr, p_1, p_2) -LSH fam. st. $p_1 \times 1, p_2 \approx 0$. How to construct DS?
- > Choose $h \sim H$ w.r.t store $h(p) \in \mathbb{R}^X$.

LSH Family \Rightarrow NNS - data structures!

- > Suppose we had a (r, cr, p_1, p_2) -LSH fam. s.t.
 $p_1 \times 1, p_2 \approx 0$. How to construct DS?
- > Choose $h \sim H$ w.r.t & store $h(p) \forall p \in X$.
- > Given $q \in U$, compute $h(q)$ & see if $\exists p \in X$ with
 $h(p) = h(q)$.

LSH Family \Rightarrow NNS - data structures!

- > Suppose we had a (r, cr, p_1, p_2) -LSH fam. s.t.
 $p_1 \times 1, p_2 \approx 0$. How to construct DS?
- > Choose $h \sim H$ w.r.t & store $h(p) \in \mathbb{P} \in \mathcal{X}$.
- > Given $q \in \mathcal{U}$, compute $h(q)$ & see if $\exists p \in \mathcal{X}$ with
 $h(p) = h(q)$. Can do this in $O(1)$ time using a hash table!

LSH Family \Rightarrow NNS - data structures!

- > Suppose we had a (r, cr, p_1, p_2) -LSH fam. st. $p_1 \times 1, p_2 \approx 0$. How to construct DS?
- > Choose $h \sim H$ hash & store $h(p) \in \mathbb{P} \in \mathcal{X}$.
- > Given $q \in U$, compute $h(q)$ & see if $\exists p \in \mathcal{X}$ with $h(p) = h(q)$. Can do this in $O(1)$ time using a hash table!
- > If there is no such pt p , then whp, there is no pt at dist. $\leq \overline{cr}$ of q .

LSH Family \Rightarrow NNS - data structures!

- > Suppose we had a (r, cr, p_1, p_2) -LSH fam. st. $p_1 \approx 1, p_2 \approx 0$. How to construct DS?
- > Choose $h \sim H$ hash & store $h(p) \in \mathbb{F}$.
- > Given $q \in U$, compute $h(q)$ & see if $\exists p \in X$ with $h(p) = h(q)$. Can do this in $O(1)$ time using a hash table!
- > If there is no such pt p , then w.h.p., there is no pt at dist. $\leq \overline{cr}$ of q .
- > So, only need to show: given (r, cr, p_1, p_2) -LSH fam with $p_1 > p_2$, we can "boost" it to get $p_1 \approx 1, p_2 \approx 0$.

LSH Parameter Boosting

> Two steps: first just make p_L small

LSH Parameter Boosting

- > Two steps: first just make ρ_L small
- > For this, take k independent hashes from \mathcal{H} & hash each PEX to

$$h(p) = (h_1(p), \dots, h_k(p)) \in \mathbb{Z}^k$$

LSH Parameter Boosting

- > Two steps: first just make ρ_L small
- > For this, take k independent hashes from \mathcal{H} & hash each $p \in X$ to
$$h(p) = (h_1(p), \dots, h_k(p)) \in \mathbb{Z}^k$$
- > Then $\forall p, p'$, $d(p, p') \geq c\tau \Rightarrow \Pr[h(p) = h(p')] \leq \rho_L^k$.

LSH Parameter Boosting

- > Two steps: first just make p_L small
- > For this, take k independent hashes from \mathcal{H} & hash each $P \in \mathcal{X}$ to

$$h(p) = (h_1(p), \dots, h_k(p)) \in \mathbb{Z}^k$$

- > Then $\forall p, p'$, $d(p, p') \geq c\tau \Rightarrow \Pr[h(p) = h(p')] \leq p_2^k$.
- > Great, but this doesn't $\uparrow p_1$! In fact, 2 "close" pts are mapped to same vector only with prob. p_1^k .

Boosting $P_1 \uparrow$

Boosting $P_1 \uparrow$

> choose l ind. topics of above k-dim hash f_1, \dots, f_L !

Boosting $P_1 \uparrow$

- > choose l ind. topics of above k -dim hash f_1, \dots, f_L !
- > we show that if pts are close, then w.h.p,
at least one of the hashes agree.

Boosting $P_1 \uparrow$

- > choose l ind. topics of above k -dim hash f_1, \dots, f_L !
- > we show that if pts are close, then w.h.p.,
at least one of the hashes agree.
- > Assume $d(p, q) \leq r$. Then:

Boosting $p_i \uparrow$

- > choose l ind. topics of above k -dim hash f_1, \dots, f_L !
- > we show that if pts are close, then w.h.p.,
at least one of the hashes agree.
- > Assume $d(p, q) \leq r$. Then:
$$\Pr[d_i | f_i(p) = f_i(q)] = 1 - \Pr[\forall i \mid f_i(p) \neq f_i(q)]$$
$$= 1 - \Pr[f_i(p) \neq f_i(q)]^l$$
$$\geq 1 - (1 - p_i^k)^l$$

Boosting $p_i \uparrow$

- > choose l ind. topics of above k -dim hash f_1, \dots, f_L !
- > we show that if pts are close, then w.h.p.,
at least one of the hashes agree.
- > Assume $d(p, q) \leq r$. Then:

$$\begin{aligned} \Pr[d_i | f_i(p) = f_i(q)] &= 1 - \Pr[\exists i \mid f_i(p) \neq f_i(q)] \\ &= 1 - \Pr[f_i(p) \neq f_i(q)]^l \\ &\geq 1 - (1 - p_i^k)^l \end{aligned}$$

gotta make sure this is small!

The Algorithm

Algorithm 1 LSH Algorithm

Preprocessing:

Choose $k \cdot \ell$, $h_{1,1}, \dots, h_{\ell,k}$ functions uniformly at random from \mathcal{H} .

Construct ℓ hash tables; for all $1 \leq i \leq \ell$ store $f_i(p) = (h_{i,1}(p), \dots, h_{i,k}(p))$ for all $p \in P$ in the i -th has table.

Query(q):

for $i = 1 \rightarrow \ell$ **do**

 Compute $f_i(q)$.

 Go over all points p where $f_i(p) = f_i(q)$. For all such points if $\text{dist}(p, q) \leq c \cdot r$, output p .

end for

Parameter Tuning

Parameter

Tuning

∴ Choose k s.t. $P_2^k = 1/n$.

Parameter Tuning

- Choose k s.t. $P_2^k = 1/n$.
- Assume $P_1 = P_2^f$ (recall $P_1 > P_2 \Rightarrow f < 1$).

Parameter Tuning

- ↪ Choose k s.t. $p_2^k = 1/n$.
- ↪ Assume $p_1 = p_2^f$ (recall $p_1 > p_2 \Rightarrow f < 1$).
- ↪ Choose $\ell = n^{-f} \cdot \ln n$

Parameter Tuning

- ↪ Choose k s.t. $p_2^k = 1/n$.
- ↪ Assume $p_1 = p_2^f$ (recall $p_1 > p_2$, $\Rightarrow f < 1$).
- ↪ Choose $\ell = n^{-f} \cdot \ln n$
- > Fix a query $q \in V$. By lin. of exp., H_i ,
 $\Pr[\exists p : d(p, q) > cr, f_i(p) = f_i(q)] = n \cdot p_2^k \leq 1$

Parameter Tuning

- ↪ Choose k s.t. $p_2^k = 1/n$.
- ↪ Assume $p_1 = p_2^f$ (recall $p_1 > p_2 \Rightarrow f < 1$).
- ↪ Choose $\ell = n^{-f} \cdot \ln n$
- ↪ Fix α every $q \in V$. By lin. of exp., $\forall i$,
 $\Pr[\exists p : d(p, q) > \alpha, f_i(p) = f_i(q)] = n \cdot p_2^k \leq 1$
- ↪ Summing over i , in exp., $O(\ell)$ far pts hash to same value as q for some i . $\Rightarrow O(\ell)$ overhead in q -time.

Parameter Tuning (contd)

Given $d(p, q) \leq \epsilon$ for some $p \in X$, then,

$$\Pr[\exists i : f_i(p) = f_i(q)] \geq 1 - ((1 - p_1^k)^l$$
$$= 1 - (1 - p_2^{8k})^l$$
$$= 1 - (1 - n^{-\beta})^l$$
$$\approx 1 - e^{ln^{-\beta}} = 1 - 1/n.$$

Parameter Tuning (contd)

- Given $d(p, q) \leq r$ for some $p \in X$, then,
- $$\Pr[\exists i : f_i(p) = f_i(q)] \geq 1 - ((1 - p_1^k)^l$$
- $$= 1 - (1 - p_2^{8k})^l$$
- $$= 1 - (1 - n^{-\beta})^l$$
- $$\approx 1 - e^{\ln(-\beta)} = 1 - 1/n.$$
- Summary: for a pt q with $d(p, q) \leq r$, our algo outputs p w/ $P \geq 1 - \frac{1}{n}$. In exp., $O(l \cdot d)$ overhead to examine $O(l)$ far pts with same hash value.

Space Analysis

Space Analysis

> Maintain $O(l)$ hash tables.

Space Analysis

- > Maintain $O(l)$ hash tables.
- > In each table, store $n = |X|$ hash values, each a k -dim vector.

Space Analysis

- > Maintain $O(l)$ hash tables.
- > In each table, store $n = |X|$ hash values, each a k -dim vector.
- > Overall, $O(l \cdot n \cdot k) = O(n^{1+\beta} \log \frac{n}{p_e})$

Query Time Analysis

Query Time Analysis

> $O(\lg l)$ time to some $f_i(q) \wedge i \in [l]$

Query Time Analysis

- > $O(\lg l)$ time to some $f_i(q) \leq i \in [l]$
- > any candidate close pt p , spend $O(d)$ time to calculate $\text{dist}(p, q)$.

Query Time Analysis

- > $O(\ell q)$ time to some $f_i(q) \cup i \in [\ell]$
- > any candidate close pt p , spend $O(d)$ time to calculate $dist(p, q)$.
- > but $|O| = \text{size of output} = \# \text{ of pts at dist} \leq cr \text{ from } q$

Query Time Analysis

- > $O(\ell q)$ time to some $f_i(q) \leq i \in [\ell]$
- > any candidate close pt p , spend $O(d)$ time to calculate $de(p, q)$.
- > let $|O| = \text{size of output} = \# \text{ of pts at dist} \leq cr \text{ from } q$
- > In exp., we examine $O(\ell)$ far pts that we don't output. So, overall, time is $O(d(\ell + |O|) + \ell \cdot k) = O(n^{\frac{d}{2}}(d + \log \frac{n}{P_2}) + O(d))$

Conclusion

> Memory $O(n^{1+5})$, q-time $O(n^5)$

Conclusion

> Memory $O(n^{1+S})$, q-time $O(n^S)$

> for the binary vector example,

recall $p_1^S = p_2 \Rightarrow$

$$g = \frac{\ln 1/p_1}{\ln 1/p_2} = \frac{r/d}{c r/d} = \frac{1}{c}$$

- > For eg if $c=2$, need $O(n^{1.5})$ space for hash tables & $O(\sqrt{n})$ q-time.
- > q-time & memory get significantly better as we relax c.

for other metrics

for other metrics

> For ℓ_1 , can 'reduce' to the
binary vector case using a
'coordinate red' idea to get
 $f = \|c\|$ as well.

For ℓ_2 :

- > Use Johnson-Lindenstrauss to reduce dimension d to $\log n$.
- > Assume all pts in X & query pts lie in S^d . Then $h(x) = \text{sgn}(\langle g, x \rangle)$ where $g = (g_i)_{i=1}^d$ is sf.
 $g_i \sim N(0, 1)$

Picture

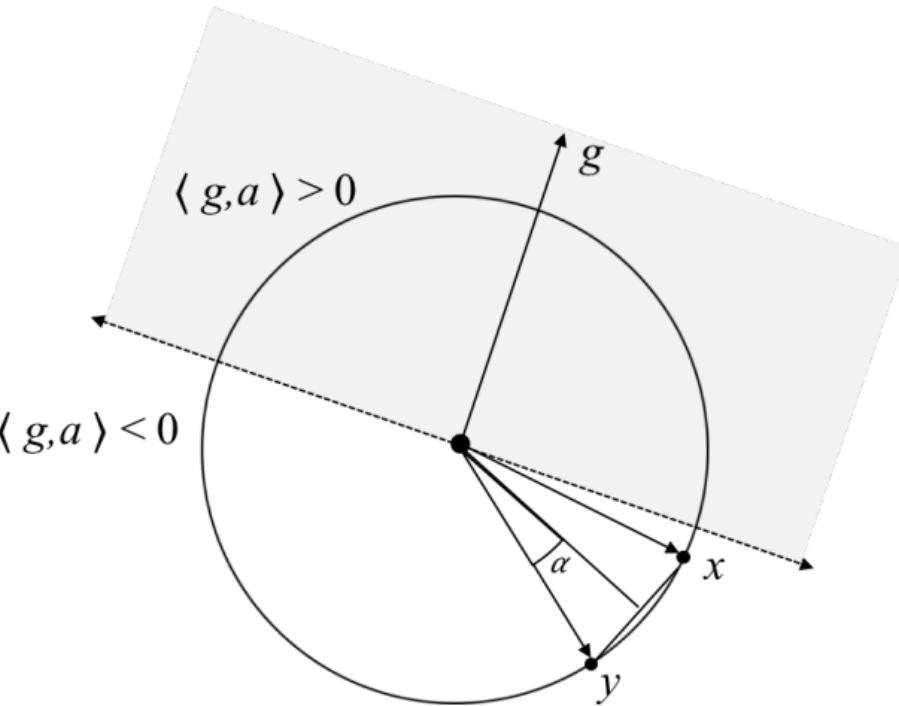


Figure 4: Our Euclidean locality sensitive hash function draws a random Gaussian vector g . Consider the hyperplane orthogonal to g and passing through $\{0\}^d$. All points above the hyperplane are hashed to 1 and all points below are hashed to -1 .

> Overall, each hash fn. g in the LSH scheme is a bit vector constructed by taking the sign of each entry in Πx , $\Pi =$ random gaussian matrix.

- > Overall, each hash fn. g in the LSH scheme is a bit vector constructed by taking the sign of each entry in Πx , $\Pi =$ random gaussian matrix.
- > this also gives a $g = I/c$ construction.

- > Overall, each hash fn. g in the LSH scheme is a bit vector constructed by taking the sign of each entry in Πx , $\Pi =$ random gaussian matrix.
- > this also gives a $\|c\|_1$ construction.
- > Possible to improve to (the optimal) $\|c\|_2$ using a more complex LSH scheme.

Recent Developments

(On board)

Thank you!