

Approximate Nearest Neighbor
using

Locality Sensitive Hashing

IM (198)

Based on Lecture Notes by Sariel Har-Peled

Presented by

K.V.N. Sreenivas, PhD Student, IISc Bengaluru

The Nearest Neighbor Problem

- Input: Set of points P in a metric Space (X, dist)
- Goal: Preprocess them such that the following queries can be efficiently solved.

Query: Given a point $q \in X$, output the point $p \in P$ such that $\text{dist}(p, q)$ is minimized.



The Nearest Neighbor Problem

Trade off between Preprocessing Space and Query time.

One extreme:

Don't preprocess (constant preprocessing space)

$|P| = n$ query time

Other extreme:

Compute Voronoi Diagram ($n^{d/2}$ space required)

Constant Query time.

Approximate Nearest Neighbor (ANN)

Input: A set of points P in a metric space $(\mathcal{X}, \text{dist})$

Goal: Preprocess set P s.t. for query $q \in \mathcal{X}$, return $p \in P$

$$\text{dist}(q, p) \leq (1 + \epsilon) \text{dist}(q, p^*)$$

where

p^* = nearest neighbor of q in P $\equiv: nn(q)$.

Approximate Nearest Neighbor (ANN)

We will assume $x \in \mathbb{R}^d$.

Naive method: nd query time

Goal : $\text{Sub}(n) \cdot d$ query time with reasonable space

This Presentation:

Query time: $\tilde{\mathcal{O}}(d n^{\frac{1}{1+\epsilon}})$ whp

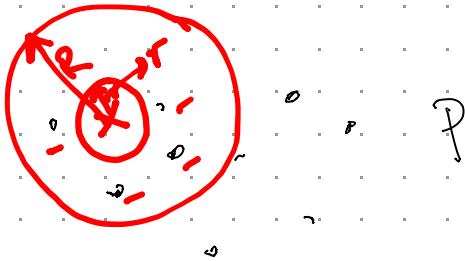
Preprocessing space: $\tilde{\mathcal{O}}(nd + n^{\frac{1}{1+\epsilon}})$

A Simpler Problem: Proximity Neighbor (PN)

PN (q, r):

For a given query point q :

- 1) If there is $p \in \text{Ball}(q, r)$, return it.
- 2) If no point in $\text{Ball}(q, r)$, say so.



Approximate Proximity Neighbor (APN)

Fix accuracy parameter ϵ .

APN(q, r):

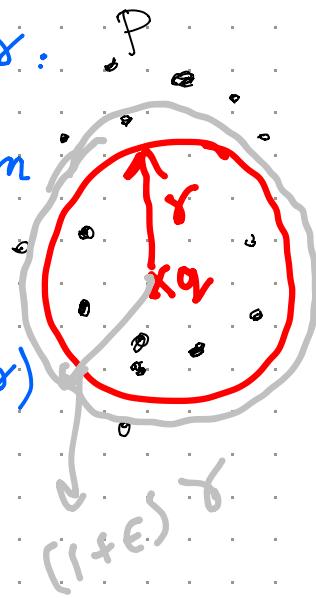
For a query point q , and a radius r :

- 1) If there is a point in $\text{Ball}(q, r)$, return some point in $\text{Ball}(q, (1+\epsilon)r)$.

- 2) If all points are outside $\text{Ball}(q, (1+\epsilon)r)$

Then report "No point in $\text{Ball}(q, r)$ ".

- 3) The intermediate case is immaterial.



$$\underline{\text{ANN}} \cong \underline{\text{APN}}$$

$$\underline{\text{APN} \leq \text{ANN}}$$

Input to APN: query q , radius γ

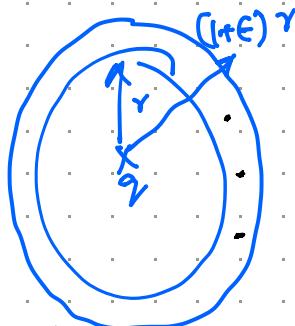
Directly solve the ANN problem with query q .

Returns p s.t. $\text{dist}(p, q) \leq (1+\epsilon) \text{dist}(\text{nn}(q), q)$

Two cases: Case 1: $\text{dist}(p, q) \leq (1+\epsilon)\gamma$. Then we return q

Case 2: $\text{dist}(p, q) > (1+\epsilon)\gamma \Rightarrow \text{dist}(\text{nn}(q), q) > r$

Return "No point inside the ball of radius r ".

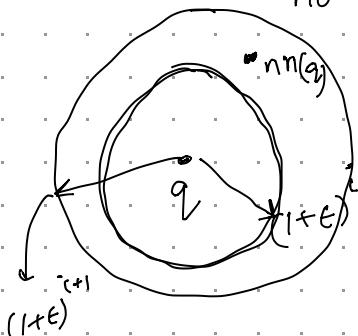
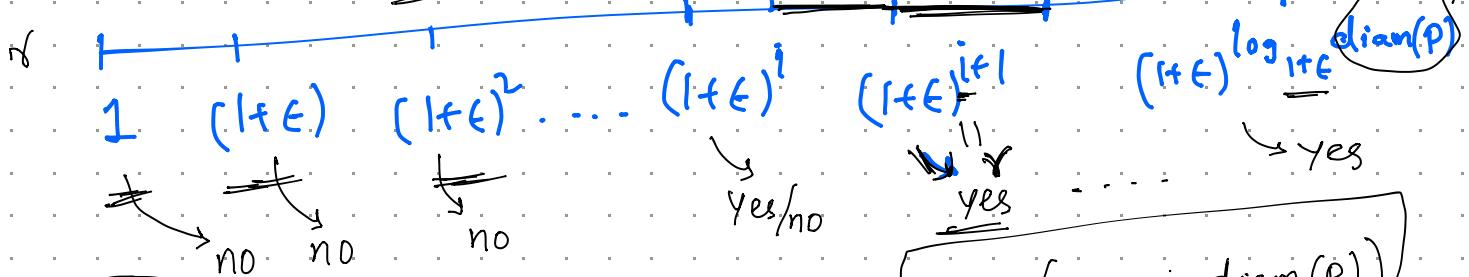


$$\underline{\underline{ANN}} \cong \underline{\underline{APN}}$$

$[1, \underline{\text{diam}(P)}]$

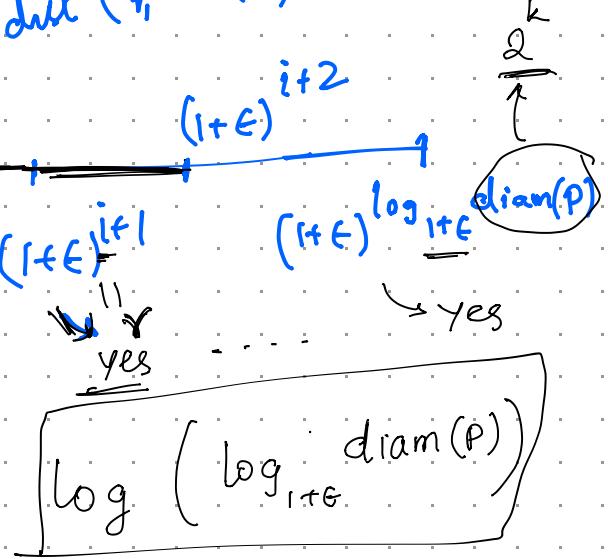
$$\underline{\underline{ANN}} \leq (\log(\log_{1+\epsilon} \text{diam}(P))) \underline{\underline{APN}}$$

$$\log(k)$$



$$\underline{\underline{\text{dist}(P, q)}} \leq (1+\epsilon)^{i+2} = \underline{\underline{(1+\epsilon)^2}} \underline{\underline{(1+\epsilon)^i}}$$

$$\leq \underline{\underline{(1+3\epsilon)}} \underline{\underline{\text{dist}(nn(q), q)}}$$



SUMMARY SO FAR

Nearest Neighbor (Hard)

Proximity Query

Approximate Nearest Neighbor - ANN

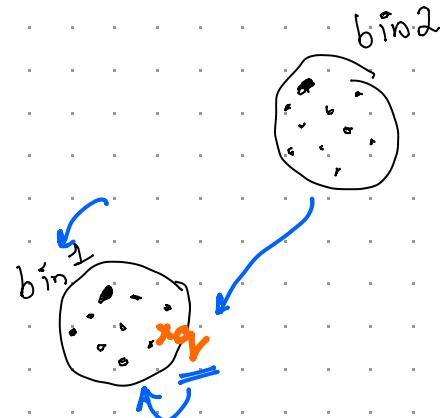
Approximate Proximity
Query - APQ

IDEAL SCENARIO

Store points in P in a hash table such that

Close points → Same bin

Far points → different bins.



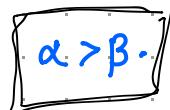
LOCALITY SENSITIVE HASHING (LSH)

A family \mathcal{F} of hash functions is said to be (r, R, α, β) -sensitive if for any $p, q \in \mathcal{U}$, we have the following:

A) If $q \in \text{Ball}(p, r)$, then $P[f(p) = f(q)] \geq \alpha$

B) If $q \notin \text{Ball}(p, R)$, then $P[f(p) = f(q)] \leq \beta$

We also require $\alpha > \beta$.



$= 1 - \beta$

AN LSH FAMILY FOR $(\{0,1\}^d, \|\cdot\|_1)$

The universe is the corners of d-dim unit hypercube.

$$f_i(p) = i^{\text{th}} \text{ coord of } p \text{ for all } i \in [d]$$

Claim: $\mathcal{F} = \{f_1, f_2, \dots, f_d\}$ is (r, R, α, β) -sensitive where

$$\alpha = 1 - \frac{r}{d} \quad \beta = 1 - \frac{R}{d} \quad (\text{for } r < R, \alpha > \beta).$$

Proof: (A) $\|p - q\|_1 \leq r$. We chose $f_i \sim \mathcal{F}$. $P[f_i(p) = f_i(q)] \geq 1 - \frac{r}{d}$

$$(B) \|p - q\|_1 \geq R. \quad P[f_i(p) = f_i(q)] \leq 1 - \frac{R}{d}$$

~~$d - R$~~

NAIVE APQ^(*) ($\exists \gamma$ which is $(r, (1+\epsilon)r, \alpha, \beta) - \text{LSH}$)

Preprocessing (The point set is P)

For every $p \in P$, we compute $f(p)$ where $f \sim \mathcal{F}$.

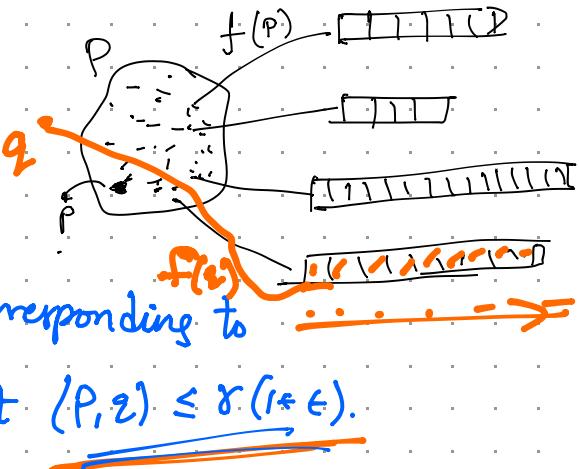
For a query q .

We compute $f(q)$

We search through the list/bucket corresponding to $f(q)$ for a point p s.t. $\text{dist}(p, q) \leq r(1+\epsilon)$.

If such point p exists, return p .

O/w, say "no point in $\text{Ball}(q, r)$ ".



NAIVE APQ (Analysis)

Suppose query q .

Expec Time to answer query q : in bucket comp. to $f(q)$

Once we encounter a point p , s.t. $\text{dist}(p, q) \leq (1 + \epsilon)\gamma$.

of points p s.t. $\text{dist}(p, q) > (1 + \epsilon)\gamma$:

the prob that p, q match to same bucket

(i.e., $f(p) = f(q)$) is at most β . (β is a constant).

So, # $\leq \underline{\beta(n)} \underline{[(\text{constant}) n]}$

NAIVE APQ

If there is p s.t

$$\text{dist}(p, q) \leq ((1+\epsilon)r,$$

then we have

a probability

of at least α $\left[\alpha = \text{constant} \right]$

$$(f_1, f_2, \dots, f_g) \sim \mathcal{F}$$

We deem that two points p, q collide iff

$$f(p) = f(q)$$

$$\# p \text{ s.t } \text{dist}(p, q) \geq (1+\epsilon)r,$$

$$\equiv \underline{\beta^n} \text{ (linear)} \quad \checkmark$$

$$f \sim \mathcal{F}$$

$$(f_1, f_2, \dots, f_k) \sim \mathcal{F}$$

We deem that two points $p \leq q$

"collide" iff

$$f_1(p) = f_1(q) \wedge f_2(p) = f_2(q) \dots$$

$$\# p \text{ s.t } \text{dist}(p, q) \geq (1+\epsilon)r \text{ & }$$

$$\Lambda(\mathcal{F}, k)(p) = \Lambda(\mathcal{F}, k)(q) \rightarrow \underline{(\beta^k n)} \quad \checkmark$$

$$\underline{f_1(p) = f_1(q)} \vee \underline{f_2(p) = f_2(q)} \vee \dots$$

Then,

Case 1 (p, q) dist is $\leq (1+\epsilon)r$

$$P[V(\pi, \tau)(p) = V(\pi, \tau)(q)]$$

$$= 1 - (1-\alpha)^T$$

It can be shown that

$$\# p \text{ s.t } V(\pi, \tau)(p) = V(\pi, \tau)(q)$$

but $\text{dist}(p, q) \geq (1+\epsilon)r$.

$$\leq \underbrace{(1 - (1-\beta)^T)}_{\text{crossed out}} n = (1 - o(1))n$$

Suppose $\text{dist}(p, q) \leq (1+\epsilon)r$,

then

$$f(p) = f_1(q) \text{ & } f_2(p) = f_2(q) \text{ & } \dots$$

↓

$$p \rightarrow \alpha^k$$

COMBINED APQ

$$\mathcal{H} = V(\Lambda(\mathcal{F}, k), \mathcal{T}) \rightarrow k, \mathcal{T} \text{ are unknowns.}$$

will be fixed later.

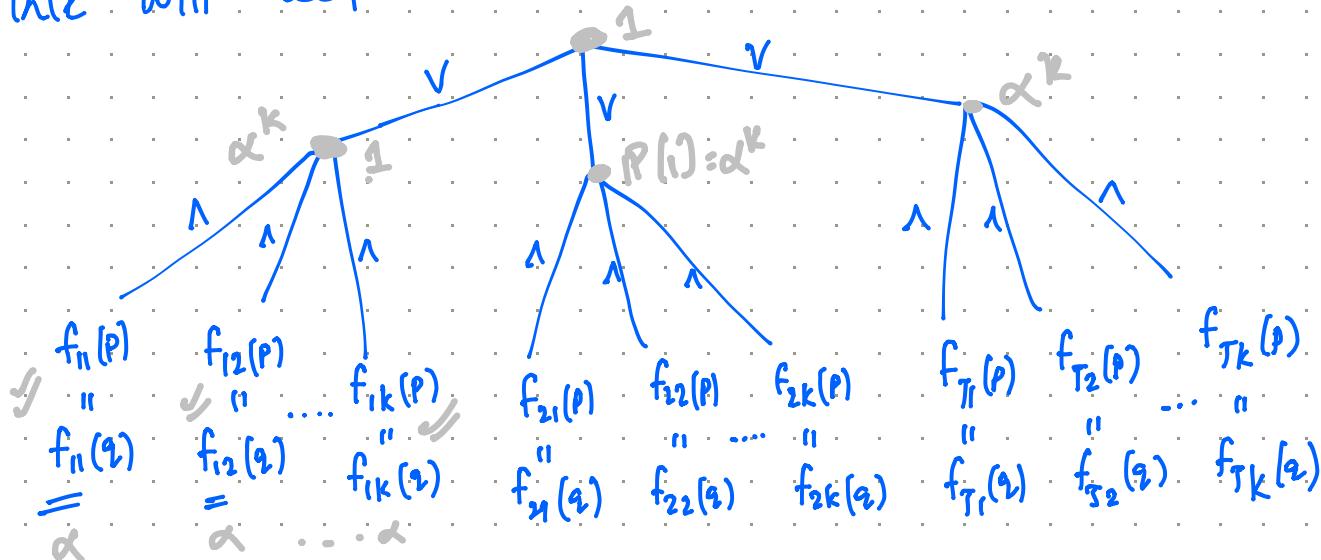
A hash function $h \in \mathcal{H}$, given as input a point P ,
computes the following matrix.

$$h(P) = \begin{bmatrix} f_{11}(P) & f_{12}(P) & \dots & f_{1k}(P) \\ f_{21}(P) & f_{22}(P) & \dots & f_{2k}(P) \\ \vdots & & & \\ f_{\mathcal{T}1}(P) & f_{\mathcal{T}2}(P) & \dots & f_{\mathcal{T}k}(P) \end{bmatrix}$$

each f_{ij} ($i \in [\mathcal{T}]$ and $j \in [k]$) is unif. rand. $\sim \mathcal{F}$.

COMBINED APO :

we will need a collision criterion to decide when $h(p) = h(q)$



Lemma:- Given an (r, R, α, β) -sensitive LSH \mathcal{F} , the family $\mathcal{H} = V(\lambda(\mathcal{F}, k), \tau)$ is a $(r, R, 1 - (1 - \alpha^k)^\tau, 1 - (1 - \beta^k)^\tau)$ -sensitive LSH.

Proof: Suppose $\text{dist}(p, q) \leq r$. We know that $\Pr_{f \in \mathcal{F}}[f(p) = f(q)] > \alpha$

$$\text{So, } \Pr_{\substack{h \in \mathcal{H} \\ h \neq f}}[h(p) = h(q)] \geq 1 - (1 - \alpha^k)^\tau$$

Similarly, if $\text{dist}(p, q) > R$, we can show that

$$\Pr_{\substack{h \in \mathcal{H} \\ h \neq f}}[h(q) = h(p)] \leq 1 - (1 - \beta^k)^\tau$$



Choosing k, T

$$\text{dist}(p, q) \leq r \implies P[h(p) = h(q)] \geq 1 - (1-\alpha^k)^T \quad // \quad \alpha > \beta$$
$$\text{dist}(p, q) > (1+\epsilon)r \implies P[h(p) = h(q)] \leq 1 - (1-\beta^k)^T$$

Impose the cond'n that $1 - (1-\alpha^k)^T = 3/4$

$$(1-\alpha^k)^T = \frac{1}{4}$$

$$\frac{1}{4} = (1-\alpha^k)^T \leq \exp(-\alpha^k T) \xrightarrow{\alpha^k > 0} e^{-4} > \frac{1}{4}$$

It suffices to take $T = \frac{4}{\alpha^k}$

$$T = \frac{4}{\alpha^k}$$

Prob. that a point in $P \setminus \text{Ball}(q, (1+\epsilon)\gamma)$ that collides with a given point q is

$$\begin{aligned} &\leq 1 - (1-\beta^k)^J = (1-(1-\beta^k)) \left(1 + \underbrace{(1-\beta^k)}_{\leq 1} + \underbrace{(1-\beta^k)^2}_{\leq 1} + \dots \right) \\ &\leq \beta^k J = \underline{\underline{4 \left(\frac{\beta}{\alpha}\right)^k}} \end{aligned}$$

Query time :- $h(q)$ and they go thru the Bucket B which contains all points \underline{p} s.t $h(q)=h(p)$.

If we find p s.t $\text{dist}(p, q) \leq (1+\epsilon)\gamma$,

If the num of points \underline{p} s.t $\text{dist}(p, q) > (1+\epsilon)\gamma$

Expec num of p s.t. $\text{dist}(p, q) > (1+\epsilon) r$ and $h(p) = h(q)$

is given by $= 4n \left(\frac{\beta}{\alpha}\right)^k$.

and to compute $h(q)$ itself, we need Tk time.

Query time $\approx \left(T_k + 4n \left(\frac{\beta}{\alpha}\right)^k \right) d$

Take

$$T_k = 4n \left(\frac{\beta}{\alpha}\right)^k$$

$$= n \beta^k T$$

\Rightarrow

$$K = n \beta^k$$

$$\left(\frac{r}{\beta}\right)^k = n$$

$$n \beta^k = 1$$

If we choose

$$k = \log \left(\frac{n}{\beta}\right) =$$

$$\text{Query time} = d \left(\tau k + 4n \left(\frac{p}{\alpha} \right)^k \right)$$

$$\leq d \left(\tau k + 4 \frac{(k)}{\alpha^k} \right)$$

$$= \underline{2d \tau k} \approx \underline{O(d \tau k)}$$

$$k = \frac{\ln n}{\ln(1/\beta)}$$

$$\tau = \frac{n}{\alpha^k}$$

$$\begin{aligned} &= \frac{4}{\alpha^{\frac{\ln n}{\ln(1/\beta)}}} = \frac{4}{(\alpha^{\ln n})^{\frac{1}{\ln(1/\beta)}}} \\ &= 4(n^{\ln(\frac{1}{\beta})})^{\frac{1}{\ln(1/\beta)}} \end{aligned}$$

$$= 4n \frac{\ln(\frac{1}{\alpha})}{\ln(\frac{1}{\beta})}$$

$$\underline{\underline{\alpha = 1 - \frac{r}{d}}}, \underline{\underline{\beta = 1 - \frac{(1+\epsilon)\gamma}{d}}} = \underline{\underline{T}}$$

Query time $\approx 2d \left(4n \frac{\ln(\frac{1}{\alpha})}{\ln(\frac{1}{\beta})} \cdot \frac{\ln(n)}{\ln(\frac{1}{\beta})} \right)$

$$\approx O \left(n \frac{\ln(\frac{1}{\alpha}) \epsilon}{\ln(\frac{1}{\beta})} (\ln(n) d) \right)$$

Prop: For $x \in (0, 1)$ & $t \geq 1$ s.t $1-tx > 0$, we have

$$\frac{\ln\left(\frac{t}{1-tx}\right)}{\ln\left(\frac{1}{1-tx}\right)} \leq \frac{1}{t}$$

A) $g(x) = \underline{(1-tx)} - \underline{(1-x)^t}$

$$g'(x) = -t + t(1-x)^{t-1}$$
$$= t \left((1-x)^{t-1} - 1 \right)$$
$$\leq 0 .$$

$$g(0) = (1-0) - (1-0)^t = 0.$$

$g(x) \leq 0$ for all $x \in (0, 1)$.

$$(1-tx) \leq (1-x)^t$$

$$\ln(1-tx) \leq t \ln(1-x)$$

$$\ln\left(\frac{1}{1-tx}\right) \geq t \ln\left(\frac{1}{1-x}\right)$$

\Rightarrow

$$\frac{\ln\left(\frac{1}{1-x}\right)}{\ln\left(\frac{1}{1-tx}\right)} \leq \frac{1}{t}.$$

Query time = $O\left(n \frac{\ln(\frac{1}{1-\gamma/d})}{\ln(\frac{1}{1-(1+\epsilon)\gamma/d})} (\ln n) d\right)$

$$\leq O\left(n^{\frac{1}{1+\epsilon}} (\ln n) d\right)$$

$n^{1/3}$

Preprocessing space :-

$$(n)(d)(\underline{J^k}) = \cancel{O}\left(n^{\frac{1}{1+\frac{1}{1+\epsilon}}} (\ln n) d\right)$$

$n^{4/3}$

Main Thm : (Assuming $(\mathcal{F}, \text{dist})$ has $(\gamma, R, \alpha, \beta)$ -LSH family)

(A^{PPG}) It is possible to preprocess a set of n points in

\mathbb{R}^d such that

Preprocessing space is

$$O\left(n^{1 + \frac{\ln(\frac{1}{\alpha})}{\ln(\frac{1}{\beta})}} (\ln n)^d\right)$$

Query time is

$$\underline{O}\left(n^{\frac{\ln(\frac{1}{\alpha})}{\ln(\frac{1}{\beta})}} (\ln n)^d\right)$$

Each query is successful with prob at least $\frac{3}{4}$.

LSH on Euclidean Metric (\mathbb{R}^d , dist)

Lemma :- Let x_1, x_2, \dots, x_d be normally distributed independent random variables. Let $v_1, \dots, v_d \in \mathbb{R}$. Then the variable $y = v_1 x_1 + v_2 x_2 + \dots + v_d x_d$ will be distributed as $N(0, \sigma^2 = \|v\|^2)$ where $\vec{v} = (v_1, v_2, \dots, v_d)$

Proof :- $\phi_x(t) = \mathbb{E}[e^{itx}]$

It can be shown that for $N(0, \sigma^2)$, the

$$\phi_{N(0, \sigma^2)}(t) = \exp\left(-\frac{\sigma^2 t^2}{2}\right)$$

$$Y = v_1 X_1 + \dots + v_d X_d$$

$$\phi_Y(t) = \mathbb{E}\left[e^{it(v_1 X_1 + v_2 X_2 + \dots + v_d X_d)}\right]$$

$$= \prod_{j=1}^d \mathbb{E}\left[e^{it(v_j X_j)}\right]$$

$$= \prod_{j=1}^d \exp\left(-\frac{v_j^2 t^2}{2}\right)$$

$$\phi_Y(t) = \exp\left(-\left(v_1^2 + v_2^2 + \dots + v_d^2\right) \frac{t^2}{2}\right)$$

$$\gamma \sim \mathcal{N}(0, \sigma^2 = \|\mathbf{v}\|^2).$$



LSH for Euclidean Metric?

Fix a vector \vec{v} . Then, let's define for a point $p \in \mathbb{R}^d$

$$h(p) = \left\lfloor \frac{\langle p, \vec{v} \rangle + t}{x} \right\rfloor \quad \text{where } t \sim \text{Unif}[0, x] \text{ and}$$

x is a parameter which we will choose shortly.

$$\mathcal{H} = \left\{ h : \vec{v} \text{ is a vector } \sim \mathcal{N}(0, I)^d \right\}$$

We need to choose α s.t \mathcal{H} is LSH.

Fix p, q, \vec{v}, α . Let $\|p - q\| =: n$ and

$$|\langle p, \vec{v} \rangle - \langle q, \vec{v} \rangle| =: \beta.$$

$$\langle p, \vec{v} \rangle - \langle q, \vec{v} \rangle = \beta \text{ (Assume)}$$



Fix some p, q, \vec{v}, α .

$$P[h(p) = h(q)]$$

$$h(p), h(q) \Leftrightarrow \left\lfloor \frac{\langle p, \vec{v} \rangle + t}{\alpha} \right\rfloor = \left\lfloor \frac{\langle q, \vec{v} \rangle + t}{\alpha} \right\rfloor$$

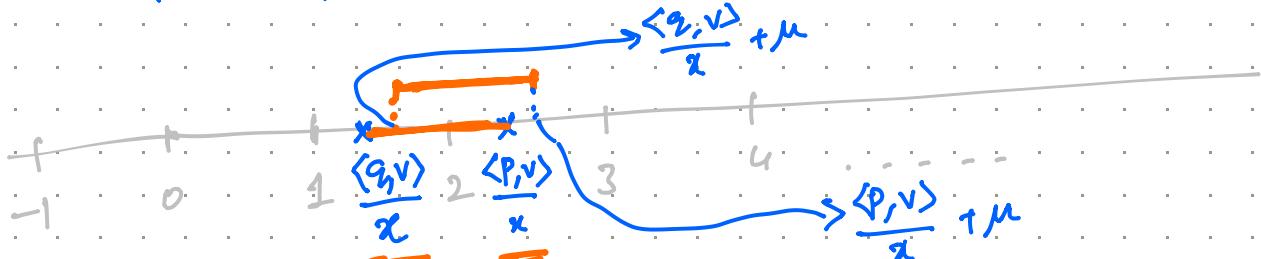
Claim : If $\beta = \langle p, v \rangle - \langle q, v \rangle > \alpha$, then $h(p) \neq h(q)$.

Pf :

$$\begin{aligned} h(p) &= \left\lfloor \frac{\langle q, v \rangle + \beta + t}{\alpha} \right\rfloor = \left\lfloor \frac{\beta}{\alpha} + \frac{\langle q, v \rangle + t}{\alpha} \right\rfloor \\ &> \left\lfloor \frac{\langle q, v \rangle + t}{\alpha} \right\rfloor = h(q) \end{aligned}$$



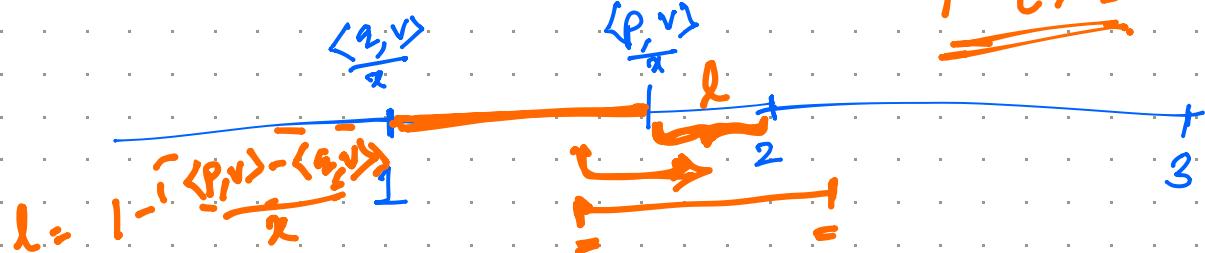
Suppose $\beta = \langle p, v \rangle - \langle q, v \rangle \leq \chi$.



$$h(p) = \left[\frac{\langle p, v \rangle}{\chi} + \mu \right] \quad \text{where } \mu \sim \text{Unif}[0, 1]$$

$$P[h(p) = h(q)] = 1 - \frac{\beta}{\chi} \quad (\beta = \langle p, v \rangle - \langle q, v \rangle)$$

$$\mu \in [0, \chi]$$



$$P[h(p) = h(q)] = \int_0^{\beta} P[|\langle p, v \rangle - \langle q, v \rangle| = \beta] \left(1 - \frac{p}{\eta}\right) dp$$

~~$\cdots \cdots \cdots$~~

$$\langle p, v \rangle - \langle q, v \rangle = \underline{\langle p-q, v \rangle} \sim N(0, \|p-q\|^2)$$

$\tilde{\eta}$

Say $\eta = \|p-q\|$.

$$\begin{aligned} P[|\langle p, v \rangle - \langle q, v \rangle| = \beta] &= P[\langle p-q, v \rangle = \beta \text{ or } \langle q-p, v \rangle = \beta] \\ &= 2 P[\langle p-q, v \rangle = \beta] \end{aligned}$$

~~$\cdots \cdots \cdots$~~

$P \left[\langle p-q, v \rangle = \beta \right] \xrightarrow[x \sim \cdot]{x} \text{the pdf of } \mathcal{N}(0, n^2)$
 at β

$$2 \left[\text{pdf of } \mathcal{N}(0, n^2) \right]^\infty = \frac{2}{\sqrt{2\pi} \eta} \exp \left(-\frac{\beta^2}{2n^2} \right)$$

$$\begin{aligned}
 P[h(p) = h(q)] &= \int_0^x \frac{2}{\sqrt{2\pi} \eta} \exp \left(-\frac{\beta^2}{2n^2} \right) \left(1 - \frac{\beta}{\eta} \right) d\beta \\
 &\quad \text{R.P.-q ||.} \quad p(n, x)
 \end{aligned}$$

$$\tilde{O} \left(n \frac{\ln(\frac{1}{\delta})}{\ln(\frac{1}{1-\beta})} \right) \quad \text{where } \alpha = P[h(p) = h(q)] \text{ given} \\ \|p-q\| \leq \gamma$$

$$\beta \geq P[h(p) = h(q)] \text{ given} \\ \|p-q\| \geq r(1+\epsilon)$$

choose α s.t

$$P \left[\ln \left(\frac{1}{P(n, \alpha)} \right) \right] \text{ is minimized}$$

$$\ln \left(\frac{1}{P(n(1+\epsilon), \alpha)} \right)$$

Datar, Immorlica, Indyk, Mirrokni (2004)

Using computations, they show that it is possible to
choose param α s.t

$$\frac{\ln \left(\frac{1}{p(\eta, x)} \right)}{\ln \left(\frac{1}{p(\eta(1+\epsilon), x)} \right)} \leq \frac{1}{1+\epsilon} \text{ for all } \eta > 0.$$

→ APQ !

$$\text{Query time} \rightarrow \tilde{O}(dn^{\frac{1}{1+\epsilon}})$$

$$\text{Preprocessing space} \rightarrow \tilde{O}\left(d(n^{1+\frac{1}{1+\epsilon}})\right).$$