CS-GY 9223 I: Lecture 4
Near neighbor search + locality sensitive
hashing

NYU Tandon School of Engineering, Prof. Christopher Musco
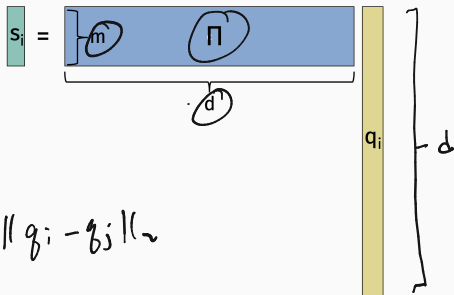
- Problem set 1.
- Reading group.

## Lemma (Johnson-Lindenstrauss, 1984)

*For any set of n data points $\mathbf{q}_1, \ldots, \mathbf{q}_n \in \mathbb{R}^d$ there exists a* <u>*linear map*</u> *$\Pi : \mathbb{R}^d \to \mathbb{R}^m$ where $m = O\left(\frac{\log n}{\epsilon^2}\right)$ such that* <u>*for all i, j,*</u>

$$(1 - \epsilon)\|\mathbf{q}_i - \mathbf{q}_j\|_2 \leq \|\mathbf{\Pi}\mathbf{q}_i - \mathbf{\Pi}\mathbf{q}_j\|_2 \leq (1 + \epsilon)\|\mathbf{q}_i - \mathbf{q}_j\|_2.$$



$$m = O\left(\frac{\log n}{\epsilon^2}\right)$$

$$\|s_i - s_j\|_2 \approx \|q_i - q_j\|_2$$

3

$\Pi \in \mathbb{R}^{k \times d}$ be chosen so that each entry equals $\frac{1}{\sqrt{m}}\mathcal{N}(0,1)$.

... or each entry equals $\frac{1}{\sqrt{m}} \pm 1$ with equal probability.



```
>> Pi = randn(m,d);
>> s = (1/sqrt(m))*Pi*q;
```



```
>> Pi = 2*randi(2,m,d)-3;
>> s = (1/sqrt(m))*Pi*q;
```

Often called "random projections". Why?

4

$q_1, \ldots, q_n \in \mathbb{R}^d$   $S_i := \Pi q_i$ for all $i$.

$\Pi : \mathbb{R}^d \to \mathbb{R}^m$

$\to S_1, \ldots, S_n \in \mathbb{R}^m$

$S = m$ dimensional subspace

When is a map $P : \mathbb{R}^d \to \mathbb{R}^m$ a projection operator? $\to$ maps each $q_i$ to closest point in $S$.

$Var[X] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$
          $0$

$m \left[ \underset{\overset{}{P}}{\overline{\underline{\hspace{2cm} d \hspace{2cm}}}} \right]$    $P P^T = $ Identity

Is $\Pi$ a projection operator?    variance

No.    $\mathbb{E}[\Pi \Pi^T] = d I$

$\mathbb{E}[[\Pi \Pi^T]_{ij}] = \mathbb{E}\left[ \sum_{\ell=1}^{d} \Pi_{i\ell} \Pi_{\ell j}^T \right] = \sum_{\ell=1}^{d} \mathbb{E}[\Pi_{i\ell} \Pi_{\ell j}^T]$ .   $d$

$\overset{i \neq j}{\nearrow} = 0$

$j \neq i$   $\mathbb{E}[[\Pi \Pi^T]_{ij}] = 0$.   $i = j$   $\sum_{\ell=1}^{d} \underbrace{\mathbb{E}[\Pi_{i\ell} \cdot \Pi_{i\ell}^T]}_{1 = var[N(0,1)]}$

$\overset{normal}{\nearrow}$

5

$k$-**means objective**: Find clusters $C_1, \ldots, C_k \subseteq \underline{\underline{\{1, \ldots, n\}}}$ to minimize:

$X_1, \ldots, X_n \in \mathbb{R}^d$

$$Cost(C_1, \ldots, C_k) = \sum_{j=1}^{k} \frac{1}{2|C_j|} \sum_{u,v \in C_j} \|X_u - X_v\|_2^2.$$

$C_1 \cup \ldots \cup C_k = \{1, \ldots, n\}$

$C_i \cap C_j = \emptyset \quad \text{for all} \quad i,j \in 1, \ldots, k$

**Approximation algorithm**: Find optimal clusters $\tilde{C}_1^*, \ldots, \tilde{C}_k^*$ for the $k$ dimension data set $\Pi X_1, \ldots, \Pi X_n$.

Want to prove:

$$\text{Cost}(\tilde{C}_1^*, \ldots, \tilde{C}_k^*) \leq (1+\epsilon) \min_{C_1, \ldots, C_n} \text{Cost}(C_1, \ldots, C_n)$$

$$k = O(\log n / \epsilon^2)$$

$\underline{\text{JL Lemma}}$ : For $\underline{\text{all}}$ $u, v,$ $(1-\epsilon) \| X_u - X_v \|_2^2 \leq \| \Pi X_u - \Pi X_v \|_2^2$
$$\leq (1+\epsilon) \| X_u - X_v \|_2^2$$
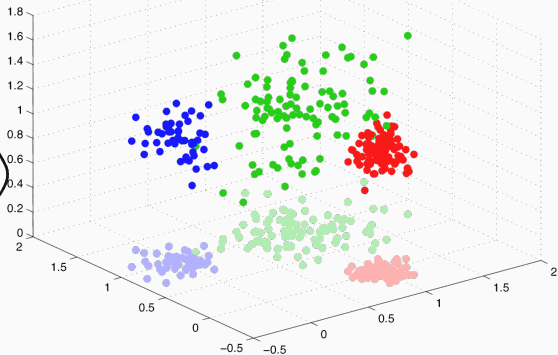
$$Cost(C_1, \ldots, C_k) = \sum_{j=1}^{k} \frac{1}{2|C_j|} \sum_{u,v \in C_j} \| X_u - X_v \|_2^2.$$

$$\widetilde{Cost}(C_1, \ldots, C_k) = \sum_{j=1}^{k} \frac{1}{2|C_j|} \sum_{u,v \in C_j} \| \Pi X_u - \Pi X_v \|_2^2.$$

$$\leq \sum_{j=1}^{k} \frac{1}{2|C_j|} \sum_{u,v \in C_j} (1+\epsilon) \| X_u - X_v \|_2^2 \quad \rightarrow \text{by JL}$$

For any $C_1, \ldots, C_k,$

$$= (1+\epsilon) \sum_{j=1}^{k} \frac{1}{2|C_j|} \sum_{u,v \in C_j} \| X_u - X_v \|_2^2 = (1+\epsilon) Cost(C_1, \ldots, C_k)$$

$$(1-\epsilon) Cost(C_1, \ldots, C_k) \leq \widetilde{Cost}(C_1, \ldots, C_k) \leq (1+\epsilon) Cost(C_1, \ldots, C_k)$$

We prove this

Proving left hand side is similar.

8

Let $C_1^*, \ldots, C_k^* = \arg\min Cost(C_1, \ldots, C_k)$ and
$\tilde{C}_1^*, \ldots, \tilde{C}_k^* = \arg\min \widetilde{Cost}(C_1, \ldots, C_k)$

**Want to prove:** $Cost(\tilde{C}_1^*, \ldots, \tilde{C}_k^*) \leq (1 + O(\epsilon)) Cost(C_1^*, \ldots, C_k^*)$

John Shin's comment:
$1/\epsilon^2$ dependence is really bad! $\epsilon = .01 \rightarrow$ $1/\epsilon^2 = 10,000$

- this is a weakness of JL. Not always good for highly accurate approximations

$$Cost(\tilde{C}_1^*, \ldots, \tilde{C}_u^a) \leq \frac{1}{1-\epsilon} \widetilde{Cost}(\tilde{C}_1^*, \ldots, \tilde{C}_u^*)$$

$$\leq \frac{1}{1-\epsilon} \widetilde{Cost}(C_1^*, \ldots, C_u^a)$$

- in some applications (including k-means)

$$\leq \frac{(1+\epsilon)}{(1-\epsilon)} Cost(C_1^*, \ldots, C_u^*)$$

It tends to do better than the theory predicts.

$\approx u(1+\epsilon)(1 + O(\epsilon))$
$\approx u(1 + O(\epsilon))$

$$\leq (1 + 4\epsilon) Cost(C_1^*, \ldots, C_u^*)$$

for any $\epsilon \leq 1/2$

If wanted $\leq (1+\epsilon)$ just set $k = O\left(\frac{\log u}{(\epsilon/4)^2}\right)$ to begin with.

9

high dimensional vector representation

| 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 |

$C$

| .45 | .68 | .10 | .92 |

sketched representation

**Goal:** Given input vectors $q$ and $y$, $C(q)$ and $C(y)$ should be similar if $q$ and $y$ are similar.

10

**Other Example:** Binary valued vectors.

> ### Definition (Jaccard Similarity)
>
> $$J(\mathbf{q}, \mathbf{y}) = \frac{|\mathbf{q} \cap \mathbf{y}|}{|\mathbf{q} \cup \mathbf{y}|} = \frac{\text{\# of non-zero entries in common}}{\text{total \# of non-zero entries}} \quad \frac{2}{4} = \frac{1}{2}$$
>
> $0 \leq J(\mathbf{q}, \mathbf{y}) \leq 1.$

for example
above
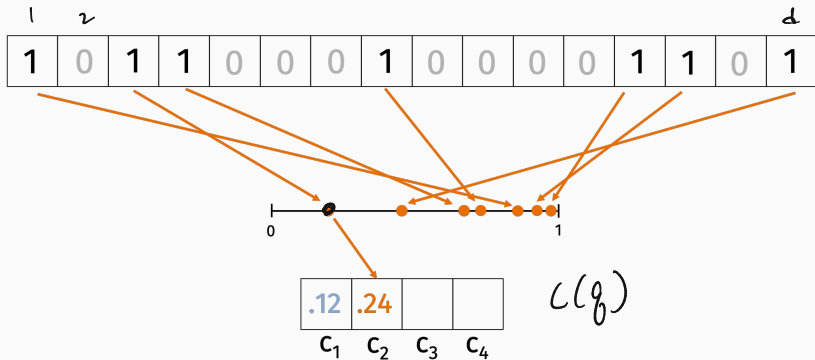
- Choose $k$ random hash functions
  $h_1, \ldots, h_k : \{1, \ldots, d\} \to [0, 1]$.
- For $i \in 1, \ldots, k$, let $c_i = \min_{j, \mathbf{q}_j = 1} h_i(j)$.
- $C(\mathbf{q}) = [c_1, \ldots, c_k]$.



$$\min\left(h_2(1), \ h_4(3), h_2(4), h_2(8) \cdots \right) = .24$$

**Example 1:** Binary valued vectors.

If $J(\mathbf{q}, \mathbf{y}) = v$ then the expected number of common entries between $C(\mathbf{q})$ and $C(\mathbf{y})$ is $v$.

> Actually, we proved:
> $$\Pr[c_i(\mathbf{q}) = c_i(\mathbf{y})] = v$$
> for all $i$.

| $C(\mathbf{q})$ | .12 | .24 | .76 | .35 |
|---|---|---|---|---|

| $C(\mathbf{y})$ | .12 | .98 | .76 | .11 |
|---|---|---|---|---|

Using a Chernoff bound, we proved that if $C$ maps to dimension $O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ we can approximate the Jaccard similarity between any two binary vectors to accuracy $\epsilon$ with probability $1 - \delta$.

$$J(Y, q) \geqslant .9$$

**Common goal:** Find all fingerprints in database $q_1, \ldots, q_n \in \mathbb{R}^d$ that are close to some input finger print $y \in \mathbb{R}^d$.
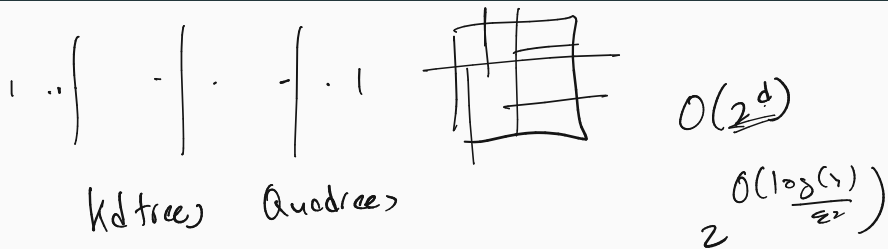
- Audio + video search.   $n = 10 \text{ million}$
- Finding duplicate or near duplicate documents.
- Seismic applications (here they want all pairs of close fingerprints).   $n = \text{millions}$

sketch dimension

<span style="color:orange">Does similarity sketching help in these applications?</span>

$m \leq d$

$O(n d)$ time

$O(n d)$ space      improvement      $O(n m)$

$O(n m)$

14

$O(2^d)$

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

kd trees    Quadrees

**New goal:** Sublinear $o(n)$ time to find near neighbors.

Can we also improve the $n$ dependence?
More important in many applications.

   — possible in low-dimensions using
      kd-trees, quad trees, etc

Let $h : \mathbb{R}^d \to \{1, \ldots, m\}$ be a random hash function.

We call $h$ <u>locality sensitive</u> if $\Pr[h(\mathbf{q}) == h(\mathbf{y})]$ is:

- Higher when $\mathbf{q}$ and $\mathbf{y}$ are more similar.
- Lower when $\mathbf{q}$ and $\mathbf{y}$ differ substantially.



16

LSH for Jaccard similarity:

- Let $c : \{0, 1\}^d \to [0, 1]$ be a single instantiation of MinHash.
- Let $g : [0, 1] \to \{1, \ldots, m\}$ be a fully random hash function.
- Let $h(\mathbf{x}) = g(c(\mathbf{x}))$.

LSH for Jaccard similarity:

- Let $c : \{0,1\}^d \to [0,1]$ be a single instantiation of MinHash.

*uniformly*

- Let $g : [0,1] \to \{1, \ldots, m\}$ be a fully random hash function.
- Let $h(\mathbf{x}) = g(c(\mathbf{x}))$.

$$g(c(q)) = g(c(y))$$
$$\text{when} \quad c(q) = c(y)$$

If $J(\mathbf{q}, \mathbf{y}) = v,$

Case 1: $c(q) = c(y)$
(happens with prob. $v$)

$$\Pr[h(q) = h(y)] = 1.$$

$$\Pr[h(\mathbf{q}) == h(\mathbf{y})] = v \cdot 1 + (1-v) \cdot \frac{1}{m}$$

$$= v + \frac{(1-v)}{m}$$

→ usually very small (think of $m = O(u)$)

Case 2: $c(q) \neq c(y)$
(happens with prob. $1-v$)

$$\Pr[h(q) = h(y)] = 1/m$$

→ why?

$h(q) = g(a)$
$h(y) = g(\beta)$  where $\alpha \neq \beta$.
Since $g$ is uniformly random, $\Pr[g(a) = g(\beta)]$ $= 1/m$

17

Basic approach for near neighbor search in a database.

**Pre-processing:**

- Select random LSH function $h : \{0,1\}^d \rightarrow 1, \ldots, m.$
- Create table $T$ with $m$ slots.    $m = O(n)$  (we won't discuss choice of $m$ rigorously)
- For $i = 1, \ldots, n$, insert $\mathbf{q}_i$ into $T(h(\mathbf{q}_i))$.

for Jaccard similarity

$$h = g(c(q_i))$$

uniform random hash

MinHash

Basic approach for near neighbor search in a database.

### Pre-processing:

- Select random LSH function $h : \{0,1\}^d \rightarrow 1, \ldots, m$.
- Create table $T$ with $m$ slots.
- For $i = 1, \ldots, n$, insert $\mathbf{q}_i$ into $T(h(\mathbf{q}_i))$.

$J(q, J)$ is large
$\downarrow$
$Pr[h(J) == h(q)]$
large.

### Query:

- Want to find near neighbors of input $\mathbf{y} \in \{0,1\}^d$.
- Linear scan through all vectors in $T(h(\mathbf{y}))$.

to find best matches

only $1/m$ buckets $\Rightarrow$ hopefully $\ll n$ elements to go through

18

T

Two main considerations:

- **False Negative** Rate: What's the probability we do not find a vector that is close to y?
- **False Positive** Rate: What's the probability we need to scan over vectors that are not close to y?

False Negative Rate
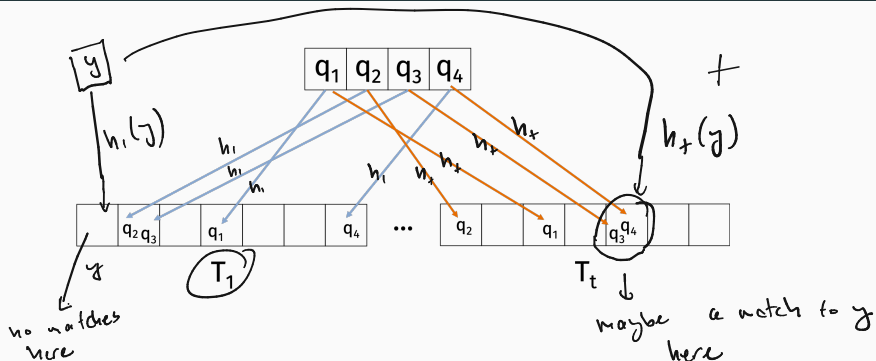
Suppose the nearest database point **q** has $J(\mathbf{y}, \mathbf{q}) = .4$.

What's the probability we do not find q?

$$Pr[\text{find } q] \approx .4$$

$$Pr[\text{don't find } q] = 1 - .4 = .6$$

$$\nearrow 60\%$$

thats terrible...

21

Pre-processing:

- Select $t$ independent LSH's $h_1, \ldots, h_t : \{0,1\}^d \to 1, \ldots, m$.
- Create tables $T_1, \ldots, T_t$, each with $m$ slots.
- For $i = 1, \ldots, n, j = 1, \ldots, t$, insert $\mathbf{q}_i$ into $T_j(h_j(\mathbf{q}_i))$.

22

Query:

- Want to find near neighbors of input $y \in \{0,1\}^d$.
- Linear scan through all vectors in
  $T_1(h_1(y)), T_2(h_2(y)), \ldots, T_t(h_t(y))$.

$t = 10$
$99\%$

$\underline{t \text{ repetitions}}$

**Query:**

- Want to find near neighbors of input $\mathbf{y} \in \{0,1\}^d$.
- Linear scan through all vectors in
  $T_1(h_1(\mathbf{y})), T_2(h_2(\mathbf{y})), \ldots, T_t(h_t(\mathbf{y}))$.

Suppose the nearest database point $\mathbf{q}$ has $J(\mathbf{y}, \mathbf{q}) = .4$.

What's the probability we find q?

$\Pr\{\text{don't find } q?\} = \Pr\{\text{don't find } q \text{ in Table 1} \underline{\text{and}}$
$\qquad\qquad\qquad\qquad\qquad \text{don't find } q \text{ in Table 2} \underline{\text{and}}$
$\Pr\{\text{find } q\} = 1 - (1 - .4)^t$
$\qquad\qquad\qquad\qquad\qquad \vdots$
$\qquad\qquad\qquad\qquad\qquad \text{don't find } q \text{ in Table T}\}$
$\boxed{\text{If } t = 10: \quad 1 - (.4)^{10} \geqslant .99} \qquad = \underline{(1 - .4)^t}$

23

t. 10
89%

Suppose there is some other database point $\mathbf{q}_j$ with
$\boxed{J(\mathbf{y}, \mathbf{q}_j) = .2?}$ What is the probability we will consider that point
in our original scheme? $\mathbb{P}\{h(y) = h(q_j)\} \approx .2$ when $J(q_j, y)$

$= .2$

why approx?

In the new scheme?

$$1 - (1 - .2)^{10} = 89\%$$

remember

$Pr\{h(y) = h(\mathbf{q})\}$

$= .2 + \underbrace{\dfrac{(1 - .2)}{m}}_{\approx 0}$

24

$h$ is uniformly random

$g([.1 \quad .6 \quad .2 \quad .93]) \to \{1, \ldots, m\}$

$\downarrow$

**Change our locality sensitive hash function**.

$\Pr(h(1) = h(2)) = 1/m$

when $1 \neq 2$.

$h(q): g(\underline{(.1)}, .2, .7)$

Tunable LSH for Jaccard similarity:

$h(y) = g(.1, \underline{.9}, .7)$

- Choose parameter $s \in \mathbb{Z}^+. = \{1, 2, \ldots\}$ positive integers

- Let $c_1, \ldots, c_s : \{0,1\}^d \to [0,1]$ be random MinHashs.

- Let $g : [0,1]^s \to \{1, \ldots, m\}$ be a fully random hash function.

- Let $h(\mathbf{x}) = g(c_1(\mathbf{x}), \ldots, c_s(\mathbf{x}))$.

If $J(\mathbf{q}, \mathbf{y}) = v$,

$\underline{\text{Case 1}}: c_1(q) = c_1(y), \ldots, c_s(q) = c_s(y)$  will depend on $s$

Happens with prob. $v^s$

$\Pr[h(q) = h(y)] = 1$.

$\Pr[h(\mathbf{q}) == h(\mathbf{y})] = v^s \cdot 1 + (1-v^s) \cdot 1/m$

$\underline{\text{Case 2}}: c_i(q) \neq c_i(y)$ for some $i$:

$\approx v^s \qquad \underbrace{\qquad}_{\approx 0}$

Happens w/ prob $(1-v^s)$  $\Pr[h(q) = h(y)] = 1/m$

25

$t = 3$

$78\%$

Parameter: $S = 1$.

Chance we find $\mathbf{q}_i$ with $J(\mathbf{y}, \mathbf{q}_i) = .8$:     $.8$

$$1 - (1 - .8)^t$$

If $t = 3 \rightarrow 1 - (1 - .8)^3 \geqslant 99\%$

Chance we need to scan $\mathbf{q}_j$ with $J(\mathbf{y}, \mathbf{q}_j) = .4$:

$$1 - (1 - .4)^3 \approx 78\%$$

Parameter: $S = 2$.

Chance we find $\mathbf{q}_i$ with $J(\mathbf{y}, \mathbf{q}_i) = .8$:

$$1 - (1 - .8^2)^t$$

$$\text{If} \quad \underline{t = 5} \longrightarrow 1 - (1 - .8^2)^5 \geqslant 99\%$$

Chance we need to scan $\mathbf{q}_j$ with $J(\mathbf{y}, \mathbf{q}_j) = .4$:

$$1 - (1 - .4^2)^5 \approx 58\%$$

$t = 12$

$12\%$

<p align="center" style="color:orange">Parameter: $S = 5$.</p>

Chance we find $\mathbf{q}_i$ with $J(\mathbf{y}, \mathbf{q}_i) = .8$:

$$1 - \left(1 - .8^5\right)^t$$

If $\quad \underline{t = 12} \rightarrow 1 - (1 - .8^5)^{12} \geqslant 99\%$

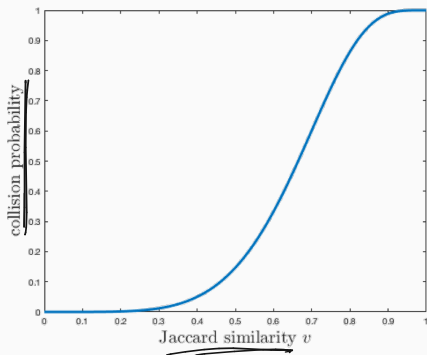Chance we need to scan $\mathbf{q}_j$ with $J(\mathbf{y}, \mathbf{q}_j) = .4$:

$$1 - (1 - .4^5)^{12} \approx \underline{\underline{12\%}}$$

reducing false positive rate a lot!

Probability we see **q** when querying **y** if $J(\mathbf{q}, \mathbf{y}) = v$:

$$\left(1 - (1 - v^s)^t\right)$$



$$s = 5, t = 5$$

Probability we see **q** when querying **y** if $J(\mathbf{q}, \mathbf{y}) = v$:

$$1 - (1 - v^s)^t$$



$s = 5, t = 40$

Increase T → Curve moves left

Probability we see **q** when querying **y** if $J(\mathbf{q}, \mathbf{y}) = v$:
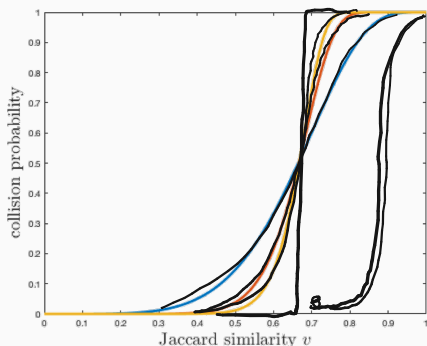
$$\approx 1 - (1 - v^s)^t$$



$$s = 40, t = 5$$

Increase $s$ → curve moves right

32

Probability we see **q** when querying **y** if $J(\mathbf{q}, \mathbf{y}) = v$:

$$1 - (1 - v^s)^t$$



Increasing both *s* and *t* gives a steeper curve. (in some place)

Better for search, but worse space complexity. 33

Use Case 1: Fixed threshold.

- Shazam wants to find match to audio clip $\mathbf{y}$ in a database of 10 million clips.
- There are 10 true matches with $J(\mathbf{y}, \mathbf{q}) > .9$.
- There are 10,000 near matches with $J(\mathbf{y}, \mathbf{q}) \in [.7, .9]$.

With $s = 25$ and $t = 40$.

- Hit probability for $J(\mathbf{y}, \mathbf{q}) > .9$ is $\gtrsim 1 - (1 - .9^{25})^{40} = .95$
- Hit probability for $J(\mathbf{y}, \mathbf{q}) \in [.7, .9]$ is $\lesssim 1 - (1 - .9^{25})^{40} = .95$
- Hit probability for $J(\mathbf{y}, \mathbf{q}) < .7$ is $\lesssim 1 - (1 - .7^{25})^{40} = .005$

$\mathbb{E}\left[\text{Total number of items scanned:}\right] \rightarrow$ using linearity of expectation

$$.95 \cdot 10 + .95 \cdot 10,000 + .005 \cdot 9,989,990 \approx 60,000 \ll 10,000,000.$$

34

traded    faster runtime
     for   more space

Space complexity: 40 hash tables $\approx 40 \cdot O(n)$.

Directly trade space for fast search.

Naive solution:   $40 \cdot n \cdot d$  space

Storing pointers to data vectors:   $40 \cdot n \cdot O(1) + nd$
                                              space

Concrete worst case result: *(no data assumptions like Shazam case study)*

### Theorem (Indyk, Motwani, 1998)

*If there exists some q with $\|\mathbf{q} - \mathbf{y}\|_0 \leq R$, return a vector $\tilde{\mathbf{q}}$ with $\underline{\|\tilde{\mathbf{q}} - \mathbf{y}\|_0} \leq C \cdot R$ in:*

If $C = 2$,

- *Time: $O\left(n^{1/C}\right)$.* $\rightarrow O(\sqrt{n})$
- *Space: $O\left(n^{1+1/C}\right)$.* $\rightarrow O\left(n^{3/2}\right)$

$\|\mathbf{q} - \mathbf{y}\|_0$ = "hamming distance" = number of elements that differ between q and y.

### Theorem (Indyk, Motwani, 1998)

*Let $q$ be the closest database vector to $y$. Return a vector $\tilde{q}$ with $\|\tilde{q} - y\|_0 \leq C \cdot \|q - y\|_0$ in:*

No $B$ parameter

- *Time: $\tilde{O}(n^{1/C})$.*
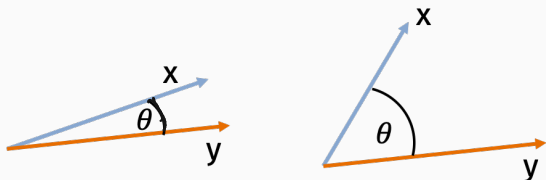- *Space: $\tilde{O}(n^{1+1/C})$.*

hides extra
log factors

Any ideas for how this is done?

Good locality sensitive hash functions exists for many other similarity measures.

Good locality sensitive hash functions exists for many other similarity measures.

**Cosine similarity** $\cos(\theta(x, y)) = \frac{\langle x, y \rangle}{\|x\|_2 \|y\|_2}$: $\quad \simeq \quad x^\top y$



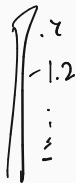$$-1 \le \cos(\theta(x, y)) \le 1.$$

$\theta = \pi \qquad\qquad\qquad\qquad \theta = 0$

Cosine similarity is natural "inverse" for Euclidean distance.

**Euclidean distance** $\underline{\|x - y\|_2^2}$:

- Suppose for simplicity that $\underline{\|x\|_2^2} = \underline{\|y\|_2^2} = 1$.

$$\|x-y\|_2^2 = \underbrace{\|x\|_2^2}_{=1} + \underbrace{\|y\|_2^2}_{=1} - 2x^T y = 2\underbrace{(1 - x^T y)}_{= 1 - \text{cosine similarity}}$$

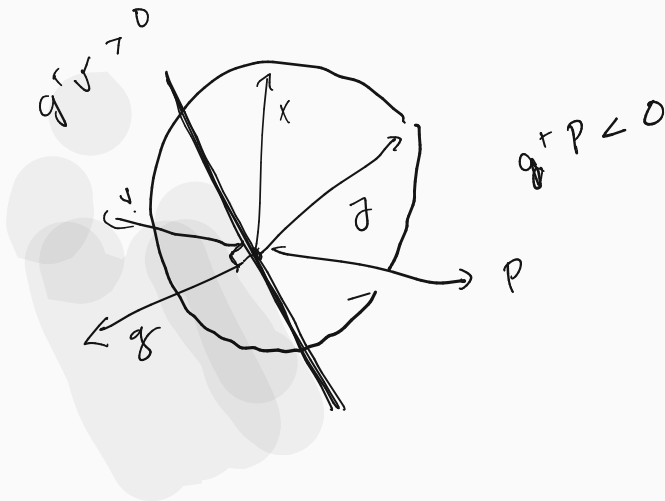$$\begin{bmatrix} .\gamma \\ -1.2 \\ \vdots \\ \_ \end{bmatrix}$$

Locality sensitive hash for cosine similarity:

- Let $\mathbf{g} \in \mathbb{R}^d$ be randomly chosen with each entry $\mathcal{N}(0,1)$.
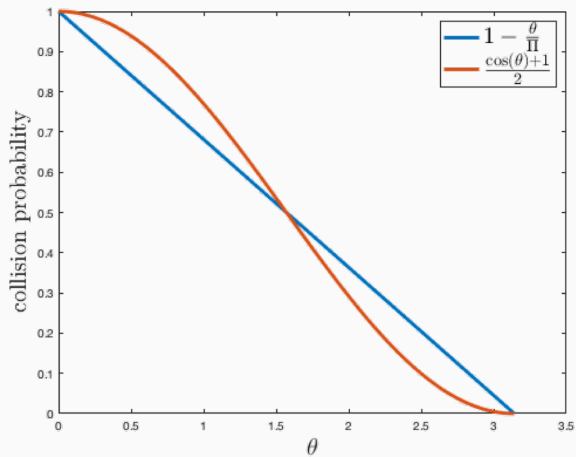- $h : \mathbb{R}^d \to \{-1, 1\}$ is definied $h(\mathbf{x}) = \text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle)$.  $\quad \text{sign}(\mathbf{g}^\top \mathbf{x})$

If $\cos(\theta(\mathbf{x}, \mathbf{y})) = v$, what is $\Pr[h(\mathbf{x}) == h(\mathbf{y})]$?

$$\cos(\theta(x,y)) = -1 \quad \Pr[h(x) = h(y)] = 0.$$
$$\cos(\theta(x,y)) = 1 \quad \Pr[h(x) = h(y)] = 1.$$

40

Inspired by Johnson-Lindenstrauss sketching

Locality sensitive hash for cosine similarity:

- Let $\mathbf{g} \in \mathbb{R}^d$ be randomly chosen with each entry $\mathcal{N}(0,1)$.
- $h : \mathbb{R}^d \rightarrow \{-1,1\}$ is definied $h(\mathbf{x}) = \text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle)$.

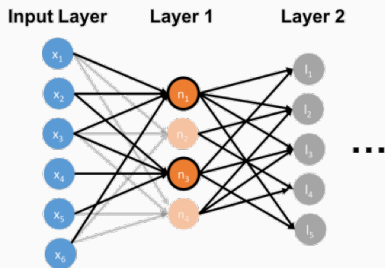If $\cos(\theta(\mathbf{x}, \mathbf{y})) = v$, what is $\Pr[h(\mathbf{x} == h(\mathbf{y})]$?

$$1 - \frac{\theta(x, z)}{\pi}$$

Work of Anshumali Shrivastava at Rice University and coauthors.



**Input Layer**    **Layer 1**    **Layer 2**
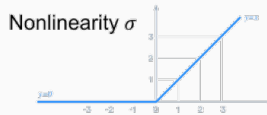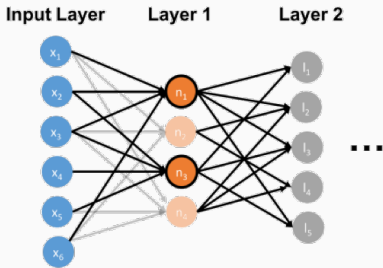
Nonlinearity $\sigma$

$$n_i = \sigma\left(\sum_{j=1}^{m} w(x_j, n_i) \cdot x_j\right) = \sigma(\langle w_i, x \rangle)$$

- Number of multiplications to evaluate $\mathcal{N}(\mathbf{x})$:
  $|\mathbf{x}| \cdot |\text{layer 1}| + |\text{layer 1}| \cdot |\text{layer 2}| + |\text{layer 2}| \cdot |\text{layer 3}| + \ldots$
- For an approximate solution, only consider neurons on each each with <u>high activation</u>.

45

Work of Anshumali Shrivastava at Rice University and coauthors.



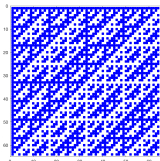$$n_i = \sigma\left(\sum_{j=1}^{m} w(x_j, n_i) \cdot x_j\right) = \sigma(\langle w_i, x\rangle)$$

- High activation $=$ large value of $\sigma(\langle \mathbf{w}_i, \mathbf{x}\rangle)$.
- Typically $\sigma(\langle \mathbf{w}_i, \mathbf{x}\rangle)$ increases as $\langle \mathbf{w}_i, \mathbf{x}\rangle$ increases.
- Use LSH/SimHash to quickly find all $\mathbf{w}_i$ for which $\langle \mathbf{w}_i, \mathbf{x}\rangle$ is large and only include these terms in the sum.
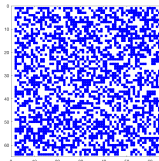
Why can't we just sample entries from vectors?

[Ailon, Chazelle, 2009 – The Fast Johnson-Lindenstrauss Transform]

[Ailon, Chazelle, 2009 – The Fast Johnson-Lindenstrauss Transform]

Deterministic Hadamard matrix.

Randomized Hadamard *HD*.

Fully random sign matrix.