# CS-GY 6763: LECTURE 5
# NEAR NEIGHBOR SEARCH IN HIGH DIMENSIONS + LOCALITY SENSITIVE HASHING

NYU Tandon School of Engineering, Prof. Rajesh Jayaram

## SIMILARITY SKETCHING

Given two length $d$ vectors $\mathbf{y}$ and $\mathbf{q}$, construct compact representations (sketches) $\tilde{\mathbf{y}}$ and $\tilde{\mathbf{q}}$ such that dist$(\mathbf{y}, \mathbf{q})$ can be estimated accurately from $\tilde{\mathbf{y}}$ and $\tilde{\mathbf{q}}$.
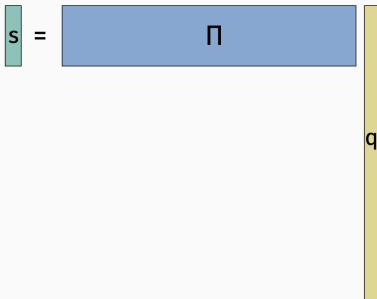
Each of $\tilde{\mathbf{y}}$ and $\tilde{\mathbf{q}}$ should require $k \ll d$ space.

## EUCLIDEAN DIMENSIONALITY REDUCTION

**Lemma (Johnson-Lindenstrauss, 1984)**

*For any two data points* $\mathbf{y}, \mathbf{q} \in \mathbb{R}^d$ *there exists a <u>linear map</u>* $\Pi : \mathbb{R}^d \to \mathbb{R}^k$ *where* $k = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ *such that with probability* $1 - \delta$,

$$(1 - \epsilon)\|\mathbf{q} - \mathbf{y}\|_2 \le \|\mathbf{\Pi q} - \mathbf{\Pi y}\|_2 \le (1 + \epsilon)\|\mathbf{q} - \mathbf{y}\|_2.$$

## EUCLIDEAN DIMENSIONALITY REDUCTION

**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> $\Pi : \mathbb{R}^d \to \mathbb{R}^k$ where $k = O\left(\frac{\log(n/\delta)}{\epsilon^2}\right)$ such that with probability $(1 - \delta)$, <u>for all $i, j$</u>,*

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

Extends to approximating all pairwise distances in a set of *n* vectors via a **union bound**.

## JACCARD SIMILARITY

Another distance measure (actually a similarity measure) between <u>binary</u> vectors in $\{0, 1\}^d$ :

**Definition (Jaccard Similarity)**

$$J(\mathbf{q}, \mathbf{y}) = \frac{|\mathbf{q} \cap \mathbf{y}|}{|\mathbf{q} \cup \mathbf{y}|} = \frac{\# \text{ of non-zero entries in common}}{\text{total } \# \text{ of non-zero entries}}$$

Natural similarity measure for binary vectors. $0 \leq J(\mathbf{q}, \mathbf{y}) \leq 1$.

## JACCARD SIMILARITY

Another distance measure (actually a similarity measure) between
<u>binary</u> vectors in $\{0,1\}^d$ :

**Definition (Jaccard Similarity)**

$$J(\mathbf{q}, \mathbf{y}) = \frac{|\mathbf{q} \cap \mathbf{y}|}{|\mathbf{q} \cup \mathbf{y}|} = \frac{\# \text{ of non-zero entries in common}}{\text{total } \# \text{ of non-zero entries}}$$

Natural similarity measure for binary vectors. $0 \leq J(\mathbf{q}, \mathbf{y}) \leq 1$.

Can be applied to any data which has a natural binary
representation (more than you might think).

$$\mathbf{y} = \begin{bmatrix} 1 & 0 & 1 & 1 & 0 & 0 \end{bmatrix}$$
$$\mathbf{q} = \begin{bmatrix} 1 & 1 & 0 & 1 & 0 & 0 \end{bmatrix}$$

# JACCARD SIMILARITY: SET DEFINITION

Jaccard similarity can also be expressed over sets.
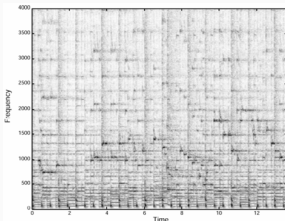
**Definition (Jaccard Similarity)**

Let $U$ be a universe of items, and $A, B \subset U$. Then

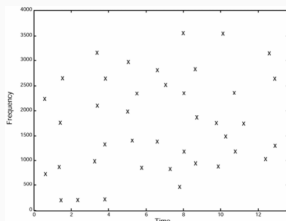$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

- Customer purchase similarities.

- Document similarity

- Similarity of sparse embeddings.

## SIMILARITY ESTIMATION

How does **Shazam** match a song clip against a library of 8 million songs (32 TB of data) in a fraction of a second?
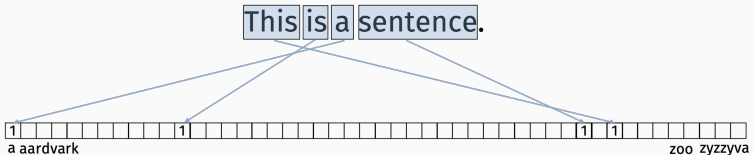


Spectrogram extracted from audio clip.



Processed spectrogram: used to construct audio "fingerprint" $\mathbf{q} \in \{0, 1\}^d$.

Each clip is represented by a high dimensional binary vector $\mathbf{q}$.

| 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**"Bag-of-words" model:**



How many words do a pair of documents have in common?

**"Bag-of-words" model:**



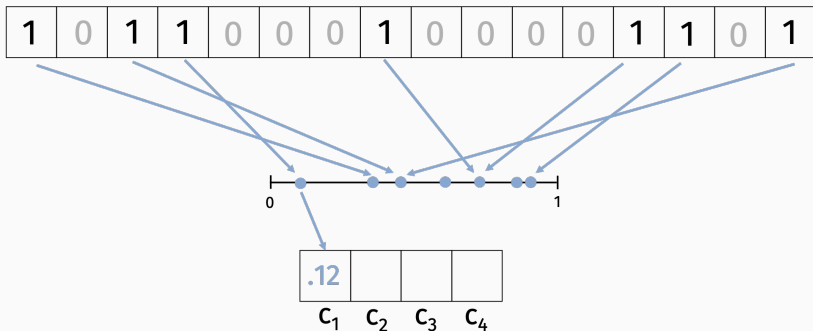How many bigrams do a pair of documents have in common?

## APPLICATIONS: DOCUMENT SIMILARITY

- Finding duplicate or new duplicate documents or webpages.
- Change detection for high-speed web caches.
- Finding near-duplicate emails or customer reviews which could indicate spam.

**MinHash (Broder, '97)**:

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

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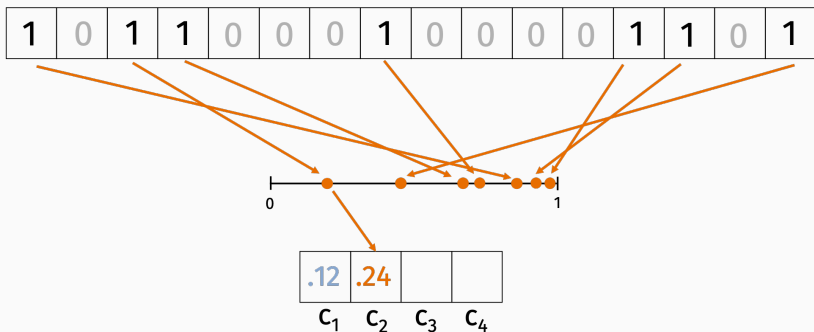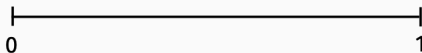- $C(\mathbf{q}) = [c_1, \ldots, c_k]$.

## MINHASH ANALYSIS

**Claim:** $\Pr[c_i(\mathbf{q}) = c_i(\mathbf{y})] = J(\mathbf{q}, \mathbf{y})$.

| $\mathbf{q}$ | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|
| $\mathbf{y}$ | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 |

$$\vdash\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\dashv$$
0                       1

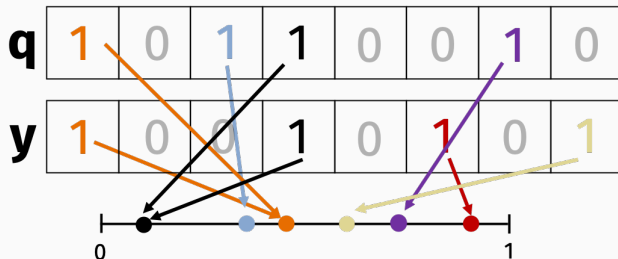**Proof:**

1. For $c_i(\mathbf{q}) = c_i(\mathbf{y})$, we need that
$\arg\min_{i \in \mathbf{q}} h(i) = \arg\min_{i \in \mathbf{y}} h(i)$.

**Claim:** $\Pr[c_i(\mathbf{q}) = c_i(\mathbf{y})] = J(\mathbf{q}, \mathbf{y})$.



2. Every non-zero index in $\mathbf{q} \cup \mathbf{y}$ is equally likely to produce the lowest hash value. $c_i(\mathbf{q}) = c_i(\mathbf{y})$ only if this index is 1 in <u>both</u> $\mathbf{q}$ and $\mathbf{y}$. There are $\mathbf{q} \cap \mathbf{y}$ such indices. So:

$$\Pr[c_i(\mathbf{q}) = c_i(\mathbf{y})] = \frac{\mathbf{q} \cap \mathbf{y}}{\mathbf{q} \cup \mathbf{y}} = J(\mathbf{q}, \mathbf{y})$$

Let $J = J(\mathbf{q}, \mathbf{y})$ denote the Jaccard similarity between $\mathbf{q}$ and $\mathbf{y}$.

**Return:** $\tilde{J} = \frac{1}{k} \sum_{i=1}^{k} \mathbb{1}[c_i(\mathbf{q}) = c_i(\mathbf{y})]$.

**Unbiased estimate for Jaccard similarity:**

$$\mathbb{E}\tilde{J} = \quad J(\mathbf{q}, y)$$

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

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

The more repetitions, the lower the variance.

$$var(\tilde{J}) = \frac{1}{k^2} \sum var(\cdot) = \frac{\cdot}{k} \cdot k \cdot J = \frac{J}{k}$$

16

$$\frac{\log(1/\delta)}{\epsilon^2}$$

Let $J = J(\mathbf{q}, \mathbf{y})$ denote the true Jaccard similarity.

**Estimator:** $\tilde{J} = \frac{1}{k} \sum_{i=1}^{k} \mathbb{1}[c_i(\mathbf{q}) = c_i(\mathbf{y})]$.

$$\text{Var}[\tilde{J}] = \frac{1}{k} \cdot J \leq \frac{1}{k}$$

Plug into Chebyshev inequality. How large does $k$ need to be so that with probability $> 1 - \delta$:

$$|J - \tilde{J}| \leq \epsilon?$$

$$P\left(|x - E(x)| \geq \epsilon\right) < \frac{\text{var}(x)}{\epsilon^2} < \frac{1}{k} \cdot \frac{1}{\epsilon^2}$$

**Chebyshev inequality:** As long as $k = O\left(\frac{1}{\epsilon^2 \delta}\right)$, then with prob. $1 - \delta$,

$$J(\mathbf{q}, \mathbf{y}) - \epsilon \leq \tilde{J}(C(\mathbf{q}), C(\mathbf{y})) \leq J(\mathbf{q}, \mathbf{y}) + \epsilon.$$
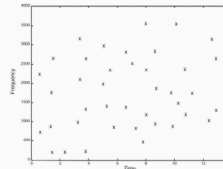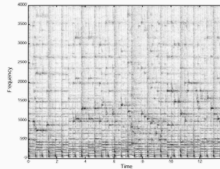
And $\tilde{J}$ only takes $O(k)$ time to compute! **Independent** of original fingerprint dimension $d$.

Can be improved to $\log(1/\delta)$ dependence. Can anyone tell me how?

input data

high dimensional vector representation

| 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

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

sketched representation

**break**

## NEAR NEIGHBOR SEARCH

**Common goal:** Find all vectors in database $\mathbf{q}_1, \ldots, \mathbf{q}_n \in \mathbb{R}^d$ that are close to some input query vector $\mathbf{y} \in \mathbb{R}^d$. I.e. find all of $\mathbf{y}$'s "nearest neighbors" in the database.

- The Shazam problem.
- Audio + video search.
- Finding duplicate or near duplicate documents.
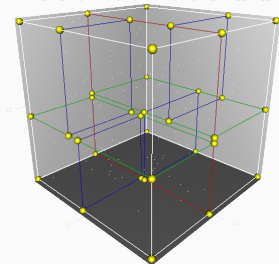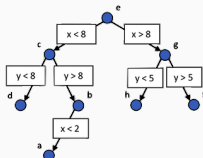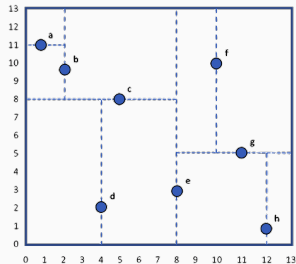- Detecting seismic events.

**How does similarity sketching help in these applications?**

- Improves runtime of "linear scan" from $O(nd)$ to $O(nk)$.
- Improves space complexity from $O(nd)$ to $O(nk)$. This can be super important – e.g. if it means the linear scan only accesses vectors in fast memory.

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

# BEYOND A LINEAR SCAN

This problem can already be solved for a small number of dimensions using space partitioning approaches (e.g. kd-tree).

## HIGH DIMENSIONAL NEAR NEIGHBOR SEARCH

Only been attacked much more recently:

- **Locality-sensitive hashing [Indyk, Motwani, 1998]**
- Spectral hashing [Weiss, Torralba, and Fergus, 2008]
- Vector quantization [Jégou, Douze, Schmid, 2009]
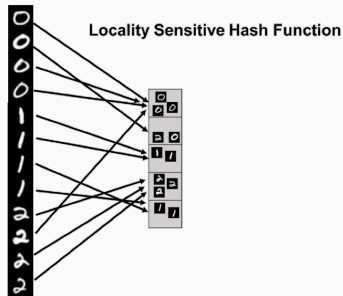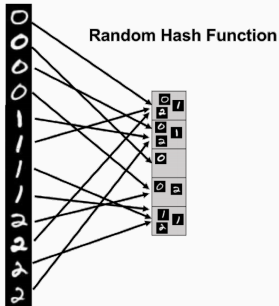    - This is most similar to the custom method e.g. Shazam uses.

**Key Insight:** Trade worse space-complexity for better time-complexity.

## LOCALITY SENSITIVE HASH FUNCTIONS

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

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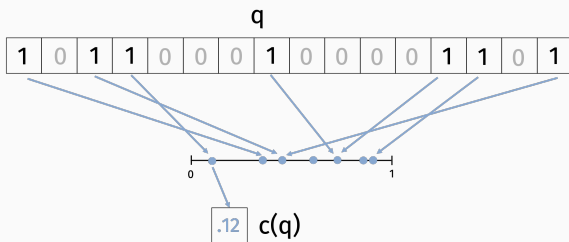$\Pr[h(\mathbf{q}) == h(\mathbf{y})]$ is:

- Higher when $\mathbf{q}$ and $\mathbf{y}$ are more similar, i.e. $s(\mathbf{q}, \mathbf{y})$ is higher.
- Lower when $\mathbf{q}$ and $\mathbf{y}$ are more dissimilar, i.e. $s(\mathbf{q}, \mathbf{y})$ is lower.

## LOCALITY SENSITIVE HASH FUNCTIONS

LSH for $s(\mathbf{q}, \mathbf{y})$ equal to Jaccard similarity:

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

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 uniform random hash function.
- Let $h(\mathbf{x}) = g(c(\mathbf{x}))$.

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

$$\Pr[h(\mathbf{q}) == h(\mathbf{y})] = J + (1-J)\frac{1}{m}$$
$$= J + O\left(\frac{1}{m}\right)$$

## NEAR NEIGHBOR SEARCH

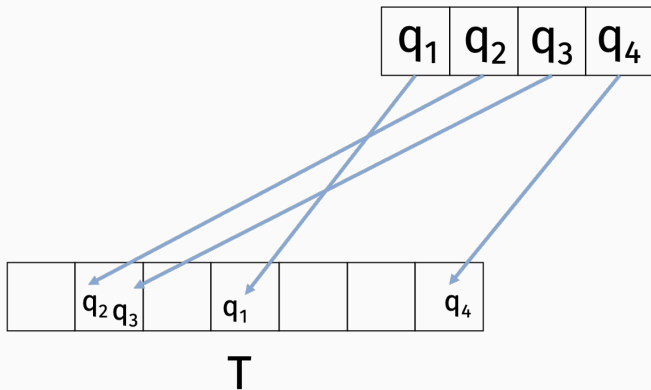Basic approach for near neighbor search in a database.

**Pre-processing:**

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

**Query:**

- Want to find near neighbors of input $\mathbf{y} \in \{0,1\}^d$.
- Linear scan through all vectors $\mathbf{q} \in T(h(\mathbf{y}))$ and return any that are close to $\mathbf{y}$. Time required is $O(d \cdot |T(h(\mathbf{y}))|)$.

---

[1]Enough to make the $O(1/m)$ term negligible.

T

## NEAR NEIGHBOR SEARCH

**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 that a vector in $T(h(\mathbf{y}))$ is not close to **y**?

A higher false negative rate means we miss near neighbors.

A higher false positive rate means increased runtime – we need to compute $J(\mathbf{q}, \mathbf{y})$ for every $\mathbf{q} \in T(h(\mathbf{y}))$ to check if it's actually close to **y**.

**Note:** The meaning of "close" and "not close" is application dependent. E.g. we might specify that we want to find anything with Jaccard similarity $> .4$, but not with Jaccard similarity $< .2$.
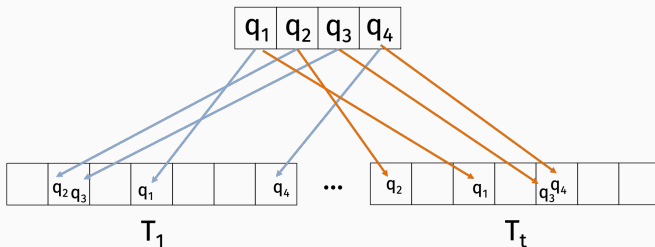
Suppose the nearest database point $\mathbf{q}$ has $J(\mathbf{y}, \mathbf{q}) = .4$.

**What's the probability we do not find q with one LSH?**

$$.6 + O\left(\frac{1}{\not{D}}\right)$$

**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 $T_j(h_j(\mathbf{q}_i))$.

**Query:**

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

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

**What's the probability we find q?**

$$\left( 1 - .6^t \right)$$

Suppose there is some other database point **z** with $J(\mathbf{y}, \mathbf{z}) = .2$.

What is the probability we will need to compute $J(\mathbf{z}, \mathbf{y})$ in our hashing scheme with one table? I.e. the probability that **y** hashes into at least one bucket containing **z**.
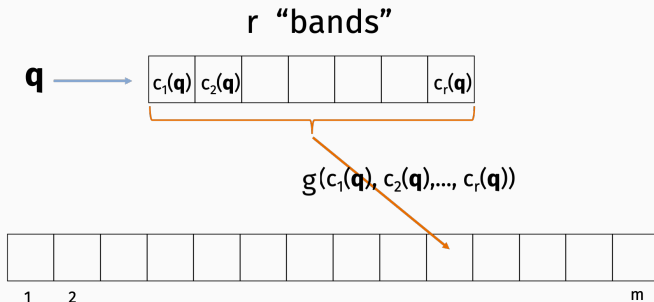
**In the new scheme with $t = 10$ tables?**

$$1 - .8^t \sim .9$$

**Change our locality sensitive hash function**.

Tunable LSH for Jaccard similarity:

- Choose parameter $r \in \mathbb{Z}^+$.

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

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

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



r "bands"

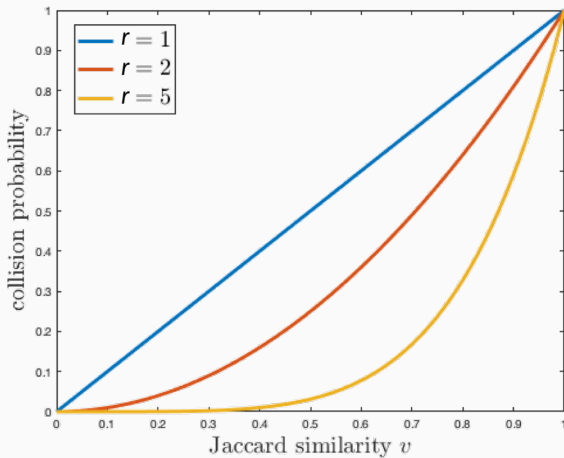$g(c_1(\mathbf{q}), c_2(\mathbf{q}),..., c_r(\mathbf{q}))$

Underline{Tunable} LSH for Jaccard similarity:

- Choose parameter $r \in \mathbb{Z}^+$.

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

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

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

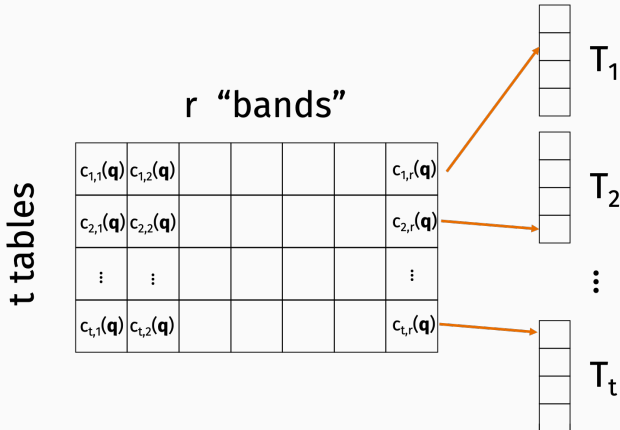If $J(\mathbf{q}, \mathbf{y}) = v$, then $\Pr[h(\mathbf{q}) == h(\mathbf{y})] = v_1^r + O\left(\frac{1}{m}\right)$

$$(1-\xi)^k \approx 1 - \xi k \quad \text{if } k \leq \frac{1}{\xi}$$

Full LSH cheme has two parameters to tune:



r "bands"

t tables

| $c_{1,1}(\mathbf{q})$ | $c_{1,2}(\mathbf{q})$ | | | | $c_{1,r}(\mathbf{q})$ |
| $c_{2,1}(\mathbf{q})$ | $c_{2,2}(\mathbf{q})$ | | | | $c_{2,r}(\mathbf{q})$ |
| $\vdots$ | $\vdots$ | | | | $\vdots$ |
| $c_{t,1}(\mathbf{q})$ | $c_{t,2}(\mathbf{q})$ | | | | $c_{t,r}(\mathbf{q})$ |

$T_1$

$T_2$

$\vdots$

$T_t$

Effect of **increasing number of tables** $t$ on:

| False Negatives | False Positives |
|---|---|
| decreasing | increases |

Effect of **increasing number of bands** $r$ on:

| False Negatives | False Positives |
|---|---|
| increasing | decrease |

Choose tables $t$ large enough so false negative rate to 1%.

**Parameter: r = 1**.

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

$$1 - .2^t > .99$$

$$t = 4$$
$$acc \rceil$$

Chance we need to check **z** with $J(\mathbf{y}, \mathbf{z}) = .4$:

$$1 - .6^t = .87$$

Choose tables $t$ large enough so false negative rate to 1%.

<div align="center">

**Parameter: r = 2**.

</div>

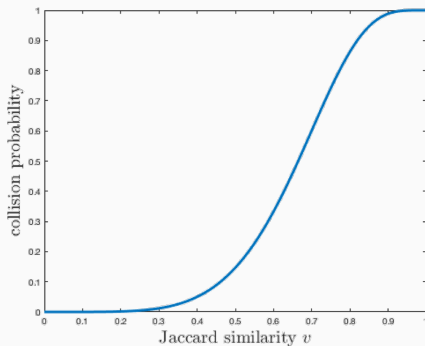Chance we find $\mathbf{q}$ with $J(\mathbf{y}, \mathbf{q}) = .8$:

$$1 - .36^t \geq .99$$

need $t = 5$

$$1 - .8^r$$

Chance we need to check $\mathbf{z}$ with $J(\mathbf{y}, \mathbf{z}) = .4$:

$$1 - .84^t \approx .58$$

Choose tables $t$ large enough so false negative rate to 1%.

**Parameter: r = 5**.

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

$$1 - .\langle 4\rangle^{\uparrow}$$

Chance we need to check **z** with $J(\mathbf{y}, \mathbf{z}) = .4$:

# S-CURVE TUNING

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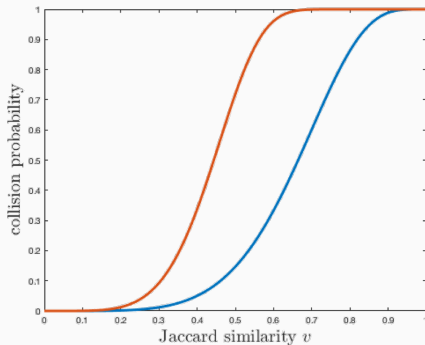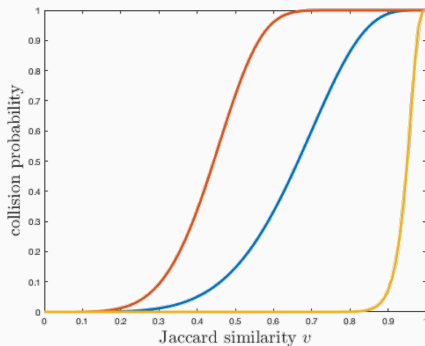
$$\approx 1 - (1 - v^r)^t$$



$$r = 5, t = 5$$

# *S*-CURVE TUNING

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

$$\approx 1 - (1 - v^r)^t$$



$r = 5, t = 40$

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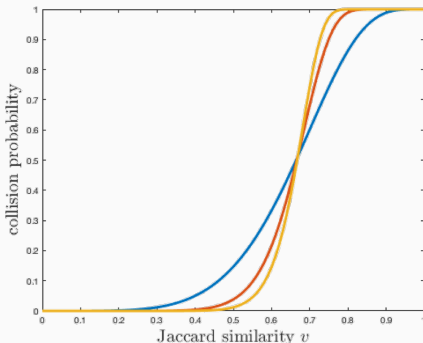
$$\approx 1 - (1 - v^r)^t$$



$$r = 40, t = 5$$

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

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



Increasing both $r$ and $t$ gives a steeper curve.

**Better for search, but worse space complexity**.

## FINDING THE BEST PARAMETERS

$t$ = number of tables, $r$ = number of "bands"

Suppose we have $y, q, z$ with

$$J(y, q) \approx p_1 = \Pr[h(y) = h(q)]$$

$$J(y, z) \approx p_2 = \Pr[h(y) = h(z)]$$

where $p_1 > p_2$. What is the probability we find $q$ when searching from $y$?

$$1 - (1 - p_1^r)^t$$

What is the probability we find $z$ when searching from $y$?

$$1 - (1 - p_2^r)^t$$

## FINDING THE BEST PARAMETERS

$t$ = number of tables, $r$ = number of "bands". Suppose we have $y, q, z$ with $p_1 > p_2$ and:

$$\Pr[\text{find } q] = 1 - (1 - p_1^r)^t, \quad \Pr[\text{find } z] = 1 - (1 - p_2^r)^t$$

False positive rate $= (1 - p_1^r)^t$, False negative rate $= 1 - (1 - p_2^r)^t$.

Suppose we want False positive $< .01$ and False negative $< .01$.

$t$ = number of tables, $r$ = number of "bands". Suppose we have $\boldsymbol{y}, \boldsymbol{q}, \boldsymbol{z}$ with $p_1 > p_2$ and:

$$\Pr[\text{find } \boldsymbol{q}] = 1 - (1 - p_1^r)^t, \quad \Pr[\text{find } \boldsymbol{z}] = 1 - (1 - p_2^r)^t$$

False positive rate = $(1 - p_1^r)^t$, False negative rate = $1 - (1 - p_2^r)^t$.

Suppose we want False positive $< .01$ and False negative $< .01$. Then we should set $t = \frac{\log(\frac{1}{100})}{\log(1 - p_1^r)} = \Theta(p_1^{-r})$.

So False negative rate $\approx 1 - (1 - p_2^r)^{p_1^{-r}} \approx \left(\frac{p_2}{p_1}\right)^r$. So

$$r = \Theta\left(\frac{1}{\log \frac{p_1}{p_2}}\right), \qquad t = \Theta(p_1^{-r})$$

**Lemma**

*Let $y, q, z$ be points with $p_1 = Pr[h(y) = h(q)]$ and $p_2 = \Pr[h(y) = h(z)]$, where $p_1 > p_2$. Then to achieve false positive and false negative rates $< .01$, it suffices to set*

$$r = \Theta\left(\frac{1}{\log \frac{p_1}{p_2}}\right), \qquad t = \Theta(p_1^{-r})$$

$t$ = *number of tables*, $r$ = *number of "bands".*

**Note:** as the gap $p_1 - p_2$ becomes smaller, need to use many more tables and bands!

## FIXED THRESHOLD

**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 <u>true matches</u> with $J(\mathbf{y}, \mathbf{q}) > .9$.
- There are 10,000 <u>near matches</u> with $J(\mathbf{y}, \mathbf{q}) \in [.7, .9]$.
- All other items have $J(\mathbf{y}, \mathbf{q}) < .7$.

With $r = 25$ 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 $\lesssim 1 - (1 - .7^{25})^{40} = .005$

**Upper bound on total number of items checked:**

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

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

**Directly trade space for fast search.**

## FIXED THRESHOLD R

**Near Neighbor Search Problem**

Concrete worst case result:

> **Theorem (Indyk, Motwani, 1998)**
>
> *If there exists some q with $\|\mathbf{q} - \mathbf{y}\|_0 \leq R$, return a vector $\tilde{\mathbf{q}}$ with $\|\tilde{\mathbf{q}} - \mathbf{y}\|_0 \leq C \cdot R$ in:*
>
> - *Time: $O\left(n^{1/C}\right)$.*
> - *Space: $O\left(n^{1+1/C}\right)$.*

$\|\mathbf{q} - \mathbf{y}\|_0 =$ "hamming distance" $=$ number of elements that differ between $\mathbf{q}$ and $\mathbf{y}$.
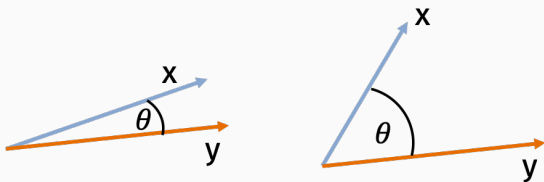
**Theorem (Indyk, Motwani, 1998)**

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

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

Good locality sensitive hash functions exists for other similarity measures.

**Cosine similarity** $\cos\left(\theta(\mathbf{x},\mathbf{y})\right) = \frac{\langle \mathbf{x},\mathbf{y}\rangle}{\|\mathbf{x}\|_2\|\mathbf{y}\|_2}$:



$$-1 \leq \cos\left(\theta(\mathbf{x},\mathbf{y})\right) \leq 1.$$

Cosine similarity is natural "inverse" for Euclidean distance.

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

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

$$|x - y|_2^2 = |x|_2^2 - 2\langle x, y \rangle + |y|_2^2$$
$$2(1 - \langle x, y \rangle)$$
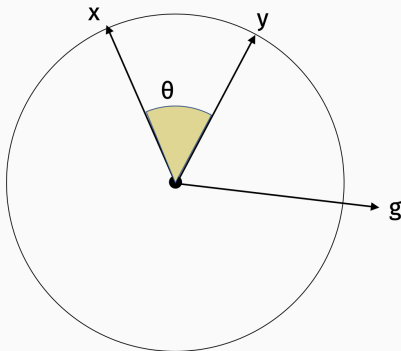$$2(1 - \cos(\theta))$$

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 defined $h(\mathbf{x}) = \text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle)$.

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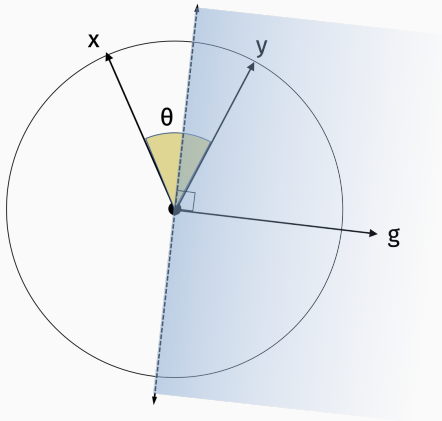
## SIMHASH ANALYSIS

**To prove:**

$\Pr[h(\mathbf{x}) == h(\mathbf{y})] = 1 - \frac{\theta}{\pi}$, where $h(\mathbf{x}) = \text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle)$.
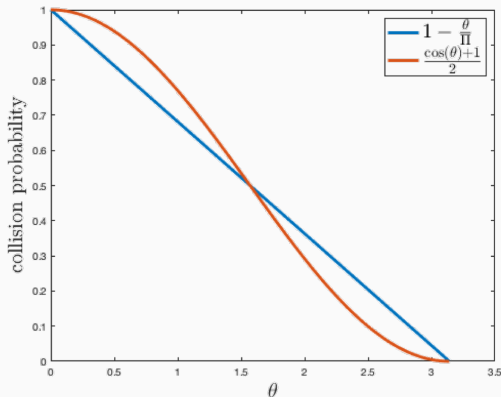
$\Pr[h(\mathbf{x}) == h(\mathbf{y})] \approx$ probability $\mathbf{x}$ and $\mathbf{y}$ are on the same side of hyperplane orthogonal to $\mathbf{g}$.

Each hyperplane is equally likely!

**Theorem:** If $\cos(\theta(\mathbf{x}, \mathbf{y})) = v$, then

$$\Pr[h(\mathbf{x}) == h(\mathbf{y})] = 1 - \frac{\theta(\mathbf{x}, \mathbf{y})}{\pi} = 1 - \frac{\cos^{-1}(v)}{\pi}$$

## SIMHASH

SimHash can be tuned, just like our MinHash based LSH function for Jaccard similarity:

- Let $\mathbf{g}_1, \ldots, \mathbf{g}_r \in \mathbb{R}^d$ be randomly chosen with each entry $\mathcal{N}(0, 1)$.
- Let $\theta = \theta(\mathbf{x}, \mathbf{y})$
- $h : \mathbb{R}^d \to \{1, -1\}$ is defined

$$h(\mathbf{x}) = [\text{sign}(\langle \mathbf{g}_1, \mathbf{x} \rangle), \ldots, \text{sign}(\langle \mathbf{g}_r, \mathbf{x} \rangle)]$$

$$\Pr[h(\mathbf{x}) == h(\mathbf{y})] = \left(1 - \frac{\theta}{\pi}\right)^r$$