

CS-GY 9223 D: Lecture 4

Near neighbor search + locality sensitive hashing

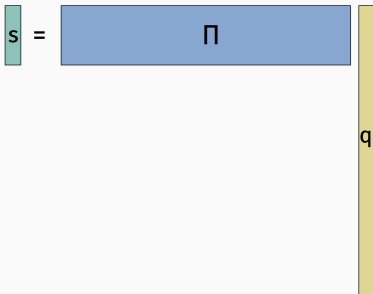
NYU Tandon School of Engineering, Prof. Christopher Musco

EUCLIDEAN DIMENSIONALITY REDUCTION

Lemma (Johnson-Lindenstrauss, 1984)

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

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



RANDOMIZED JL CONSTRUCTIONS

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

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

-2.1384	2.9888	-0.3538	0.0229	0.5201	-0.2938	-1.3320	-1.3617	-0.1952
-0.8396	0.8252	-0.8236	-0.2620	-0.0200	-0.8479	-2.3299	0.4550	-0.2176
1.3546	1.3798	-1.5771	-1.7582	-0.0348	-1.1281	-1.4491	-0.8487	-0.3831
-1.0722	-1.0582	0.5000	-0.2857	-0.7982	2.5268	0.3335	-0.3349	0.0338
0.9610	-0.4686	0.2820	-0.8314	1.0187	1.6555	0.3914	0.5528	0.0513
0.1240	-0.2725	0.0335	-0.9792	-0.1332	0.3075	0.4517	1.0391	0.8261
1.4367	1.0984	-1.3337	-1.1564	-0.7145	-1.2571	-0.1303	-1.1176	1.5278
-1.9609	-0.2779	1.1275	-0.5336	1.3514	-0.8655	0.1837	1.2607	0.4669
-0.1977	0.7815	0.3502	-2.0026	-0.2248	-0.1765	-0.4762	0.6601	-0.2897
-1.2878	-2.0518	-0.2991	0.9642	-0.5898	0.7914	0.8628	-0.0679	0.6252

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

1	1	-1	-1	-1	-1	-1	-1	1	-1	-1	1	-1	-1	1	1	-1
1	1	1	-1	1	-1	-1	-1	1	1	1	1	-1	1	-1	-1	-1
1	1	-1	-1	-1	1	-1	-1	1	1	-1	1	-1	1	-1	1	-1
-1	-1	-1	1	1	-1	-1	-1	-1	-1	-1	-1	-1	-1	1	1	1
1	-1	1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	1	1	-1
1	-1	-1	1	-1	1	1	-1	-1	1	-1	-1	-1	-1	-1	1	1
1	1	-1	1	1	-1	1	-1	1	-1	1	-1	1	1	1	-1	-1
-1	-1	-1	-1	-1	-1	1	-1	1	1	-1	-1	1	-1	-1	-1	1
-1	-1	1	1	1	1	-1	-1	1	-1	1	1	1	1	-1	1	-1
-1	1	-1	1	-1	1	1	-1	-1	1	-1	1	-1	-1	1	-1	1

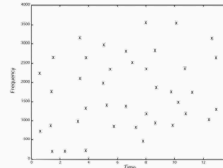
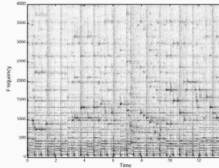
```
>> Pi = 2*randi(2,m,d)-3;  
>> s = (1/sqrt(m))*Pi*q;
```

Often called “random projections”.

SIMILARITY SKETCHING



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

Definition (Jaccard Similarity)

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

$$0 \leq J(\mathbf{q}, \mathbf{y}) \leq 1.$$

Similar result to JL: Given a MinHash sketch with $k = O\left(\frac{\log n}{\epsilon^2}\right)$ dimensions, we can estimate the Jaccard similarity between all pairs $\mathbf{x}_1, \dots, \mathbf{x}_n$ with high probability.

NEAR NEIGHBOR SEARCH

Common goal: Find all vectors in database $\mathbf{q}_1, \dots, \mathbf{q}_n \in \mathbb{R}^d$ 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.

- 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 $O(\underline{nd})$ to $O(\underline{nk})$
- 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.

Similarity

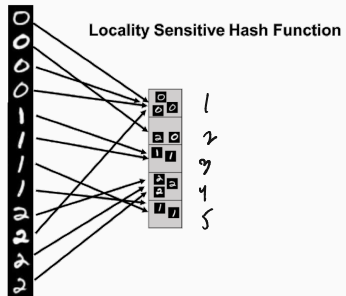
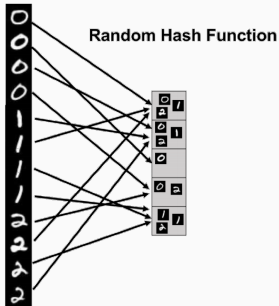
New goal: Sublinear $o(n)$ time to find near neighbors.

LOCALITY SENSITIVE HASH FUNCTIONS

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

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

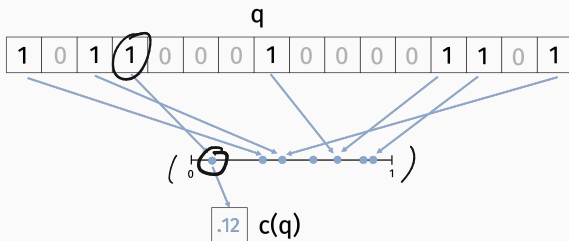
- Higher when \mathbf{q} and \mathbf{y} are more similar, i.e. $s(\mathbf{q}, \mathbf{y})$ is higher.
- Lower when \mathbf{q} and \mathbf{y} are more dissimilar, i.e. $s(\mathbf{q}, \mathbf{y})$ is lower.



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



LOCALITY SENSITIVE HASH FUNCTIONS

LSH for Jaccard similarity:

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

If $J(q, y) = v$,

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

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

2. $c(q) \neq c(y)$, then $\Pr[h(q) = h(y)] = 1/m$

NEAR NEIGHBOR SEARCH

Basic approach for near neighbor search in a database.

Pre-processing:

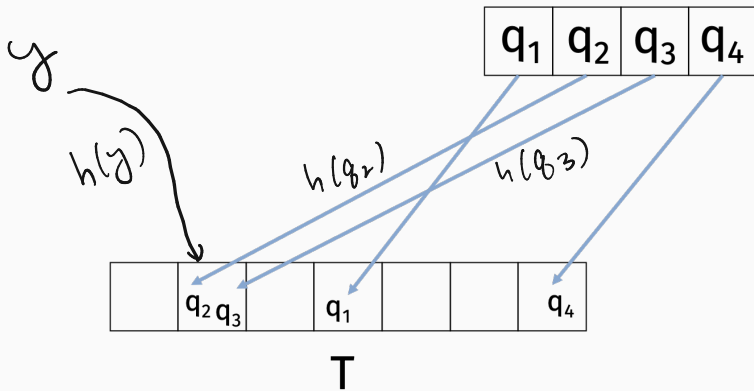
$\nearrow O(n)$

- Select random LSH function $h : \{0, 1\}^d \rightarrow 1, \dots, m$.
- Create table T with $m = O(n)$ slots.
- For $i = 1, \dots, n$, 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 $\mathbf{q} \in \underline{T(h(\mathbf{y}))}$ and return any that are close to \mathbf{y} . Time required is $O(d \cdot |T(h(\mathbf{y}))|)$.

NEAR NEIGHBOR SEARCH



Two main considerations:

- **False Negative Rate:** What's the probability we do not find a vector that is close to \mathbf{y} ?
- **False Positive Rate:** What's the probability that a vector in $T(h(\mathbf{y}))$ is not close to \mathbf{y} ?

A higher false negative rate means we miss near neighbors.

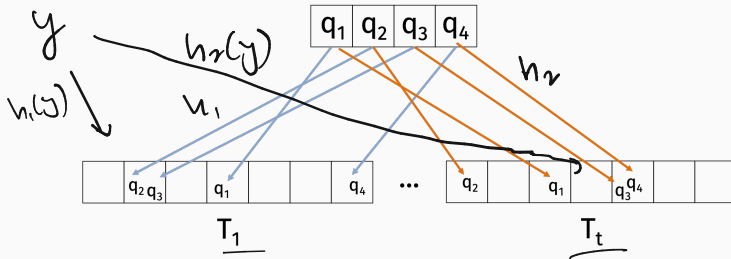
A higher false positive rate means increased runtime – we need to compute $J(\mathbf{q}, \mathbf{y})$ for every $\mathbf{q} \in T(h(\mathbf{y}))$ to check if it's actually close to \mathbf{y} .

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

What's the probability we do not find \mathbf{q} ?

$$1 - .4 = 60\%$$

REDUCING FALSE NEGATIVE RATE



Pre-processing:

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

Query:

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

$$O(d \cdot |T(h(\mathbf{y}))|)$$

$$O(d \cdot \sum_{j=1}^t |T_j(h(\mathbf{y}))|)$$

REDUCING FALSE NEGATIVE RATE

400%

$\neq 10$

99%

Query:

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

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

What's the probability we find \mathbf{q} ?

1 - (probability don't find \mathbf{q})

$$1 - .6^t$$

$$t = 10$$

$$1 - .6^{10} \approx 99\%$$

WHAT HAPPENS TO FALSE POSITIVES?

Suppose there is some other database point z with $J(y, z) = .2$.
What is the probability we will need to compute $J(z, y)$ in our hashing scheme with one table?

$$1 - .6^t$$

In the new scheme with $t = 10$ tables?

$$1 - .2 = .8$$

Prob. false positive: $\underbrace{1 - .8^t}_{\approx 89\%}$

$t=10$
89%

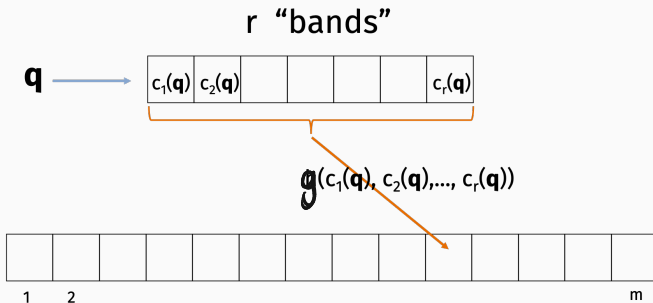
REDUCING FALSE POSITIVES

Change our locality sensitive hash function.

Tunable LSH for Jaccard similarity:

$$r = 1$$

- Choose parameter $\underline{r} \in \mathbb{Z}^+$.
- Let $c_1, \dots, c_r: \{0, 1\}^d \rightarrow [0, 1]$ be random MinHash.
- Let $g: [0, 1]^r \rightarrow \{1, \dots, m\}$ be a fully random hash function.
- Let $h(x) = g(\underline{c_1(x)}, \dots, \underline{c_r(x)})$.



REDUCING FALSE POSITIVES

Tunable LSH for Jaccard similarity:

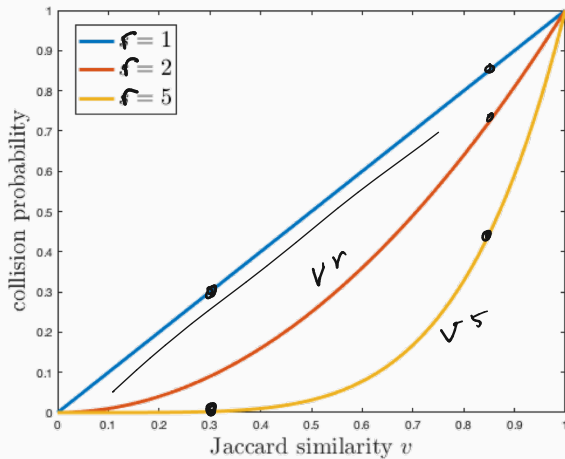
- Choose parameter $r \in \mathbb{Z}^+$.
- Let $c_1, \dots, c_r: \{0, 1\}^d \rightarrow [0, 1]$ be random MinHash.
- Let $g: [0, 1]^r \rightarrow \{1, \dots, m\}$ be a fully random hash function.
- Let $h(x) = g(c_1(x), \dots, c_r(x))$.

$$\text{If } J(\underline{q}, \underline{y}) = v, \text{ then } \Pr[\underline{h(q)} == \underline{h(y)}] = v^r \cdot 1 + \frac{1-v^r}{m} \approx v^r$$

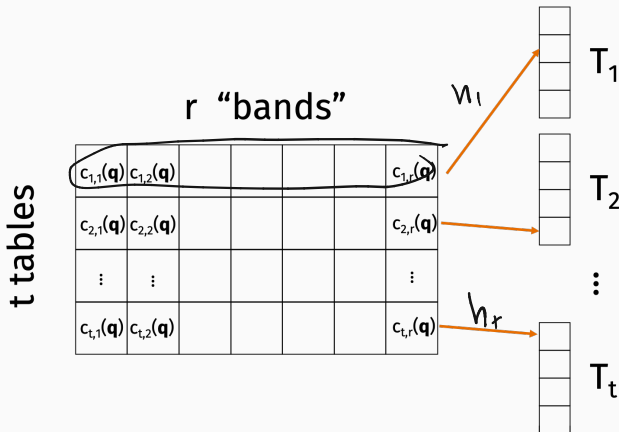
$$1. \ c_i(q) = c_i(y) \text{ for } \underline{\text{all}} \ i \in 1 \dots r$$

2. otherwise

TUNABLE LSH



Full LSH scheme has two parameters to tune:



Effect of increasing number of tables t on:

False Negatives

Decrease

False Positives

Increase

Effect of increasing number of bands r on:

False Negatives

Increase

False Positives

Decrease

SOME EXAMPLES

$t = 3$

78%

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

Parameter: $r = 1$.

Chance we find q with $J(y, q) = .8$:

$$1 - .2^t \quad t = 3 \quad 1 - .2^3 > 99\%$$

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

$$1 - .6^t \geq 78\%$$

SOME EXAMPLES

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

$$t = 5 \\ 58\%$$

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

Parameter: $r = 2$.

Chance we find q with $J(y, q) = .8$:

$$1 - (1 - .8^2)^t \rightarrow \underline{\underline{99\%}}$$
$$1 - .32^t$$

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

$$1 - (1 - .4^2)^t \quad \text{if } t = 5$$
$$1 - .84^t = 58\%$$

SOME EXAMPLES

$t = 12$
12%

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

Parameter: $r = 5$.

Chance we find q with $J(y, q) = .8$:

$$1 - (1 - .8^5)^t \geq \underline{99\%}$$

$t = 12$

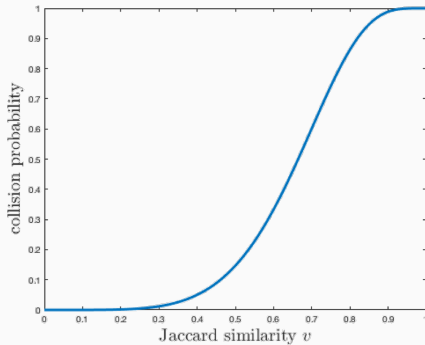
Chance we need to check z with $J(y, z) = .4$:

$$1 - (1 - .4^5)^t \leq 12\%$$

S-CURVE TUNING

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

$$\approx 1 - (1 - \underline{v}^r)^t$$

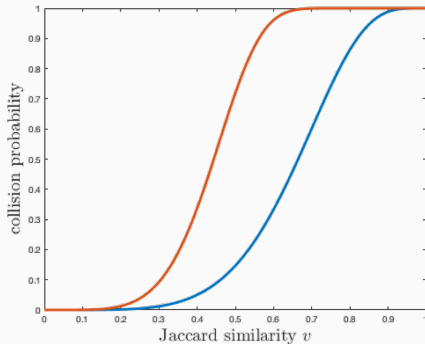


$$r = 5, t = 5$$

S-CURVE TUNING

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

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

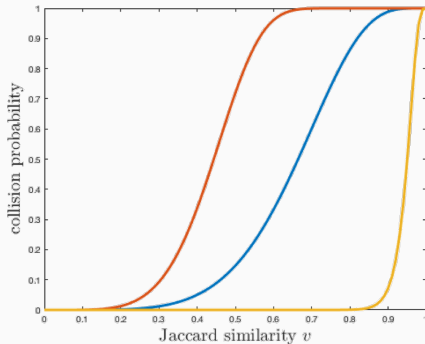


$$r = 5, t = 40$$

S-CURVE TUNING

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

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

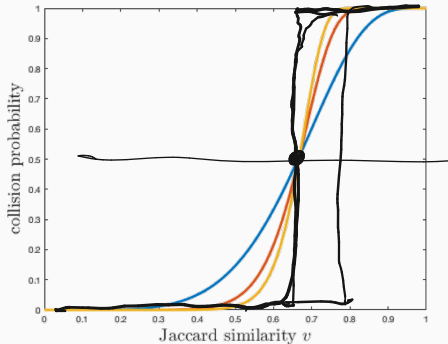


$$r = 40, t = 5$$

S-CURVE TUNING

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.

FIXED THRESHOLD

Use Case 1: Fixed threshold.

$$O(n) + F \cdot O(1) \quad F \cdot O(n)$$

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

With $r = 25$ and $t = 40$,

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

Expected 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.

Nearest Neighbors

Concrete worst case result:

Theorem (Indyk, Motwani, 1998)

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

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

$\|q - y\|_0$ = "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:

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

~~Any ideas for how this is done?~~

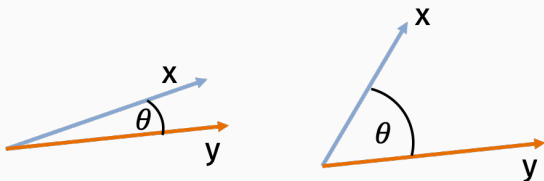
OTHER LSH FUNCTIONS

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

$$\sqrt{|A| |B|}$$

Cosine similarity $\cos(\theta(x, y)) = \frac{\langle x, y \rangle}{\|x\|_2 \|y\|_2}$

$$|A \cup B|$$



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

COSINE SIMILARITY

Cosine similarity is natural “inverse” for Euclidean distance.

$$\text{Euclidean distance } \|x - y\|_2^2 = (x - y)^T (x - y) = x^T x + y^T y - 2x^T y$$

• Suppose for simplicity that $\|x\|_2^2 = \|y\|_2^2 = 1$.

$$\|x\|_2^2 = 1 \quad \|y\|_2^2 = 1$$

$$= \|x\|_2^2 + \|y\|_2^2 - 2x^T y$$

$$2 - 2\cos(\theta(x, y))$$

$$x = \underline{y} \quad \cos(x, y) = 1 \quad \text{sign}(\langle x, y \rangle) = \text{sign}(\langle y, y \rangle) = \text{sign}(\langle \underline{x}, y \rangle)$$

Locality sensitive hash for **cosine similarity**:

- Let $\mathbf{g} \in \mathbb{R}^d$ be randomly chosen with each entry $\mathcal{N}(0, 1)$.
- Let $f: \{-1, 1\} \rightarrow \{1, \dots, m\}$ be a uniformly random hash function.
- $h: \mathbb{R}^d \rightarrow \{1, \dots, m\}$ is defined $h(\mathbf{x}) = \underbrace{f(\text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle))}_{=}$.

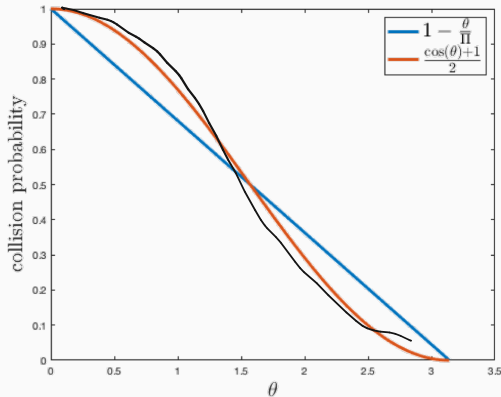
If $\cos(\theta(\mathbf{x}, \mathbf{y})) = v$, what is $\Pr[h(\mathbf{x}) == h(\mathbf{y})]$?

$$\cos(\mathbf{x}, \mathbf{y}) = 1 \quad \text{if } \mathbf{x} = \mathbf{y} \quad \Pr[h(\mathbf{x}) == h(\mathbf{y})] =$$

$$\text{If } \text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle) = \text{sign}(\langle \mathbf{g}, \mathbf{y} \rangle)$$

Theorem: If $\cos(\theta(x, y)) = v$, then

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



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

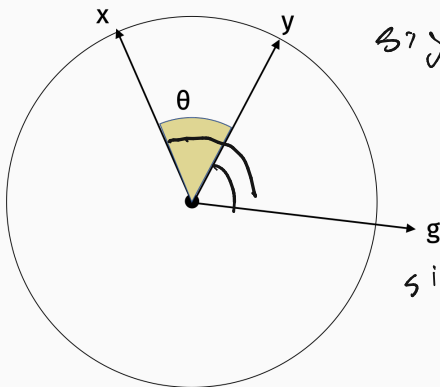
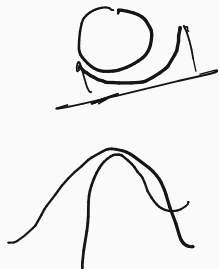
- Let $\mathbf{g}_1, \dots, \mathbf{g}_r \in \mathbb{R}^d$ be randomly chosen with each entry $\mathcal{N}(0, 1)$.
- Let $f: \{-1, 1\}^r \rightarrow \{1, \dots, m\}$ be a uniformly random hash function.
- $h: \mathbb{R}^d \rightarrow \{1, \dots, m\}$ is defined

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

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



SIMHASH ANALYSIS



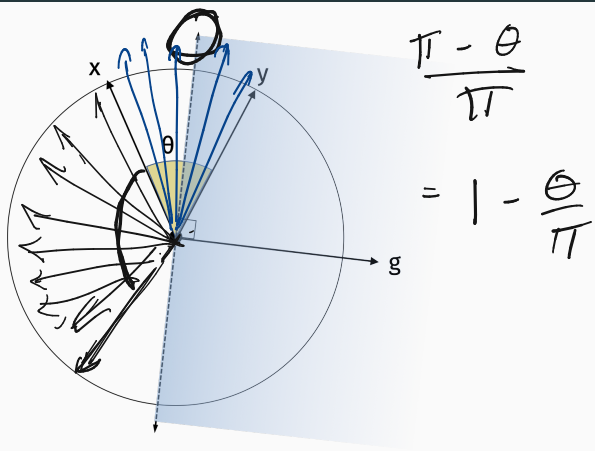
$$\text{sign}(\langle g, y \rangle) = +1$$

$$\text{sign}(\langle g, x \rangle) = -1$$

$$h(x) = f(\text{sign}(\langle g, x \rangle))$$

$$\Pr[h(x) == h(y)] = v + \frac{1-v}{m} \approx v.$$

where $v = \Pr[\text{sign}(\langle g, x \rangle) == \text{sign}(\langle g, y \rangle)]$



$\Pr[h(x) == h(y)] \approx$ probability x and y are on the same side of hyperplane orthogonal to g .