

# CS-GY 6763: Lecture 5

## Dimensionality reduction, near neighbor search in high dimensions

---

NYU Tandon School of Engineering, Prof. Christopher Musco

Despite all our warning from last class that low-dimensional space looks nothing like high-dimensional space, next we are going to learn about how to **compress high dimensional vectors to low dimensions**.

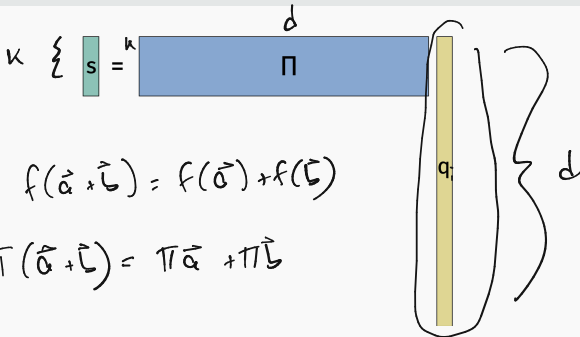
We will be very careful not to compress things too far. An extremely simple method known as (Johnson-Lindenstrauss Random Projection) pushes right up to the edge of how much compression is possible.

# 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  $k = \tilde{O}\left(\frac{\log n}{\epsilon^2}\right)$  such that for all  $i, j$ ,

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



This is equivalent to:

**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  $k = \underline{O}\left(\frac{\log n}{\epsilon^2}\right)$  such that for all  $i, j$ ,

$$(1 - \epsilon) \|\mathbf{q}_i - \mathbf{q}_j\|_2^2 \leq \|\Pi \mathbf{q}_i - \Pi \mathbf{q}_j\|_2^2 \leq (1 + \epsilon) \|\mathbf{q}_i - \mathbf{q}_j\|_2^2.$$

because for small  $\epsilon$ ,  $(1 + \epsilon)^2 = 1 + O(\epsilon)$  and  $(1 - \epsilon)^2 = 1 - O(\epsilon)$ .

$$(1 - \epsilon)^2 \sim$$

$$(1 + \epsilon)^2 \sim$$

$$(1 - \epsilon)^2 \geq 1 - 2\epsilon$$

$$(1 + \epsilon)^2 \leq 1 + 3\epsilon$$



Make pretty much any computation involving vectors faster and more space efficient.

- (Faster vector search (used in image search, AI-based web search, Retrieval Augmented Generation (RAG), etc.).
- (Faster machine learning (today we will see an application to speeding up clustering).
- Faster numerical linear algebra.

Only useful if we can explicitly construct a JL map  $\Pi$  and apply efficiently to vectors.

# EUCLIDEAN DIMENSIONALITY REDUCTION

Remarkably,  $\Pi$  can be chosen (completely at random!)

One possible construction: Random Gaussian.

$$\Pi_{i,j} = \frac{1}{\sqrt{k}} \mathcal{N}(0, 1)$$

# of rows  
(dimension we reduce to)

The map  $\Pi$  is **oblivious to the data set**. This stands in contrast to) to other vector compression methods you might know like PCA.

( [Indyk, Motwani 1998] [Arriaga, Vempala 1999] [Achlioptas 2001]  
[Dasgupta, Gupta 2003].

Many other possible choices suffice – you can use random  $\{+1, -1\}$  variables, sparse random matrices, pseudorandom  $\Pi$ . Each with different advantages.

# RANDOMIZED JL CONSTRUCTIONS

22 x 1000

Let  $\mathbf{\Pi} \in \mathbb{R}^{k \times d}$  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	-1.3328	-1.3617	-0.1952
-0.8396	0.8252	-0.8236	-0.2620	-0.0200	-0.8479	-2.3299	0.4558	-0.2176
1.3546	1.3798	-1.5771	-1.7582	-0.0348	-1.1201	-1.4491	-0.8487	-0.3831
-1.0722	-1.0582	0.5888	-0.2857	-0.7982	2.5268	0.3335	-0.3349	0.0238
0.9610	-0.4686	0.2820	-0.8314	1.0187	1.6555	0.3914	0.5528	0.0513
0.1240	-0.2725	0.0335	-0.0792	-0.1332	0.3875	0.4517	1.0391	0.8261
1.4367	1.0984	-1.3337	-1.1564	-0.7145	-1.2571	-0.1383	-1.1176	1.5278
-1.9689	-0.2779	1.1275	-0.5336	1.3514	-0.8655	0.1837	1.2607	0.4669
-0.1977	0.7015	0.3582	-2.0026	-0.2248	-0.1765	-0.4762	0.6501	-0.2897
-1.2078	-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

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

A random orthogonal matrix  $\mathbf{Q}$  also works. I.e. with  $\mathbf{Q}\mathbf{Q}^T = \mathbf{I}_{k \times k}$ .

For this reason, the JL operation is often called a “random projection”, even though it technically is not a projection when

$\mathbf{\Pi}$ 's entries are i.i.d.

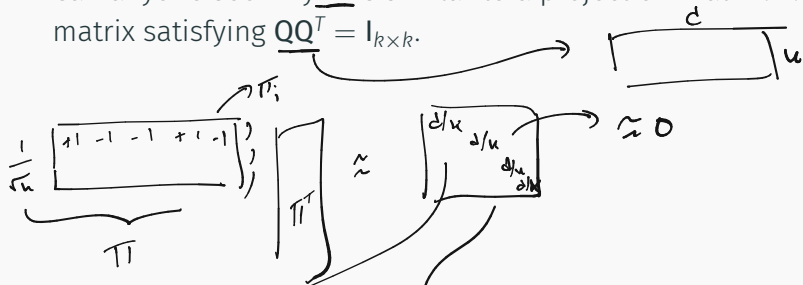
requirement  $\rightarrow$  on entries:

1) mean 0 2) variance 1 3) “sub-gaussian r.v.”

→ Vershynin Book

# RANDOM PROJECTION

Can anyone see why  $\Pi$  is similar to a projection matrix? I.e., a matrix satisfying  $\Pi\Pi^T = I_{k \times k}$ .



$$\langle \pi_i, \pi_i \rangle = \left\langle \frac{1}{\sqrt{k}} [ +1 \ -1 \ -1 \dots ], \frac{1}{\sqrt{k}} [ +1 \ -1 \ -1 \dots ] \right\rangle$$

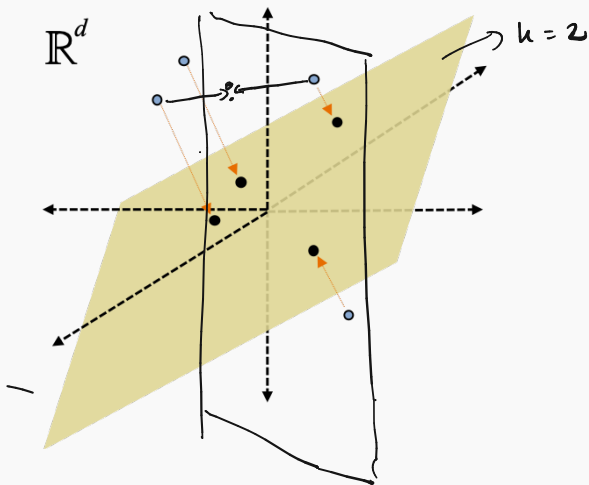
$$= \frac{1}{k} = d/k$$

$$\mathbb{E}(\langle \pi_i, \pi_j \rangle) = 0 \quad \frac{1}{k} \langle [ +1 \ -1 \ +1 \dots ], [ -1 \ +1 \ +1 \dots ] \rangle$$

$\downarrow$   $\downarrow$   
 1st row  $\downarrow$   $\downarrow$   $\downarrow$   $\downarrow$   
 1st row  $\downarrow$   $\downarrow$   $\downarrow$   $\downarrow$

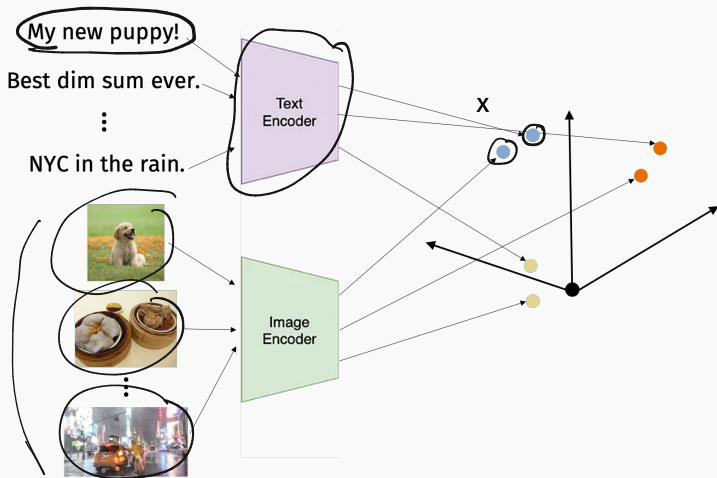
$$= \sum_{l=1}^d \pi_i[l] \pi_j[l]$$

## RANDOM PROJECTION



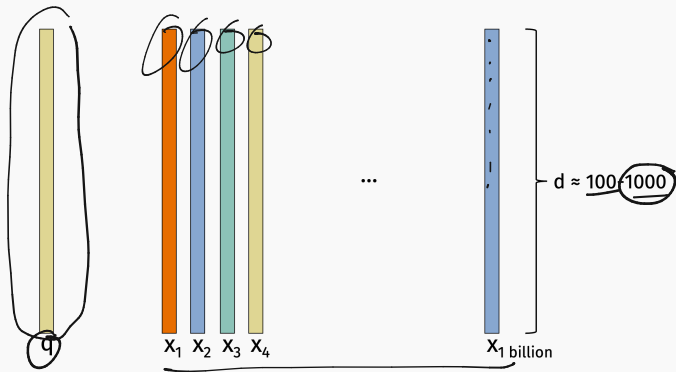
**Intuition:** Multiplying by a random matrix mimics the process of projecting onto a random  $k$  dimensional subspace in  $d$  dimensions.

## APPLICATION: THE NEW PARADIGM FOR SEARCH



Use neural network (BERT, CLIP, etc.) to convert documents, images, etc. to high dimensional vectors. Results matching search should have similar vector embeddings.

## APPLICATION: THE NEW PARADIGM FOR SEARCH



Finding results for a query reduces to finding the nearest vector in a vector database, with similarity typically measured by Euclidean distance. **This is a massive algorithmic challenge!**

$$\min_i \|q - x_i\|_2 \quad \max_i \langle q, x_i \rangle$$

## ANOTHER EXAMPLE OF VECTOR SEARCH

**Shazam** can match a song clip against a library of 8 million songs (32 TB of data) in a fraction of a second. Whole system based on vector embeddings + search.

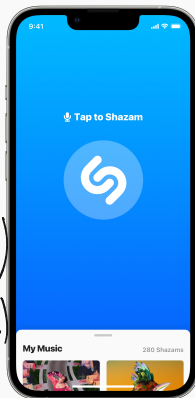
$$x_1, \dots, x_n \in \mathbb{R}^d$$

$$s_1, \dots, s_m \in \mathbb{R}^k$$

$$\min_{s_1, \dots, s_m} \sum_{i,j} (\langle x_i, x_j \rangle - \langle s_i, s_j \rangle)^2$$

$$\min_{s_1, \dots, s_m} \max_{i,j} (\langle x_i, x_j \rangle - \langle s_i, s_j \rangle)^2$$

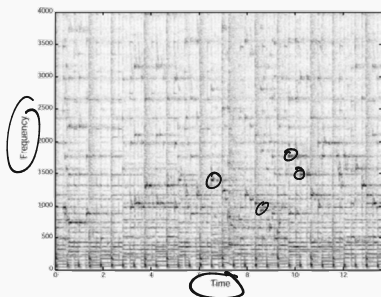
$$\min_{i,j} \max_{i,j} \left| \|x_i - x_j\|^2 - \|s_i - s_j\|^2 \right|$$



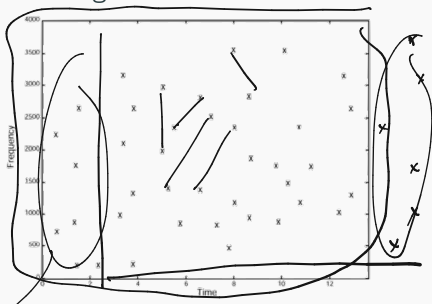


## ANOTHER EXAMPLE OF VECTOR SEARCH

**Shazam** can match a song clip against a library of 8 million songs (32 TB of data) in a fraction of a second. Whole system based on vector embeddings + search.



Spectrogram extracted from audio clip.



Processed spectrogram: used to construct audio “fingerprint”  
 $\mathbf{x} \in \mathbb{R}^d$ .

[...]

]

## VECTOR SEARCH

Tons of new startups in the space (offering managed vector databases) and all major tech companies are frantically working on speeding up vector search.

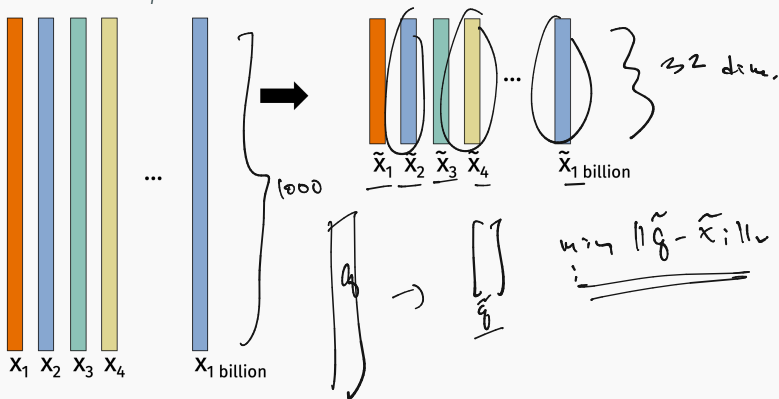


Two main ingredients:

- (1. Vector indexing methods (second half of lecture))
- (2. Vector compression methods (like Johnson-Lindenstrauss).

## APPLICATION: THE NEW PARADIGM FOR SEARCH

Main computational cost is repeatedly computing  $\|\mathbf{q} - \mathbf{x}_i\|_2$  for candidate result  $\mathbf{x}_i$ .



Vector compression leads to faster distance computations. Not only is computational complexity reduced, but we can fit more database vectors in memory.

## EUCLIDEAN DIMENSIONALITY REDUCTION

Lemma (Johnson-Lindenstrauss, 1984)

Let  $\mathbf{\Pi} \in \mathbb{R}^{k \times d}$  be chosen so that each entry equals  $\frac{1}{\sqrt{k}} \mathcal{N}(0, 1)$ , where  $\mathcal{N}(0, 1)$  denotes a standard Gaussian random variable.

If we choose  $k = O\left(\frac{\log(n)}{\epsilon^2}\right)$ , then with probability  $99/100$  for all  $i, j$ ,

$$(1 - \epsilon) \|\mathbf{q}_i - \mathbf{q}_j\|_2^2 \leq \|\mathbf{\Pi} \mathbf{q}_i - \mathbf{\Pi} \mathbf{q}_j\|_2^2 \leq (1 + \epsilon) \|\mathbf{q}_i - \mathbf{q}_j\|_2^2$$

$$\underline{1-\epsilon} \rightarrow O\left(\frac{\log(n/\epsilon)}{\epsilon^2}\right)$$

Intermediate result:

**Lemma (Distributional JL Lemma)**

Let  $\Pi \in \mathbb{R}^{k \times d}$  be chosen so that each entry equals  $\frac{1}{\sqrt{k}} \mathcal{N}(0, 1)$ , where  $\mathcal{N}(0, 1)$  denotes a standard Gaussian random variable.

If we choose  $k = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ , then for any vector  $x$ , with probability  $(1 - \delta)$ :

$$(1 - \epsilon) \|x\|_2^2 \leq \|\Pi x\|_2^2 \leq (1 + \epsilon) \|x\|_2^2$$

$$x = g_i - g_j \quad (1 - \epsilon) \|g_i - g_j\|_2^2 \leq \|\Pi(g_i - g_j)\|_2^2 \leq (1 + \epsilon) \|g_i - g_j\|_2^2$$

$$\|\Pi g_i\|_2^2 \approx \|\Pi g_j\|_2^2$$

Given this lemma, how do we prove the traditional Johnson-Lindenstrauss lemma?

Upper bound

## JL FROM DISTRIBUTIONAL JL

We have a set of vectors  $\underline{\mathbf{q}}_1, \dots, \underline{\mathbf{q}}_n$ . Fix  $\underline{i}, \underline{j} \in 1, \dots, n$ .

Let  $\mathbf{x} = \underline{\mathbf{q}}_i - \underline{\mathbf{q}}_j$ . By linearity,  $\mathbf{\Pi}\mathbf{x} = \mathbf{\Pi}(\underline{\mathbf{q}}_i - \underline{\mathbf{q}}_j) = \underline{\mathbf{\Pi}\mathbf{q}}_i - \underline{\mathbf{\Pi}\mathbf{q}}_j$ .

By the Distributional JL Lemma, with probability  $1 - \delta$ ,

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

Finally, set  $\delta = \frac{1}{100n^2}$ . Since there are  $< n^2$  total  $i, j$  pairs, by a union bound we have that with probability 99/100, the above will hold for all  $i, j$ , as long as we compress to:  $\leq \delta \cdot n^2 = 1/100$

$$k = O\left(\frac{\log(1/(1/100n^2))}{\epsilon^2}\right) = O\left(\frac{\log n}{\epsilon^2}\right) \text{ dimensions. } \square$$

$\nearrow \log(1/\delta)$

$\swarrow \log(100n^2) = O(\log(n))$

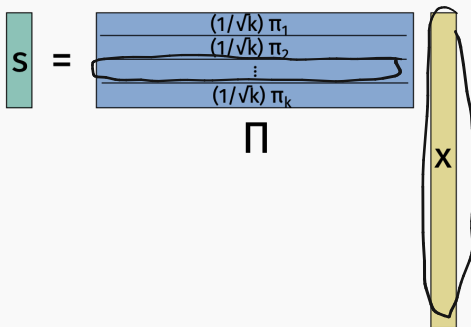
## PROOF OF DISTRIBUTIONAL JL

Want to argue that, with probability  $(1 - \delta)$ ,

$$(1 - \epsilon) \|\underline{\mathbf{x}}\|_2^2 \leq \|\Pi \mathbf{x}\|_2^2 \leq (1 + \epsilon) \|\underline{\mathbf{x}}\|_2^2$$

Claim:  $\mathbb{E} \|\Pi \mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2$ .

Some notation:



So each  $\pi_i$  contains  $\mathcal{N}(0, 1)$  entries.

## PROOF OF DISTRIBUTIONAL JL

**Intermediate Claim:** Let  $\pi$  be a length  $d$  vector with  $\mathcal{N}(0, 1)$  entries.

$$\mathbb{E} [\|\Pi x\|_2^2] = \mathbb{E} [(\langle \pi, x \rangle)^2]$$

$$\begin{aligned} \sum_{i=1}^n (\pi_i x_i)^2 &= \sum_{i=1}^n \langle \frac{1}{\sqrt{n}} \pi_i, x \rangle^2 = \sum_{i=1}^n \frac{1}{n} \langle \pi_i, x \rangle^2 \\ \mathbb{E} \{ \|\Pi x\|_2^2 \} &= \sum_{i=1}^n \frac{1}{n} \mathbb{E} \{ \langle \pi_i, x \rangle^2 \} = \sum_{i=1}^n \frac{1}{n} \mathbb{E} \{ \langle \pi, x \rangle^2 \} \\ &= \mathbb{E} \{ \langle \pi, x \rangle^2 \}. \end{aligned}$$

**Goal:** Prove  $\mathbb{E} \|\Pi x\|_2^2 = \|x\|_2^2$ .



# PROOF OF DISTRIBUTIONAL JL

$$z_1, \dots, z_d$$

$$\langle \pi, x \rangle = Z_1 \cdot x[1] + Z_2 \cdot x[2] + \dots + Z_d \cdot x[d]$$

where each  $Z_1, \dots, Z_d$  is a standard normal  $\mathcal{N}(0, 1)$ .

We have that  $Z_i \cdot x[i]$  is a normal  $\mathcal{N}(0, x[i]^2)$  random variable.

$$\mathbb{E}[\langle \pi, x \rangle^2] = \|x\|_2^2$$

$$\text{Var}(\langle \pi, x \rangle) = \text{Var}\left(\sum_{i=1}^d \mathcal{N}(0, x[i]^2)\right) \stackrel{\text{non-random scalar}}{=} \sum_{i=1}^d x[i]^2 = \|x\|_2^2$$

independence of  $z_1, \dots, z_d$

$$= \mathbb{E}[\langle \pi, x \rangle^2] - \mathbb{E}[\langle \pi, x \rangle]^2$$

Goal: Prove  $\mathbb{E}\|\Pi x\|_2^2 = \|x\|_2^2$ . Established:  $\mathbb{E}\|\Pi x\|_2^2 = \mathbb{E}[(\langle \pi, x \rangle)^2]$

What type of random variable is  $\langle \underline{\pi}, \underline{x} \rangle$ ?

Fact (Stability of Gaussian random variables)

$$\underline{\mathcal{N}}(\mu_1, \sigma_1^2) + \underline{\mathcal{N}}(\mu_2, \sigma_2^2) = \underline{\mathcal{N}}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

$$\begin{aligned} \langle \underline{\pi}, \underline{x} \rangle &= \underline{\mathcal{N}}(0, x[1]^2) + \underline{\mathcal{N}}(0, x[2]^2) + \dots + \underline{\mathcal{N}}(0, x[d]^2) \\ &= \underline{\mathcal{N}}(0, \underline{\|\mathbf{x}\|_2^2}). \end{aligned}$$

So  $\underline{\mathbb{E}} \|\underline{\pi} \mathbf{x}\|_2^2 = \mathbb{E} [(\langle \underline{\pi}, \underline{x} \rangle)^2] = \mathbb{E} [\underline{\mathcal{N}}(0, \underline{\|\mathbf{x}\|_2^2})^2] = \underline{\underline{\|\mathbf{x}\|_2^2}}$ , as desired.

## PROOF OF DISTRIBUTIONAL JL

$$k = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$$

Want to argue that, with probability  $(1 - \delta)$ ,

$$\left( (1 - \epsilon) \|\mathbf{x}\|_2^2 \leq \|\Pi \mathbf{x}\|_2^2 \leq (1 + \epsilon) \|\mathbf{x}\|_2^2 \right)$$

1.  $\mathbb{E} \|\Pi \mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2$ .

2. Need to use a concentration bound.

$\pi_1, \dots, \pi_n$  are i.i.d  
Gaussian random vectors.

$$\|\Pi \mathbf{x}\|_2^2 = \frac{1}{k} \sum_{i=1}^k (\langle \pi_i, \mathbf{x} \rangle)^2 = \frac{1}{k} \sum_{i=1}^k \mathcal{N}(0, \|\mathbf{x}\|_2^2)^2$$

“Chi-squared random variable with  $k$  degrees of freedom.”

# CONCENTRATION OF CHI-SQUARED RANDOM VARIABLES

## Lemma

Let  $H$  be a Chi-squared random variable with  $k$  degrees of freedom.

$$\Pr[|\mathbb{E}H - H| \geq \epsilon \mathbb{E}H] \leq 2e^{-k\epsilon^2/8}$$

$$\Pr[|\|x\|_2^2 - \|\pi x\|_2^2| \geq \epsilon \|x\|_2^2] \leq 2e^{-k\epsilon^2/8}$$

$\underbrace{\hspace{10em}}_{\text{want} = \delta}$

$$2e^{-k\epsilon^2/8} = \delta$$

$$-k\epsilon^2/8 = \log(\delta/2)$$

$$k\epsilon^2/8 = \log(2/\delta)$$

$$k = \frac{\log(2/\delta) \cdot 8}{\epsilon^2}$$

**Goal:** Prove  $\|\pi x\|_2^2$  concentrates within  $1 \pm \epsilon$  of its expectation, which equals  $\|x\|_2^2$ .

If high dimensional geometry is so different from low-dimensional geometry, why is dimensionality reduction possible?

Doesn't Johnson-Lindenstrauss tell us that high-dimensional geometry can be approximated in low dimensions?

## CONNECTION TO DIMENSIONALITY REDUCTION

Hard case:  $\underline{x}_1, \dots, \underline{x}_n \in \mathbb{R}^d$  are all <sup>nearly orthogonal</sup> mutually orthogonal unit  <sup>$2^d$</sup>  vectors:  $\langle \underline{x}_i, \underline{x}_j \rangle \leq \epsilon$

$$\|\underline{x}_i - \underline{x}_j\|_2^2 = 2 - 2\epsilon \quad \text{for all } i, j.$$

When we reduce to  $k$  dimensions with JL, we still expect these vectors to be nearly orthogonal. Why?

$$\langle \Pi \underline{x}_i, \Pi \underline{x}_j \rangle = \frac{1}{2} \left( \overset{\approx 1}{\|\Pi \underline{x}_i\|_2^2} + \overset{\approx 1}{\|\Pi \underline{x}_j\|_2^2} - \overset{\approx \|\underline{x}_i - \underline{x}_j\|_2^2 = 2 - 2\epsilon}{\|\Pi \underline{x}_i - \Pi \underline{x}_j\|_2^2} \right) \leq \underline{\underline{O(\epsilon)}}.$$

$$\|\Pi \underline{x}_i - \Pi \underline{x}_j\|_2^2 = \|\Pi \underline{x}_i\|_2^2 + \|\Pi \underline{x}_j\|_2^2 - 2 \langle \Pi \underline{x}_i, \Pi \underline{x}_j \rangle$$

## CONNECTION TO DIMENSIONALITY REDUCTION

**Hard case:**  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  are all mutually orthogonal unit vectors:

$$\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 = 2 \quad \text{for all } i, j.$$

From our result last class, in  $O(\log n / \epsilon^2)$  dimensions, there exists  $2^{O(\epsilon^2 \cdot \log n / \epsilon^2)} \geq n$  unit vectors that are close to mutually orthogonal.  $O(\log n / \epsilon^2)$  = just enough dimensions.

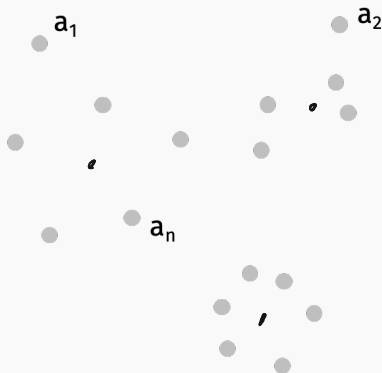
$$2^{O(k \cdot \epsilon^2)} = \# \text{ of } \epsilon\text{-nearly orthogonal unit vectors in } k \text{ dimensional space.}$$

$$k = O(\log(n) / \epsilon^2) \quad 2^{O(\log(n))} \geq n.$$

## SECOND APPLICATION

**k-means clustering:** Give data points  $\underline{a_1}, \dots, \underline{a_n} \in \mathbb{R}^d$ , find centers  $\underline{\mu_1}, \dots, \underline{\mu_k} \in \mathbb{R}^d$  to minimize:

$$\text{Cost}(\mu_1, \dots, \mu_k) = \sum_{i=1}^n \min_{j=1, \dots, k} \|\mu_j - a_i\|_2^2$$

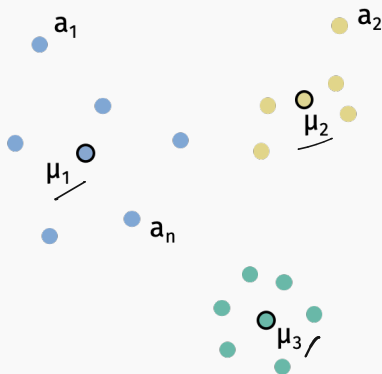




## SAMPLE APPLICATION

**k-means clustering:** Give data points  $\mathbf{a}_1, \dots, \mathbf{a}_n \in \mathbb{R}^d$ , find centers  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k \in \mathbb{R}^d$  to minimize:

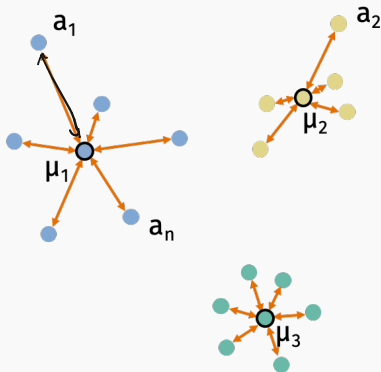
$$\text{Cost}(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k) = \sum_{i=1}^n \min_{j=1, \dots, k} \|\boldsymbol{\mu}_j - \mathbf{a}_i\|_2^2$$



## SAMPLE APPLICATION

k-means clustering: Give data points  $\underline{a}_1, \dots, \underline{a}_n \in \mathbb{R}^d$ , find centers  $\underline{\mu}_1, \dots, \underline{\mu}_k \in \mathbb{R}^d$  to minimize:

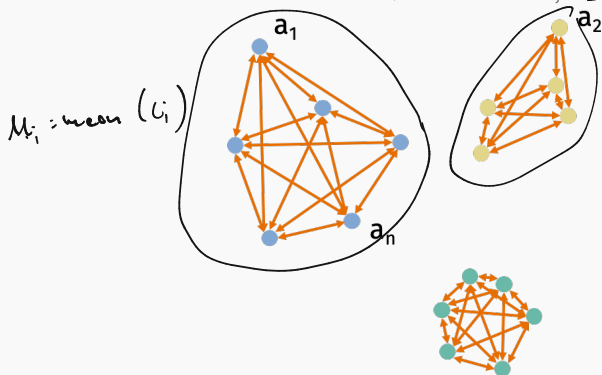
$$\text{Cost}(\underline{\mu}_1, \dots, \underline{\mu}_k) = \sum_{i=1}^n \min_{j=1, \dots, k} \|\underline{\mu}_j - \underline{a}_i\|_2^2$$



# K-MEANS CLUSTERING

**Equivalent form:** Find clusters  $C_1, \dots, C_k \subseteq \{1, \dots, n\}$  to minimize:

$$\text{Cost}(\underline{C_1}, \dots, \underline{C_k}) = \sum_{j=1}^k \left( \frac{1}{2|C_j|} \right) \sum_{u,v \in C_j} \|a_u - a_v\|_2^2.$$

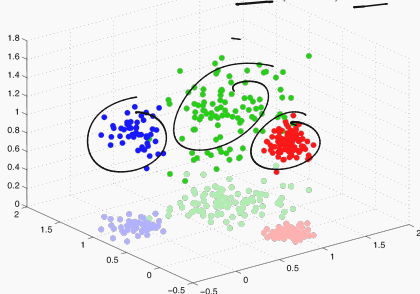


**Exercise:** Prove this to your self.)

# K-MEANS CLUSTERING

NP-hard to solve exactly, but there are many good approximation algorithms. All depend at least linearly on the dimension  $d$ .

**Approximation scheme:** Find clusters  $\tilde{C}_1, \dots, \tilde{C}_k$  for the  $k = O\left(\frac{\log n}{\epsilon^2}\right)$  dimension data set  $\underline{\mathbf{a}}_1, \dots, \underline{\mathbf{a}}_n$ .



Argue these clusters are near optimal for  $\mathbf{a}_1, \dots, \mathbf{a}_n$ .

# K-MEANS CLUSTERING

$$\underline{\text{Cost}(C_1, \dots, C_k)} = \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u,v \in C_j} \underbrace{\|a_u - a_v\|_2^2}_{\text{circled}}$$

$$\underline{\widetilde{\text{Cost}}(C_1, \dots, C_k)} = \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u,v \in C_j} \underbrace{\|\Pi a_u - \Pi a_v\|_2^2}_{\text{underlined}}$$

**Claim:** For any clusters  $\underline{C_1}, \dots, \underline{C_k}$ :

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

$$\frac{1}{1 - \epsilon}$$

$$\sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u,v \in C_j} \underbrace{\|\Pi a_u - \Pi a_v\|_2^2}_{\leq (1+\epsilon) \|a_u - a_v\|_2^2} \leq (1+\epsilon) \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u,v \in C_j} \|a_u - a_v\|_2^2$$

# K-MEANS CLUSTERING

Assume

Suppose we use an approximation algorithm to find clusters

$B_1, \dots, B_k$  such that:

$$\widetilde{\text{Cost}}(\underline{B_1}, \dots, \underline{B_k}) \leq (1 + \alpha) \widetilde{\text{Cost}}^*$$

approx. function

$\widetilde{\text{Cost}}(B_1, \dots, B_k)$

Then:

$$\begin{aligned} \text{Cost}(B_1, \dots, B_k) &\leq \frac{1}{1 - \epsilon} \widetilde{\text{Cost}}(B_1, \dots, B_k) \\ &\leq (1 + O(\epsilon))(1 + \alpha) \widetilde{\text{Cost}}^* \\ &\leq (1 + O(\epsilon))(1 + \alpha)(1 + \epsilon) \text{Cost}^* \\ &= \underline{(1 + O(\alpha + \epsilon)) \text{Cost}^*} \end{aligned}$$

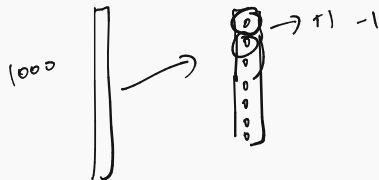
Break until

3:55 pm.

$$\text{Cost}^* = \min_{C_1, \dots, C_k} \text{Cost}(C_1, \dots, C_k) \text{ and}$$

$$\widetilde{\text{Cost}}^* = \min_{C_1, \dots, C_k} \widetilde{\text{Cost}}(C_1, \dots, C_k)$$

## DIMENSIONALITY REDUCTION



Product Quantization

The Johnson-Lindenstrauss Lemma let us sketch vectors and preserve their  $\ell_2$  **Euclidean distance**.

We also have dimensionality reduction techniques that preserve alternative measures of similarity.

## JACCARD SIMILARITY

$$\begin{aligned} q &= \{0, 0, 1, 0, 0, 0\} \\ y &= \{1, 1, 1, 0, 0, 0\} \end{aligned} \quad \frac{1}{4}$$

Often vector embeddings used in semantic search are binary.  
For such vectors, Jaccard similarity is often used instead of Euclidean distance or inner product to compute similarity.

### Definition (Jaccard Similarity)

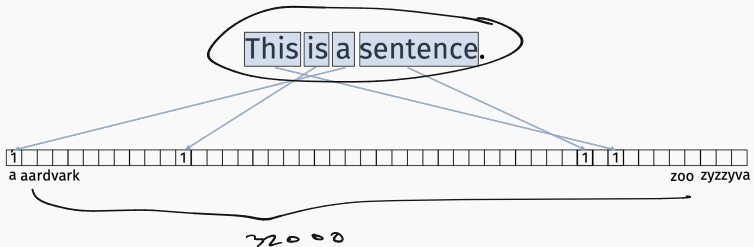
$$J(\underline{q}, \underline{y}) = \frac{|\underline{q} \cap \underline{y}|}{|\underline{q} \cup \underline{y}|} = \frac{\text{\# of non-zero entries in common}}{\text{total \# of non-zero entries}}$$

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



# JACCARD SIMILARITY FOR DOCUMENT COMPARISON

“Bag-of-words” model:

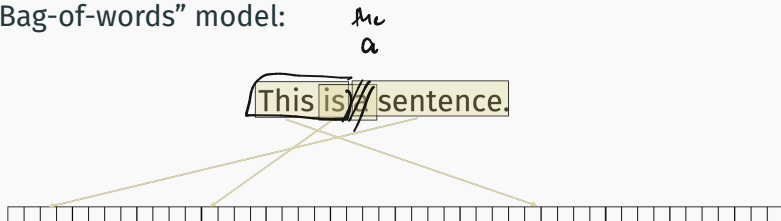


How many words do a pair of documents have in common?

## JACCARD SIMILARITY FOR DOCUMENT COMPARISON

SLADE

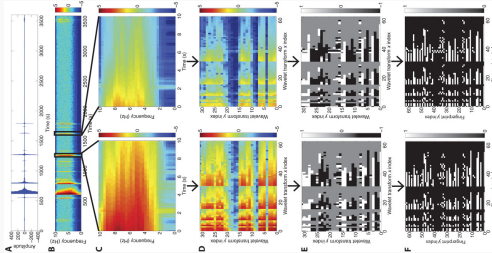
“Bag-of-words” model:



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

- ① 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.

# JACCARD SIMILARITY FOR SEISMIC DATA

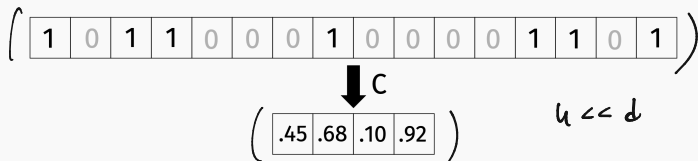


Feature extract pipeline for earthquake data.

(see paper by Rong et al. posted on course website)

## SIMILARITY ESTIMATION

**Goal:** Design a compact sketch  $C : \{0, 1\}^d \rightarrow \mathbb{R}^k$ :



Want to use  $C(\mathbf{q}), C(\mathbf{y})$  to approximately compute the Jaccard similarity  $J(\mathbf{q}, \mathbf{y}) = \frac{|\mathbf{q} \cap \mathbf{y}|}{|\mathbf{q} \cup \mathbf{y}|}$ .

MinHash (Broder, '97): *→ sketch dimension*

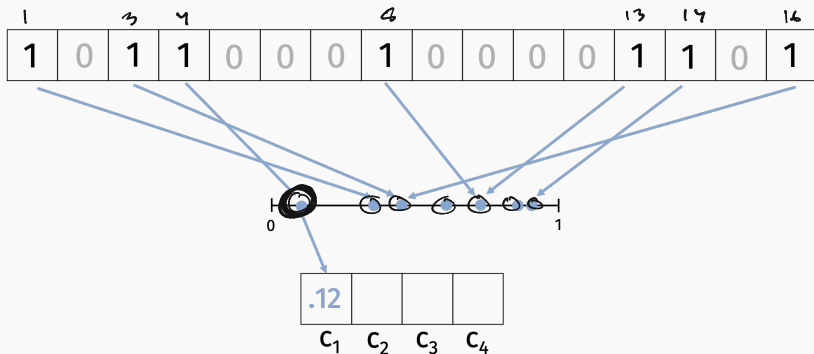
- Choose  $k$  random hash functions

$\underbrace{h_1}, \dots, \underbrace{h_k} : \{1, \dots, \underbrace{d}\} \rightarrow \underbrace{[0, 1]}.$

- For  $i \in 1, \dots, k,$

- Let  $c_i = \min_{j, q_j=1} h_i(j).$

- $C(q) = [c_1, \dots, c_k].$

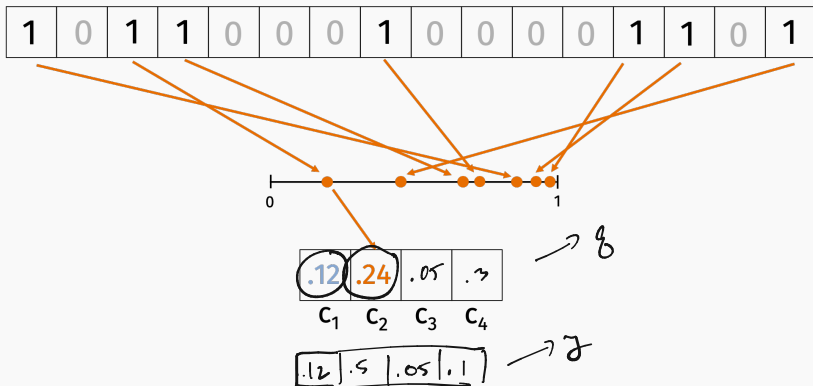


# MINHASH

- Choose  $k$  random hash functions

$$h_1, \dots, h_k : \{1, \dots, n\} \rightarrow [0, 1].$$

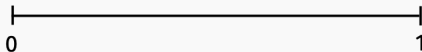
- For  $i \in 1, \dots, k$ ,
  - Let  $c_i = \min_{j, q_j=1} h_i(j)$ .
- $C(\mathbf{q}) = [c_1, \dots, c_k]$ .



## MINHASH ANALYSIS

Claim: For all  $i$ ,  $\Pr[\underline{c_i(q)} = \underline{c_i(y)}] = \underline{J(q, y)} = \frac{|q \cap y|}{|q \cup y|}$ .

<b>q</b>	1	0	1	1	0	0	1	0
<b>y</b>	1	0	0	1	0	1	0	1



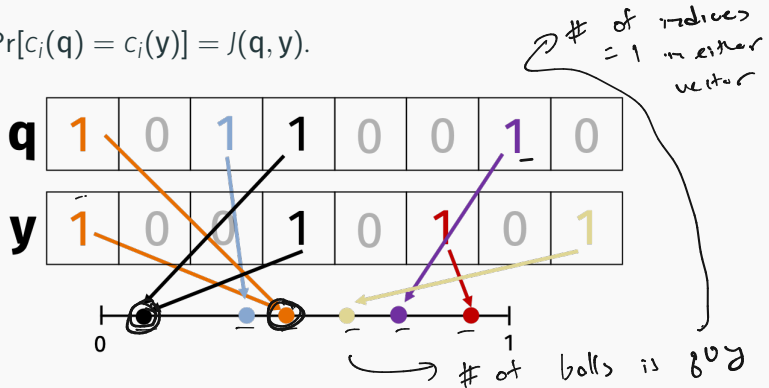
Proof:

1. For  $c_i(q) = c_i(y)$ , we need that  $\arg \min_{i, b_i \neq 0} h(i) = \arg \min_{j, b_j \neq 0} h(i)$ .



## MINHASH ANALYSIS

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



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

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

## MINHASH ANALYSIS

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:

$$\mathbb{E}\tilde{J} = \frac{1}{k} \sum_{i=1}^k \mathbb{E}[\mathbb{1}(c_i(\mathbf{q}) = c_i(\mathbf{y}))] = \frac{1}{k} \cdot k \cdot J(\mathbf{q}, \mathbf{y}) = J(\mathbf{q}, \mathbf{y}).$$

$C(\mathbf{q})$	.12	.24	.76	.35
$C(\mathbf{y})$	.12	.98	.76	.11

$$\tilde{J} = 1/2$$

The more repetitions, the lower the variance.

# MINHASH ANALYSIS

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

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

$$\text{Var}[\tilde{J}] = \frac{1}{k^2} \sum_{i=1}^k \text{Var}(\mathbb{1}[c_i(\mathbf{q}) = c_i(\mathbf{y})]) = \frac{1}{k^2} \sum_{i=1}^k J - J^2 \leq 1/k.$$

Plug into Chebyshev inequality. How large does  $k$  need to be so that with probability  $> 1 - \delta$ ,  $\underline{\underline{|J - \tilde{J}| \leq \epsilon}}$

$$\Pr(|J - \tilde{J}| \geq \alpha \epsilon) \leq \frac{1/\alpha^2}{\epsilon^2} = \delta$$

$$\alpha = 1/\sqrt{\delta}$$

$$k = \left( \frac{\log(1/\delta)}{\epsilon^2} \right)$$

$$\Pr(|J - \tilde{J}| \geq \underbrace{\frac{1}{\sqrt{\delta}} \epsilon}_{1/\sqrt{k}}) \leq \delta$$

$$\frac{1}{\sqrt{\delta} \sqrt{k}} = \epsilon \quad \sqrt{k} = \frac{1}{\sqrt{\delta} \epsilon}$$

$$k = \frac{1}{\epsilon^2 \delta}$$

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

$$J(\mathbf{q}, \mathbf{y}) - \epsilon \leq \underline{\tilde{J}(C(\mathbf{q}), C(\mathbf{y}))} \leq J(\mathbf{q}, \mathbf{y}) + \epsilon.$$

And  $\tilde{J}$  only takes  $O(k)$  time to compute! **Independent** of original vector dimension,  $d$ .

Can be improved to  $\log(1/\delta)$  dependence?

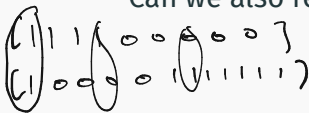
## VECTOR SEARCH / NEAR NEIGHBOR SEARCH

**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.

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.

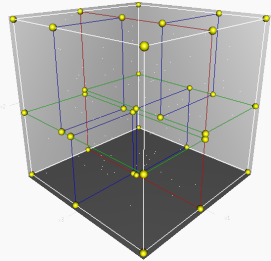
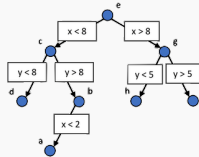
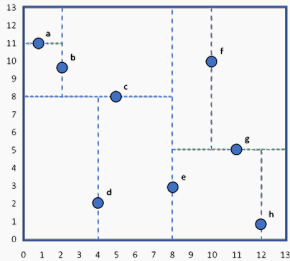
Can we also reduce the dependence on  $n$ ?



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

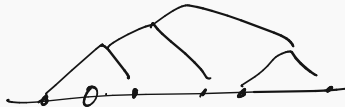
# BEYOND A LINEAR SCAN

This problem can already be solved in low-dimensions using space partitioning approaches (e.g. kd-tree).



Runtime is roughly  ~~$O(n \cdot \min(n, d))$~~   $O(\log(n))$ , which is only sublinear for  $d = o(\log n)$ .

$$O(\log(n) \cdot 2^d)$$



$$O(\log(n)) \text{ time}$$

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]
- Graph-based vector search [Malkov, Yashunin, 2016, Subramanya et al., 2019]

Key ideas behind all of these methods:

1. Trade worse space-complexity + preprocessing time for better time-complexity. I.e., preprocess database in data structure that uses  $\Omega(n)$  space.
2. Allow for approximation.

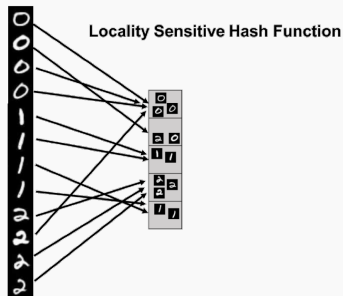
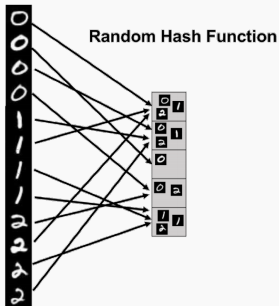


## LOCALITY SENSITIVE HASH FUNCTIONS

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

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

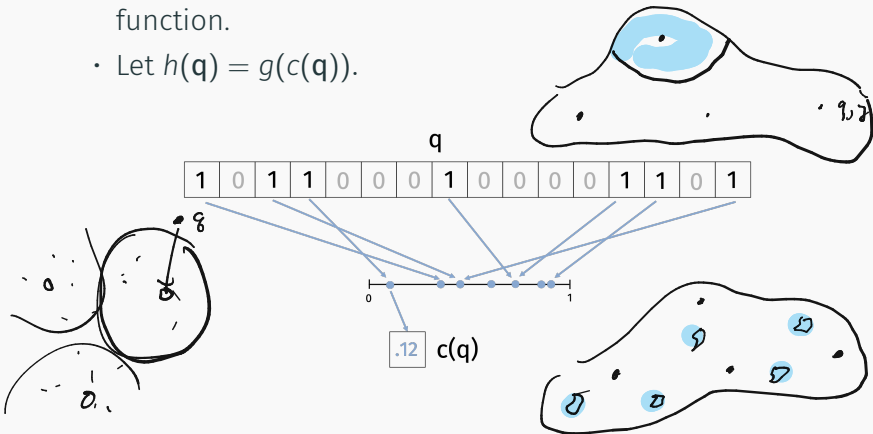
- Higher when  $\mathbf{q}$  and  $\mathbf{y}$  are more similar, i.e.  $s(\mathbf{q}, \mathbf{y})$  is higher.
- Lower when  $\mathbf{q}$  and  $\mathbf{y}$  are more dissimilar, i.e.  $s(\mathbf{q}, \mathbf{y})$  is lower.



## LOCALITY SENSITIVE HASH FUNCTIONS

LSH for  $s(q, y)$  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 uniform random hash function.
- Let  $h(q) = g(c(q))$ .



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

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

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

Basic approach for LSH-based near neighbor search in a database.

### Pre-processing:

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

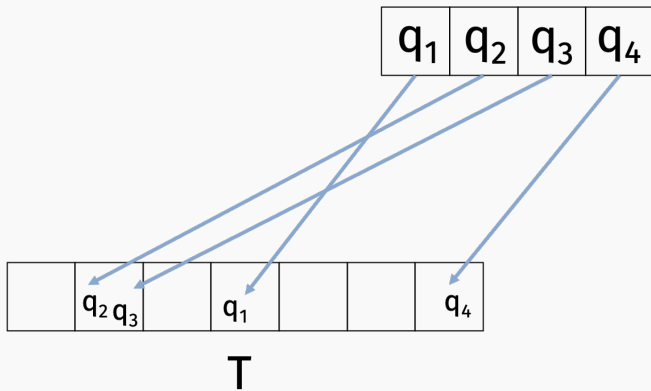
### Query:

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

---

<sup>1</sup>Enough to make the  $O(1/m)$  term negligible.

## 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  $S(\mathbf{q}, \mathbf{y})$  for every  $\mathbf{q} \in T(h(\mathbf{y}))$  to check if it's actually close to  $\mathbf{y}$ .

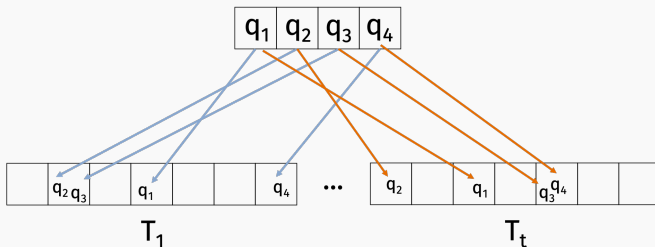
**Note:** The meaning of “close” and “not close” is application dependent. E.g. we might specify that we want to find anything with Jaccard similarity  $> .4$ , but not with Jaccard similarity  $< .2$ .

## REDUCING FALSE NEGATIVE RATE

Let's use Jaccard similarity as a running example. We will discuss LSH for inner product/Euclidean distance as well. Suppose the nearest database point  $\mathbf{q}$  has  $J(\mathbf{y}, \mathbf{q}) = .4$ .

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

## 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  $T_1(h_1(\mathbf{y})) \cup T_2(h_2(\mathbf{y})) \cup \dots, T_t(h_t(\mathbf{y}))$ .

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

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

(10, 99%)

## WHAT HAPPENS TO FALSE POSITIVES?

Suppose there is some other database point  $\mathbf{z}$  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  $\mathbf{y}$  hashes into at least one bucket containing  $\mathbf{z}$ .

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

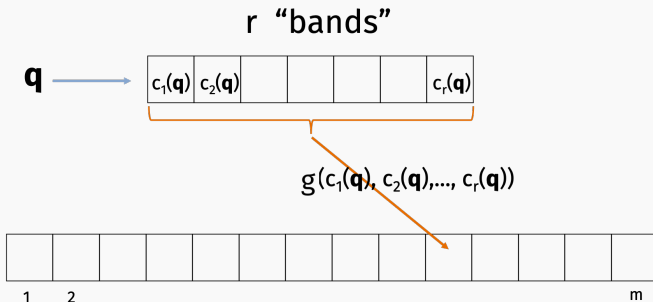
(89%)

## REDUCING FALSE POSITIVES

Change our locality sensitive hash function.

Tunable LSH for Jaccard similarity:

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

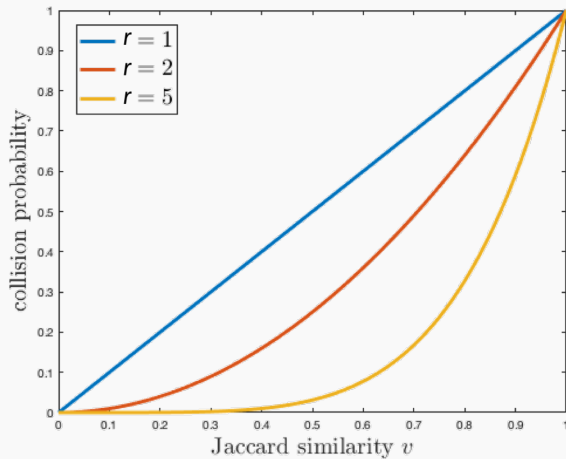


## REDUCING FALSE POSITIVES

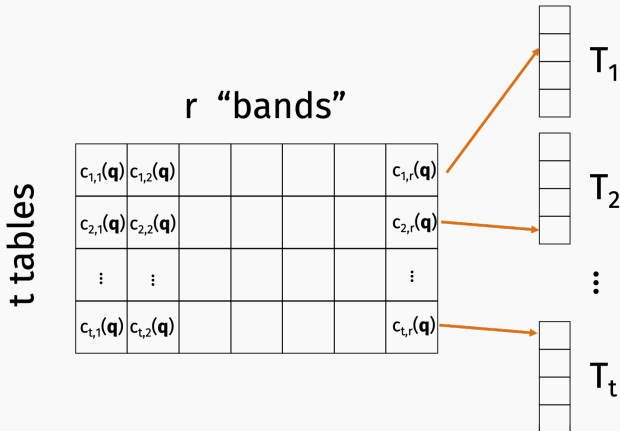
Tunable LSH for Jaccard similarity:

- 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 uniform random hash function.
- Let  $h(\mathbf{x}) = g(c_1(\mathbf{x}), \dots, c_r(\mathbf{x}))$ .

If  $J(\mathbf{q}, \mathbf{y}) = v$ , then  $\Pr[h(\mathbf{q}) == h(\mathbf{y})] =$



Full LSH scheme has two parameters to tune:



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

False Negatives

False Positives

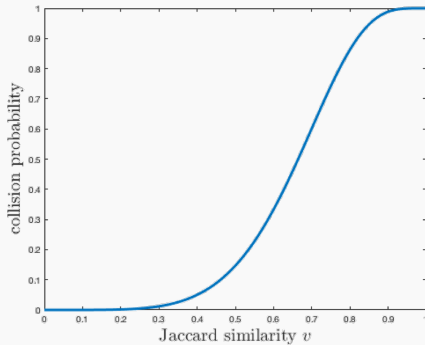
Effect of **increasing number of bands  $r$**  on:

False Negatives

False Positives

## S-CURVE TUNING

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



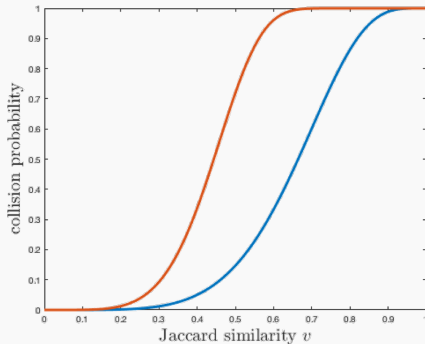
$$r = 5, t = 5$$



## S-CURVE TUNING

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

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

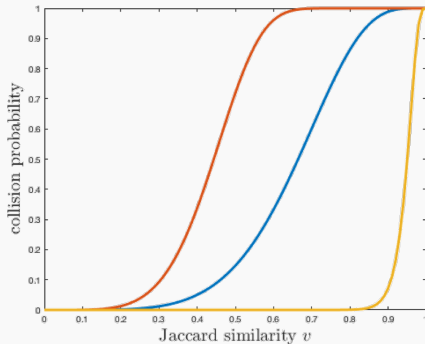


$$r = 5, t = 40$$

## S-CURVE TUNING

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

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

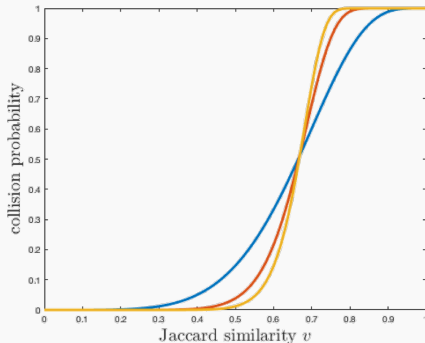


$$r = 40, t = 5$$

## 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.

- 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  $J(\mathbf{y}, \mathbf{q}) \in [.7, .9]$ .
- 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:

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

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

Directly trade space for fast search.

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

$\|\mathbf{q} - \mathbf{y}\|_0$  = “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  $y$ . Return a vector  $\tilde{q}$  with  $\|\tilde{q} - y\|_0 \leq C \cdot \|q - y\|_0$  in:

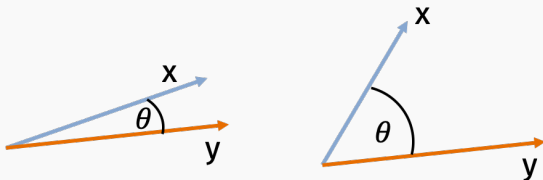
- Time:  $\tilde{O}(n^{1/C})$ .
- Space:  $\tilde{O}(n^{1+1/C})$ .

Similar results can be proven for other metrics, including Euclidean distance. But you need a good LSH function.

## OTHER LSH FUNCTIONS

Good locality sensitive hash functions exist for other similarity measures.

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



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



Cosine similarity is natural “inverse” for Euclidean distance.

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

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

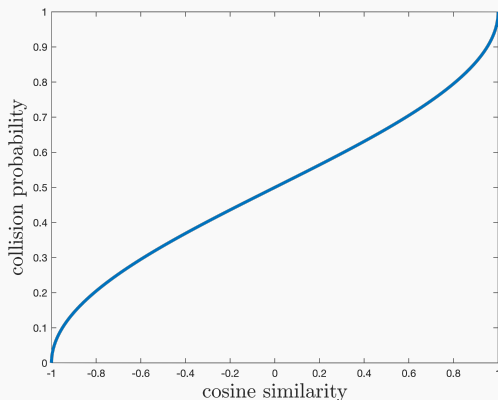
Locality sensitive hash for **cosine similarity**:

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

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

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

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



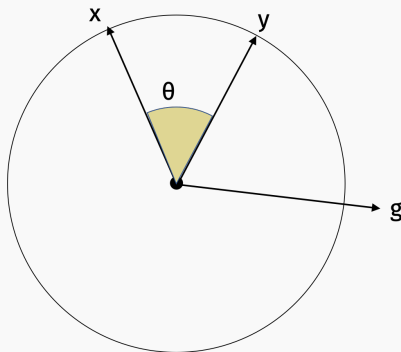
SimHash can be banded, just like our MinHash based LSH function for Jaccard similarity:

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

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

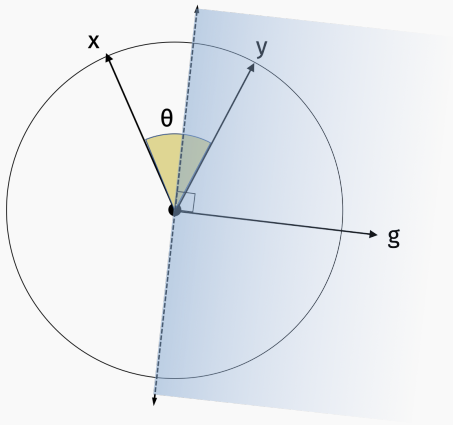
## SIMHASH ANALYSIS IN 2D

To prove:  $\Pr[h(x) == h(y)] \approx 1 - \frac{\theta}{\pi}$ , where  $h(x) = f(\text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle))$  and  $f$  is uniformly random hash function.

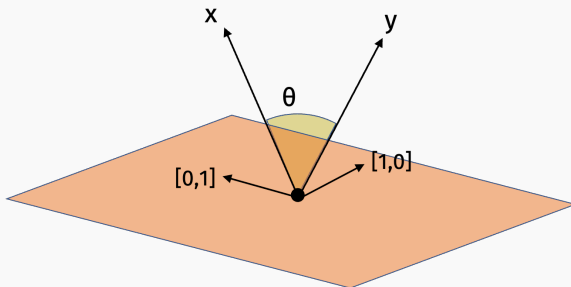


$$\Pr[h(x) == h(y)] = z + \frac{1 - z}{m} \approx z.$$

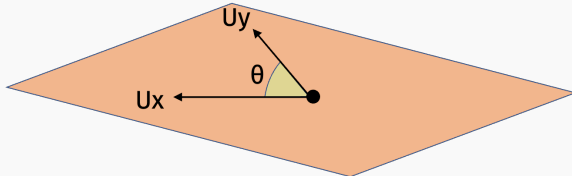
where  $z = \Pr[\text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle) == \text{sign}(\langle \mathbf{g}, \mathbf{y} \rangle)]$



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



There is always some rotation matrix  $\mathbf{U}$  such that  $\mathbf{Ux}, \mathbf{Uy}$  are spanned by the first two standard basis vectors and have the same cosine similarity as  $\mathbf{x}$  and  $\mathbf{y}$ .



There is always some rotation matrix  $\mathbf{U}$  such that  $\mathbf{x}, \mathbf{y}$  are spanned by the first two-standard basis vectors.

**Note:** A rotation matrix  $\mathbf{U}$  has the property that  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ . I.e.,  $\mathbf{U}^T$  is a rotation matrix itself, which reverses the rotation of  $\mathbf{U}$ .



Claim:

$$\begin{aligned}\Pr[\text{sign}(\langle \mathbf{g}, \mathbf{x} \rangle) == \text{sign}(\langle \mathbf{g}, \mathbf{y} \rangle)] &= \Pr[\text{sign}(\langle \mathbf{g}, \mathbf{Ux} \rangle) == \text{sign}(\langle \mathbf{g}, \mathbf{Uy} \rangle)] \\ &= \Pr[\text{sign}(\langle \mathbf{g}[1, 2], (\mathbf{Ux})[1, 2] \rangle) == \text{sign}(\langle \mathbf{g}[1, 2], (\mathbf{Uy})[1, 2] \rangle)] \\ &= 1 - \frac{\theta}{\pi}.\end{aligned}$$

The first step is the trickiest here. Why does it hold?