

Plan

Logistics

Review

Johnson-Lindenstrauss Lemma

Games!

Tea!

Proposal

## Review

$\|x_i\|_2^2 = 1 \quad \forall i$   
 $x_1, \dots, x_t$  nearly orthogonal

if  $|\langle x_i, x_j \rangle| < \epsilon$  for  $i \neq j$

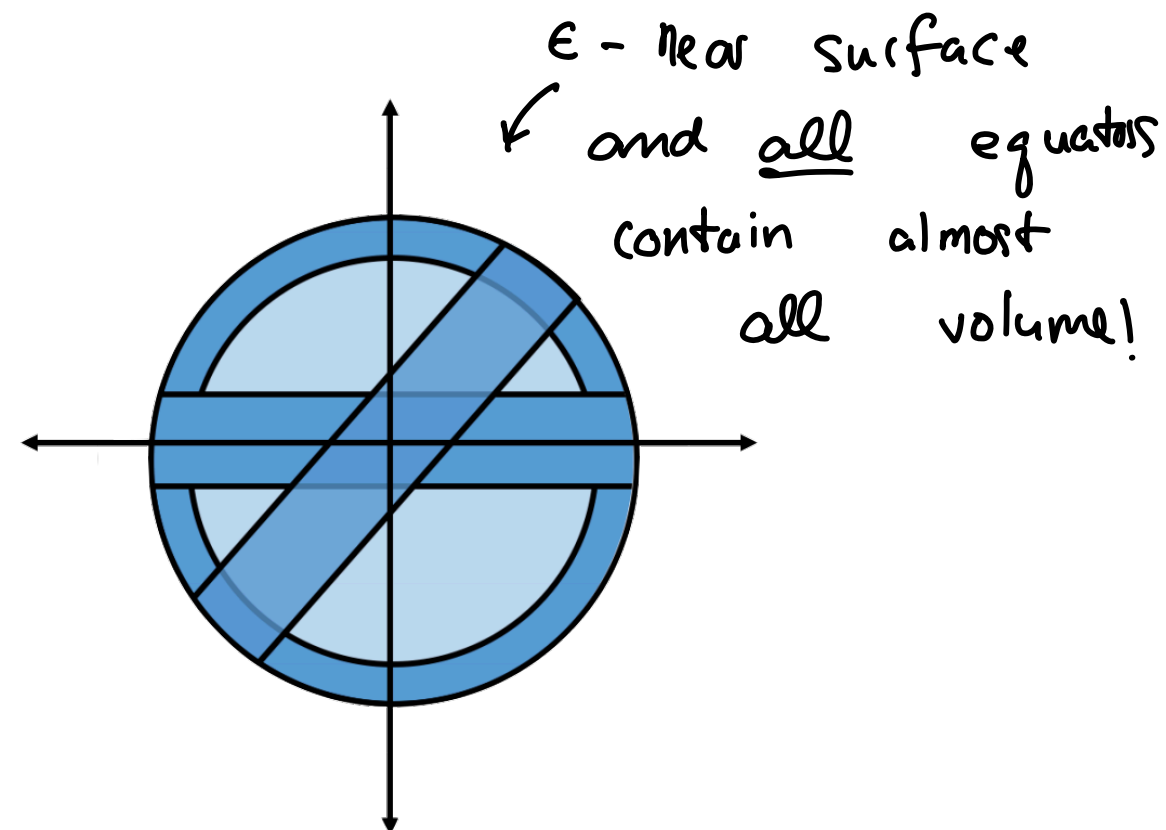
Probabilistic method:

$\Pr(x_1, \dots, x_t \text{ nearly ortho}) > 0$

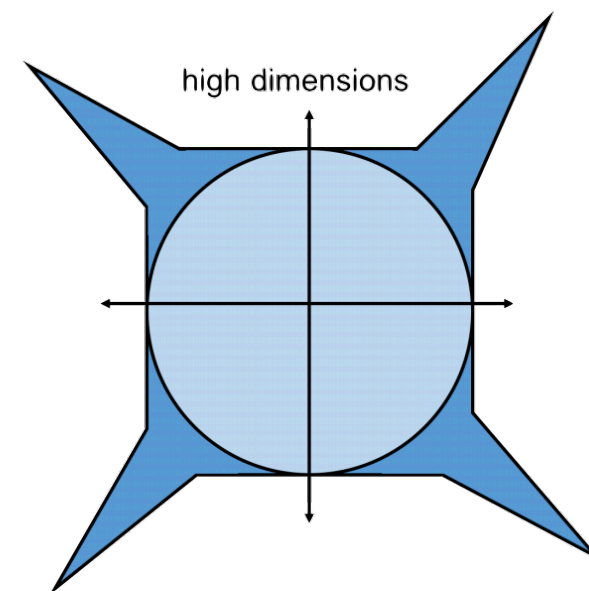
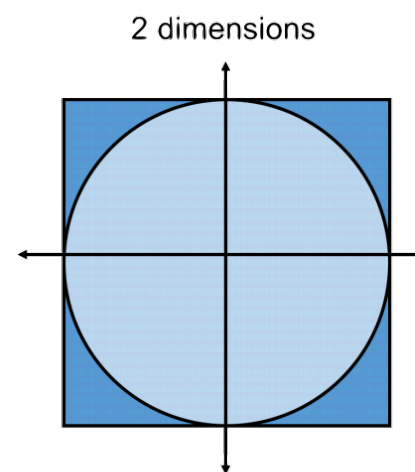
$\Rightarrow \exists$  nearly ortho  $x_1, \dots, x_t$

Proved when  $t = 2^{\epsilon^2 d}$

exponential!



$$\frac{\text{Vol}(\text{cube})}{\text{Vol}(\text{sphere})} \approx d^d$$



High-dimensional geometry  
is weird but we want to  
work with it...

How do we represent  
data using less space  
while approximately preserving  
structure?

Johnson-Lindenstrauss Lemma

↳ Lemma in math paper

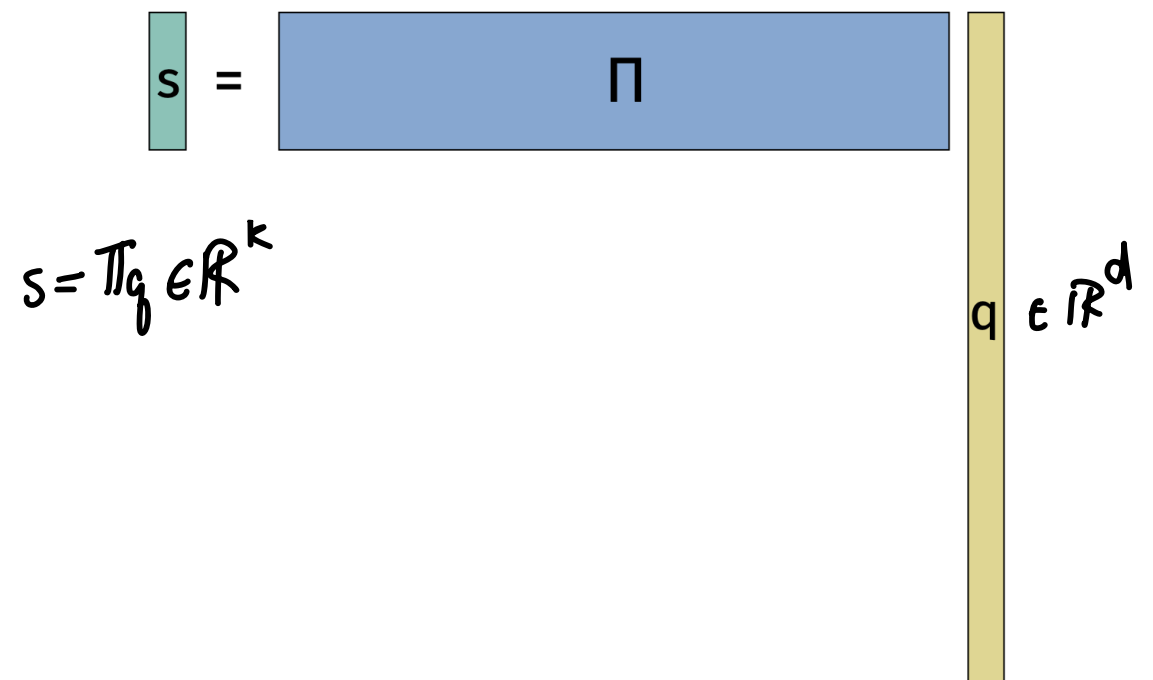
↳ Years for CS community to find

JL Lemma

$$q_1, \dots, q_n \in \mathbb{R}^d$$

There exists  $\Pi: \mathbb{R}^d \rightarrow \mathbb{R}^k$   
for  $k = O\left(\frac{\log n}{\epsilon^2}\right)$  ← independent of  $d$  so that

$$(1-\epsilon) \|q_i - q_j\|_2 \leq \|\Pi q_i - \Pi q_j\|_2 \\ \leq (1+\epsilon) \|q_i - q_j\|_2$$



What about squared norm?

$$(1-\epsilon) \|q_i - q_j\|_2 \leq \|\pi q_i - \pi q_j\|_2 \leq (1+\epsilon) \|q_i - q_j\|_2 \quad \text{with } k = O\left(\frac{\log n}{\epsilon^2}\right)$$

↖  $(1+\epsilon)^2 = 1+2\epsilon+\epsilon^2 = 1+\epsilon \cdot \text{constant}$  for small  $\epsilon$

↘  $(1-\epsilon)^2 \|q_i - q_j\|_2^2 \leq \|\pi q_i - \pi q_j\|_2^2 \leq (1+\epsilon)^2 \|q_i - q_j\|_2^2$  with  $k = O\left(\frac{\log n}{\epsilon^2}\right)$

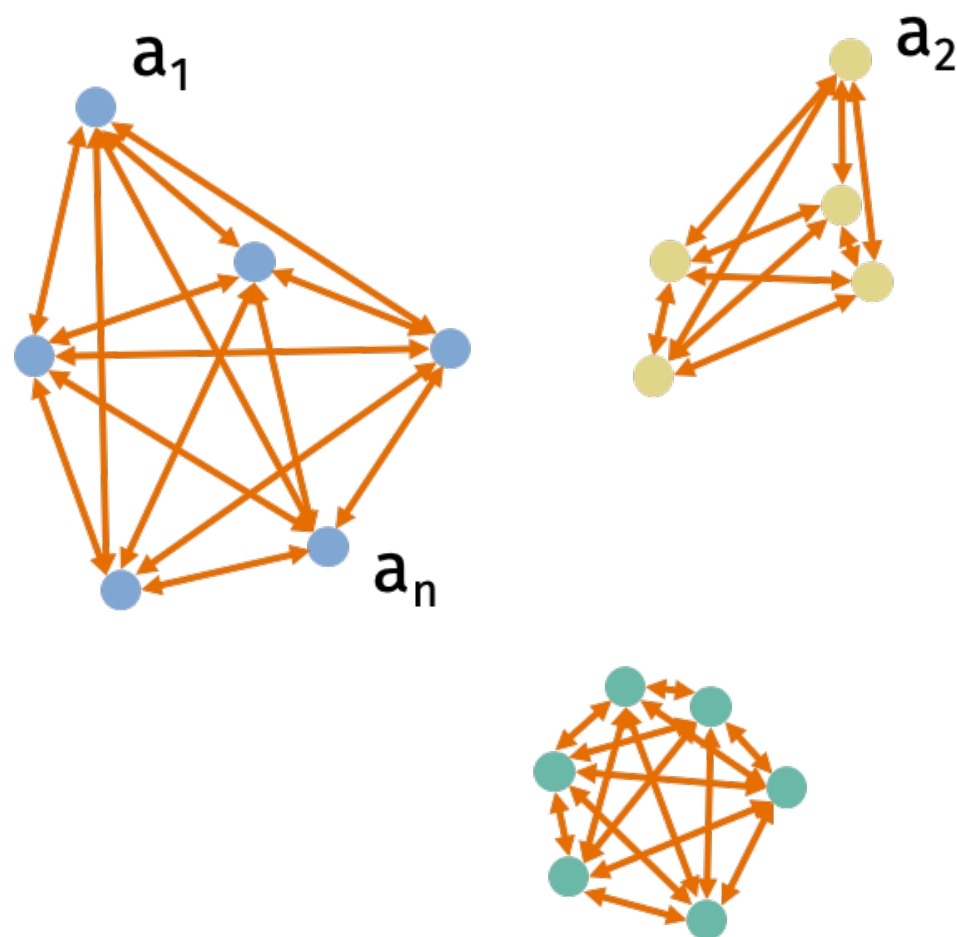
$$(1-\epsilon') \|q_i - q_j\|_2^2 \leq \|\pi q_i - \pi q_j\|_2^2 \leq (1+\epsilon) \|q_i - q_j\|_2^2 \quad \text{with } k = O\left(\frac{\log n}{(\epsilon')^2}\right)$$

$$\begin{aligned} \epsilon' &= 2\epsilon + \epsilon^2 \\ &= \epsilon \cdot \text{constant} \end{aligned}$$

# Clustering with k-means

Problem: Group points  $a_1, \dots, a_n \in \mathbb{R}^d$  into  $k$  clusters  $\{C_1, \dots, C_k\} = C$

$$\text{cost}(C) = \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u,v \in C_j} \|a_u - a_v\|_2^2$$



NP-hard but we can approximate in time depending on  $d$

Idea:

- ① Compress data (Approximately preserve distance)
- ② Cluster on compressed (Approximate solution)
- ③ Return cluster

faster! →

## JL Lemma

What is  $\Pi$ ?

Can we efficiently compute  $\Pi$ ?

$\Pi \in \mathbb{R}^{k \times d}$  is random matrix

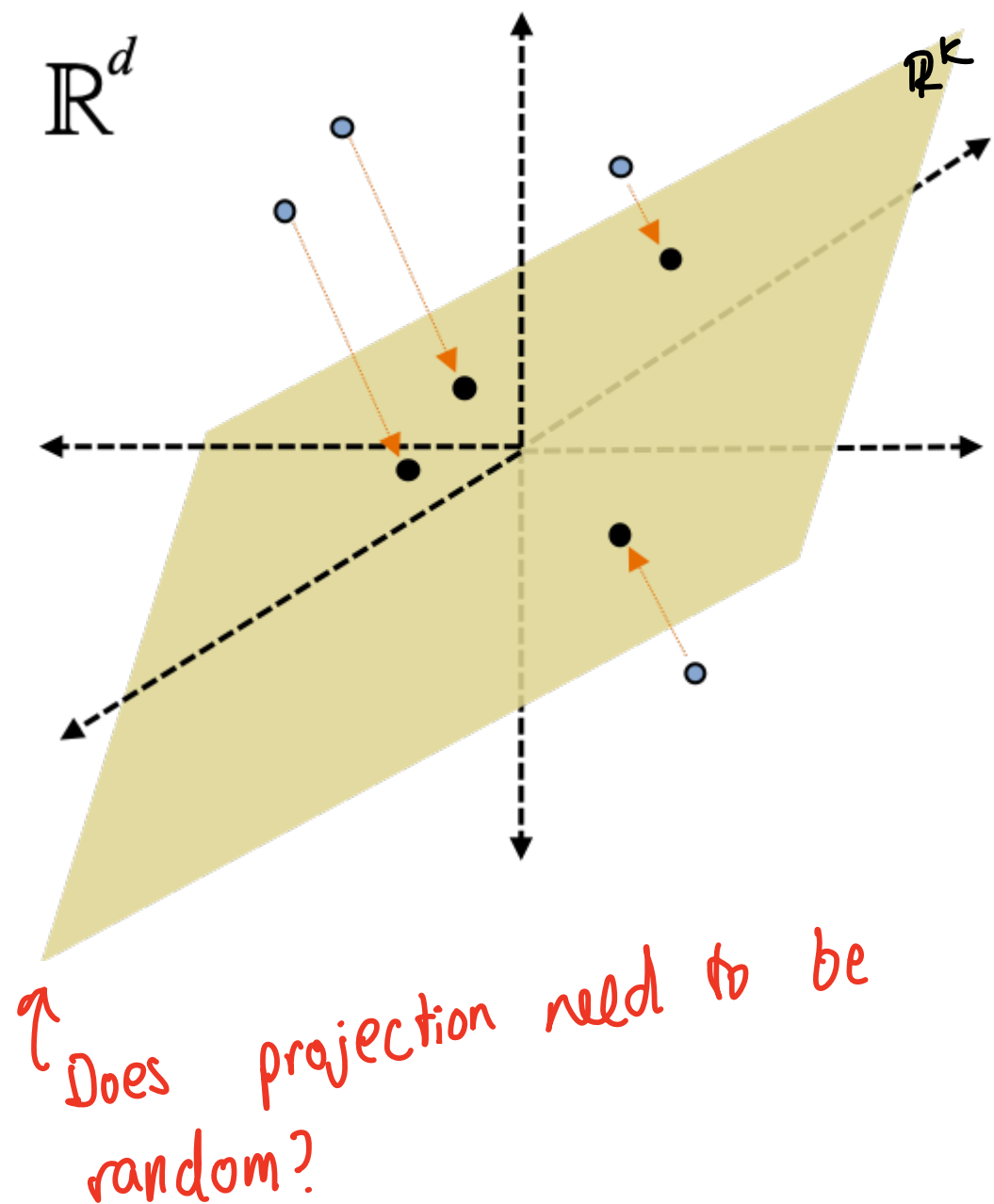
$$\Pi_{ij} \sim \mathcal{N}(0, 1) \cdot \frac{1}{\sqrt{k}} \quad \leftarrow \text{preserve norm}$$

Other random matrices work, too!

↳ binary

↳ sparse

↳ pseudorandom



## Distributional JL Lemma

$\Pi \in \mathbb{R}^{k \times d}$  random scaled normal

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

Then w.p.  $1 - \delta$

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

How do we prove JL with distributional JL?

Proving Distributional JL

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

$\Leftrightarrow$

$$|\|\pi x\|_2^2 - \|x\|_2^2| \leq \epsilon \|x\|_2^2$$

Concentration!

$i^{\text{th}}$  row of  $\pi$

$$\mathbb{E}[\|\pi x\|_2^2] = \sum_{i=1}^k \frac{1}{k} \mathbb{E}[\langle \pi_i, x \rangle^2]$$

$$= \frac{1}{k} \sum_{i=1}^k \mathbb{E} \left[ \sum_{j=1}^d (\pi_i[j] x[j])^2 \right]$$

linearity  
of variance

$$= \frac{1}{k} \sum_{i=1}^k \mathbb{E}[Z_i^2] = \frac{1}{k} \sum_{i=1}^k \|x\|_2^2 = \|x\|_2^2$$

Stability of Gaussians

$$X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$$

$$X_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$$

$$X_1 + X_2 \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

$$Z_i = \sum_{j=1}^d \pi_i[j] x[j]$$

$$\pi_i[j] \sim \mathcal{N}(0, 1)$$

$$\text{Var}(c x) = c^2 \text{Var}(x)$$

$$\pi_i[j] \cdot x[j] \sim \mathcal{N}(0, x[j]^2)$$

$$Z_i \sim \mathcal{N}(0, \|x\|_2^2)$$



Chernoff?  $\therefore$

We actually know this distribution!

$Z$  is chi-squared r.v. with  $k$  degrees of freedom (sum of  $k$  squared normals)

$$\Pr(|Z - \mathbb{E}[Z]| > \epsilon \mathbb{E}[Z]) \leq 2e^{-\epsilon^2 k/8}$$

$$Z = \frac{1}{k} \sum_{i=1}^k z_i^2 = \frac{1}{k} \sum_{i=1}^k (\langle \pi_i, x \rangle)^2 = \|\pi x\|_2^2$$

$$\mathbb{E}[Z] = \|x\|_2^2$$

$$\Pr(\| \pi x \|_2^2 - \|x\|_2^2 > \epsilon \|x\|_2^2) \leq 2e^{-\epsilon^2 k/8} \stackrel{\text{want}}{=} \delta$$

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

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

$$k = \frac{8 \log^2 \frac{2}{\delta}}{\epsilon^2}$$

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

Can we hope for fewer dimensions?

$$x_1, \dots, x_n \in \mathbb{R}^d \quad \text{orthogonal normal}$$

Then  $\|x_i - x_j\|_2^2 = \langle x_i - x_j, x_i - x_j \rangle$

$$= \langle x_i, x_i - x_j \rangle - \langle x_j, x_i - x_j \rangle$$

$$= \|x_i\|_2^2 - \langle x_i, x_j \rangle - \langle x_j, x_i \rangle + \|x_j\|_2^2$$

of the normal  $^1$   
= 2

JL says we can preserve to  $(1 \pm \epsilon)$  in  $k = O\left(\frac{\log n}{\epsilon^2}\right)$  dimensions

From nearly orthogonal,  
we know there are  $2^{O(\epsilon^2 d)}$  nearly orthogonal vectors in  $d$

In  $k$ , there are  $2^{O(\epsilon^2 \cdot \frac{\log n}{\epsilon^2})} \approx n$  nearly orthogonal vectors so we can't put any more in (w/o constants)!