

Plan

Logistics

Review

Johnson-Lindenstrauss Lemma

Games tonight @ 6 tonight here

Tea Time! @ 2 Friday Bihall

Project Proposal

Gradescope

Recommend: Finish problem
in class (then write
up on your own)

Review

$$\|x_i\|_2^2 = 1$$

$x_1, \dots, x_t \in \mathbb{R}^d$ nearly orthogonal

$$|\langle x_i, x_j \rangle| \leq \epsilon \quad i \neq j$$

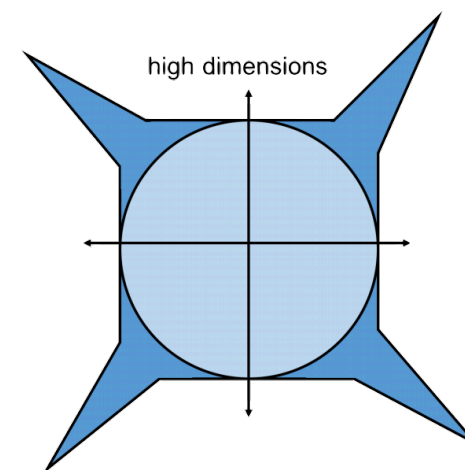
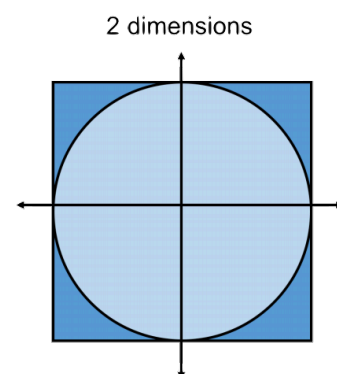
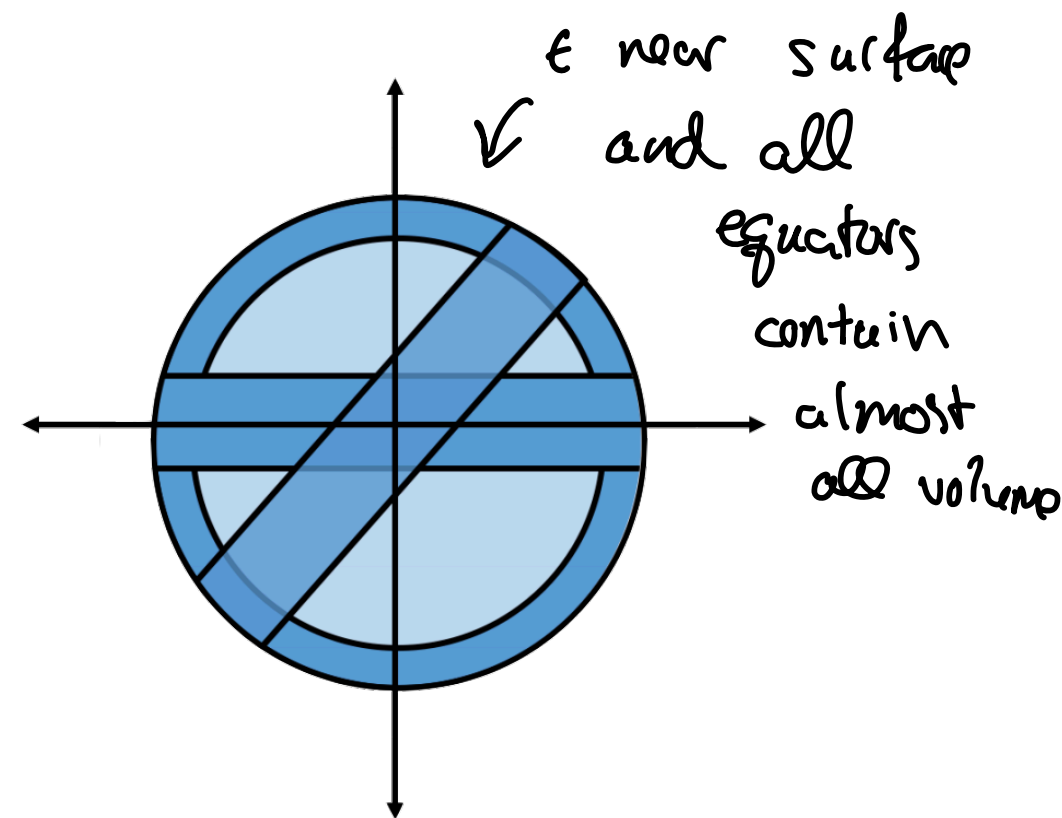
Probabilistic Method

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

$\Rightarrow \exists$ at least one set
of nearly ortho vectors

$$t = 2^{c \cdot \epsilon^2 d}$$

\uparrow exponential in d



$$\frac{\text{Vol}(C_d)}{\text{Vol}(B_d)} \approx d^d$$

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

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

Johnson-Lindenstrauss Lemma

↳ Lemma in math paper

↳ Took CS people years to find

JL Lemma

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

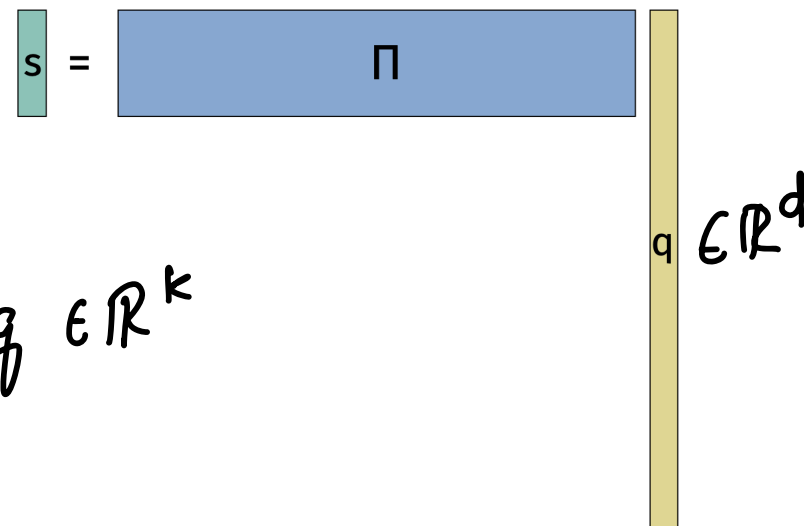
There exists (random) linear map

$$\Pi: \mathbb{R}^d \rightarrow \mathbb{R}^k \quad k = O\left(\frac{\log n}{\epsilon^2}\right)$$

no dependence on d

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

for all $i \neq j$ w.p. $9/10$



$$\|x\|_2^2 = \sum_{i=1}^d (x[i])^2 \quad \|x\|_2 = \sqrt{\sum_{i=1}^d (x[i])^2}$$

$$\epsilon < 1/2$$

$$\epsilon^2 < \frac{1}{2}\epsilon$$

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

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

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

$$(1-\epsilon)^2 = 1 - 2\epsilon + \epsilon^2$$

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

$$(1-\epsilon)^2 = 1 - \epsilon \cdot \text{constant}$$

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

$$\epsilon' = \epsilon \cdot \text{constant}$$

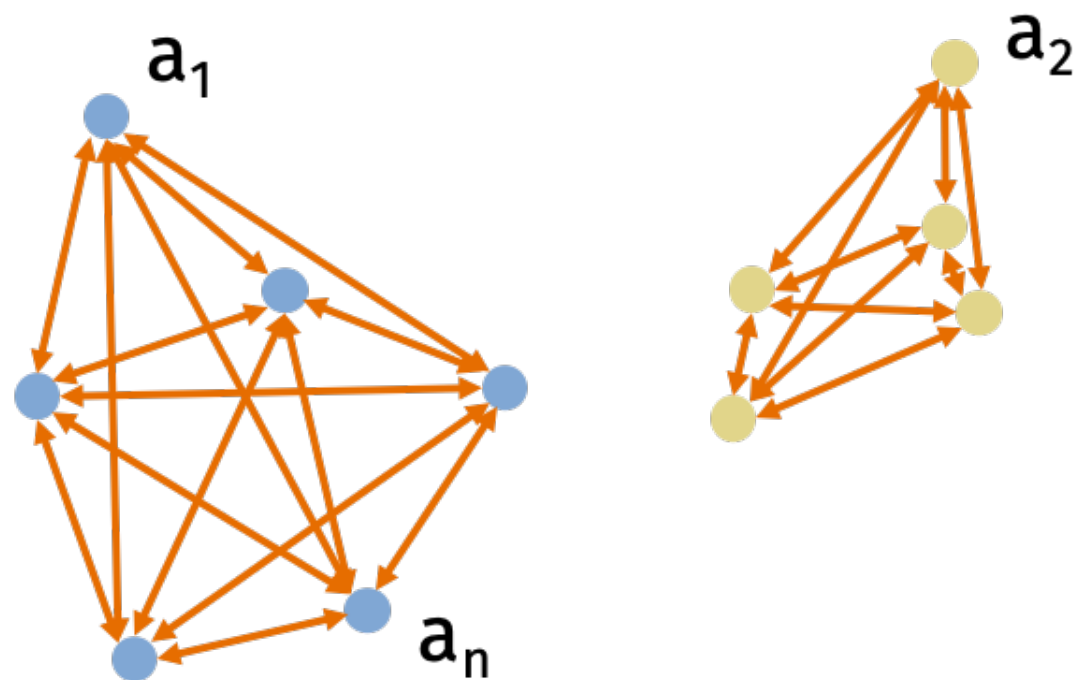
$$(1-\epsilon')\|q_i - q_j\|_2 \leq \|\pi q_i - \pi q_j\|_2 \leq (1+\epsilon')\|q_i - q_j\|_2$$

$$\text{when } k = O\left(\frac{\log n}{\epsilon'^2}\right)$$

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

Clustering

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

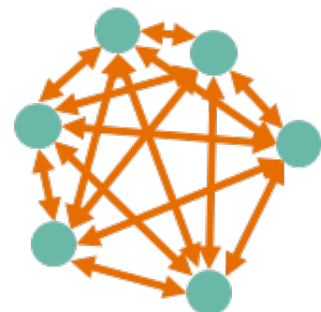


$$\text{Cost}(C) = \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{u, v \in C_j} \|a_u - a_v\|_2^2$$

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

Idea:

- ① compress data (approximately preserves distance)
- ② Cluster on compressed (approx solution)
- ③ return cluster



JL Lemma

What is π ?

Can we efficiently compute π ?

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

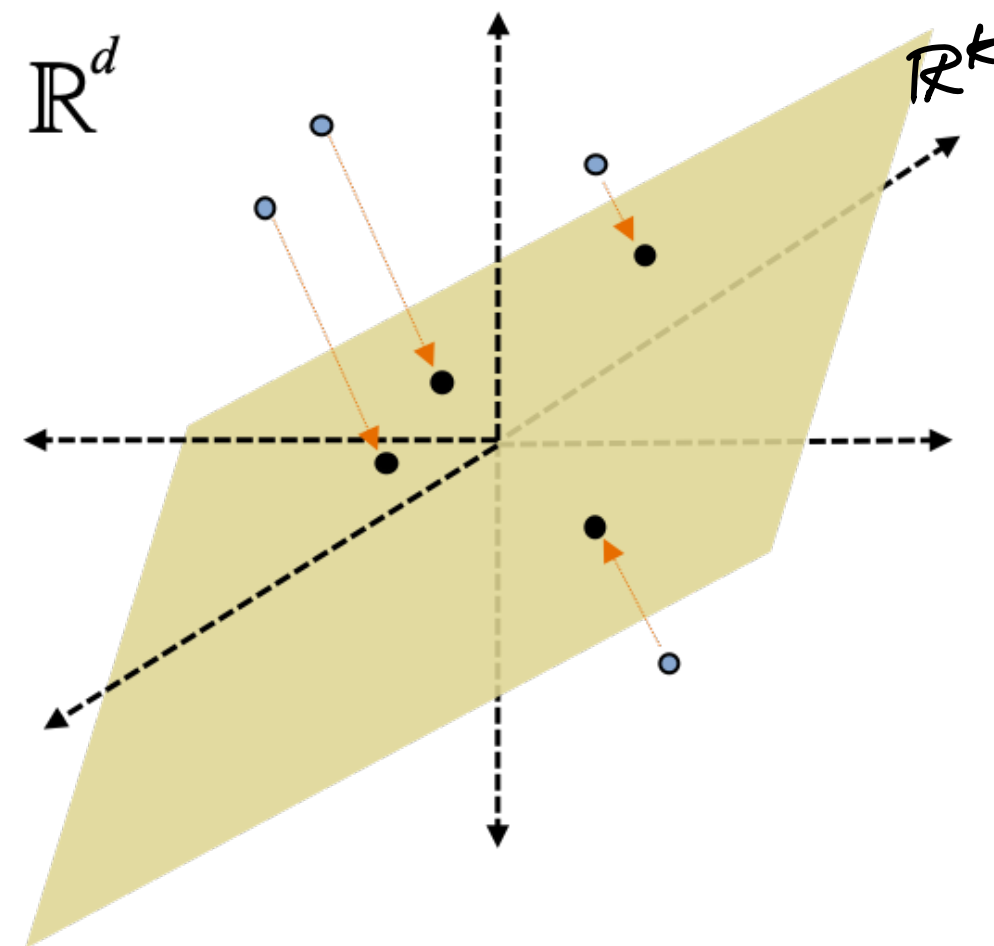
$$\pi_{i,j} \sim \mathcal{N}(0, 1) \cdot \frac{1}{\sqrt{k}} \quad \leftarrow \text{preserving norm}$$

Other random matrices work, too!

↳ binary

↳ sparse

↳ pseudorandom



Does projection need to be random? Why?

Distributional JL Lemma

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

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

Then w.p. $1 - \delta$

$$x = q_i - q_j$$

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

$$(1 - \epsilon) \|q_i - q_j\|_2^2 \leq \|\Pi q_i - \Pi q_j\|_2^2 \leq (1 + \epsilon) \|q_i - q_j\|_2^2 \quad \text{w.p. } 1 - \delta$$

$$\# \text{ pairs } i, j = \binom{n}{2} \leq n^2$$

How do we prove JL lemma using distributional JL lemma?

$$O\left(\frac{\log(1/\delta)}{\epsilon^2}\right) = O\left(\frac{\log 10n^2}{\epsilon^2}\right) = O\left(\frac{\log n}{\epsilon^2}\right)$$

Pr(fails on any)

$$\leq \binom{n}{2} \text{Pr(fails on one)}$$

$$\leq n^2 \underbrace{\text{Pr(fails on one)}}_{\delta} \stackrel{\text{want}}{\leq} \frac{1}{10}$$

$$\delta = \frac{1}{10n^2}$$

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!

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

$$= \frac{1}{k} \sum_{i=1}^k \sum_{j=1}^d \mathbb{E}[(\pi_{i,j} x_j)^2]$$

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

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

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

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

$$Z_i = \sum_{j=1}^d \pi_{i,j} x_j$$

$$\mathbb{E}[cY] = c \mathbb{E}[Y]$$

$$\pi_{i,j} \sim \mathcal{N}(0, 1)$$

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

$$\pi_{i,j} x_j \sim \mathcal{N}(0, x_j^2)$$

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

Chernoff? :-)

We actually know this distribution!

Z be chi-squared r.v. with k degrees of freedom (sum of k squared normal)

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

$$Z = \frac{1}{k} \sum_{i=1}^k z_i^2 = \| \pi x \|_2^2$$

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

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

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

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

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

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

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

Is JL tight?

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

$$\begin{aligned}\|x_i - x_j\|_2^2 &= \langle x_i - x_j, x_i - x_j \rangle = \langle x_i, x_i - x_j \rangle - \langle x_j, x_i - x_j \rangle \\ &= \|x_i\|_2^2 - \langle x_i, x_j \rangle - \langle x_j, x_i \rangle + \|x_j\|_2^2 \\ &= \|x_i\|_2^2 - 2\langle x_i, x_j \rangle + \|x_j\|_2^2 \stackrel{\text{orthonormal}}{=} 2\end{aligned}$$

JL says we can compress 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 ortho in d

When $d = k$, there $2^{O(\epsilon^2 \cdot \frac{\log n}{\epsilon^2})} \approx n$ nearly orthogonal ^{vectors} in k