

Tuesday, Feb 24

- Missed quiz policy

1. Reach out before class

2. Take quiz at first OH after

Goal: Check understanding

Incentivize class attendance

- Late assignment policy

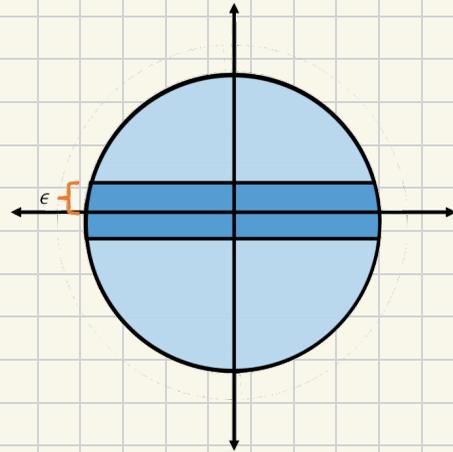
↳ 24 hours late window (life happens)

↳ hard deadline after that

Plan

- High dimensional geometry is weird

- JL projection



Q: What fraction of volume is  $\epsilon$  close to equator?

A: All but  $\frac{1}{2}$ ced fraction of volume

$\epsilon$ -close to any equator

Draw random points from unit ball

Goal: Show that  $x \sim B_d$  has  $|x_i| \leq \epsilon$  w.p.  $1 - \frac{1}{2^c \epsilon^d}$

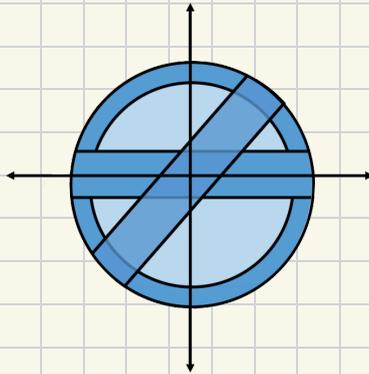
Given  $x$  from interior of unit ball,  $w = \frac{x}{\|x\|_2}$  from surface

If  $|w_i| \leq \epsilon \Rightarrow |x_i| < \epsilon$  since  $\|x\|_2 \leq 1$

New Goal: Show that  $w$  from surface has  $|x_i| \leq \epsilon$  w.p.  $1 - \frac{1}{2^c \epsilon^d}$

Let  $g \sim \mathcal{N}(0, I)$ .  $w = \frac{g}{\|g\|_2}$  from surface by rotational invariance

$$\mathbb{E} \|g\|_2^2 = \dots$$



If:

$$\textcircled{1} \|g\|_2 \geq \sqrt{d}/2$$

$$\textcircled{2} |g_1| \leq \epsilon \sqrt{d}/2$$

Then

$$|w_1| = \frac{|g_1|}{\|g\|_2} \leq \frac{\epsilon \sqrt{d}/2}{\sqrt{d}/2} = \epsilon$$

$$\begin{aligned} \Pr(|w_1| \leq \epsilon) &\geq \Pr(\textcircled{1} \text{ and } \textcircled{2}) \\ &\geq 1 - \Pr(\textcircled{1}^c) - \Pr(\textcircled{2}^c) \end{aligned}$$

$$\begin{aligned} \Pr(\textcircled{1} \text{ and } \textcircled{2}) &= \Pr(\textcircled{1}) + \Pr(\textcircled{2}) - \Pr(\textcircled{1} \text{ or } \textcircled{2}) \\ &\geq \Pr(\textcircled{1}) + \Pr(\textcircled{2}) - 1 \\ &= (1 - \Pr(\textcircled{1}^c)) + (1 - \Pr(\textcircled{2}^c)) - 1 \\ &= 1 - \Pr(\textcircled{1}^c) - \Pr(\textcircled{2}^c) \end{aligned}$$

$$\Pr(\textcircled{1}^c) = \Pr(\|g\|_2 < \sqrt{d}/2) \leq \frac{1}{2cd} \text{ by Johnson-Lindenstrauss Lemma}$$

$$\begin{aligned} \Pr(\textcircled{2}^c) &= \Pr(|g_1| > \epsilon \sqrt{d}/2) \leq \frac{1}{2(c\epsilon \sqrt{d}/2)^2} \text{ by Gaussian tail bound} \\ &\geq 1 - \frac{1}{2c^2\epsilon^2 d/2} - \frac{1}{2cd} \end{aligned}$$

larger for small  $\epsilon$

Despite warnings, let's play with high-dimensional data!

Goal: Compress to smaller dimension while preserving structure

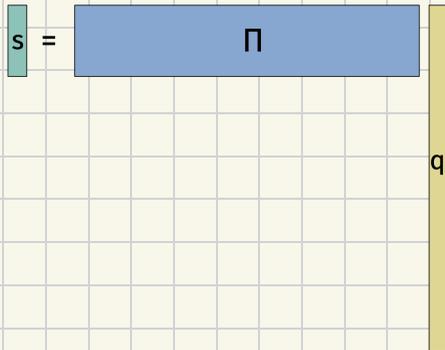
Johnson-Lindenstrauss Lemma:  $q_1, \dots, q_n \in \mathbb{R}^d$   $\exists \Pi: \mathbb{R}^d \rightarrow \mathbb{R}^k$

for  $k = O\left(\frac{\log n}{\epsilon^2}\right)$  s.t., for all  $i, j \in [n]$ , w.p.  $9/10$ ,

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

"Lemma" as a stepping stone to another result, immensely useful

$(1 \pm \epsilon)^2 \approx (1 \pm \epsilon)$   
for small  $\epsilon$



# Clustering Application

Points  $a_1, \dots, a_n \in \mathbb{R}^d$

Partition  $[n]$  into  $m$  clusters

$$C = \{C_1, \dots, C_m\}$$

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

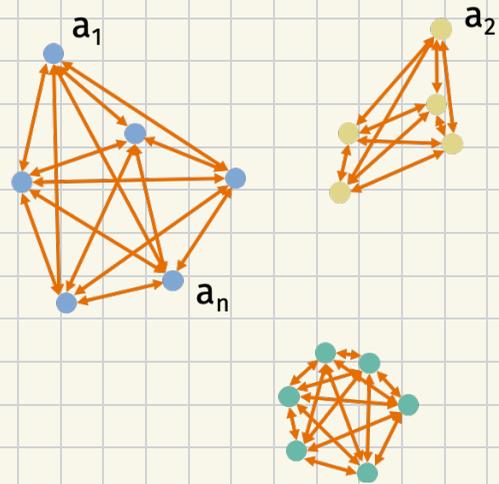
Exact solution is NP-hard

Approximate efficiently with polynomial dependence on  $d$ .

Goal: Decrease  $d$

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

is cost on projected data



By JL,

$$(1-\epsilon)\text{Cost}(c) \leq \tilde{\text{Cost}}(c) \leq (1+\epsilon)\text{Cost}(c)$$

With an approx algorithm, find  $\leftarrow$  optimal on projected

$$\tilde{\text{Cost}}(c) \leq (1+\alpha)\tilde{\text{Cost}}(\tilde{c}^*)$$

$\leftarrow$  optimal on original

$$\text{Claim: } \text{Cost}(c) \leq (1+o(\alpha+\epsilon))\text{Cost}(c^*)$$

$$\text{Hint: } \frac{1}{1-\epsilon} \approx 1+\epsilon \text{ for small } \epsilon$$

Thursday, Feb 26

- Exam 3/5

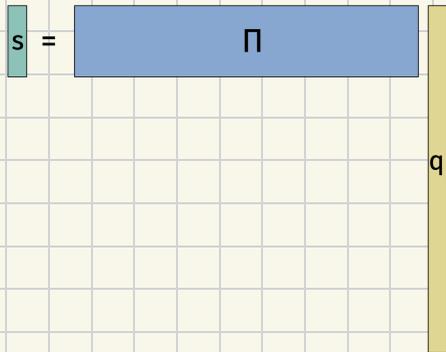
Plan

- Continue JL

Johnson-Lindenstrauss Lemma:  $q_1, \dots, q_n \in \mathbb{R}^d$   $\exists \Pi: \mathbb{R}^d \rightarrow \mathbb{R}^k$   
for  $k = O\left(\frac{\log n}{\epsilon^2}\right)$  s.t., for all  $i, j \in [n]$ , w.p.  $9/10$ ,

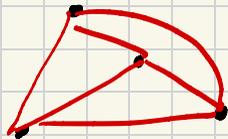
$$(1 - \epsilon) \|q_i - q_j\|_2^2 \leq \|\Pi q_i - \Pi q_j\|_2^2 \leq (1 + \epsilon) \|q_i - q_j\|_2^2$$

"Lemma" as a stepping stone to another result, immensely useful



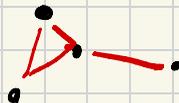
# Clustering Application

$C^*$  ← optimal in original



$$C = \{C_1, \dots, C_m\}$$

↔



$\tilde{C}^*$  ← optimal in projected



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

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

$$\text{JL: } (1-\epsilon) \text{cost}(C) \leq \tilde{\text{cost}}(C) \leq (1+\epsilon) \text{cost}(C) \quad \epsilon \ll 1$$

$$\text{Approx: Find } C \text{ s.t. } \tilde{\text{cost}}(C) \leq (1+\alpha) \text{cost}(\tilde{C}^*)$$

$$\begin{aligned} \text{cost}(C) &\leq \frac{1}{1-\epsilon} \tilde{\text{cost}}(C) \leq (1+2\epsilon) \tilde{\text{cost}}(C) \leq (1+2\epsilon)(1+\alpha) \tilde{\text{cost}}(\tilde{C}^*) \\ &\leq (1+4\epsilon+2\alpha) \tilde{\text{cost}}(\tilde{C}^*) \leq (1+4\epsilon+2\alpha) \text{cost}(\tilde{C}^*) \\ &\leq (1+4\epsilon+2\alpha)(1+\epsilon) \text{cost}(C^*) \end{aligned}$$

## Distributional JL Lemma

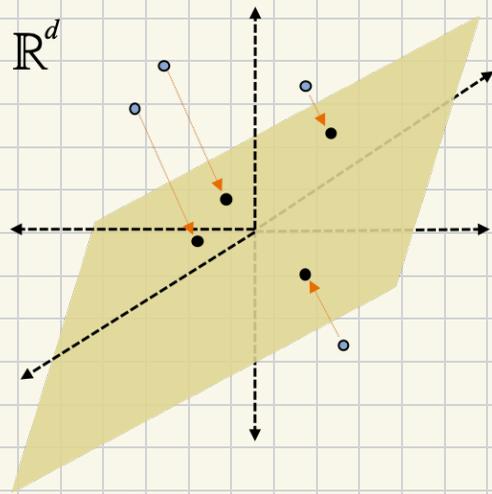
Q: Can we efficiently compute  $\Pi \in \mathbb{R}^{k \times d}$ ?

Easiest to analyze:

$$[\Pi]_{i,j} = g/\sqrt{k} \quad \text{where } g \sim \mathcal{N}(0, 1)$$

Let  $k = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ . For fixed  $x \in \mathbb{R}^d$ , w.p.  $1-\delta$ ,

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



close  $\leftrightarrow$  close  
far  $\leftrightarrow$  far

Prove JL Lemma with  $x = q_i - q_j$  and Union Bound

$$g \in \mathbb{R}^k \quad g \sim \mathcal{N}(0, I)$$

Using dist. JL, prove  $\Pr(\|g\|_2^2 \leq \frac{1}{2} \mathbb{E} \|g\|_2^2) \leq \frac{1}{2ck}$  for constant  $c$

## Proof of dist. JL

Goal: Show  $\|\Pi x\|_2^2$  concentration

$$\mathbb{E} \|\Pi x\|_2^2 = \sum_{i=1}^k \frac{1}{k} \mathbb{E} [\langle \pi_i, x \rangle^2] = \frac{1}{k} \sum_{i=1}^k \mathbb{E} \left[ \left( \sum_{j=1}^d \pi_i c_j \right) x_{[j]} \right)^2$$

$z_i$

Fact: Stability of Gaussians. For  $x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$  and  $x_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$   
 $x_1 + x_2 \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$

Claim:  $z_i \sim \mathcal{N}(0, \|x\|_2^2)$ .

$$\text{Then } \mathbb{E} \|\Pi x\|_2^2 = \frac{1}{k} \sum_{i=1}^k \mathbb{E} [z_i^2] = \|x\|_2^2$$

$$Z = \frac{1}{k} \sum_{i=1}^k z_i$$

↖ scaled chi-squared rv with  $k$  degrees of freedom

### Chi-squared Concentration

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

Our case:  $Z = \|\Pi x\|_2^2 = \frac{1}{k} \sum_{i=1}^k z_i^2$

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

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

JL tells us we can preserve high dimensional behavior...

But wait, isn't high-dimensional geometry weird?

Reconcile with "hard" case. Consider orthonormal  $x_1, \dots, x_n \in \mathbb{R}^d$ .

① JL says we can compress to  $k = O(\log n / \epsilon^2)$  dimensions to retain  $\epsilon$ -approx

② Probabilistic method tells us there are  $2^{c \epsilon^2 k}$  nearly orthogonal vectors in  $k$  i.e.,  $2^{c \epsilon^2 k} = 2^{c \cdot \log n / \epsilon^2} = 2^{c \log n} \approx n$  (between friends)

$\Rightarrow$  two sides of the same coin!

Alternate View: To have random points be "close" in  $d$ , we need  $n \approx 2^d$  of them.

In this case, JL doesn't help since  $\log n \approx d$