

Tuesday, Feb 24

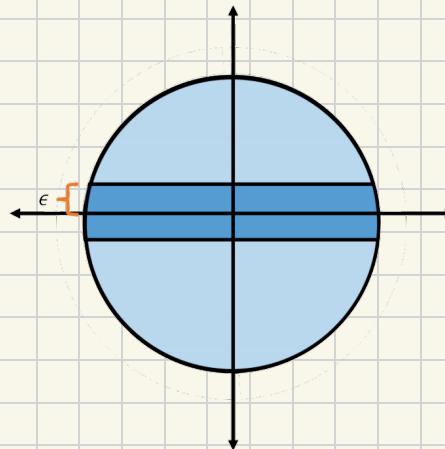
- Missed quiz policy
  1. Reach out before class
  2. Take quiz at first OH after

Goal: Check understanding

Incentivize class attendance

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.  $1 - \frac{1}{2^{c\epsilon d}}$

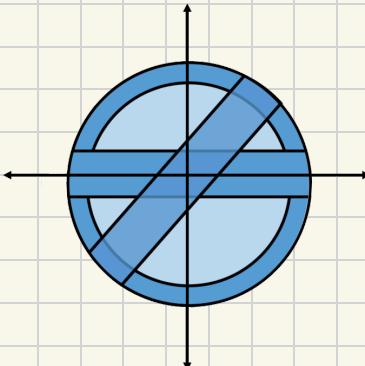
Given  $x$  from interior of unit ball,  $w = \frac{x}{\|x\|_2}$  from surface

If  $|w_i| \leq \epsilon \Rightarrow |x_i| \leq \epsilon$  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  $g \sim N(0, I)$ .  $w = \frac{g}{\|g\|_2}$  from surface by rotational invariance

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



If :

Then

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

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

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

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

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

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

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

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

Goal: Compress to smaller dimension while preserving structure

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

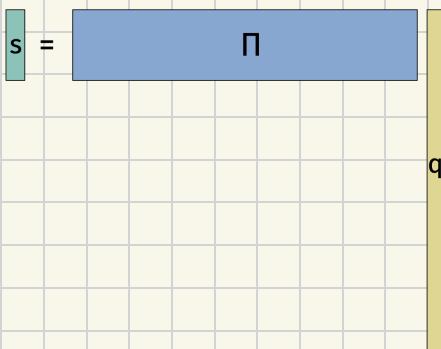
for  $k = O\left(\frac{\log n}{\epsilon^2}\right)$  s.t., for all  $i, j \in [n]$ , w.p.  $9/10$ ,

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

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

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

for small  $\epsilon$



## Clustering Application

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

Partition  $[n]$  into  $m$  clusters

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

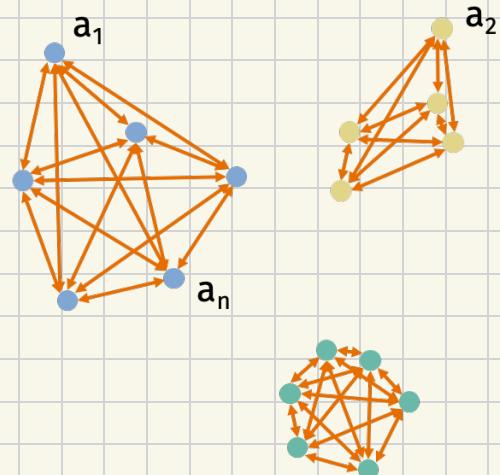
$$\text{Cost}(C) = \sum_{j=1}^m \frac{1}{2|C_j|} \sum_{u,v \in C_j} \|a_u - a_v\|_2^2$$

Exact solution is NP-hard.

Approximate efficiently with polynomial dependence on  $d$ .

Goal: Decrease  $d$

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



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

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

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

## Distributional JL Lemma

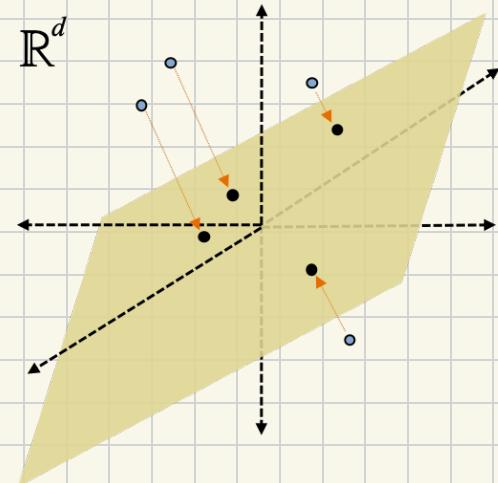
Q: Can we efficiently compute  $\Pi \in \mathbb{R}^{K \times d}$ ?

Easiest to analyze:

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

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  $x = g_i - g_j$  and Union Bound

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

Proof of dist. JL

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

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

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:  $z_i \sim \mathcal{N}(0, \|x\|_2^2)$ . Then  $\mathbb{E} \|\pi x\|_2^2 = \|x\|_2^2$