

Announcement:

Recitation

Friday

3-4 pm

Remote

Feb 15 2022

## Dimension Reduction

Data in high dimensions.

E.g. Documents, Movie or product ratings by users,  
gene expression data

Curse of ↓ dimensionality  
(high)

E.g. Nearest-neighbor search

Dimen. Redun. : Transform the vectors into lower  
dimension while retaining useful properties

① pair-wise distances

② variance of the data points

# 1. Johnson-Lindenstrauss Transform (JL)

- Linear transform

$$x \rightarrow Ax$$

Linear transformation  $m \times n$ .

# of data points

- maintain pairwise distances (L2 norm) b/w data points

Set of points

$$X = \{x_1, x_2, \dots, x_n\}, \quad n \in \mathbb{R}^D$$

Each  $x_i$

initial dim

Want Linear Transform

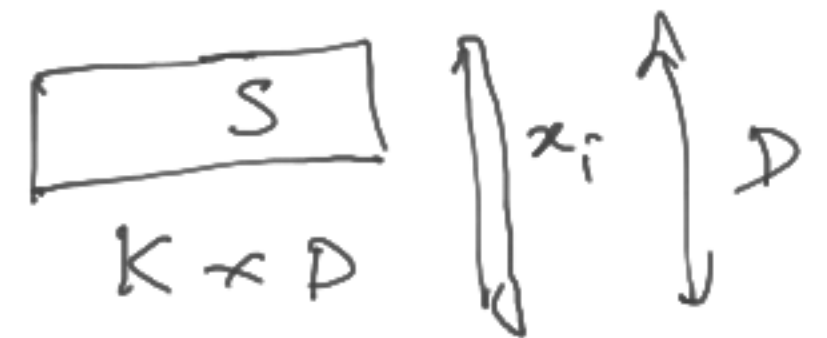
$$S: \mathbb{R}^D \rightarrow \mathbb{R}^K$$

final dim

$$K < D$$

$$\text{Any vector } x_i \rightarrow Sx_i$$

$K \times D$



# JL Lemma

Let  $\varepsilon \in (0, 1/2)$ .

$X = \{x_1, x_2, \dots, x_n\}$  each  $x_i$  in  $\mathbb{R}^D$

There exists a map (transformation)  $S: \mathbb{R}^D \rightarrow \mathbb{R}^k$

with  $k = O\left(\frac{\log n}{\varepsilon^2}\right)$  s.t.

$$(1-\varepsilon)\|x_i - x_j\|^2 \leq \|Sx_i - Sx_j\|^2 \leq (1+\varepsilon)\|x_i - x_j\|^2$$

$$v = [v_1, \dots, v_d]$$

$$\|v\|^2 = \sum_{i=1}^d v_i^2$$

Obs:

①.  $k$  is indep. of  $D$

②. Log factor!

E.g. 10 billion points

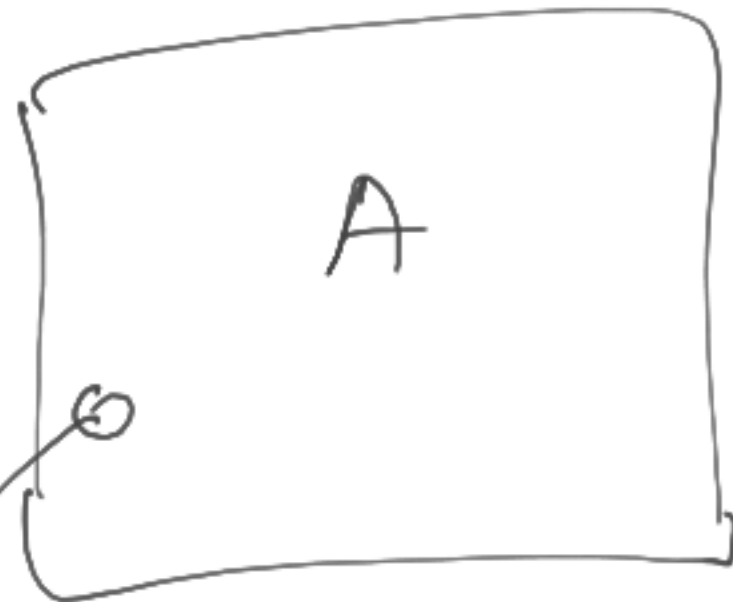
$\sim 10^{10}$  pts

$\sim n \sim 2^{30}$

$\Rightarrow k = O\left(\frac{30}{\varepsilon^2}\right)$

Construction:

Matrix  $A$



each  
entry  
of  
the matrix

$K \times D$

Transformation  
matrix:

$$S = \frac{1}{\sqrt{K}} A$$

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

(Standard Normal /  
Gaussian distribution)

Gaussian R.V.s:

Cont. R.V.

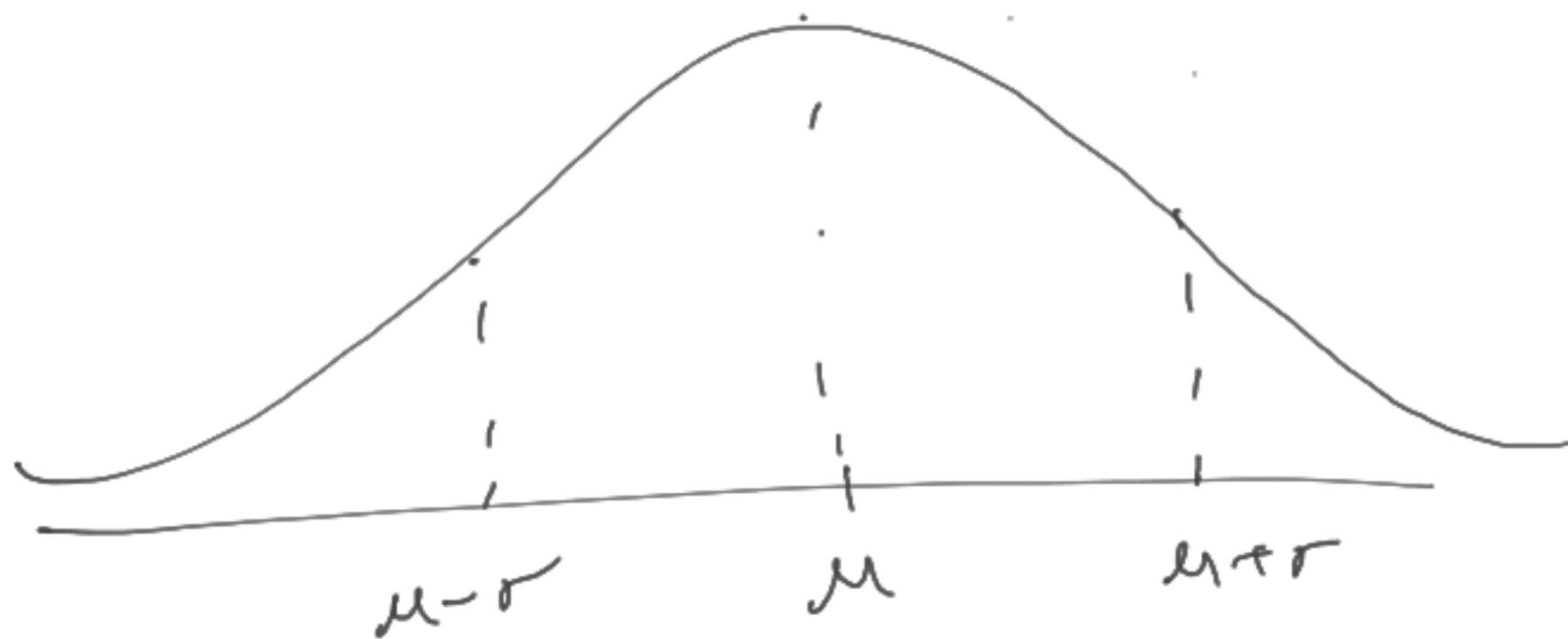
Notation

pdf:

$$\mathcal{N}(\mu, \sigma^2)$$

↑ mean                      ↑ variance

$$f(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



- exponential tapering off the mean

$$e^{-\frac{y^2}{2\sigma^2}}$$

properties:

1.  $X_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$  & indep.

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

2.  $X \sim \mathcal{N}(\mu, \sigma^2) \Rightarrow aX \sim \mathcal{N}(\underline{a\mu}, \underline{a^2\sigma^2})$

3.  $X_i \sim \mathcal{N}(\underline{0}, \underline{1}) \rightarrow$  std. normal

$$\sum_{i=1}^d a_i X_i \sim \mathcal{N}\left(0, \sum_{i=1}^d a_i^2\right)$$

$$= \mathcal{N}\left(0, \|a\|^2\right)$$

---

$$a = (a_1, \dots, a_d)$$

... JL Transform

Fact 1  
 $x \in \mathbb{R}^D$

$A =$   
 $K \times D$

$$\begin{bmatrix} - A_1 - \\ - A_2 - \\ \vdots \\ - A_K - \end{bmatrix}$$

$\longleftarrow D$

Transformations

$$Sx = \frac{1}{\sqrt{K}} A \cdot x$$

$$= \begin{bmatrix} \frac{1}{\sqrt{K}} \langle A_1, x \rangle \\ \vdots \\ \frac{1}{\sqrt{K}} \langle A_K, x \rangle \end{bmatrix}$$

$\longleftarrow K$

$\longleftarrow$  given vector

$$\sum A_{i,j} x_j$$

$\nwarrow$  scalar

$\nwarrow$   $\mathcal{N}(0,1)$

$$s_i = \begin{bmatrix} \mathcal{N}\left(0, \frac{\|x\|^2}{K}\right) \\ \vdots \\ \mathcal{N}\left(0, \frac{\|x\|^2}{K}\right) \end{bmatrix}$$

$\longleftarrow K$

Recall:

$$x \rightarrow \underline{Sx}$$

$$\downarrow$$

$$\frac{1}{\sqrt{K}} Ax$$

$$\uparrow$$

iid  $\mathcal{N}(0,1)$  entries

Fact 2

For any  $Y \in \mathbb{R}^D$

$$E[\|SY\|^2]$$

$$= E\left[\sum_{i=1}^K \frac{1}{K} \langle A_i, Y \rangle^2\right]$$

$$= \left(\sum_{i=1}^K \frac{1}{K}\right) E[\langle AX, Y \rangle^2]$$

$$= \|Y\|^2$$

$$\Rightarrow E[\|SY\|^2] = \|Y\|^2$$

$$\Rightarrow E[\|Sx_1 - Sx_2\|^2] = \|x_1 - x_2\|^2 \Rightarrow$$

result holds in expectation.

$$Sx_i - Sx_j$$

↓

$$S(x_i - x_j)$$

$$SY$$

Variance

if mean = 0,

$$E[X^2] = \text{variance}$$



Lemma (Conc. bound for squared - Gaussians)

Let  $U_1, \dots, U_k$  be i.i.d  $\mathcal{N}(0, \sigma^2)$

$$\text{Let } Z = \sum_{i=1}^k U_i^2$$

$$E[Z] = k\sigma^2$$

Then for  $\varepsilon \in (0, 1/2)$ , some constant  $c$

$$P(|Z - E[Z]| \geq \varepsilon E[Z]) \leq e^{-\frac{k\varepsilon^2}{c}}$$

Conc. bound

$$P(|Z - E[Z]| \geq \alpha) \leq \underline{\hspace{2cm}}$$

... JL Transform

As before  $Y = X_1 - X_2$

$$P \left[ \left| \|SY\|^2 - E[\|SY\|^2] \right| \geq \epsilon \|Y\|^2 \right] \leq ?$$

LHS

$$\|SY\|^2 = \sum_{i=1}^k v_i^2 \quad \text{where } v_i \sim \mathcal{N}\left(0, \frac{\|Y\|^2}{k}\right)$$

From Cone. bound Lemma

$$\text{L.H.S.} \leq e^{-\frac{K\epsilon^2}{c}}$$

Choose  $K = \frac{2c \log n}{\epsilon^2}$

want

$$\leq \frac{1}{n^2}$$

all pairwise

dist =

$$\binom{n}{2}$$

$$\sim n^2$$

to use union bound over all pair-wise distances.

By union bound, over all  $\binom{n}{2}$  pair-wise distances.  
then squared lengths maintained w.p. at least

$$1 - \binom{n}{2} \frac{1}{n^2} \geq \frac{1}{2}$$

$\rightarrow$  want  $\frac{1}{n}$

make this  $\frac{1}{n^3}$  gives

$\hookrightarrow$  Choose  $k = 3c \dots$

## Discussions:

1. Runtime.

2. Applications: ① fast (but approximate) mat. mul.

② Compressive Sensing

$$y = A x \quad (Sparse)$$

E.g. MRI imaging, cameras