#### Lecture 9

# Dimensionality Reduction: Johnson-Lindenstrauss Lemma

Course: Algorithms for Big Data

Instructor: Hossein Jowhari

Department of Computer Science and Statistics Faculty of Mathematics K. N. Toosi University of Technology

Spring 2021

Distance-Preserving Dimensionality Reduction

Given n vectors  $A = \{ \boldsymbol{x}_1, \dots, \boldsymbol{x}_n \}$  in  $\mathbb{R}^d$ ,

we want a fast transformation  $f: \mathbb{R}^d \rightarrow \mathbb{R}^t$  so that

• (Distances are approximately preserved):

For all pairs  $\boldsymbol{x}, \boldsymbol{y} \in A$ , we have  $\|f(\boldsymbol{x}) - f(\boldsymbol{y})\| \approx \|\boldsymbol{x} - \boldsymbol{y}\|$ 

• (The dimension is reduced considerably)

 $t \ll d$ 

Such a transformation will be very useful in practice (when the input dimension is large)

## Two applications

Faster clustering algorithms: As an example, recall that for the k-center problem we had a 2-factor approximation algorithm with running time O(nkd) (when the points lie in  $\mathbb{R}^d$ .)

Using the transformation  $f:\mathbb{R}^d\to\mathbb{R}^t,$  we first compute f(x) for all  $x\in A$ 

Then we run the k-center alg. on  $A' = \{f(x_1), \ldots, f(x_n)\}$ 

Running time =  $\underbrace{\text{time of transformation}}_{O(nkt)}$  +  $\underbrace{\text{time of clustering } A'}_{O(nkt)}$ 

Approximation quality  $\approx 2$ 

Approximate nearest neighbor queries: Most exact nearest neighbor data structures have time complexity  $n^{O(d)}$ 

#### How to reduce the dimension?

 A bad idea: Randomly select a small subset of dimensions. Restrict every *x* to the selected dimensions. In other words, choose S ⊆ {1,...,d} randomly. f(*x*) = *x*<sub>S</sub>

Example:  

$$\mathbf{x}_1 = \overbrace{(1,0,\ldots,0)}^d, \mathbf{x}_2 = \overbrace{(0,1,0,\ldots,0)}^d, ||x - y|| = 1$$
  
 $S = \{3, 5, 20, 25\}$   
 $f(\mathbf{x}_1) = (0,0,0,0), f(\mathbf{x}_2) = (0,0,0,0), ||f(x) - f(y)|| = 0$ 

 Dimensionality reduction methods such as PCA do not preserve the pairwise distances.

## Johnson-Lindenstrauss Lemma

JL Lemma (existential formulation): Let  $\epsilon \in (0, \frac{1}{2})$ . Given a set of n vectors  $A = \{x_1, \ldots, x_n\} \in \mathbb{R}^d$ , there is a mapping  $f : \mathbb{R}^d \to \mathbb{R}^t$  where  $t = O(\frac{\log n}{\epsilon^2})$  where

$$\forall x, y \in A, \ (1-\epsilon) \| \boldsymbol{x} - \boldsymbol{y} \| \le \| f(\boldsymbol{x}) - f(\boldsymbol{y}) \| \le (1+\epsilon) \| \boldsymbol{x} - \boldsymbol{y} \|$$

JL lemma is essentially tight with respect to the target dimension t:

Noga Alon has shown that any such mapping requires  $t = \Omega(\frac{\log n}{\epsilon^2 \log \frac{1}{\epsilon}})$ 

Larsen and Nelson have shown that any linear mapping requires  $t = \Omega\bigl(\frac{\log n}{\epsilon^2}\bigr)$ 

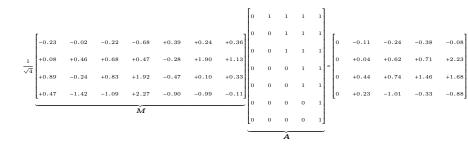
Construction of the mapping f: Let M be a  $t \times d$  matrix where every entry  $M_{ij}$  is an independent random sample from the normal standard distribution N(0,1). In other words, each  $M_{ij} \sim N(0,1)$ 

$$\frac{1}{\sqrt{t}} \underbrace{\begin{bmatrix} M_{11} & \dots & M_{1d} \\ M_{21} & \dots & M_{2d} \\ \vdots & \vdots & \vdots \\ M_{t1} & \dots & M_{td} \end{bmatrix}}_{M} \underbrace{\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \dots & \vdots \\ x_{d1} & x_{d2} & \dots & x_{dn} \end{bmatrix}}_{A} = \begin{bmatrix} \vdots & \vdots & \vdots \\ f(\boldsymbol{x}_1) & \vdots & f(\boldsymbol{x}_n) \\ \vdots & \vdots & \vdots \end{bmatrix}$$

Let  $M^{(i)}$  be the *i*-th row of M.

$$f(\boldsymbol{x}) = \frac{1}{\sqrt{t}} \boldsymbol{M} \boldsymbol{x} = (\frac{1}{\sqrt{t}} \boldsymbol{M}^{(1)} \boldsymbol{.} \boldsymbol{x}, \ldots, \frac{1}{\sqrt{t}} \boldsymbol{M}^{(t)} \boldsymbol{.} \boldsymbol{x})$$

An example: n = 5, d = 7, t = 4



Lemma 1: If  $t \ge \frac{c \log n}{\epsilon^2}$  for a large enough constant c then with probability at least 3/4 for all pairs  $\boldsymbol{x}, \boldsymbol{y} \in A$ , we have

$$\|f(\boldsymbol{x}) - f(\boldsymbol{y})\| \in [(1-\epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|, (1+\epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|]$$

$$\|f(\boldsymbol{x}) - f(\boldsymbol{y})\| \approx_{\epsilon} \|\boldsymbol{x} - \boldsymbol{y}\|$$

Lets consider a special case:

 $A = \{\mathbf{0}, \mathbf{x}\}$  where  $\mathbf{x}$  is a unit vector in  $\mathbb{R}^n$ .  $\|\mathbf{x}\| = 1$ .

We want  $||f(\boldsymbol{x}) - f(\boldsymbol{0})|| = ||f(\boldsymbol{x})|| \approx_{\epsilon} 1.$ 

Lemma 2: Given  $\boldsymbol{x} \in \mathbb{R}^d$  where  $\|\boldsymbol{x}\| = 1$ , assuming  $t \ge \frac{c \log(\frac{1}{\delta})}{\epsilon^2}$ when c is a large enough constant then we have  $Pr(\|f(\boldsymbol{x})\|^2 \in [1-\epsilon, 1+\epsilon]) \ge 1-\delta$ 

Before proving Lemma 2, we show Lemma 1 is a consequence of Lemma 2.

Observation 1: Since f(x) - f(y) = f(x - y) (the mapping f is linear) then it is enough to show that for any arbitrary vector  $z \in \mathbb{R}^d$  we have

$$\|f(oldsymbol{z})\|pprox_\epsilon\|oldsymbol{z}\|$$

Observation 2: Let  $z' = \frac{z}{\|z\|}$ . The vector z' is a unit vector. If we have  $\|f(z')\| \approx_{\epsilon} 1$  then (by linearity of f) we have

$$||f(z)|| = ||f(||z||z')|| = |||z||f(z')|| = ||z|||f(z')|| \approx_{\epsilon} ||z||$$

Observation 3: There are  $\binom{n}{2}$  pair of vectors in A. From Lemma 2 and the above observations we have for a pair  $\boldsymbol{x}, \boldsymbol{y} \in A$ , if  $t \ge \frac{c \log(\frac{1}{\delta})}{\epsilon^2}$  then  $\|f(\boldsymbol{x}) - f(\boldsymbol{y})\| \approx_{\epsilon} \|\boldsymbol{x} - \boldsymbol{y}\|$  with probability  $1 - \delta$ .

Setting  $\delta = \frac{1}{4n^2}$ , from the union bound, the statement is true for all pairs in A with probability at least  $1 - \binom{n}{2} \frac{1}{4n^2} > 3/4$ .

Therefore we get the statement of Lemma 1.

## Proof of Lemma 2

Basic facts regarding Gaussian distribution:

The Gaussian distribution N(μ, σ<sup>2</sup>) with mean μ and variance σ<sup>2</sup> has the following probability density function

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

• If 
$$X \sim N(\mu_1, \sigma_1^2)$$
 and  $Y \sim N(\mu_2, \sigma_2^2)$  then

$$cX \sim N(c\mu_1, c^2\sigma_1^2)$$

$$X + Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

Consider  $\boldsymbol{x} \in \mathbb{R}^d$  where  $\|\boldsymbol{x}\| = 1$ .

Let  $Y_i$  be the *i*-th oordinate of Mx. Note that  $f(x) = \frac{1}{\sqrt{t}}Mx$ .

**Observation:**  $Y_i = (G_1, \ldots, G_d) \cdot \boldsymbol{x}$  where each  $G_i$  is an independent sample from N(0, 1). In other words,

$$Y = G_1 x_1 + \ldots + G_d x_d$$

 $Y_i$  is a linear combination of independent Gaussians. Therefore

$$Y_i \sim N(0, x_1^2 + \ldots + x_d^2) = N(0, 1)$$

We need to analyze  $Y_i^2$  since  $Y = ||f(x)||^2 = \frac{1}{t}(Y_1^2 + ... + Y_t^2)$ 

$$E[Y_i^2] = Var[Y_i] + E^2[Y_i] = 1 \implies E[Y] = 1$$

So, in expectation, Y = ||f(x)|| is exactly 1. Very good but not enough.

We need to bound the probability  $Pr(Y > 1 + \epsilon)$ .

 $Y = \frac{1}{t} \sum_{i}^{t} Y_{i}^{2}$  is the sum of independent random variables. We could use Chernoff but unfortunately  $Y_{i}$  is not bounded. Still similar ideas that were used in the proof of Chernoff are helpful here. For any r > 0

$$Pr(Y > 1+\epsilon) = Pr(e^{trY} > e^{tr(1+\epsilon)}) \le \underbrace{\frac{E[e^{trY}]}{e^{tr(1+\epsilon)}}}_{E[e^{rY_i^2}]} = \prod_{i=1}^t \frac{E[e^{rY_i^2}]}{e^{r(1+\epsilon)}}$$
$$E[e^{rY_i^2}] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{ry^2} e^{-\frac{y^2}{2}} dy = \frac{1}{\sqrt{1-2r}} \quad \text{when } r \le \frac{1}{2}$$

Therefore for  $0 < r < \frac{1}{2}$ , we have

$$Pr(Y > 1 + \epsilon) \le \left(\frac{1}{e^{r(1+\epsilon)}\sqrt{1-2r}}\right)^t$$

Also one can show that 
$$\frac{1}{e^{r(1+\epsilon)}\sqrt{1-2r}} \leq e^{\frac{r^2}{1-2r}}$$

Therefore we have

$$Pr(Y > 1 + \epsilon) \le e^{\frac{tr^2}{1 - 2r}}$$

We set 
$$r = \frac{\epsilon}{4}$$
. Using  $1 - 2r \ge \frac{1}{2}$  when  $\epsilon \le \frac{1}{2}$ , we get  
 $Pr(Y > 1 + \epsilon) \le e^{-\frac{t\epsilon^2}{8}} \le \frac{\delta}{2} \implies t \ge \frac{8}{\epsilon^2} \ln(\frac{2}{\delta})$ 

Similarly we can show

$$Pr(Y < 1 - \epsilon) \le e^{-\frac{t\epsilon^2}{8}} \le \frac{\delta}{2}$$

Therefore having  $t \geq \frac{8}{\epsilon^2} \ln(\frac{2}{\delta})$  we get

$$Pr(1 - \epsilon \le Y \le 1 + \epsilon) \ge 1 - \delta$$

Good news: JL lemma still holds when the Gaussian distribution N(0,1) is replaced with random -1,+1 coefficients.