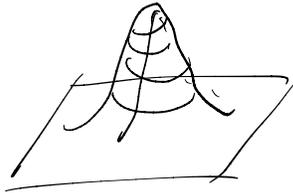


LECTURE 15

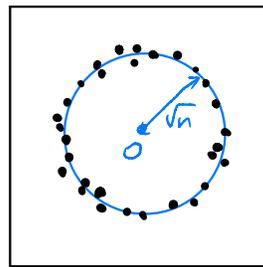
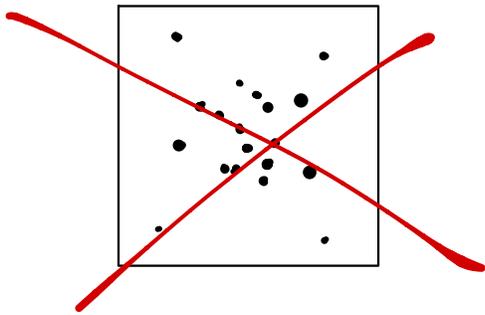
THE THIN SHELL PHENOMENON

- Gaussian random vector $g \sim N(0, I_n)$: $g = (g_1, \dots, g_n)$, $g_i \sim N(0, 1)$ iid
- pdf: $f(x) = f(x_1) \dots f(x_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-x_i^2/2} = \frac{1}{(2\pi)^{n/2}} e^{-\|x\|_2^2/2}$

Rotation invariant :



- Sample from $N(0, I_n)$ for large n :



Thm (Thin shell)

$$\text{For } g \sim N(0, I_n) : \mathbb{P}\{0.99\sqrt{n} \leq \|g\|_2 \leq 1.01\sqrt{n}\} \geq 1 - 2e^{-cn}$$

↑
exponentially close to 1!

Proof :

$$\|g\|_2^2 - n = \sum_{i=1}^n (g_i^2 - 1)$$

↑
iid, mean zero, subexponential :

$$\begin{aligned} \|g_i^2 - 1\|_{\psi_1} &\leq \|g_i^2\|_{\psi_1} + 1 \quad (\Delta \text{ inequality}) \\ &= \|g_i\|_{\psi_2}^2 + 1 \leq C \quad (\text{abs. constant}) \end{aligned}$$

Bernstein's inequality \Rightarrow

$$\mathbb{P}\{\|g\|_2^2 - n \geq \underbrace{0.01n}_t\} \leq 2 \exp\left(-c \cdot \min\left(\frac{n^2}{\sigma^2}, \frac{n}{K}\right)\right) \quad (\leq)$$

where $\sigma^2 = \sum_{i=1}^n \|g_i^2 - 1\|_{\Psi_1}^2 \leq Cn$; $K = \max_i \|g_i^2 - 1\|_{\Psi_1} \leq C$.

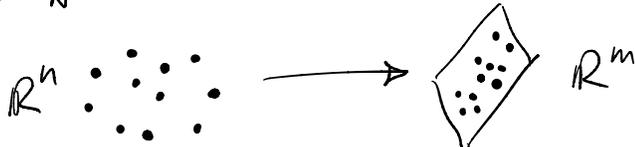
$$\leq 2\exp(-c'n).$$

\Rightarrow with prob. $\geq 1 - 2\exp(-c'n)$,

$$0.99n \leq \|g\|_2^2 \leq 1.01n. \quad \text{QED.}$$

APPLICATION: DIMENSION REDUCTION

• Data: $x_1, \dots, x_N \in \mathbb{R}^m \xrightarrow{?} \mathbb{R}^n, m \ll n ?$



Data compression to save on storage, speed.

• Possible with $n = O(\log N)$; geometry of data preserved pairwise distances.

THM (Johnson-Lindenstrauss lemma '1984) $\forall x_1, \dots, x_N \in \mathbb{R}^d$
 \exists linear map $T: \mathbb{R}^d \rightarrow \mathbb{R}^n$ such that $n \leq C \log N$ and
 $0.99 \|x_i - x_j\| \leq \|T(x_i) - T(x_j)\|_2 \leq 1.01 \|x_i - x_j\| \quad \forall i, j = 1, \dots, N.$

Proof. Probabilistic Method: choose random map $n \times \begin{matrix} d \\ G \end{matrix} \begin{matrix} \\ \\ \\ \end{matrix} = \begin{matrix} \\ \\ \\ \end{matrix} n$
 $G := n \times d$ Gaussian random matrix $G_{ij} \sim N(0, 1)$ iid.

Claim: \forall fixed $z \in \mathbb{R}^d, \|z\|_2 = 1$: $Gz \sim N(0, I_n)$

$$\left[(Gz)_i = \sum_{j=1}^d G_{ij} z_j \sim N(0, \sum_{j=1}^d z_j^2) = N(0, 1) \text{ indep.} \right]$$

$\underbrace{\hspace{2cm}}_{N(0, z_j^2) \text{ indep.}}$

Thin Shell Thm $\Rightarrow \mathbb{P}\{0.99\sqrt{n} \leq \|Gz\|_2 \leq 1.01\sqrt{n}\} \geq 0.99 \quad (*)$

Proof of JL ① Fix (i, j) , use (*) for $z = \frac{x_i - x_j}{\|x_i - x_j\|_2} \Rightarrow$

$$P \left\{ 0.99\sqrt{n} \leq \frac{\|G(x_i - x_j)\|_2}{\|x_i - x_j\|_2} \leq 1.01\sqrt{n} \right\} \geq 1 - 2e^{-cn}$$

\Rightarrow for $T := \frac{1}{\sqrt{n}}G$,

$$P \left\{ 0.99 \|x_i - x_j\|_2 \leq \|Tx_i - Tx_j\| \leq 1.01 \|x_i - x_j\|_2 \right\} \geq 1 - 2e^{-cn}$$

② Union Bound: $\exists N^2$ pairs $(i, j) \Rightarrow$

$$P \left\{ \forall i, j: E_{ij} \text{ holds} \right\} \geq 1 - N^2 \cdot 2e^{-cn} = 1 - 2\exp(2\log N - cn) \quad (\geq)$$

Choose n s.t. $\boxed{cn = 4\log N} \Rightarrow$

$$(\geq) 1 - 2\exp(-2\log N) = 1 - \frac{2}{N^2} > 0 \quad \text{if } N \geq 2.$$

(and if $N=1$, the thm is trivial).

\Rightarrow such T exists.

QED.

REMARKS 1. Why does not JL lift the "curse of high dimensionality"?

Preprocess HD data by JL \Rightarrow low D data.

2. Fast JL transforms. [Ailon-Chazelle '2009]

3. JL without dependence on $N = \# \text{data}$:

Prop 9.32 of Book [Liu-Mehrabian-Plan-V' 2017]

HW: JL for Ber
[Achlioptas]