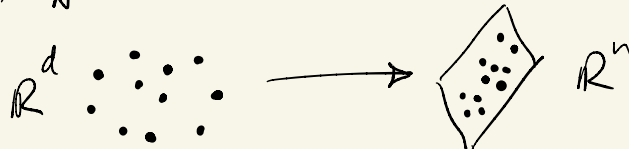


# LECTURE 12

## DIMENSION REDUCTION

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



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  $n = C \log N$  and

$$0.99 \|x_i - x_j\|_2 \leq \|T(x_i) - T(x_j)\|_2 \leq 1.01 \|x_i - x_j\|_2 \quad \forall i, j = 1, \dots, N.$$

Proof. ↑ T = random projection. Probabilistic Method: choose a random linear map  $n \times d$   $G$   $\int_z = \int_n$

$G := n \times d$  Gaussian random matrix  $G_{ij} \sim N(0, 1)$  iid.

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

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

$\underbrace{\hspace{10em}}_{N(0, z_j^2) \text{ indep.}}$

Thin Shell Thm  $\Rightarrow$

$$\mathbb{P}\{0.99\sqrt{n} \leq \|Gz\|_2 \leq 1.01\sqrt{n}\} \geq 1 - 2e^{-cn} \quad (*)$$

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

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

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.

Q.E.D.

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]

# COMBINATORIAL OPTIMIZATION

- Problem A: Given  $(a_i)_{i=1}^n$ , find  $\max_{x_i \in \{\pm 1\}} \sum_{i=1}^n a_i x_i$

Although exhaustive search is hard ( $2^n$  configurations of  $(x_i)$ ),  
 $\exists$  direct solution:  $x_i = \text{sign}(a_i)$

- Problem B: Given  $(a_{ij})_{i,j=1}^n$ , find

$$\max_{x_i \in \{\pm 1\}} \sum_{i,j=1}^n a_{ij} x_i x_j \quad (*)$$

NP-hard. Integer quadratic program.

## EXAMPLES:

- ① Ising model of magnetism:

$x_i = \text{spin of atom } i \ (\pm 1)$   
 $a_{ij} = \text{strength of interaction between atoms } i, j$   
 e.g.  $a_{ij} = \begin{cases} 1 & \text{if } i, j \text{ are neighbors} \\ 0 & \text{if not} \end{cases}$

$$H(x) = \sum_{i,j} a_{ij} x_i x_j = \text{"Kauzmannian"}$$

of the system (free energy)

- Probability that the system is in state  $x = (x_1, \dots, x_n) \in \{\pm 1\}^n$  is

$$P(x) := \frac{1}{Z} \exp\left(-\frac{H(x)}{T}\right)$$

"Gibbs measure"

Here  $T = \text{"temperature"}$ ,  $Z = \sum_{x \in \{\pm 1\}^n} \exp(-H(x)/T)$  is a normalizing const ("partition function")

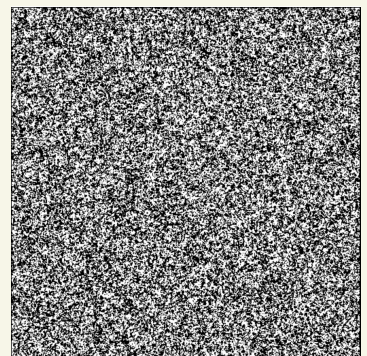
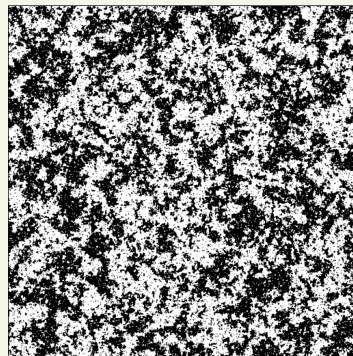
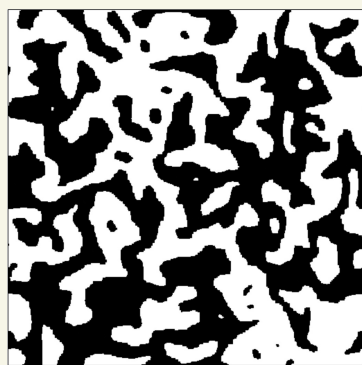
- $T \rightarrow \infty \Rightarrow P(x) \rightarrow \text{uniform on } \{\pm 1\}^n$  (chaos)

$T \rightarrow 0 \Rightarrow P(x) \rightarrow \begin{cases} 1 & \text{for } x = x_0 \\ 0 & \text{elsewhere} \end{cases}$  where

lowest-energy state

$$H(x_0) = \min_{x \in \{\pm 1\}^n} H(x) = ?$$

equivalent to Problem B.



→ Increasing the temperature T →