

JOHNSON-LINDENSTRAUSS THEOREM

* MAKE NOTE OF CONNECTION BETWEEN FAST KNN AND APPROX.

METHODS BASED ON SEARCH OVER LOW-DIMENSIONAL

* KNN CONNECT ISOMAP / LLE W/ CLUSTERING

GIVEN N VECTORS IN \mathbb{R}^d (d, N LARGE) WE CAN DEFINE

THE RANDOM PROJECTION

$$\begin{aligned} f_k: \mathbb{R}^d &\longrightarrow \mathbb{R}^k \\ \vec{x} &\longmapsto \begin{bmatrix} \vec{u}_1^T \vec{x} \\ \vdots \\ \vec{u}_k^T \vec{x} \end{bmatrix} = \mathbf{U}^T \vec{x} \end{aligned} \quad \mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_k]$$

WHERE THE COORDINATES OF $\vec{u}_1, \dots, \vec{u}_k$ ARE INDEPENDENT $N(0, 1)$

RANDOM VARIABLES. WHAT HAPPENS TO DISTANCES $\|\vec{x}_i - \vec{x}_j\|$

AFTER THE PROJECTION $\|f(\vec{x}_i) - f(\vec{x}_j)\|$?

JOHNSON-LINDENSTRAUSS LEMMA

DERIVED LATER

FOR ANY $0 < \epsilon < 1$ AND ANY INTEGER N , LET $k \geq \frac{3}{c\epsilon^2} \log N$. ($c = 1/96$)

FOR ANY SET OF N POINTS $\vec{x}_1, \dots, \vec{x}_N$ IN \mathbb{R}^d , THE RANDOM PROJECTION $f_k: \mathbb{R}^d \rightarrow \mathbb{R}^k$ (ABOVE)

SATISFIES

$$P \left((1-\epsilon)\sqrt{k} \|\vec{x}_i - \vec{x}_j\| \leq \|f(\vec{x}_i) - f(\vec{x}_j)\| \leq (1+\epsilon)\sqrt{k} \|\vec{x}_i - \vec{x}_j\|, 1 \leq i < j \leq N \right) \geq 1 - \frac{3}{2N}$$

\approx

$$P \left(\left| \frac{\frac{1}{k} \|f(\vec{x}_i) - f(\vec{x}_j)\|^2 - \|\vec{x}_i - \vec{x}_j\|^2}{\|\vec{x}_i - \vec{x}_j\|^2} \right| \leq 2\epsilon, \forall 1 \leq i < j \leq N \right)$$

NOTE:

- 1) THE THEOREM HOLDS FOR ALL x_i, x_j PAIRS NOT JUST MOST OF THEM
- 2) FOR K -NEAREST NEIGHBOR SCHEMES CAN PROJECT TO $k \propto \frac{\log d}{\epsilon^2}$ - DIMENSIONAL SPACE AND FIND THE CORRECT NEIGHBORS W/ HIGH PROBABILITY

- MAKE NO ASSUMPTION ON THE DISTRIBUTION OF $\vec{x}_1, \dots, \vec{x}_n$

- DOMINANT TERM IS $1/\epsilon^2$ IN APPLICATIONS,

- $\frac{1}{\sqrt{k}}$ IS UNIFORM SCALING OF ALL PROJECTED DISTANCES

- ϵ IS MAXIMUM TOLERABLE RELATIVE ERROR BETWEEN

$$\underbrace{\sqrt{k} \|\vec{x}_i - \vec{x}_j\|}_{\text{RESCALED ORIGINAL EUCLIDEAN DISTANCES}} \approx \underbrace{\|f(\vec{x}_i) - f(\vec{x}_j)\|}_{\text{PROJECTED DISTANCES}}$$

3) OTHER VERSIONS OF THEM

- 1) REQUIRE THE COMPONENTS OF $\vec{u}_1, \dots, \vec{u}_k$ ARE UNCORRELATED VARIANCE 1, MEAN 0 RVs
 - GAUSSIAN ONES PROOF

- 2) PROVE THE EXISTENCE OF MAP f WHICH IS LIPSCHITZ (AND CAN BE FOUND IN POLYNOMIAL TIME) W/ SLIGHTLY BETTER GUARANTEES ON ERROR.

OUR GOAL

PROVE THE LEMMA, USING PROPERTIES AND BOUNDS ON THE THE IID $N(0,1)$ COMPONENTS OF $\vec{u}_1, \dots, \vec{u}_k$!

INTUITION: IF k IS STILL LARGE $\|\vec{u}_j\|^2 \approx k$ $\left. \begin{array}{l} \vec{u}_i \cdot \vec{u}_j \approx k \delta_{ij} \end{array} \right\} \text{ BY LLN}$

PROBABILITY REVIEW

RECALL,

MARKOV'S INEQUALITY

Let Z be a random variable w/ $Z \geq 0$ and $EZ = \mu_Z < \infty$

$$P(Z \geq a) \leq \frac{EZ}{a}$$

$$\left(P(|Z - EZ| \geq a) \leq \frac{E|Z - EZ|}{a} \leq \frac{2EZ}{a} \right)$$

Chebyshev's Inequality

Let Z be a random variable w/ $EZ = \mu_Z < \infty$

$$\text{Var}(Z) = \sigma_Z^2 < \infty$$

$$P(|Z - \mu_Z| \geq a) = \frac{\sigma_Z^2}{a^2}$$

Others: BERNSTEIN, BENNET, Hoeffding

Let Z_1, \dots, Z_d be IID random variables drawn from the same distribution as Z w/ mean 0 and variance σ^2

$$P(|Z_1 + \dots + Z_d| \geq a) \leq \frac{d\sigma^2}{a^2} \quad (\text{CHEBYSHEV})$$

BUT WE'LL NEED TO DO BETTER.

MASTER TAIL BOUND (BLUM)

LET Z_1, \dots, Z_d BE MUTUALLY IND. RANDOM VARIABLES W/ MEAN 0
AND VARIANCE AT MOST σ^2 . SUPPOSE $a \in [0, \sqrt{2} d \sigma^2]$ AND
 $s \leq d \sigma^2 / 2$ IS A POSITIVE EVEN INTEGER AND $|E(Z_i^r)| \leq \sigma^2 r!$ FOR
 $r=2, 3, \dots, s$ THEN

$$P(|Z_1 + \dots + Z_d| \geq a) \leq \left(\frac{2sd\sigma^2}{a^2} \right)^{s/2}.$$

FURTHER, IF $s \geq a^2 / (4d\sigma^2)$, THEN WE ALSO HAVE

$$s \in \left[\frac{a^2}{4d\sigma^2}, \frac{d\sigma^2}{2} \right]$$

$$P(|Z_1 + \dots + Z_d| \geq a) \leq 3e^{-\frac{a^2}{12n\sigma^2}}$$

PROOF:

• FIRST FOR $a \geq 0$ AND EVEN INTEGERS r

$$P(|Z_1 + \dots + Z_d| \geq a) = P((Z_1 + \dots + Z_d)^r \geq a^r) \leq \frac{E((Z_1 + \dots + Z_d)^r)}{a^r}$$

SO WE NEED TO WORK ON BOUNDS OF $E(|Z_1 + \dots + Z_d|^r)$!

$$\begin{aligned} (Z_1 + \dots + Z_d)^r &= \sum_{r_1 + \dots + r_d = r} \binom{r}{r_1, \dots, r_d} Z_1^{r_1} \dots Z_d^{r_d} \\ &= \sum \frac{r!}{r_1! \dots r_d!} Z_1^{r_1} \dots Z_d^{r_d} \end{aligned}$$

BY INDEPENDENCE,

$$E((Z_1 + \dots + Z_d)^r) = \sum \left(\frac{r!}{r_1! \dots r_d!} \right) E(Z_1^{r_1}) \dots E(Z_d^{r_d})$$

If $r_i = 1$ THEN $E z_i^{r_i} = E z_i = 0$ BY ASSUMPTION. So

WE FOCUS ON (r_1, \dots, r_d) RUNS OVER SETS W/

$r_i = 0$ OR $r_i \geq 2$ $i=1, \dots, d$

* THERE ARE ONLY $\frac{d}{2}$ SUCH SETS!

• SINCE $E(z_i^{r_i}) \leq r_i! \sigma^2$

$$E((z_1 + \dots + z_d)^r) \leq \sum_{r_1 + \dots + r_d = r} \frac{r!}{r_1! \dots r_d!} |E(z_1^{r_1})| \dots |E(z_d^{r_d})|$$

$$\leq r! \sum_{r_1 + \dots + r_d = r} \sigma^{2(\# \text{ NONZERO } r_i)}$$

$$\leq r! \sum_{t=1}^{d/2} \binom{d}{t} \binom{r-t-1}{t-1} \sigma^{2t}$$

OF SUBSETS OF $\{1, \dots, d\}$ W/ CARDINALITY t
OF WAYS TO ALLOCATE r TO t NONZERO r_i 'S W/ $r_i \geq 2$!

$$\binom{d}{t} = \frac{d(d-1) \dots (d-t+1)}{t!} \leq \frac{d^t}{t!}$$

$$\binom{r-t-1}{t-1} \leq 2^{r-t-1} \quad (\text{TRIVIAL BOUNDS ON BINOMIAL COEFFICIENTS})$$

$$\Rightarrow E[(z_1 + \dots + z_d)^r] \leq r! \sum_{t=1}^{d/2} h(t) \quad , \quad h(t) = \frac{d^t}{t!} 2^{r-t-1} \sigma^{2t}$$

$$\frac{h(t)}{h(t-1)} = \frac{d^t 2^{r-t-1}}{t!} \sigma^{2t} \bigg/ \left(\frac{d^{t-1} 2^{r-t}}{(t-1)!} \sigma^{2(t-1)} \right) = \frac{d \sigma^2}{2t} \quad t \leq \frac{r}{2} \left(\leq \frac{d \sigma^2}{4} \right)$$

$$\geq 2$$

$$\Rightarrow E[(z_1 + \dots + z_d)^r] \leq r! h(r/2) \left(1 + \frac{1}{2} + \dots \right) = \frac{r!}{(r/2)!} 2^{r/2} (d \sigma^2)^{r/2} \quad \sqrt{5}$$

B1 MARKOV'S INEQUALITY

$$\begin{aligned} P(|z_1 + \dots + z_d| \geq a) &= P((z_1 + \dots + z_d)^2 \geq a^2) \leq \frac{E[(z_1 + \dots + z_d)^2]}{a^2} \leq \frac{r! (2d\sigma^2)^{r/2}}{(r/2)! a^r} = g(r) \\ &\leq \left(\frac{2rd\sigma^2}{a^2} \right)^{r/2}, \quad \frac{r!}{(r/2)!} \leq r^{r/2} \end{aligned}$$

WHICH HELDS FOR $r \leq S$, r EVEN, AND APPLYING IT TO $r=S$ GIVES THE FIRST INEQUALITY!

$$\begin{aligned} \text{FOR EVEN } r, \quad g(r) &= \frac{r!}{(r/2)!} \left(\frac{2d\sigma^2}{a^2} \right)^{r/2} \\ \frac{g(r)}{g(r-2)} &= \frac{r! (2d\sigma^2)^{r/2}}{(r/2)! a^r} \cdot \frac{(r-2)! a^{r-2}}{(r-2)! (2d\sigma^2)^{(r-2)/2}} = \frac{r(r-1) (2d\sigma^2)}{\frac{r}{2}} \\ &= \frac{4(r-1)d\sigma^2}{a^2} \leq 1 \quad \left(\Leftrightarrow r \leq \frac{a^2}{4d\sigma^2} + 1 \right) \end{aligned}$$

$$\Rightarrow g(r) \text{ IS DECREASING SO LONG AS } r-1 \leq \frac{a^2}{4d\sigma^2}$$

IF $\left\lceil \frac{a^2}{4d\sigma^2} \right\rceil \geq r-1$ IS EVEN, WE MUST TAKE

$$r^* = \text{LARGEST EVEN INTEGER w/ } \frac{a^2}{6d\sigma^2} - 2 \leq r^* \leq \frac{a^2}{6d\sigma^2}$$

$$\begin{aligned} \Rightarrow P(z_1 + \dots + z_d \geq a) &\leq \left(\frac{2r^* d\sigma^2}{a^2} \right)^{r^*/2} \\ &\leq \left(\frac{1}{3} \right)^{r^*/2} \leq \left(\frac{1}{e} \right)^{r^*/2} = e^{-r^*/2} \\ &\leq e^{-\left(\frac{a^2}{6d\sigma^2} - 2 \right)/2} = e^{-\frac{a^2}{12d\sigma^2}} \\ &\leq \boxed{3 e^{-\frac{a^2}{12d\sigma^2}}} \end{aligned}$$

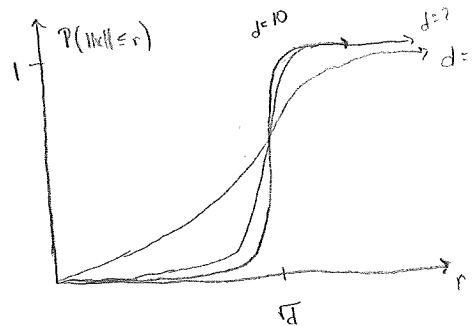
$$\frac{a^2}{6d\sigma^2} - 2 \leq r^*$$

□

• Consider $\vec{x} \sim \mathcal{N}(\vec{0}_d, I_d)$

$$\begin{aligned}
 P(\|\vec{x}\| \leq r) &= \int_{\|\vec{x}\| \leq r} \frac{1}{(2\pi)^{d/2}} e^{-\frac{\|\vec{x}\|^2}{2\pi}} d\vec{x} \\
 &= \frac{2\pi^{d/2}}{\Gamma(\frac{d}{2})} \cdot \frac{1}{(2\pi)^{d/2}} \int_0^r y^{d-1} e^{-y^2/2} dy \\
 &\quad \uparrow \text{SPHERICAL COORDINATES} \\
 &= \frac{2^{\frac{-d}{2}+1}}{\Gamma(d/2)} \int_0^r y^{d-1} e^{-y^2/2} dy
 \end{aligned}$$

ALL PROBABILITY MASS
LIES NEAR SURFACE OF
UNIT SPHERE w/ RADIUS \sqrt{d}



GAUSSIAN ANNULUS THM

For $\vec{u} \sim \mathcal{N}(\vec{0}_d, I_d)$ and $\beta \leq \sqrt{d}$

$$P(\sqrt{d} - \beta \leq \|\vec{u}\| \leq \sqrt{d} + \beta) \geq 1 - e^{-c\beta^2}$$

FOR SOME POSITIVE CONSTANT c .

PROOF:

$$\begin{aligned}
 P(\|\vec{u}\| \notin [\sqrt{d} - \beta, \sqrt{d} + \beta]) &\leq P(|\|\vec{u}\|^2 - d| \geq \beta\sqrt{d}) \\
 &= P\left(\left|(u_1^2 - 1) + \dots + (u_d^2 - 1)\right| \geq \beta\sqrt{d}\right) \\
 &= P\left(\left|\underbrace{\frac{u_1^2 - 1}{2}}_{z_1} + \dots + \underbrace{\frac{u_d^2 - 1}{2}}_{z_d}\right| \geq \frac{\beta\sqrt{d}}{2}\right)
 \end{aligned}$$

OBSERVE

$$u_i \sim N(0,1) \Rightarrow E\left(\frac{u_i^2-1}{2}\right) = 0$$

$$\begin{aligned} \text{Var}(z_i) &= \text{Var}\left(\frac{u_i^2-1}{2}\right) = \frac{1}{4} \text{Var}(u_i^2) \leq 2 \\ E(z_i)^r &\leq 2 r! \quad r=1,2,\dots \end{aligned} \left\{ \begin{array}{l} \text{USING PROPERTIES} \\ \text{OF MOMENTS} \\ \text{OF GAUSSIANS} \end{array} \right.$$

$$\Rightarrow \text{BY MARKOV TAIL BOUND w/ } \sigma^2=2$$

$$P\left(\left|\frac{u_1^2-1}{2} + \dots + \frac{u_d^2-1}{2}\right| \geq \frac{\beta \sqrt{d}}{2}\right) \leq 3 e^{-\frac{(\beta \sqrt{d}/2)^2}{\text{Var}(z)}} = 3 e^{-\frac{\beta^2}{46}}$$

$$\Rightarrow P\left(\sqrt{d}-\beta \leq \|u\| \leq \sqrt{d}+\beta\right) \geq 1 - 3 e^{-\frac{\beta^2}{46}}$$

□

RANDOM PROJECTION THM

LET $\vec{x} \in \mathbb{R}^d$ AND f BE THE RANDOM PROJECTION TO \mathbb{R}^k

$$\begin{aligned} P\left(\left|\|f(\vec{x})\| - \sqrt{k} \|\vec{x}\|\right| \geq \varepsilon \sqrt{k} \|\vec{x}\|\right) &\leq 3 e^{-\frac{k \varepsilon^2}{96}} \\ P\left(\|f(\vec{x})\| \notin \left[(1-\varepsilon)\sqrt{k} \|\vec{x}\|, (1+\varepsilon)\sqrt{k} \|\vec{x}\|\right]\right) &\leq 6 e^{-\frac{k \varepsilon^2}{96}} \end{aligned}$$

PROOF: WE MAY ASSUME $\|\vec{x}\|=1$ ABOVE.

$$f(\vec{x}) = \begin{bmatrix} \vec{u}_1^T \\ \vdots \\ \vec{u}_k^T \end{bmatrix} \vec{x} \quad \vec{u}_1, \dots, \vec{u}_k \stackrel{i.i.d.}{\sim} N(0, I_d)$$

$$E(\vec{u}_i^T \vec{x}) = E[u_i^T] \vec{x} = 0$$

$$\text{Var}(\vec{u}_i^T \vec{x}) \stackrel{i.i.d.}{=} \sum_{j=1}^d \text{Var}(u_{ij} x_j) = \sum_{j=1}^d x_j^2 \text{Var}(u_{ij}) = \sum_{j=1}^d x_j^2 = \|\vec{x}\|^2 = 1$$

$$\Rightarrow f(\vec{x}) \sim N(0, I_k) \Rightarrow P\left(\left|\|f(\vec{x})\| - \sqrt{k} \|\vec{x}\|\right| \geq \varepsilon \sqrt{k} \|\vec{x}\|\right) \leq 3 e^{-\frac{k \varepsilon^2}{96}}$$

□

PROOF (JOHANSSON - LINDENSTRAUSS)

$$P \left((1-\epsilon) \sqrt{k} \|\vec{x}_i - \vec{x}_j\| \leq \underbrace{\|f(\vec{x}_i) - f(\vec{x}_j)\|}_{f(\vec{x}_i - \vec{x}_j)} \leq (1+\epsilon) \sqrt{k} \|\vec{x}_i - \vec{x}_j\|, 1 \leq i < j \leq N \right)$$

$$= 1 - P \left(\|f(\vec{x}_i) - f(\vec{x}_j)\| \notin \left[(1-\epsilon) \sqrt{k} \|\vec{x}_i - \vec{x}_j\|, (1+\epsilon) \sqrt{k} \|\vec{x}_i - \vec{x}_j\| \right] \right. \\ \left. \text{FOR SOME } i, j \right)$$

$$\geq 1 - \sum_{i=1}^{N-1} \sum_{j=i+1}^N P \left(\|f(\vec{x}_i) - f(\vec{x}_j)\| \notin \left[(1-\epsilon) \sqrt{k} \|\vec{x}_i - \vec{x}_j\|, (1+\epsilon) \sqrt{k} \|\vec{x}_i - \vec{x}_j\| \right] \right)$$

$$\geq 1 - \sum_{i=1}^{N-1} \sum_{j=i+1}^N 3e^{-\frac{k\epsilon^2}{96}} = 1 - 3 \binom{N}{2} e^{-\frac{k\epsilon^2}{96}}$$

$$\geq 1 - 3 \binom{N}{2} e^{-\frac{3 \cdot 96 \epsilon^2}{96 \epsilon^2} \log N} \quad k \geq \frac{3 \cdot 96}{\epsilon^2} \log N$$

$$= 1 - 3 \binom{N}{2} \left(\frac{1}{N^3} \right) \geq 1 - 3 \left(\frac{N^2}{2} \right) \left(\frac{1}{N^3} \right) = 1 - \frac{3}{2N}$$

□