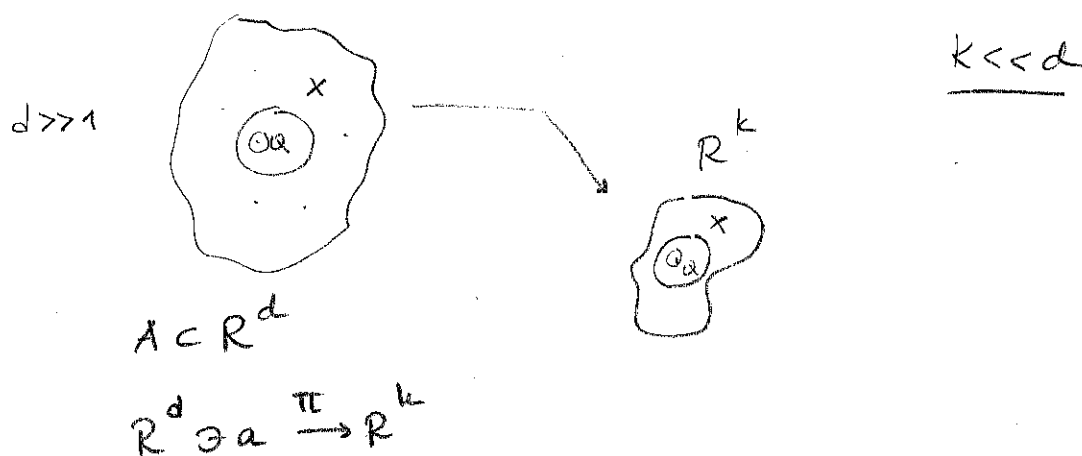


14.02.2020

1

Johnson - Lindenstrauss Theorem



Isometry: L_2 -norm

Given a $\varepsilon \in (0, 1)$ a map $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$ is called an

ε -isometry if $\forall a, a' \in A$

$$(1 - \varepsilon) \|a - a'\|^2 \leq \|f(a) - f(a')\|^2 \leq (1 + \varepsilon) \|a - a'\|^2$$

$$\begin{cases} \ell(c_1 a + c_2 b) = c_1 \ell(a) + c_2 \ell(b) \\ \pi(v_1 + v_2) = \pi(v_1) + \pi(v_2) \end{cases}$$

* Which k ?

- choose vectors randomly.

- concentration estimates: that my points mainly concentrate around the mean.

If I assign a gaussian measure on my entire space and

choose

- I take a big space, in each coordinate I assign gaussian measure

* unit variance > gaussian measure
* zero mean

- each.

iid. X_{ij} $i=1, \dots, k$ and $j=1, \dots, d$.

2

$$E[X_{ij}] = 0 \quad \text{Var}(X_{ij}) = 1 \quad \text{definition.}$$

Vector $v \in \mathbb{R}^d$, $\pi_i(v) \equiv \sum_{j=1}^d v_j X_{ij}$ $i=1, \dots, k$
 \parallel
 (v_1, v_2, \dots, v_d)

$$\pi(v) \equiv (\pi_1(v), \dots, \pi_k(v))$$

$$E[\|\pi(v)\|^2] = E\left[\sum_{i=1}^k \|\pi_i(v)\|^2\right] = E\left[\sum_{i=1}^k \sum_{j=1}^d v_j^2 X_{ij}^2\right] = k \|v\|^2$$

$$= \sum_{i=1}^k \|\pi_i(v)\|^2 = \sum_{i=1}^k \sum_{j=1}^d v_j^2 X_{ij}^2 + 2 \sum_{i=1}^k \sum_{\substack{\dim=1 \\ e \neq m}} v_e v_m X_{ie} X_{im}$$

physical coordinates $= 0$

- Distribution -

which gives the concentration around the mean \Rightarrow Gaussian.

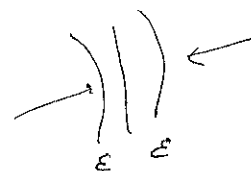
$$\pi_i(v) = \frac{1}{\sqrt{k}} \sum_{j=1}^d v_j X_{ij} \quad i=1, \dots, k$$

$$v_1 - v_2$$

$$E \|\pi(v)\|^2 = \frac{1}{k} E \left[\sum_{i=1}^k \|\pi_i(v)\|^2 \right]$$

$$= \frac{1}{k} E \left[\sum_{i=1}^k \|\pi_i(v)\|^2 \right] = \frac{1}{k} E \left[\sum_{i=1}^k \sum_{j=1}^d v_j^2 X_{ij}^2 \right] = \|v\|^2$$

$v_1 - v_2$



Gaussian (Centered) — Same machinery can be used.

3

$$M_X(t) = e^{t^2/2}, \text{ Sub-gaussian } M_X(t) \leq e^{t^2/2}$$

* Gaussian Annulus Theorem:

For a d -dim'l. spherical Gaussian with unit variance in each direction
for $\beta \leq \sqrt{d}$, all but at most $3e^{-\beta^2}$ of the prob. mass is in.

the annulus $\sqrt{d} - \beta \leq \|x\| \leq \sqrt{d} + \beta$, c is a fixed true number

$$Y = \|X\| \quad \sqrt{d} - \beta \leq Y \leq \sqrt{d} + \beta \quad |Y - \sqrt{d}| \geq \beta$$

$$|Y^2 - d| \geq \beta(Y + \sqrt{d}) \geq \beta\sqrt{d}$$

$$Y_i = X_i^2 - 1, \quad Y^2 - d = Y_1 + Y_2 + \dots + Y_d$$

$$|Y_1 + Y_2 + \dots + Y_d| \geq \beta\sqrt{d} \quad E[Y_i] = 0$$

$$\begin{aligned} \text{For } |X_i| \leq 1 \quad |Y_i|^s &\leq 1 \\ |X_i| \geq 1 \quad |Y_i|^s &\leq |X_i|^{2s} \end{aligned} \Rightarrow |E[Y_i^s]| = E[|Y_i|^s] \leq E[1 + |X_i|^{2s}]$$

$$= 1 + \sqrt{\frac{2}{\pi}} \int_0^\infty x^{2s} e^{-x^2/2} dx$$

$$E[|Y_i|^s] \leq 2^s \cdot s!$$

$$E(Y_i) = 0$$

$$\text{Var}(Y_i) \leq 2^2 \cdot 2! = 8$$

$$W_i = \frac{Y_i}{2}, \quad \text{Var}(W_i) \leq 2, \quad E[W_i^s] \leq 2^s!$$

$$|w_1, w_2, \dots, w_d| \geq \frac{\beta \sqrt{d}}{2}$$

$$\text{Var}(w) = \sigma^2 = 2 \quad n=d$$

$$3e^{-\frac{\beta^2 d}{|x|^2 2}} \longrightarrow ?$$

$$|\sqrt{d}|, E[|X|^2] = d.$$

* The order ^{tail estimates.} will not change.

$$\pi: \mathbb{R}^d \rightarrow \mathbb{R}^k \quad u_1, \dots, u_k \in \mathbb{R}^d$$

$$\pi(u) = (u_1 u, \dots, u_k u) \quad u_i = (X_{i1}, X_{i2}, \dots, X_{id})$$

Theorem: (Random Projection Theorem) (R.P. Theorem)

Fix $v \in \mathbb{R}^d$. Define π as above $\exists c > 0$ s.t.

$$\forall \epsilon \in (0, 1), \mathbb{P}[|\pi(v)| - \sqrt{k} \|v\| \geq \epsilon \sqrt{k} \|v\|] \leq 3e^{-c k \epsilon^2}$$

— $\pi(v)$ is a random spherical Gaussian with unit variance $k=d$.

(?) Johnson-Lindenstrauss Theorem:

For $0 < \epsilon < 1$ and any integer $d \gg 1$ with $k \gg \frac{3}{c \epsilon^2} \ln N$.

N is the number of points

for $\pi: \mathbb{R}^d \mapsto \mathbb{R}^k$ w.p. $1 - \frac{3}{2N}$

$$(1-\epsilon) \sqrt{k} \|v_i - v_j\| \leq \|\pi(v_i) - \pi(v_j)\| \leq (1+\epsilon) \sqrt{k} \|v_i - v_j\| \quad \forall v_i, v_j \in \mathbb{R}^d$$

\hookleftarrow distance between query and data is smaller than ϵ in lower-dim'l space.

$$P[| \Pi(v_i - v_j) | \notin [(1-\epsilon)\sqrt{k} |v_i - v_j|, (1+\epsilon)\sqrt{k} |v_i - v_j|]] \leq 3e^{-ck\epsilon^2}$$

we have N -points
 $\binom{N}{2}$ pairs

$$\binom{N}{2} O(N^2) \rightarrow \frac{3}{2} N^2 \cdot \frac{1}{N^3} = \frac{3}{2N}$$

$$\frac{3}{N^3} \cdot \underbrace{\left(\frac{N^2}{2}\right)}_{\sim \binom{N}{2}}$$

$$e^{-ck\epsilon^2 \cdot \frac{3}{c\epsilon^2} \ln N}$$

$$= e^{\ln N^{-3}}$$

$$= \frac{1}{N^3}$$

- query & data lie in any dimension. (even in ∞ -dim.)
- ∞ -dimension version of Euclidean norm. = Hilbert's space
- Square-integrable function. $\rightarrow ?$