

Random Projections

MDS d_{ij} $i, j = 1, \dots, n$
 $Y_i \in \mathbb{R}^k$

$$\min_{Y \in \mathbb{R}^{kn}} \sum_{i,j} (\|Y_i - Y_j\|^2 - d_{ij}^2)^2$$

Total square error
 "Average" distortion

$$d_{ij}(1-\varepsilon) \leq \|Y_i - Y_j\| \leq d_{ij}(1+\varepsilon) \quad \forall i, j$$

uniform distortion

① SDP $\|Y_i - Y_j\|^2 \rightarrow$ SD. Relaxation

② Universal basis 依赖数据

Random Projection

$$X_i \in \mathbb{R}^p, \quad d_{ij} = \|X_i - X_j\| \quad i, j = 1, \dots, n.$$

$$Y_i = f(X_i) \in \mathbb{R}^k \quad k = O(\varepsilon^{-2} \log n)$$

s.t. $1-\varepsilon \leq \frac{\|Y_i - Y_j\|}{\|X_i - X_j\|} \leq 1+\varepsilon$ with probability

f : random proj.

$$\geq 1 - n^{-\alpha}, \quad \alpha > 0$$

\rightarrow

Almost Isometry Embedding!

Johnson - Lindenstrauss Lemma

John, Lindenstrauss Copyright by 极书

2001 Computer Science 2003

Lipschitz Extension

概率性方法证明存在性

Approximate NN.

Sanjoy Dasgupta, Anupam Gupta

Dimitris Achlioptas '03

Thm Given $\epsilon > 0$, n, ∞ let.

$$k = C(\alpha, \epsilon) \log n = (4 + 2\alpha) \left(\frac{\epsilon^2}{2} - \frac{\epsilon^3}{8} \right)^{-1} \log n.$$

Then for any n points $x_i \in \mathbb{R}^d$ ($i=1, \dots, n$), there exists a map

$f: \mathbb{R}^d \rightarrow \mathbb{R}^k$ such that $\forall x_i, x_j$

$$1 - \epsilon \leq \frac{\|f(x_i) - f(x_j)\|^2}{\|x_i - x_j\|^2} \leq 1 + \epsilon. \quad (*)$$

f can be founded in randomized polynomial time. $(*)$ holds with probability at least $1 - n^{-\alpha}$, $\alpha > 0$.

Ex: f can be constructed by random matrices / projections

$$f(x_i) = R x_i$$

$$X = [x_1 \dots x_n]^{d \times n}$$

$$\Rightarrow Y = R X$$

$$Y = [y_1 \dots y_n]^{k \times n}$$

• $R = [r_1 \dots r_k]^T$ $r_i \in S^{d-1}$, eg. $r_1 = \frac{(a_1^i \dots a_d^i)}{\|a^i\|}$ $a_d^i \sim N(0,1)$

• $R = A / \sqrt{k}$ $A_{ij} \sim N(0,1)$, Variance = 1

• $R = A / \sqrt{k}$ $A_{ij} = \begin{cases} 1 & p = \frac{1}{2} \\ -1 & (1-p) = \frac{1}{2} \end{cases}$, $A_{ij} = \begin{cases} 1 & p = \frac{1}{6} \\ -1 & p = \frac{5}{6} \end{cases}$

Human Immigration

640k

640k.

5k.

100k

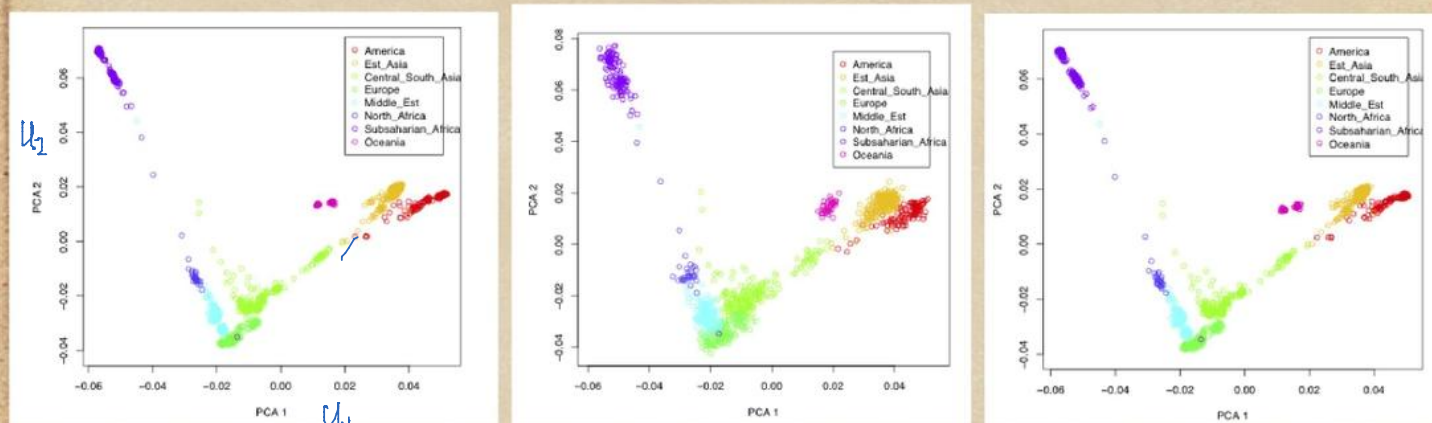


FIGURE 1. (Left) Projection of 1043 individuals on the top 2 MDS principal components. (Middle) MDS computed from 5,000 random projections. (Right) MDS computed from 100,000 random projections. Pictures are due to Qing Wang.

14年10月7日星期二

SNPs. 分析 (Single Nucleotide Polymorphisms)

600k.

Human Genome Diversity Project

n=1064 individual

p=644,258 SNPs

 $X^{n \times p}$

0: AA 1: AC 2: CC

9: Missing value.

-21 h.

 $n \times p$

n=1043, p=644,258

X

SVD.

$$\tilde{X} = HX$$

$$H = (I - \frac{1}{n} 44^T)$$

$$\tilde{X} = U \Sigma V^T \quad \underline{p \times q}$$

 $\|x\|_1, \dots, \|x\|_k$

PCA.

n=1043

\mathbb{R}^d uniformly projected to k -subspace

$\xleftrightarrow{\text{dist}}$ random vector on S^{d-1} restricted onto top k -coord.

$$\underbrace{(a_1^i \dots a_k^i)}_k \quad \|a^i\|_2 = 1$$

$$(\underbrace{a_1^i \dots a_k^i}_k, 0)$$

$$x_i \sim \mathcal{N}(0, 1), \quad i=1, \dots, d.$$

$$X = (x_1, \dots, x_d) \quad Y = \frac{X}{\|X\|} \in S^{d-1} \quad \text{uniformly distributed.}$$

$$Z = (x_1, \dots, x_k, 0) \in \mathbb{R}^d \quad k\text{-subspace}$$

$$L = \|Z\|^2 \quad \mathbb{E}[L] = \frac{k}{d} = \mu$$

Lemma Concentration Inequality $k < d$.

(a) $\beta < 1$ lower bound.

$$\text{Prob}[L \leq \beta \mu] \leq \beta^{k/2} \left(1 + \frac{(1-\beta)k}{d-k} \right)^{d-k/2} \leq \exp\left(\frac{k}{2}(1-\beta + \ln \beta)\right)$$

(b) $\beta > 1$ upper bound

$$\text{Prob}[L \geq \beta \mu] \leq \beta^{k/2} \left(1 + \frac{(1-\beta)k}{d-k} \right)^{d-k/2} \leq \exp\left(\frac{k}{2}(1-\beta + \ln \beta)\right)$$

$$\mathbb{E}L = \mu. \quad \text{exponentially} \sim \mu \quad \downarrow$$

Proof of JL Lemma

$d < k$, trivial

$d > k$.

k -subspace $X_i \in \mathbb{R}^d \rightarrow Y_i \in \mathbb{R}^k$ $i=1, \dots, n$

$$L = \sum_{i,j} \|Y_i - Y_j\|^2 \quad \mu = \mathbb{E} L = \frac{k}{d} \sum_{i,j} \|X_i - X_j\|^2 \quad Y_i = (x_1^i, \dots, x_k^i, 0, \dots, 0)$$

$$\text{Prob} [L \leq (1-\varepsilon)\mu] \leq \exp\left(\frac{k}{2} (1 - (1-\varepsilon) + \ln(1-\varepsilon))\right)$$

\uparrow

$$\leq \exp\left[\frac{k}{2} \left(\varepsilon - \left(\varepsilon + \frac{\varepsilon^2}{2}\right)\right)\right] = \exp\left(-\frac{k\varepsilon^2}{4}\right)$$

$$\leq \exp(-(k+\alpha)\ln n)$$

$$= n^{-(k+\alpha)} \rightarrow 0$$

$$k \geq 4(1+\alpha/2)(\varepsilon^2/2)^{-1} \ln n$$

$$\text{Prob} [L \geq (1+\varepsilon)\mu] \leq \exp\left(\frac{k}{2} (1 - (1+\varepsilon) + \ln(1+\varepsilon))\right)$$

$$\leq \exp\left[\frac{k}{2} (\varepsilon^2/2 - \varepsilon^3/3)\right],$$

$$\ln(1+\varepsilon) \leq \varepsilon - \frac{\varepsilon^2}{2} + \frac{\varepsilon^3}{3}, \text{ so}$$

$$\leq \exp(-(k+\alpha)\ln n) = n^{-(k+\alpha)} \rightarrow 0$$

$$k \geq 4(1+\alpha/2)(\varepsilon^2/2 - \varepsilon^3/3)^{-1} \ln n$$

Given pair (i, j) .

$$\text{Prob} \left\{ 1-\varepsilon \leq \frac{\|Y_i - Y_j\|^2}{\|X_i - X_j\|^2} \leq 1+\varepsilon \right\} \leq 1 - \frac{2}{n^{k+\alpha}}$$

$\forall (i, j)$ from $\{i=1, \dots, n; j=1, \dots, n\}$ of $\binom{n}{2}$

$$\text{Prob} \left\{ \forall (i, j), 1-\varepsilon \leq \frac{\|Y_i - Y_j\|^2}{\|X_i - X_j\|^2} \leq 1+\varepsilon \right\} \leq 1 - \binom{n}{2} \frac{2}{n^{k+\alpha}} \leq 1 - n^{-\alpha}$$

$$\binom{n}{2} = n(n-1)/2$$