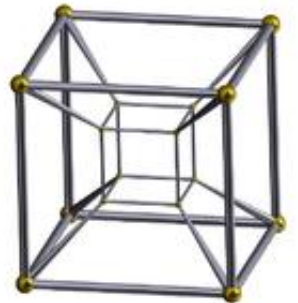
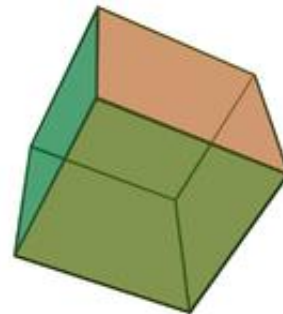
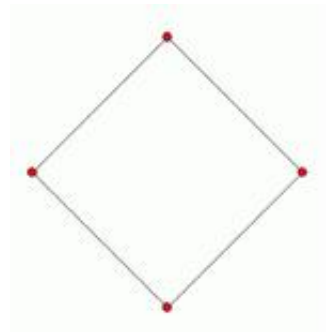


High Dimensional Space

longhuan@sjtu.edu.cn



Word Vector Model

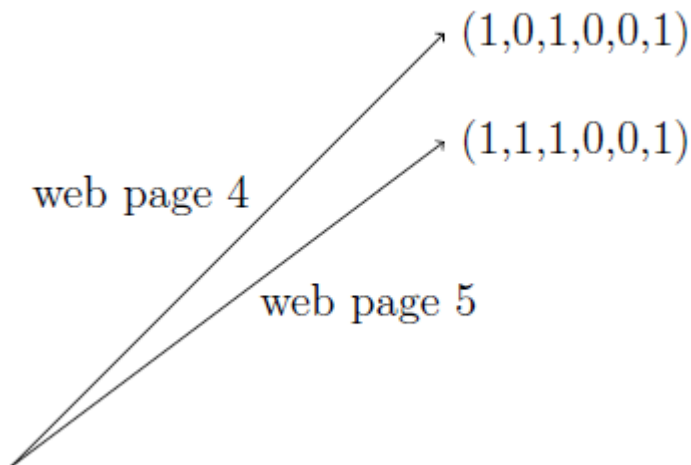
稀疏



抽象成 高维向量

高维矩阵

Web Page Model



- Nearest neighbor query
- Information retrieval 检索
- Web page rank
- Online recommendation
-

The law of Large numbers

Properties of High-Dimensional space,
unit ball

Generating points uniformly at random
from a ball

Gaussians in High Dimension

Random Projection and Johnson-
Lindenstrauss Lemma

Seperating Gaussians

Normal distribution (Gauss Distribution)

$X \sim N(\mu, \sigma^2)$, with density function:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad -\infty < x < +\infty$$

Variance $Var(X) = E((X - E[X])^2)$

$$\begin{aligned} &= E(X^2 + E[X]^2 - 2XE[X]) \\ &= E(X^2 - E[X]^2) \\ &= E[X^2] - E[X]^2 \end{aligned}$$

Chebyshev's Inequality

$$\forall a > 0, \Pr(|X - E(X)| \geq a) \leq \frac{Var[X]}{a^2}$$

Law of Large Numbers

- In probability theory, the **law of large numbers (LLN)** is a theorem that describes the result of performing the same experiment a large number of times.
- According to the law, **the average of the results obtained from a large number of trials should be close to the expected value**, and will tend to become closer as more trials are performed.

Law of large numbers

Let x_1, x_2, \dots, x_n be n independent samples of a random variable x , then 独立同分布

$$\Pr\left(\left|\frac{x_1 + x_2 + \dots + x_n}{n} - E(x)\right| \geq \epsilon\right) \leq \frac{\text{Var}(x)}{n\epsilon^2}$$

Proof. (Chebychev's Inequality)

$$\begin{aligned}\Pr\left(\left|\frac{x_1 + x_2 + \dots + x_n}{n} - E(x)\right| \geq \epsilon\right) &\leq \frac{\text{Var}\left(\frac{x_1 + x_2 + \dots + x_n}{n}\right)}{\epsilon^2} \\ &= \frac{\text{Var}(x_1 + x_2 + \dots + x_n)}{n^2 \epsilon^2} \\ &= \frac{\text{Var}(x)}{n \epsilon^2}\end{aligned}$$

Application

- \mathbf{x} be a d –dimensional random point whose coordinates are each selected from $N\left(0, \frac{1}{2\pi}\right)$, 独立同分布
- i.e. $\mathbf{x} = [x_1, x_2, \dots, x_d]$ with $x_i \sim N\left(0, \frac{1}{2\pi}\right)$
- By LLN: $|\mathbf{x}|^2 = \sum_{i=1}^d x_i^2 = \frac{d}{2\pi} = \Theta(d)$ with high probability. \rightarrow 向量长度 与 1 的概率接近
- The probability that point \mathbf{x} lie in the unit ball is *vanishingly small*.

Application

- $\mathbf{x}, \mathbf{y} : [z_1, z_2, \dots, z_d]$ with $z_i \sim N(0, 1)$
- $|\mathbf{x}|^2 \approx d, |\mathbf{y}|^2 \approx d,$
- $|\mathbf{x} - \mathbf{y}|^2 \approx ?$

Application

- $\mathbf{x}, \mathbf{y} : [z_1, z_2, \dots, z_d]$ with $z_i \sim N(0, 1)$
- $|\mathbf{x}|^2 \approx d, |\mathbf{y}|^2 \approx d,$
- $|\mathbf{x} - \mathbf{y}|^2 = \sum_{i=1}^d (x_i - y_i)^2$
$$E(x_i - y_i)^2 = E(x_i^2) + E(y_i^2) - 2E(x_i y_i)$$
$$= 1 + 1 - 2E(x_i)E(y_i) = 2.$$

Application

- $\mathbf{x}, \mathbf{y} : [z_1, z_2, \dots, z_d]$ with $z_i \sim N(0, 1)$
- $|\mathbf{x}|^2 \approx d, |\mathbf{y}|^2 \approx d,$
- $|\mathbf{x} - \mathbf{y}|^2 = \sum_{i=1}^d (x_i - y_i)^2 = 2d$
$$E(x_i - y_i)^2 = E(x_i^2) + E(y_i^2) - 2E(x_i y_i)$$
$$= 1 + 1 - 2E(x_i)E(y_i) = 2.$$
- $|\mathbf{x} - \mathbf{y}|^2 \approx |\mathbf{x}|^2 + |\mathbf{y}|^2$
- **Pythagorean theorem** \Rightarrow random d –dimensional \mathbf{x}, \mathbf{y} are **approximately orthogonal**.

x, y 正交, 内积为 0

Application

- $\mathbf{x}, \mathbf{y} : [z_1, z_2, \dots, z_d]$ with $z_i \sim N(0, 1)$
- **Pythagorean theorem** \Rightarrow random d –dimensional \mathbf{x}, \mathbf{y} are **approximately orthogonal**.

- If we scale these random points to be unit length and call \mathbf{x} the **North Pole**, *much of the surface area of the unit ball must lie near the equator.*

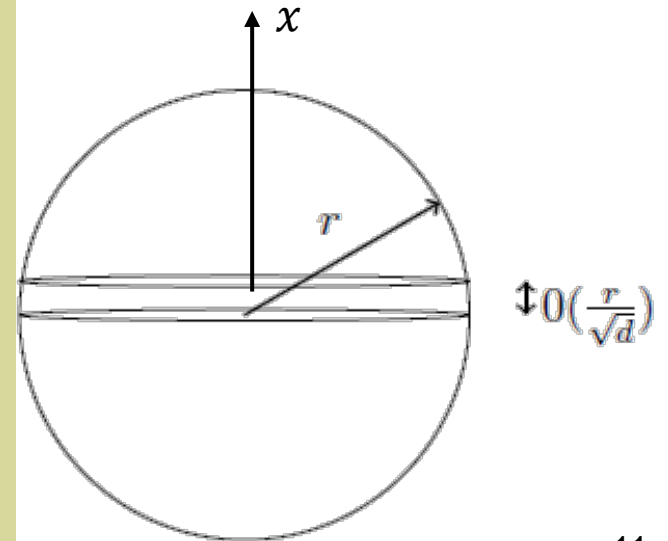


Table of tail bounds

	Condition	Tail bound
Markov	$x \geq 0$	$\text{Prob}(x \geq a) \leq \frac{E(x)}{a}$
Chebychev	Any x	$\text{Prob}(x - E(x) \geq a) \leq \frac{\text{Var}(x)}{a^2}$
Chernoff	$x = x_1 + x_2 + \dots + x_n$ $x_i \in [0, 1]$ i.i.d. Bernoulli;	$\text{Prob}(x - E(x) \geq \varepsilon E(x)) \leq 3e^{-c\varepsilon^2 E(x)}$
Higher Moments	r positive even integer	$\text{Prob}(x \geq a) \leq E(x^r)/a^r$
Gaussian Annulus	$x = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$ $x_i \sim N(0, 1); \beta \leq \sqrt{n}$ indep.	$\text{Prob}(x - \sqrt{n} \geq \beta) \leq 3e^{-c\beta^2}$
Power Law for x_i ; order $k \geq 4$	$x = x_1 + x_2 + \dots + x_n$ x_i i.i.d ; $\varepsilon \leq 1/k^2$	$\text{Prob}(x - E(x) \geq \varepsilon E(x)) \leq (4/\varepsilon^2 kn)^{(k-3)/2}$

有时 r 有较小的界...

The law of Large numbers

Properties of High-Dimensional space,
unit ball

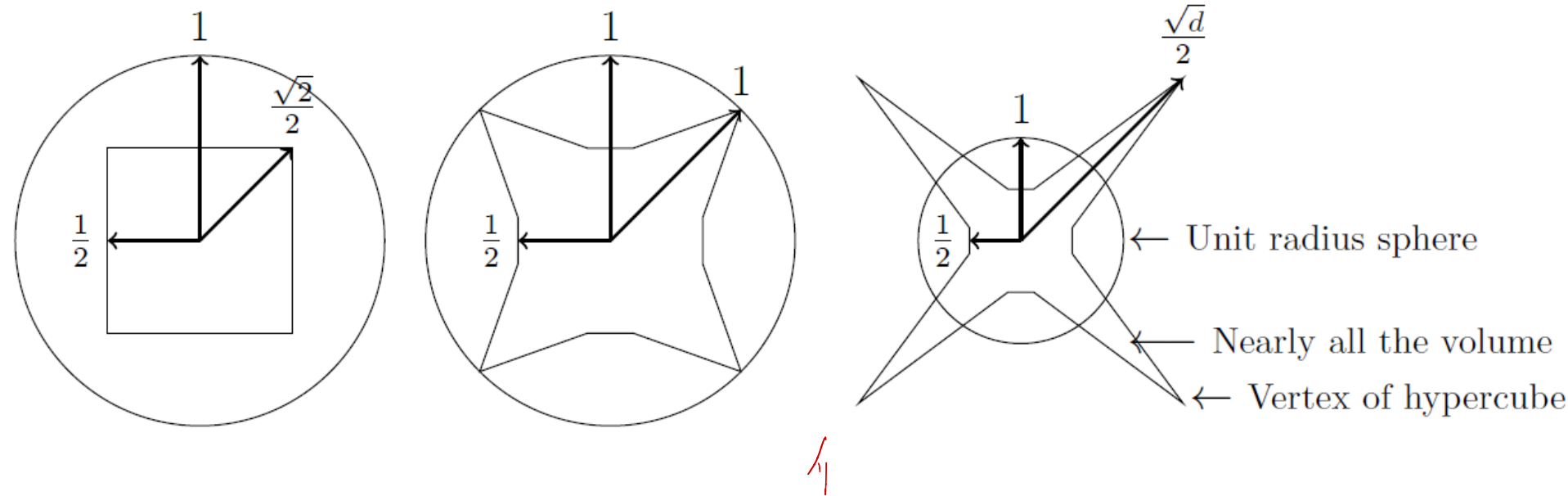
Generating points uniformly at random
from a ball

Gaussians in High Dimension

Random Projection and Johnson-
Lindenstrauss Lemma

Seperating Gaussians

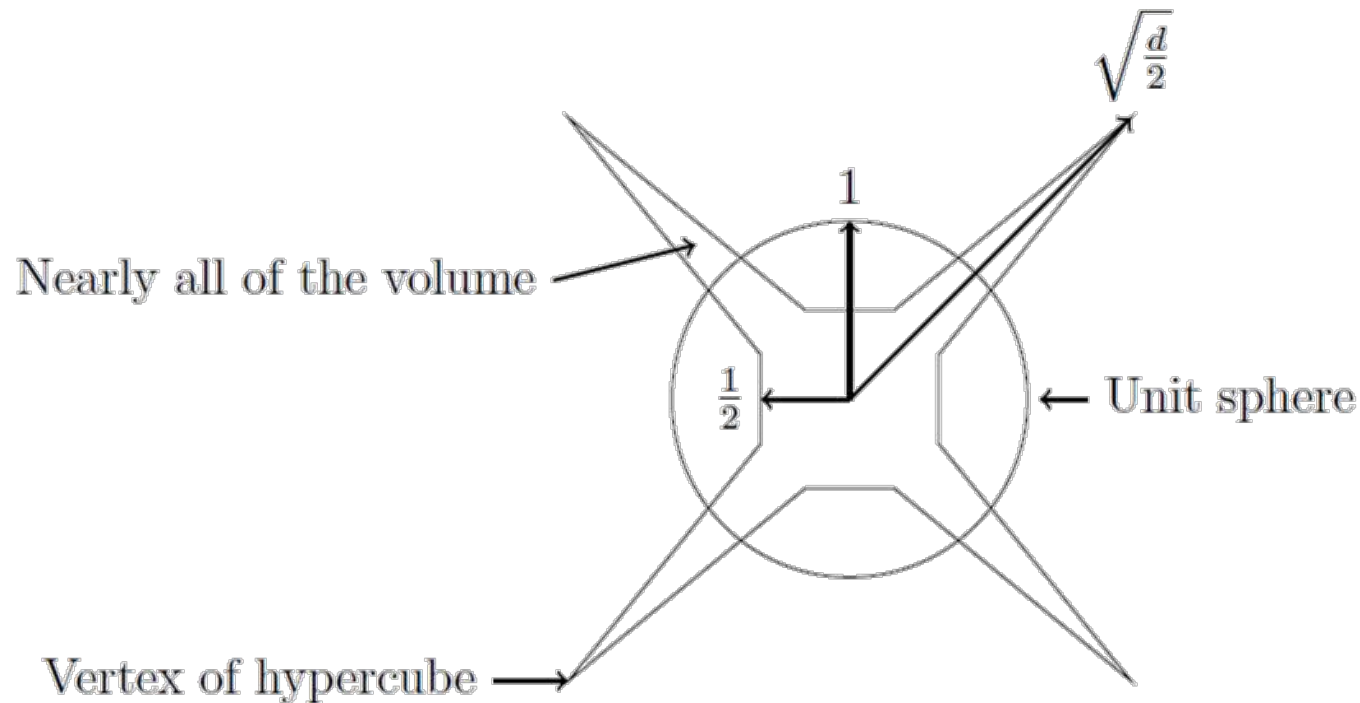
Relationship between the sphere and cube



The difference between the volume of a **cube** with unit-length sides and the volume of a unit-radius **sphere** at the dimensions: 2, 4 and d .

單位球內の手必 沒有樣

Conceptual drawing of a sphere and a cube



For large d , almost all the volume of the cube is located outside the sphere.

Geometry of High Dimensions

- Most of the volume of the high-dimensional objects is near the surface:

高维球缩-点就没东西}

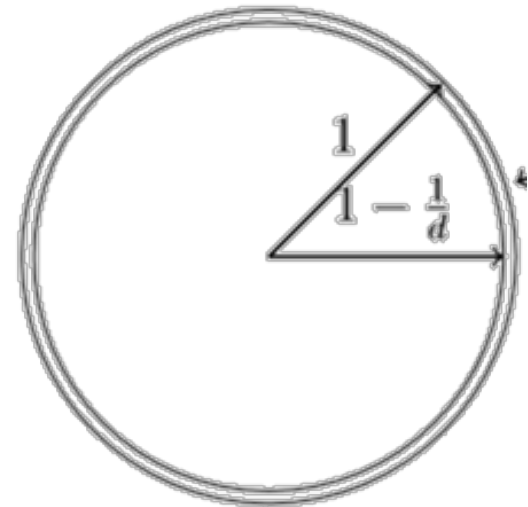
$$\frac{\text{Volume}((1 - \epsilon)A)}{\text{Volume}(A)} = (1 - \epsilon)^d \leq e^{-\epsilon d}$$

Fix ϵ and letting $d \rightarrow \infty$, the above quantity rapidly approaches zero.

- Application:**

S be the unit ball in d –dimensions (i.e., the set of points within distance 1 of the origin). Then $1 - e^{-\epsilon d}$ fraction of the volume is in $S \setminus (1 - \epsilon)S$.

Especially, consider $\epsilon = \frac{1}{d}$.

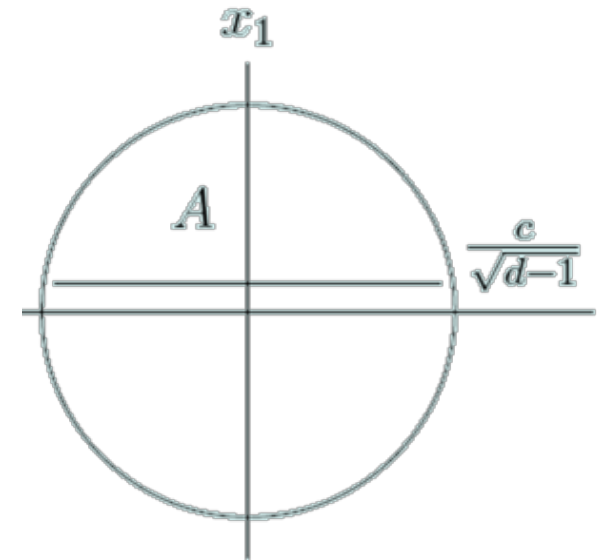


Unit ball in d – dimensions

- **Surface:** $A(d) = \frac{2\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2})}$, **Volume:** $V(d) = \frac{2}{d} \frac{\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2})}$.
面积函数 体积函数
- $V(2) = \pi, V(3) = \frac{4}{3}\pi$, $\lim_{n \rightarrow \infty} V(d) = 0$.
- **Most of the volume** of a **unit ball** in high dimensions is concentrated **near its equator** no matter which direction is defined to be the North Pole.

Theorem: For $c \geq 1$ and $d \geq 3$, at least a $1 - \frac{2}{c} e^{-c^2/2}$ fraction of the volume of the d –dimensional unit ball has $|x_1| \leq \frac{c}{\sqrt{d-1}}$.

坐标趋向于0
x1轴



The law of Large numbers

Properties of High-Dimensional space,
unit ball

Generating points uniformly at random
from a ball

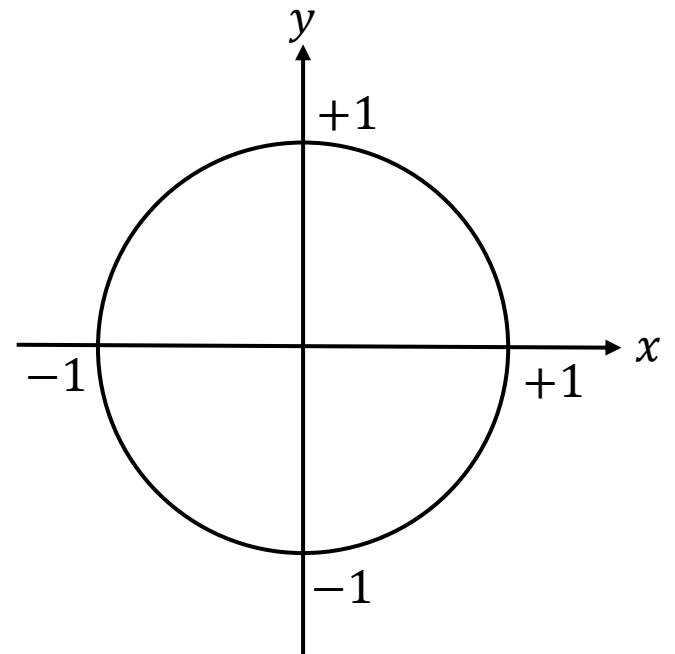
Gaussians in High Dimension

Random Projection and Johnson-
Lindenstrauss Lemma

Seperating Gaussians

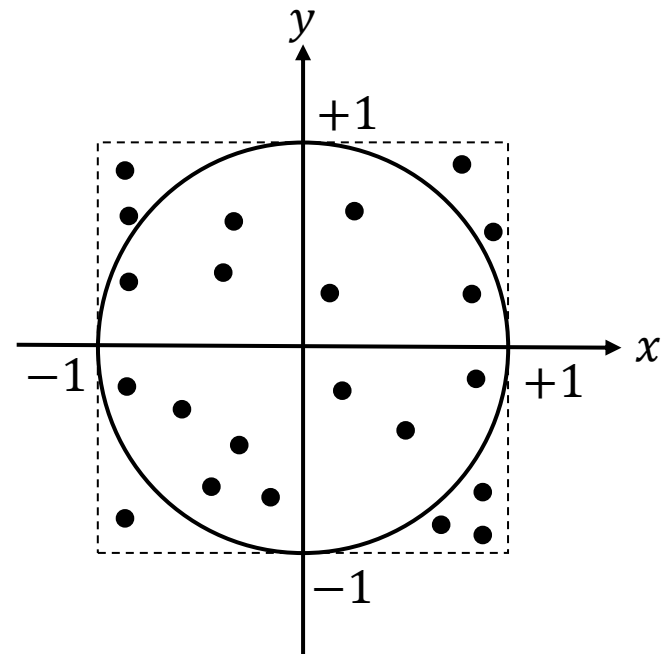
Generating points uniformly at random on the surface of the unit ball

- $d = 2$



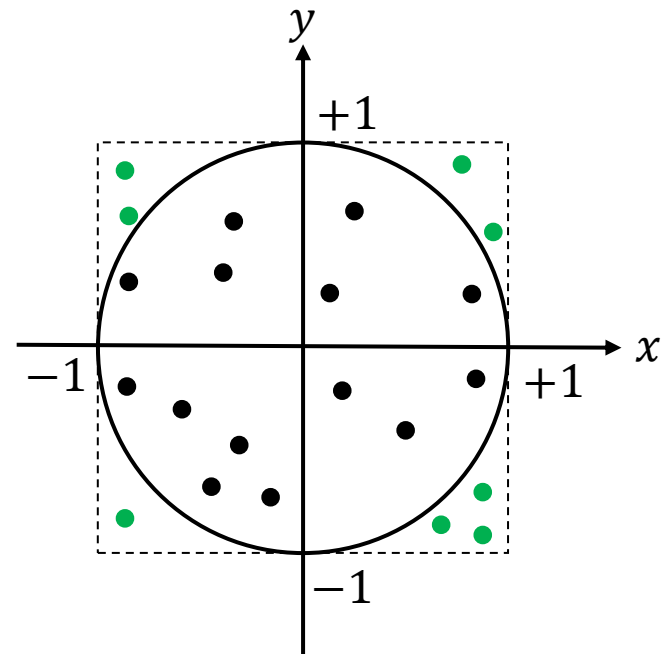
Generating points uniformly at random on the surface of the unit ball

- $d = 2$
 - Generate x_i, y_i u.a.r from the interval $[-1,1]$;



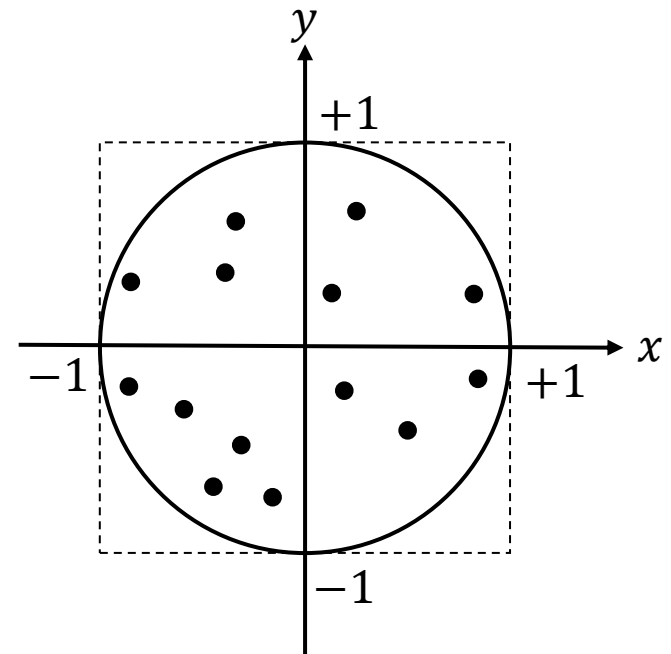
Generating points uniformly at random on the surface of the unit ball

- $d = 2$
 - Generate x_i, y_i u.a.r from the interval $[-1,1]$;



Generating points uniformly at random on the surface of the unit ball

- $d = 2$
 - Generate x_i, y_i u.a.r from the interval $[-1, 1]$;
 - Discard the points outside the unit circle;

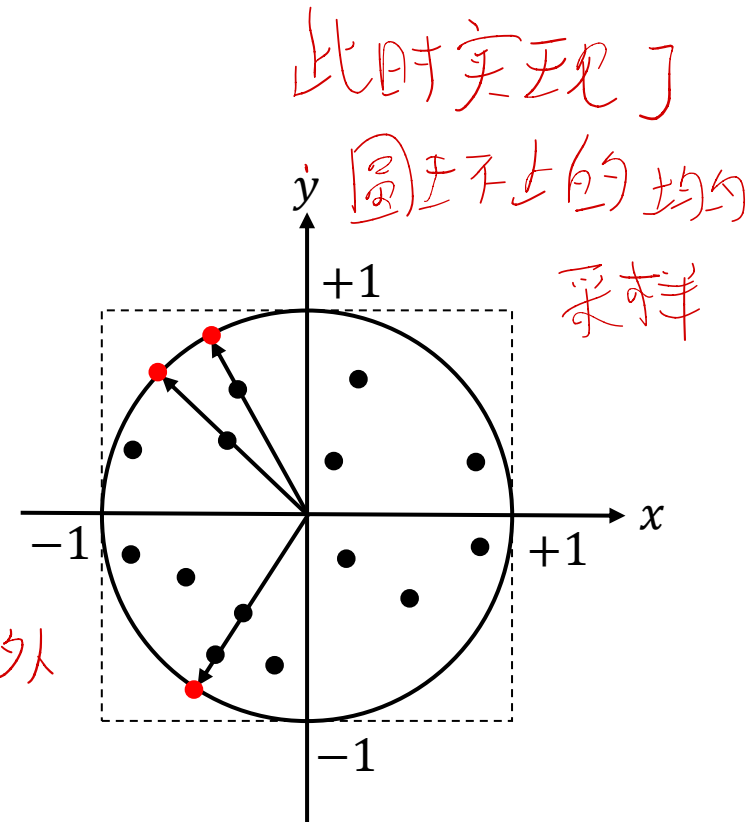


Generating points uniformly at random on the surface of the unit ball

- $d = 2$
 - Generate x_i, y_i u.a.r from the interval $[-1,1]$;
 - Discard the points outside the unit circle;
 - Project the remaining points onto the circle.

- How about d is large?
 - The above strategy would fail. (why?)

- ① **Surface:** Spherical normal distribution + Normalizing.
- ② **Surface+interior:** Scale the point on the surface.



The law of Large numbers

Properties of High-Dimensional space,
unit ball

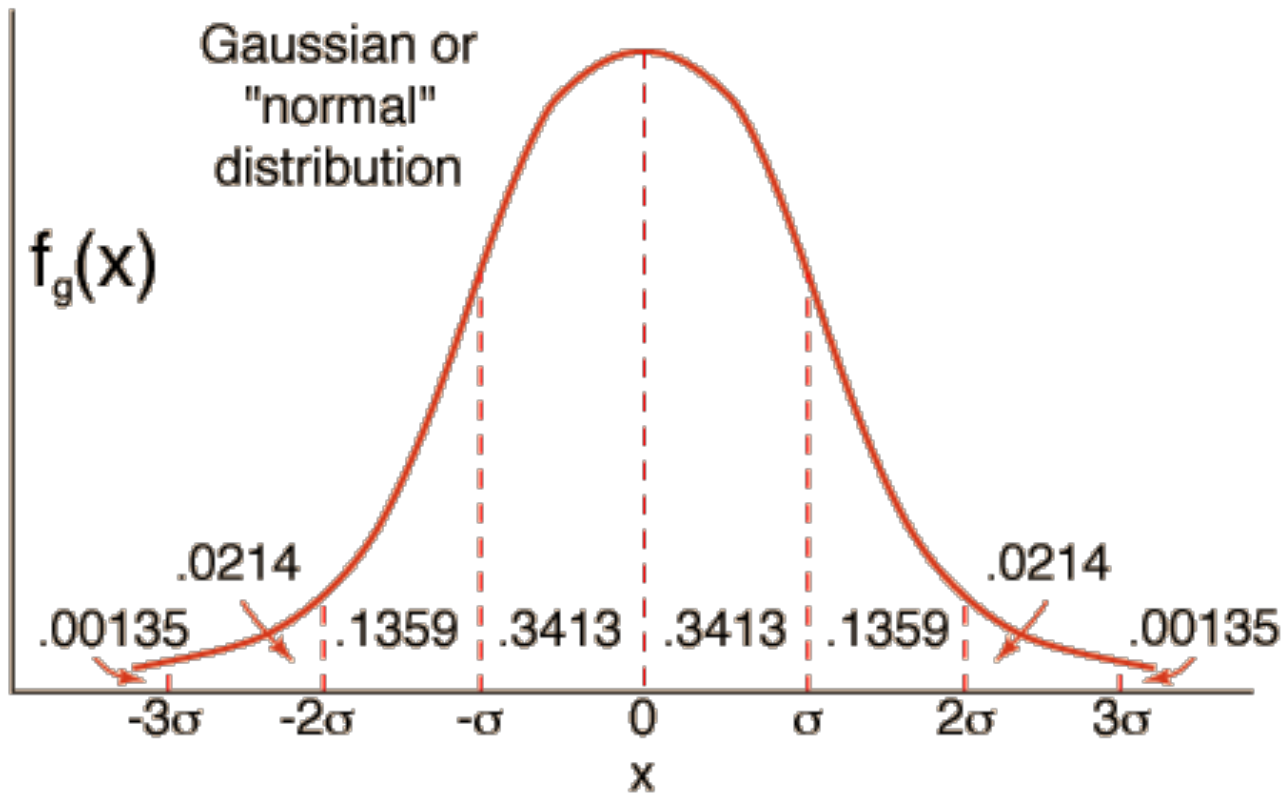
Generating points uniformly at random
from a ball

Gaussians in High Dimension

Random Projection and Johnson-
Lindenstrauss Lemma

Seperating Gaussians

- 1-dimensional Gaussian



- d –dimensional spherical Gaussian with 0 means and variance σ^2 in each coordinate has density function: 几个高斯分布的联合分布

$$p(x) = \frac{1}{(2\pi)^{d/2} \sigma^d} \exp\left(-\frac{|x|^2}{2\sigma^2}\right)$$

- Integrate the PDF over a **unit ball** centered at the origin will cover **almost 0 mass**, for the volume of such a ball is negligible.
- The radius of the ball need to be nearly \sqrt{d} before there is a **significant volume** and hence significant probability mass.

Gaussian Annulus Theorem

- For a d –dimensional spherical Gaussian with unit variance in each direction ,for any $\beta \leq \sqrt{d}$, all but at most $3e^{-c\beta^2}$ of the probability mass lies within the annulus

$$\sqrt{d} - \beta \leq |x| \leq \sqrt{d} + \beta$$

where c is a fixed positive constant.

不在范围内

The law of Large numbers

Properties of High-Dimensional space,
unit ball

Generating points uniformly at random
from a ball

Gaussians in High Dimension

Random Projection and Johnson-
Lindenstrauss Lemma

Separating Gaussians

Database query: Nearest neighbor search

n points from R^d :
$$\begin{bmatrix} v_{11} & v_{21} & \vdots & v_{n1} \\ v_{12} & v_{22} & \vdots & v_{n2} \\ \vdots & \vdots & \vdots & \vdots \\ v_{1d} & v_{2d} & \vdots & v_{nd} \end{bmatrix}$$

- **Nearest neighbor search**: find the nearest or approximately nearest database point to the query point.
- When d is large, it could cost more than expected.
- **Dimension reduction** : *Project* the database points to a k dimensional space with $k \ll d$. It will work so long as the relative distances between points are approximately preserved.

Projection function

- Pick k vectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$, independently from the Gaussian distribution

$\frac{1}{(2\pi)^{d/2} \sigma^d} \exp\left(-\frac{|\mathbf{x}|^2}{2\sigma^2}\right)$, for any vector \mathbf{v} , the projection $f: R^d \rightarrow R^k$ is:

$$f(\mathbf{v}) = (\mathbf{u}_1 \cdot \mathbf{v}, \mathbf{u}_2 \cdot \mathbf{v}, \dots, \mathbf{u}_k \cdot \mathbf{v})$$

Projection function

Pick k vectors

$\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$, independent
ly from the Gaussian
distribution

$\frac{1}{(2\pi)^{d/2} \sigma^d} \exp\left(-\frac{|\mathbf{x}|^2}{2\sigma^2}\right)$, for
any vector \mathbf{v} , the
projection $f: R^d \rightarrow R^k$ is:

$$f(\mathbf{v}) = (\mathbf{u}_1 \cdot \mathbf{v}, \mathbf{u}_2 \cdot \mathbf{v}, \dots, \mathbf{u}_k \cdot \mathbf{v})$$

- $f(\mathbf{v}_1 - \mathbf{v}_2) = f(\mathbf{v}_1) - f(\mathbf{v}_2)$
- $|f(\mathbf{v})| \approx \sqrt{k}|\mathbf{v}|$ w.h.p.
- To estimate $|\mathbf{v}_1 - \mathbf{v}_2|$, it suffices to compute $|f(\mathbf{v}_1) - f(\mathbf{v}_2)|$

Random Projection Theorem

- Let v be a fixed vector in R^d and let f be defined as above. Then there exists constant $c > 0$ such that for $\epsilon \in (0,1)$

$$\Pr \left(\left| \|f(v) - \sqrt{k}|v|\right| \geq \epsilon \sqrt{k}|v| \right) \leq 3e^{-ck\epsilon^2}$$

Johnson-Lindenstrass Lemma

- For any $0 < \epsilon < 1$ and any integer n , let $k \geq \frac{3}{c\epsilon^2} \ln n$ for c as in the Gaussian Annulus theorem, for any set of n points in R^d , the random projection f defined above has the property that for all pairs of points v_i and v_j , with probability at least $1 - \frac{1.5}{n}$
$$(1 - \epsilon)\sqrt{k}|v_i - v_j| \leq |f(v_i) - f(v_j)| \leq (1 + \epsilon)\sqrt{k}|v_i - v_j|.$$

Comments

- JL lemma works for all pairs of points,
- k depends on $\ln n$,
- To the database, JL Lemma says the algorithm will yield the right answer with high probability whatever the query is.

The law of Large numbers

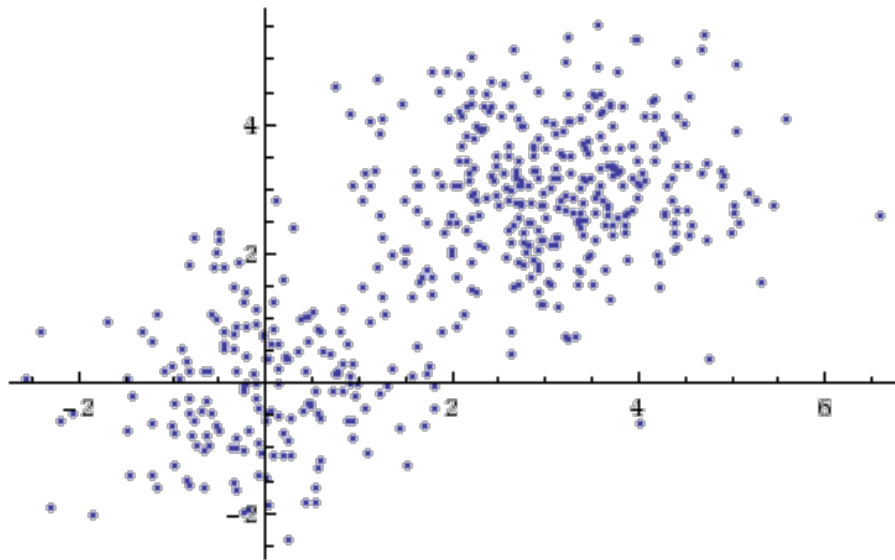
Properties of High-Dimensional space,
unit ball

Generating points uniformly at random
from a ball

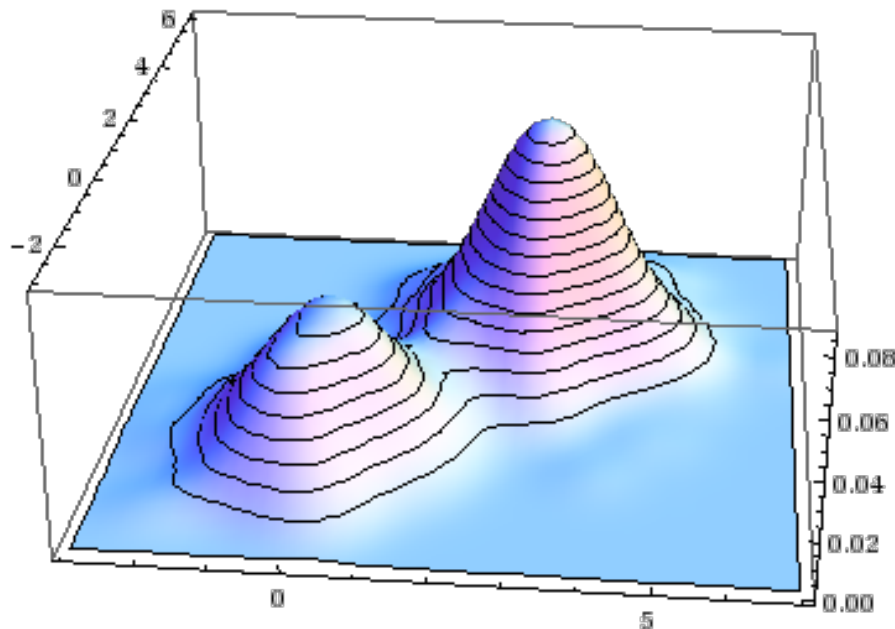
Gaussians in High Dimension

Random Projection and Johnson-
Lindenstrauss Lemma

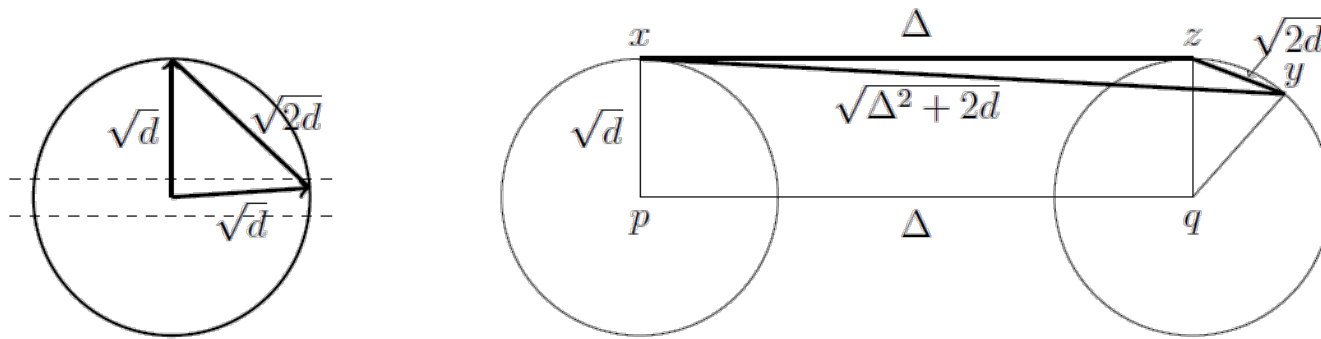
Separating Gaussians



- Mixtures of Gaussians
- Parameter estimation problem



- When $\Delta \in \omega(d^{1/4})$



- Algorithm for separating points from two Gaussians:** Calculate all pairwise distance between points. The cluster of smallest pairwise distances must come from a single Gaussian. Remove these points. The remaining points come from the second Gaussian.