

High-dimensional spaces

Consider x_1, \dots, x_d independent, $E(x_i) = 0$, $\text{Var}(x_i) = 1$

$$\text{object} = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix} =: X$$

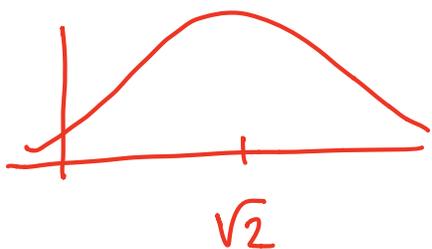
① Norms of high-dim vectors are concentrated:

$$E(\|X\|^2) = E\left(\sum_{i=1}^d x_i^2\right) = \sum_{i=1}^d \underbrace{E x_i^2}_1 = d$$

By Bernstein inequality one can prove:

$$P(|\|X\| - \sqrt{d}| > t) \leq 2 \exp(-ct^2)$$

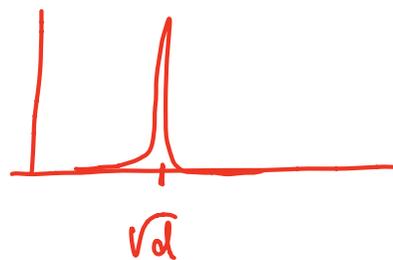
2-dim



$$X \sim N(0, I_2)$$



d-dim



$$X \sim N(0, I_d)$$



②

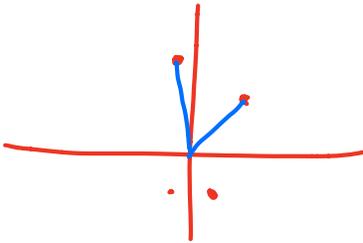
Two random vectors are nearly orthogonal

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}, \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_d \end{pmatrix}$$

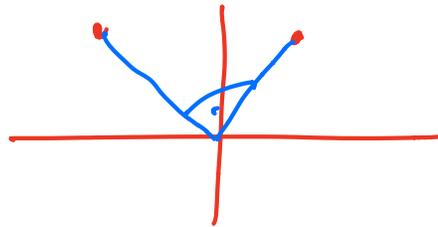
$$E(\langle x, y \rangle) = 1/d.$$

But again: strong concentration.

2 dim normal



d-dim



③

Distances in high-dim spaces are dominated by noise

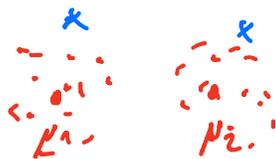
$$x \sim N(\mu_1, \sigma_1^2 I_d), \quad y \sim N(\mu_2, \sigma_2^2 I_d)$$

$$E(\|x - y\|^2) = E\left(\sum_{i=1}^d |x_i - y_i|^2\right)$$

$$= \sum_{i=1}^d (\text{Var}(x_i - y_i) + (E(x_i - y_i))^2)$$

$$= \underbrace{d \cdot (\sigma_1^2 + \sigma_2^2)}_{\text{noise}} + \underbrace{\|\mu_1 - \mu_2\|^2}_{\text{signal}}$$

2 dim



d-dim

noise \gg signal

$$\text{dist}(x, y) =$$

$\|\mu_1 - \mu_2\| + \text{small noise}$

Even worse: $x, x' \sim \mathcal{N}(\mu_1, \sigma_1^2 \mathbb{I}_d)$, $y, y' \sim \mathcal{N}(\mu_2, \sigma_2^2 \mathbb{I}_d)$

$$E(\|x - x'\|^2) = 2d \sigma_1^2 \quad \text{within-class distance}$$

$$E(\|x - y\|^2) = \|\mu_1 - \mu_2\|^2 + d(\sigma_1^2 + \sigma_2^2) \quad \text{between class distance}$$

x x'

μ_1

y

y'

μ_2

μ_2

$$E(\|y - y'\|^2) = 2d \sigma_2^2 \quad \text{within-class distance in second class}$$

If $\sigma_1 \ll \sigma_2$, then

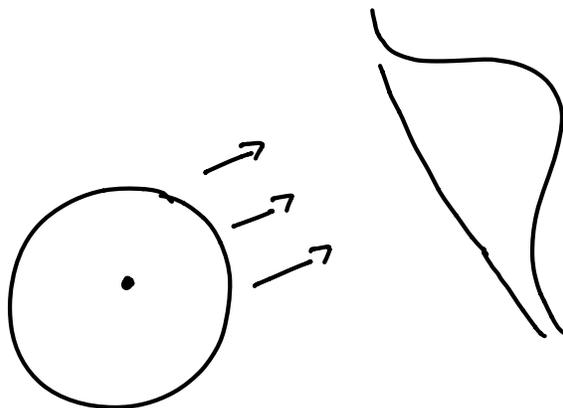
$$\underbrace{\sigma_1^2 + \sigma_2^2}_{\text{between}} \ll \underbrace{\sigma_2^2}_{\text{within-class 2}}$$

④

Projections always look like Gaussians

$$X \sim \text{Unif}(\mathbb{S}_{d-1} \sqrt{d})$$

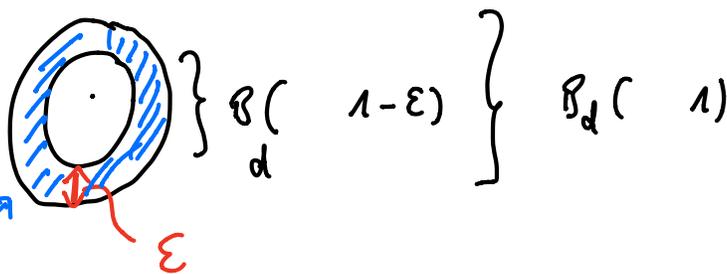
Project X on a fixed vector v :



$$\langle X, v \rangle \rightarrow N(0, 1) \text{ in distribution}$$

⑤

Volume is concentrated along surface:



$$\text{vol}(B_d(1) - B_d(1-\epsilon)) =$$

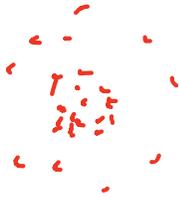
$$= 1^d \cdot \text{vol}(B_d(1)) - (1-\epsilon)^d \text{vol}(B_d(1))$$

$$= \underbrace{(1 - (1-\epsilon)^d)}_{\rightarrow 1 \text{ as } d \rightarrow \infty} \text{vol}(B_d(1))$$

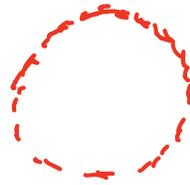
$\rightarrow 1$ as $d \rightarrow \infty$

In particular, if you sample from a normal distribution:

2-dim



d-dim

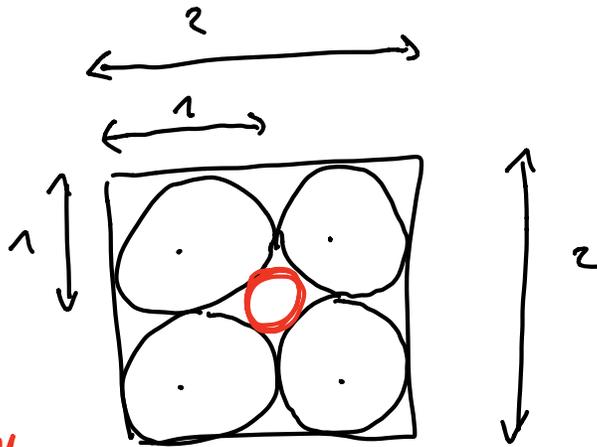


Intuition is often wrong

Riddle

2-dim:

- Square of side length 2
- Circles of diam 1
- *Small circle in the middle*



Radius of small circle:

$$\frac{1}{2} (\text{diagonal of square} - 2 \text{ diam circles})$$

$$\frac{1}{2} (2\sqrt{2} - 2) = \sqrt{2} - 1$$

d-dim

outer circle $\text{diam} = 2$

diam cube $2\sqrt{d}$

inner circle $\sqrt{d} - 1$