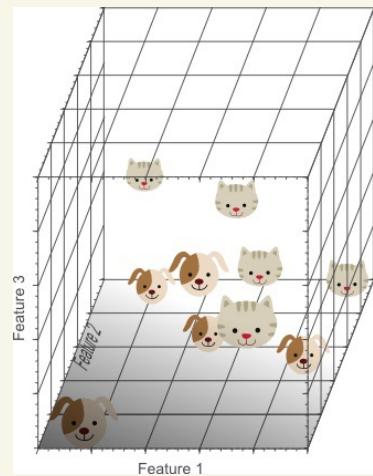
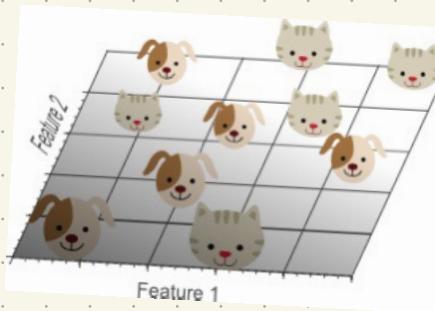
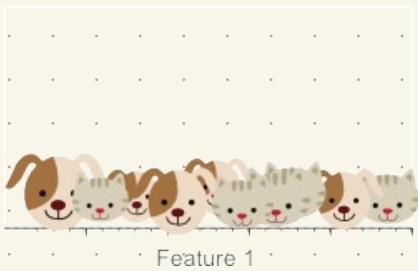


Curse of dimensionality

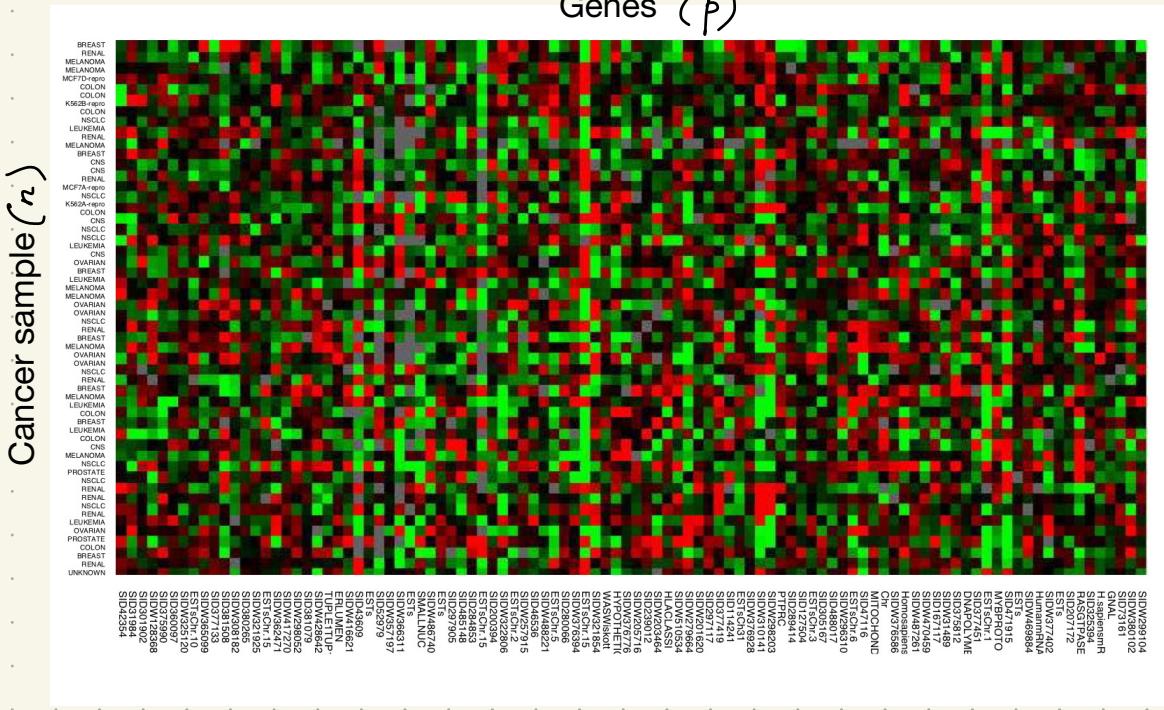


<https://medium.com/@prasantNv/the-curse-of-dimensionality-in-data-analysis-b16ea6903611>

Very often $p(\text{features}) \gg n(\text{observations})$

Example 1: DNA microarray data

Genes (p)



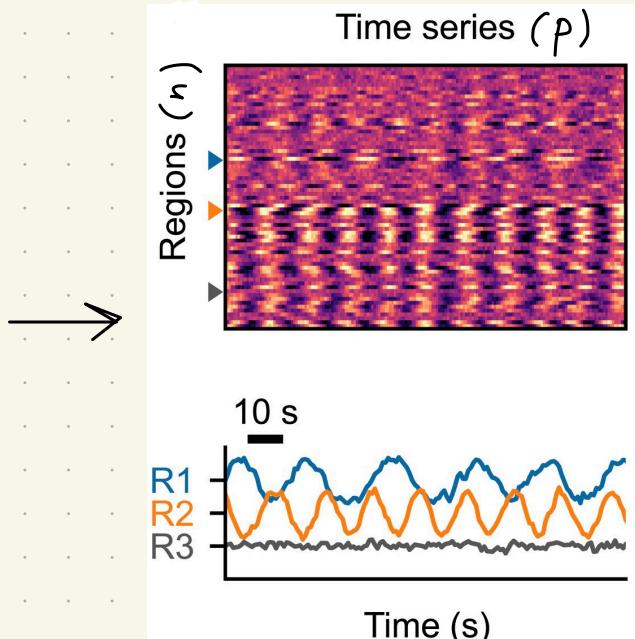
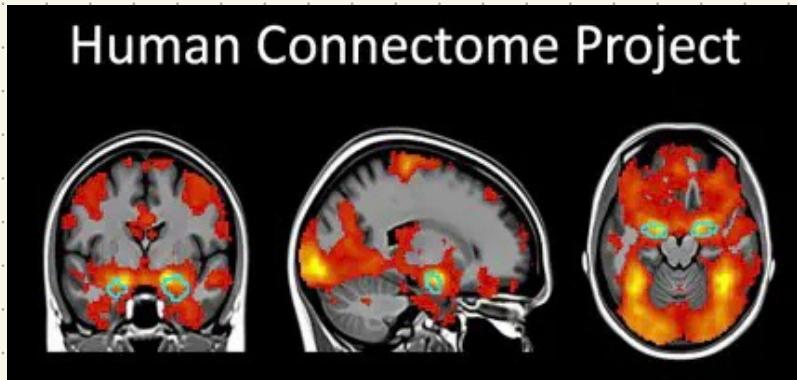
green : under-expressed

red : overexpressed

grey: missing

Very often $p(\text{features}) \gg n(\text{observations})$

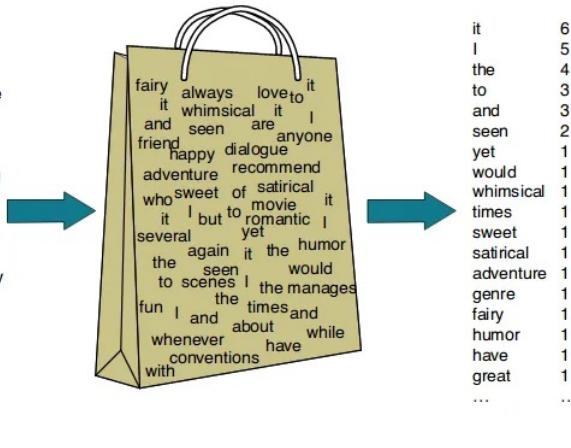
Example 2: fMRI data



Very often $p(\text{features}) \gg n(\text{observations})$

Example 3 : text data

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



Movie reviews	
Review 1	This movie is good.
Review 2	The movie is not good.
Review 3	I love this movie. Watch, you will love it too.

Bag of words (Bow) representation

Reviews (n)

	This	Movie	Is	The	Good	Of	Times	Not	I	Love	Watch	You	Will	It	too
Review 1	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0
Review 2	0	1	1	0	1	0	0	1	0	0	0	0	0	0	0
Review 3	1	1	0	0	0	0	0	0	1	2	1	1	1	1	1

Vocabulary (p)

Why do we care?

- interpretability and visualization
- noise level
- computational issues
- combinatorics

Suppose we have $n = 1000$ patients and P symptoms
(categorical: "mild", "moderate", "severe")

$x_1, \dots, x_n \in \{0, 1, 2\}^P$, each x_i takes one of 3^P values.

$\textcircled{P=3}$ $1000/27 \approx 37$ patients / symptom combination

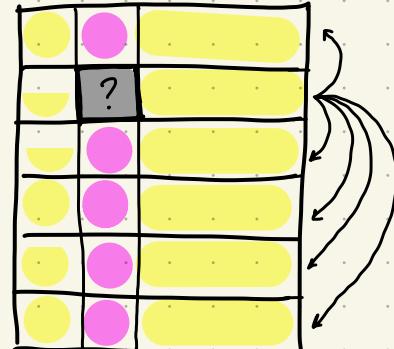
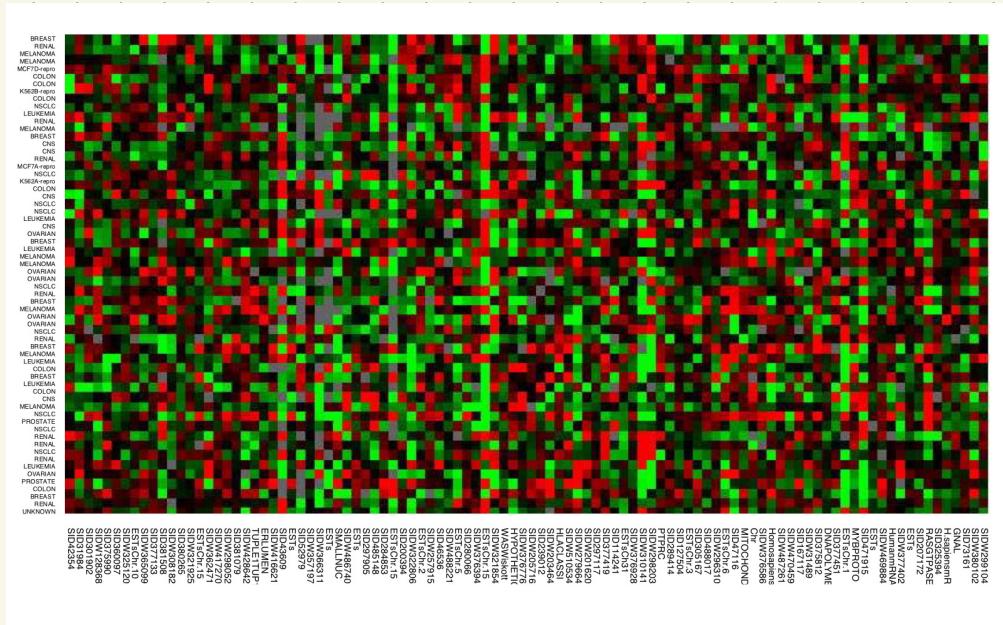
$\textcircled{P=6}$ $1000/729 \approx 1.4$ patients / symptom combination

We need more observations to draw conclusions.

What else?

Many ML / stats methods are based on **distances** between observations.

Example: impute missing values in the DNA data with the nearest neighbor.



Distance behaviour is wierd in high dimensions.

- ① In high-dimensional space nobody can hear you scream

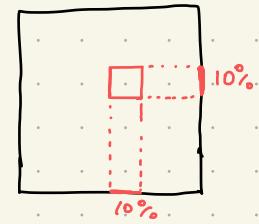


Assume $n=100$ points live in a cube $[0, 1]^P$

- ① each observation takes $1/100 = 0.01$ of the length
i.e. 1% of the segment.



- ② each observation takes $1/100 = 0.01$ of the area
i.e. a square of size $\sqrt{0.01} = 0.1 = 10\%$

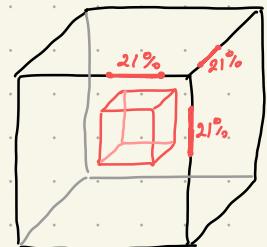


$p=3$ each observation takes $\frac{1}{100} = 0.01$ of the volume

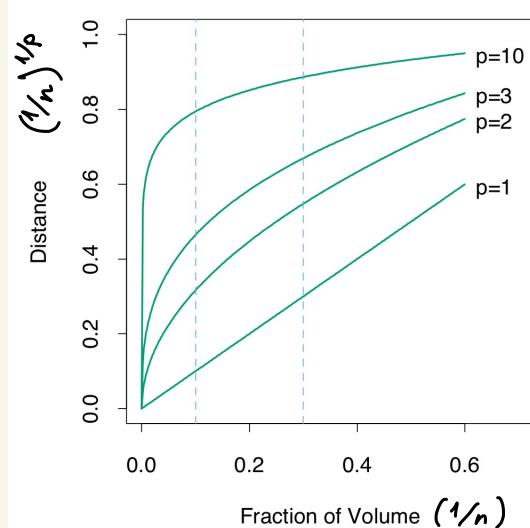
i.e. a cube of size $\sqrt[3]{0.01} = 0.21 = 21\%$

In general, to capture 1% of the volume we need a hypercube

of size $\sqrt[p]{0.01}$.



$p > 7$ The hypercube side is $> 50\%$.



② Orange peel

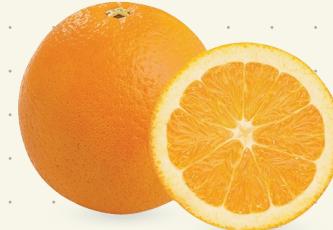
The ball of radius r is

$$B_p(r) = \{x \in \mathbb{R}^P : \|x\|_2 \leq r\}$$

Suppose $x_1, \dots, x_n \stackrel{\text{iid}}{\sim} \text{Unif}(B_p(1))$

Then the median distance from the origin to the closest data point is

$$d(p, n) = \left(1 - \left(\frac{1}{2}\right)^{\frac{1}{n}}\right)^{\frac{1}{p}}$$



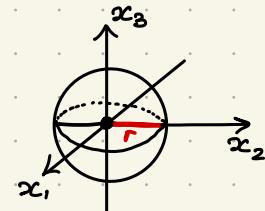
We want to find r s.t.

$$P(\min_{i=1 \dots n} \|x_i\|_2 > r) = \frac{1}{2}$$

$$P(\|x_i\|_2 < r) = \frac{\text{Vol}(B_p(r))}{\text{Vol}(B_p(1))} \quad \text{where } \text{Vol}(B_p(r)) = \frac{\pi^{\frac{P}{2}} r^P}{\Gamma(\frac{P}{2} + 1)}$$

$$P(\|x_i\|_2 < r) = r^P$$

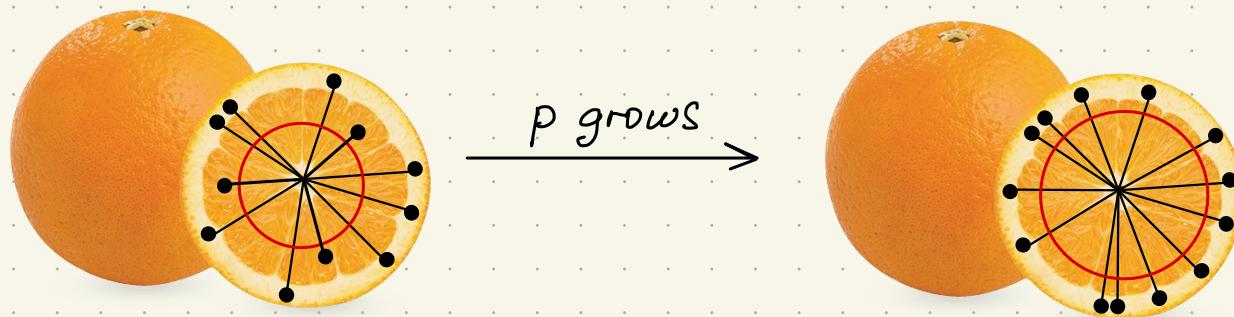
$$P(\min_{i=1 \dots n} \|x_i\|_2 > r) = \prod_{i=1}^n P(\|x_i\|_2 > r) = (1 - r^P)^n = \frac{1}{2}$$



If $n=100$ then $d(p, n) > 0.5$ for $p > 7$.

Thus

- most points are close to the boundary
- the points in high-dim. space are isolated



③ In high dimensions all distances are similar

$$x \sim N_p(0, I)$$

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$$

$$\bullet E(\|x\|^2) = P$$

$$\left| E(\|x\|^2) = \sum_{i=1}^p E(x_i^2) = P \cdot E(x_i^2) = P \right.$$

$$\bullet \text{Var}(\|x\|^2) = 2P$$

$$\left| \text{Var}(\|x\|^2) = P \cdot \text{Var}(x_i^2) = P \cdot \left(\underbrace{E(x_i^4)}_3 - \underbrace{E(x_i^2)}_1 \right) \right.$$

We can also show that:

$$|E(\|x\|) - \sqrt{P}| \leq \frac{1}{\sqrt{P}} \quad \text{and} \quad \text{Var}(\|x\|) \leq 2$$

That is, $\|x\| \rightarrow \sqrt{P}$ and the spread is bounded

If $x, y \stackrel{i.i.d.}{\sim} N_p(0, I)$ then $E(\|x-y\|^2) = 2p$

$$E(\|x-y\|^2) = \sum_{i=1}^p E(x_i^2) - 2E(x_i y_i) + E(y_i^2) = 2p$$

Moreover,

$$|E(\|x-y\|) - \sqrt{2p}| \leq \frac{1}{\sqrt{2p}} \quad \text{and} \quad \text{Var}(\|x-y\|) \leq 3$$

More formally,

$$P(|\|x\| - \sqrt{p}| \geq \epsilon) \leq 2e^{-c\epsilon^2} \quad \forall \epsilon \in [0, \sqrt{p}]$$



E.g. When $p \geq 100$ then $\|x\| \in \sqrt{p} \pm 10$
with probability 0.99.

④ In high dimensions two random vectors are orthogonal.

Let $x, y \sim N_p(0, I)$ and $a \in \mathbb{R}^p$ is constant.

- $E(\langle x, a \rangle) = 0$ and $\text{Var}(\langle x, a \rangle) = \|a\|^2$

$$E(\langle x, a \rangle) = \sum_{i=1}^p a_i E(x_i) = 0$$

$$\text{Var}(\langle x, a \rangle) = \sum_{i=1}^p a_i^2 \text{Var}(x_i) = \|a\|^2$$

- $E(\langle x, y \rangle) = 0$ and $\text{Var}(\langle x, y \rangle) = p$

$$\langle x, y \rangle = \frac{1}{2} (\|x\|^2 + \|y\|^2 - \|x-y\|^2) \quad \text{so}$$

$$E(\langle x, y \rangle) = \frac{1}{2} (p+p-2p) = 0$$

$$\text{Var}(\langle x, y \rangle) = p \cdot \text{Var}(x_i y_i) = p \text{Var}(x_i) \text{Var}(y_i) = p$$

Combining ③ and ④

$$\|x\|^2 \simeq p \pm \sqrt{2p} \quad \|y\|^2 = p \pm \sqrt{2p}$$
$$\langle x, y \rangle \simeq 0 \pm \sqrt{p}$$

The cosine between x and y is

$$\frac{\langle x, y \rangle}{\|x\| \|y\|} \simeq \frac{0 \pm \sqrt{p}}{(\sqrt{p \pm \sqrt{2p}})^2} = \frac{0 \pm \sqrt{p}}{p \pm \sqrt{2p}} \longrightarrow 0 \quad \text{for } p \rightarrow \infty$$

Formally, $\forall \varepsilon > 0$ and $p \geq 1$

$$P\left(|\frac{\langle x, y \rangle}{\|x\| \|y\|}| \geq \varepsilon\right) \leq \frac{2/\varepsilon + 7}{\sqrt{p}}$$



If $x \sim N_p(0, I)$ what is the distribution of $\frac{x}{\|x\|}$?

$$P=1$$

then $\frac{x}{\|x\|} = \text{Sign}(x) \sim \text{Unif}(-1, 1)$

$$P=2$$



$\frac{x}{\|x\|} \sim \text{Unif}(\text{circle})$

For general P , if $x \sim N_p(0, I)$ the density $f(x) \propto e^{-\frac{1}{2}\|x\|^2}$ depend on the $\|x\|$ but not the direction.

Moreover, for any Q orthogonal $Qx \sim N_p(0, I)$ thus $N_p(0, I)$ is rotation invariant and x has the same probability to point in any direction.

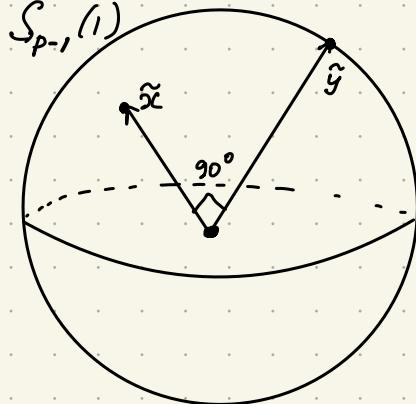
Thus $\frac{x}{\|x\|} \sim \text{Unif}(S_{p-1}(1))$ where $S_{p-1}(r) = \{x \in \mathbb{R}^P : \|x\| = r\}$

$\tilde{x} = \frac{x}{\|x\|}$, $\tilde{y} = \frac{y}{\|y\|}$ are uniform on a sphere $S_{p-1}(1)$

where $S_{p-1}(r) = \{x \in \mathbb{R}^p : \|x\| = r\}$

$\frac{\langle x, y \rangle}{\|x\| \|y\|} = 0$ means $\tilde{x} \perp \tilde{y}$

In other words, in high dimensions
two random x, y on a sphere are
Orthogonal



	$\ x\ $	$\ x-y\ $	$\langle x, y \rangle$
$x, y \sim N_p(0, I)$	$\simeq \sqrt{p}$	$\simeq \sqrt{2p}$	$\simeq 0$
$x, y \sim \text{Unif}(S_{p-1}(1))$	$= 1$	$\simeq \sqrt{2}$	$\simeq 0$