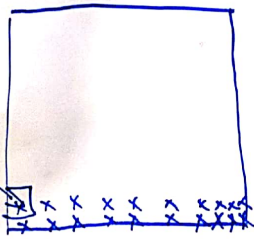


Curse of dimensionality

exponential growth with number of dimensions
needed

2^d

how big does the box
have to be



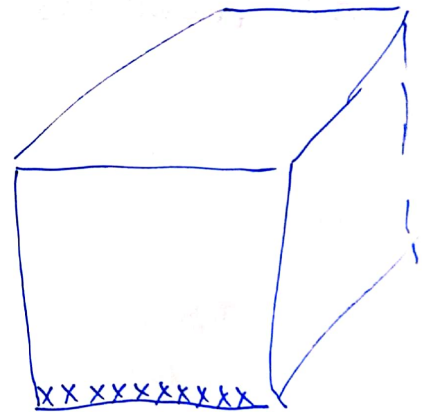
100 points 1 cm^2

- In order to cover
the same average
distance



10 points 1 cm

- average distance to
neighbour: 0.05 cm



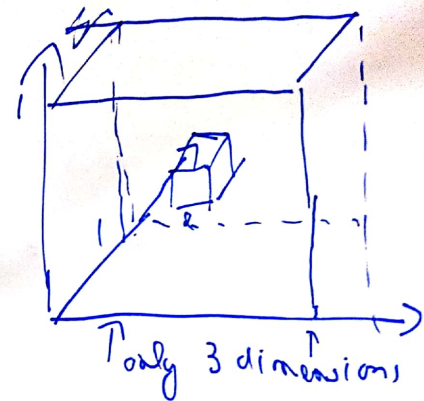
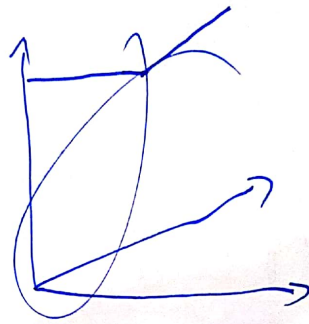
1000 1 cm^3

- In order to ~~cover~~ have
the same average distance



d dimensionality
data drawn from hypercube

edges length 1
all with distance 1



①
 \mathbb{R}^d

- Draw data uniformly at random from this hypercube
- What is the smallest little cube length l that encapsulates the K nearest neighbours of this point (to find K nearest neighbours draw tiny box around it)
- How big does the little box have to be

volume of the little box $l^d \approx \frac{K}{n}$

$$l \approx \left(\frac{K}{n}\right)^{1/d}$$

solve for this

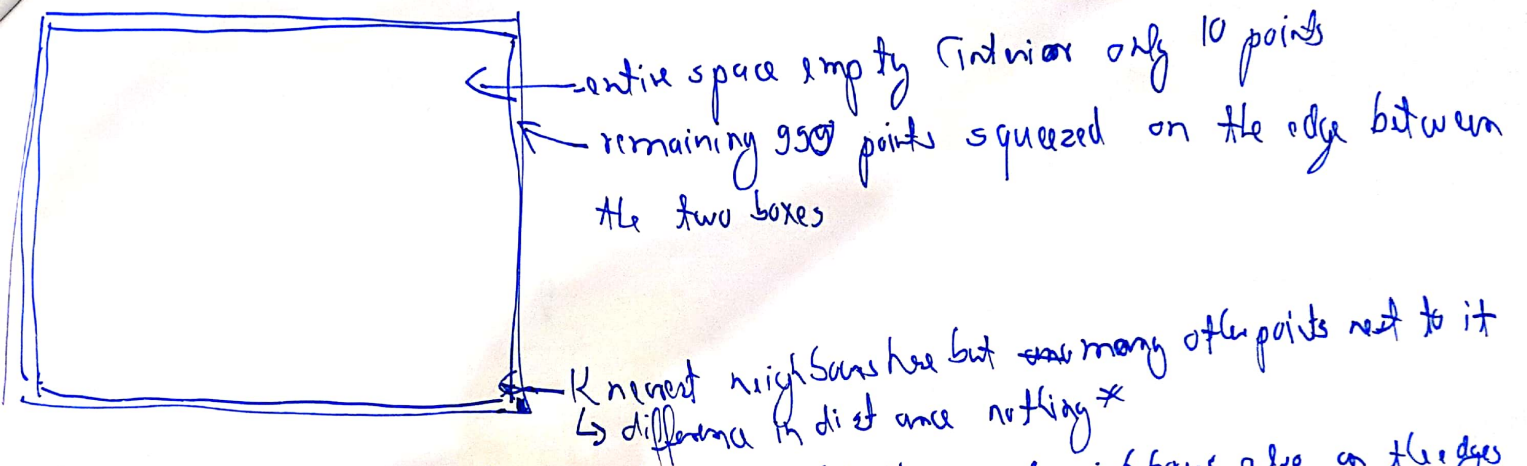
$l^d \leftarrow$ contains K points out of n if since they are uniformly distributed

$K \equiv K$ nearest neighbours

$n \equiv$ total number of points

- volume contains K out of n points
- volume same as ratio

does it mean?



- Fact that little box is so big can only mean that K nearest neighbors also on the edges
⇒ otherwise could draw smaller box
- Points not K nearest neighbors between the two boxes (right next to K nearest neighbors)
- K nearest neighbors not close at all - no notion of closeness, really far away on the edges
- All the points are really far away at the edges
- All points roughly same distance from each other
- Counterintuitive ⇒ brain made for 3 dimensional spaces

* unreasonable to say one point should have same label as another
- not significantly further from point

for $l \approx \left(\frac{K}{n}\right)^{\frac{1}{d}}$, $K=10$ $n=1000$

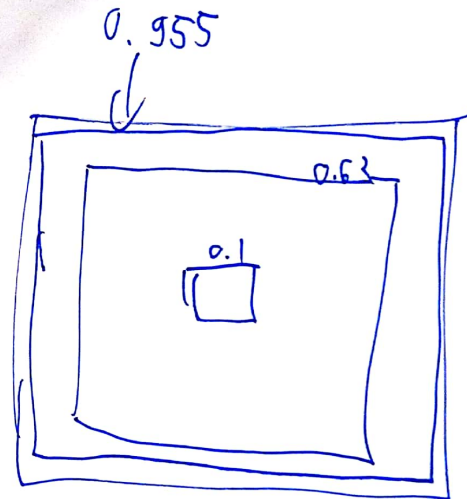
$$\approx \left(\frac{10}{1000}\right)^{\frac{1}{d}}$$

$$\approx \left(\frac{1}{100}\right)^{\frac{1}{d}}$$

Q

\Rightarrow Solve for several values of d

d	l
2	0.1
10	0.63
100	0.955
1000	0.9956



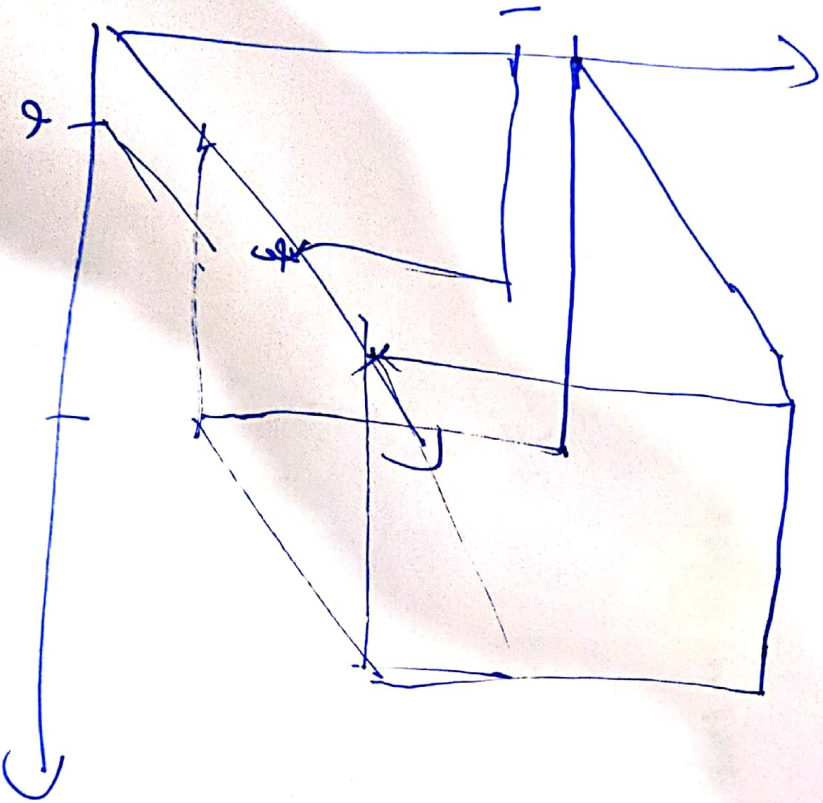
100 dimensions: have the box from which the data is sampled

- For any given point if we just look at K nearest neighbors how much space
- how much space does this box contain containing only my K nearest neighbors \rightarrow essentially same size as total box

- ⇒ in high dimensional spaces if you draw points randomly would be in interior
- ⇒ in high dimensional spaces never the center → interior is empty & never anything there
- ⇒ everything far away from each other
- ⇒ Assumption of K nearest neighbors that nearby points have the same label means
⇒ nothing nearby, everything is about same distance from each other

Question why everything is at the edges

- Drawing a point at random ^{in the cube} could mean drawing every ~~we~~ could just draw every single coordinate independently

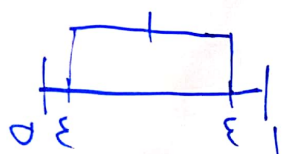


— length of 1 | some number between 0 and 1 randomly

- In order for x in the interval $[0, 1]$ cannot be at the edge

Probability of not being at the edge

Interior $1 - 2\epsilon \leftarrow$ probability of not being on the edge



$< \epsilon$ away from the edge = at the edge
 ~~ϵ per~~

- ~~Interior in~~ Probability of interior in every single dimension

$$(1 - 2\epsilon)^d$$

$$< 1$$



- in every single dimension I cannot be at the edge.

- If there is a single dimension at the edge makes you an edge point

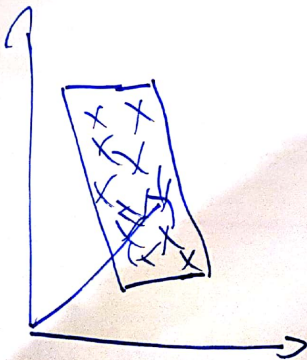
goes to 0 quickly with large power

\Rightarrow probability of hitting interior in high dimensional space is basically 0

How can this work eg. on images?

- pictures are not uniformly distributed
- no assumptions about the space, high dimensional data algorithms e.g. K nearest neighbours do not work
- BUT - High dimensional data, but within this space is a subspace that is much smaller
 - Data only lies on that subspace

↑ not uniformly distributed

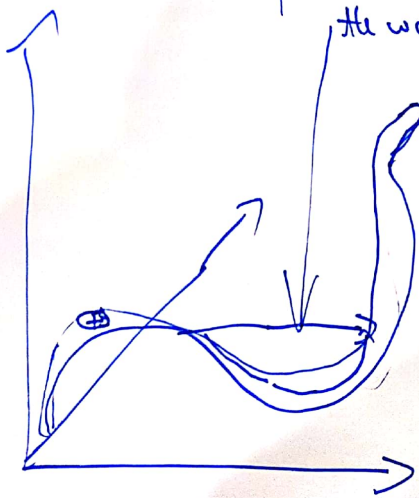


\mathbb{R}^d

- never draw data out of that subspace
- Assumption - low intrinsic dimensionality
 - ↳ Low dimensional sub-space
 - ↳ Low dimensional sub-manifold

dimensional sub manifold.

points close, but if you would want to go on the manifold would have to traverse all the way



Surface - eg. 10 dimensional in 1000 dimensional space

- can be curled up, so really need 1000 dimensional space to represent data

- Surface itself explores 1000 dimensions

- But data itself never leaves underlying many fold which is low dimensional

Manifold: property 1: locally euclidean, if you live in a tiny area of the subspace and move around. would have no idea space is curved \Rightarrow looks flat

- globally not euclidean

- K nearest neighbors only works locally

