

② k-NN

Let $N_k(\vec{x}) = \{\text{indices of } k\text{-NN of } \vec{x} \text{ in } D\}$

k-NN predictor

$$\rightarrow \text{Regression } \hat{y} = g(\vec{x}) = \frac{1}{k} \sum_{i \in N_k(\vec{x})} y_i$$

Predict unweighted average of neighbours

$$\rightarrow \text{Classification } \hat{y} = g(\vec{x}) = \underset{\text{class } c}{\text{argmax}} \#(y_i = c)_{i \in N_k(\vec{x})}$$

$$= \underset{c}{\text{argmax}} \sum_{i \in N_k(\vec{x})} I(y_i = c)$$

unweighted majority vote

③ Distances

→ Most common

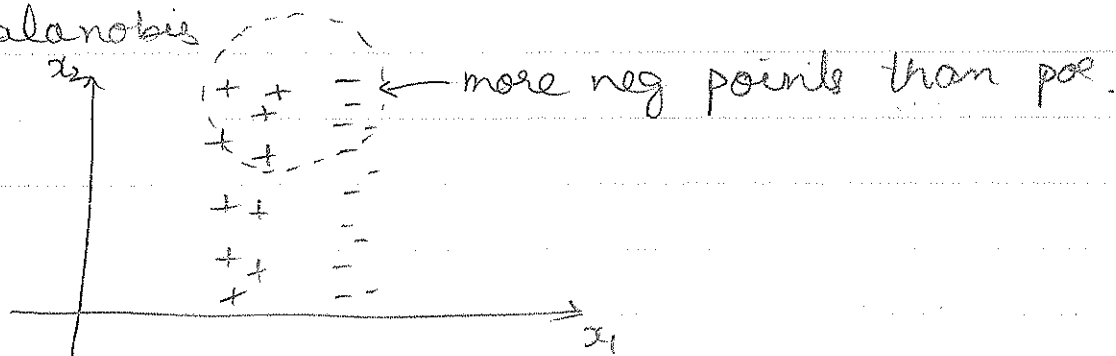
Euclidean Distance / L_2 -norm of difference

$$\vec{x}, \vec{z} \in \mathbb{R}^d$$

$$d(\vec{x}, \vec{z}) = \left[\sum_{i=1}^d (x_i - z_i)^2 \right]^{1/2}$$

→ let's generalize this in 2 ways

① Mahalanobis



New definition

$$d^2(\vec{x}, \vec{z}) = 10(x_1 - z_1)^2 + (x_2 - z_2)^2$$

↑
deviations in dim 1 should be penalized more

In general
$$d^2(\vec{x}, \vec{z}) = \sum_{i=1}^d \sigma_i^2 (x_i - z_i)^2$$

$$= (\vec{x} - \vec{z})^T \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_d^2 \end{bmatrix} (\vec{x} - \vec{z})$$

More generally,

$$d^2(\vec{x}, \vec{z}) = (\vec{x} - \vec{z})^T A (\vec{x} - \vec{z})$$

Set $A = I_{d \times d} \Rightarrow$ Euc. dist

Note $A \succeq 0$

↑
positive semi-definite

Definition: $A = A^T$ symmetric
& $\vec{x}^T A \vec{x} \geq 0 \quad \forall \vec{x} \in \mathbb{R}^d$

←
Other generalization

Minkowski-distance / L_p -norm of difference

$$d(\vec{x}, \vec{z}) = \left[\sum_{i=1}^d |x_i - z_i|^p \right]^{1/p}$$

$p=2$ = Euc. dist

$p=1$ = Manhattan distance

$$= \sum_{i=1}^d |x_i - z_i|$$

$p \rightarrow \infty$ = Max-distance

$$= \max_i |x_i - z_i| \quad 1 \leq i \leq d$$

Why? Simple proof.

$$\lim_{p \rightarrow \infty} \left[\sum_{i=1}^d |x_i - z_i|^p \right]^{1/p}$$

Let j = index of max-difference
= $\arg \max_{i=1, \dots, d} |x_i - z_i|$

[For simplicity, assume unique $\arg \max$]

$$= \lim_{p \rightarrow \infty} \left[|x_j - z_j|^p + \sum_{i \neq j} |x_i - z_i|^p \right]^{1/p}$$

$$= \lim_{p \rightarrow \infty} |x_j - z_j|^{p/p} \left[1 + \sum_{i \neq j} \left(\frac{|x_i - z_i|^p}{|x_j - z_j|^p} \right)^{1/p} \right]^{1/p}$$

$$\underbrace{\left(< 1 \right)^p}_{\text{as } p \rightarrow \infty} \rightarrow 0$$

$$\left[1 + \sum_{i \neq j} (\rightarrow 0) \right]$$

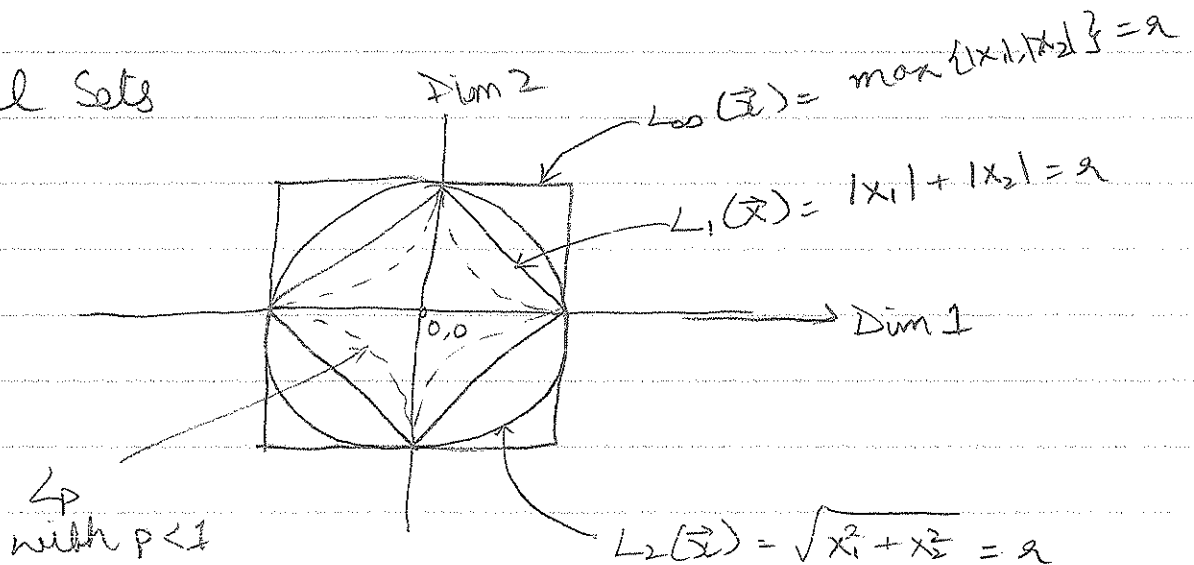
$$= |x_j - z_j|$$

□

Similarly $p=0$

$d(\vec{x}, \vec{z}) = \# \text{ dims where } x_i \neq z_i$

Level Sets

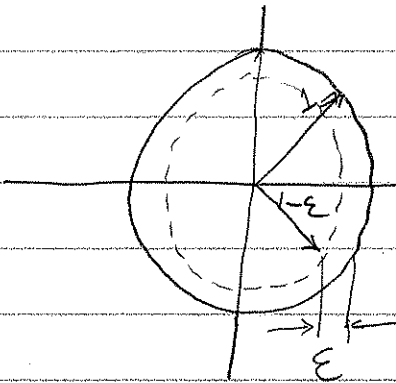


③ Curse of Dimensionality

Learning in high-dimensional space $\equiv d$ -large is difficult.

In particular, NN "shouldn't work". Why?
°° distances/neighbours become meaningless.

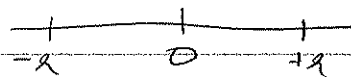
→ Example #1: Consider sphere in \mathbb{R}^d centred at $\vec{0}$
radius $r=1$



What is volume of outer ϵ -shell?

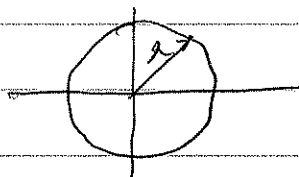
Well, what is volume of sphere?

1D



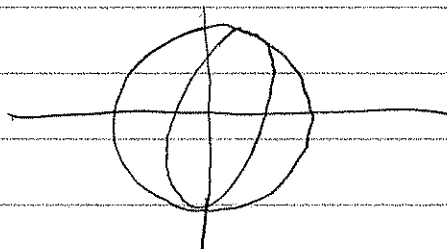
$2r$

2D



πr^2

3D



$\frac{4}{3} \pi r^3$

$\equiv k_d r^d$

(3)

$$\text{Now, } \frac{\text{Volume (Shell)}}{\text{Volume (Sphere)}} = \frac{k_d (1)^d - k_d (1-\epsilon)^d}{k_d 1^d}$$

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

$$\rightarrow 1 \quad \text{as } d \rightarrow \infty$$

⇒ Nearly all volume lies in outer ϵ -shell
 Assume uniform density of data [Hint: Problem!]

⇒ Nearly all mass lies in shell

⇒ Nearly all data-points lie in shell

⇒ All neighbours are equally apart!



Example 2: $\vec{x} = (x_1, \dots, x_d)$

Assume x_1, \dots, x_d are I.I.D random vars
 [Hint: Problem!]

Consider Normalized distance² to origin:

$$D = \frac{1}{d} \|\vec{x} - \vec{0}\|_2^2 = \frac{1}{d} \sum_{i=1}^d x_i^2$$

Recall, Central Limit Theorem

If z_1, \dots, z_n are I.I.D RVs with

$$E[z_i] = \mu \quad \forall i$$

$$\text{Var}(z_i) = \sigma^2 \quad \forall i$$

then $\frac{1}{n} \sum z_i \rightarrow N\left(\mu, \frac{\sigma^2}{n}\right)$

as $n \rightarrow \infty$

So in our case $z_i = x_i^2$

$$D = \frac{1}{d} \sum x_i^2 \rightarrow N\left(E[x_i^2], \frac{\text{Var}(x_i^2)}{d}\right)$$

$$\text{Var}(x_i^2) \equiv \text{constant}$$

$$\Rightarrow \frac{\text{Var}(x_i^2)}{d} \rightarrow 0 \text{ as } d \rightarrow \infty$$

$\Rightarrow D$ is nearly a constant.

Problem!