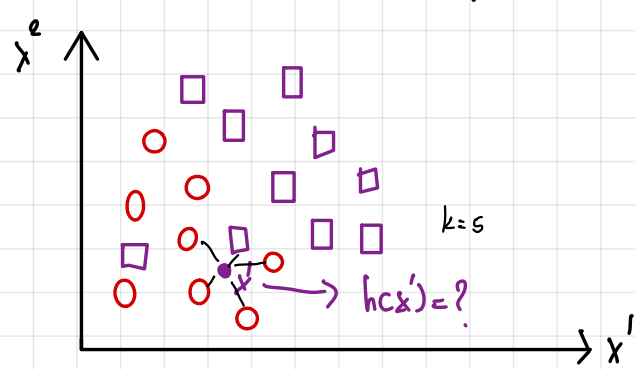


k-Nearest Neighbor (k-NN)

- feature vector in point in space

e.g. $x \in \mathbb{R}^2 \rightarrow \vec{x}$ is a 2D-point



$$D = \{ (x_1, y_1), \dots, (x_n, y_n) \}$$

$$x_i \in \mathbb{R}^2, y_i \in \{ \square, \circ \}$$

ค่า k

การหาค่า $hc(x')$ ได้ต้องหาเพื่อนบ้านที่ใกล้ที่สุด (k)

- ควรเป็น odd

- ถ้าเป็น k คู่ less ดี

- Training: store all data points

- Testing: For a test point x' , the prediction $hc(x')$ is determined by the majority of its nearest neighbors. ^(k)

- k-NN's assumption: close points should have similar labels

closeness \Rightarrow similarity

Formal definition of k-NN

- given a test point \vec{x}

- Denote $S_{\vec{x}}$ the set of k-NN of \vec{x}

$S_{\vec{x}} \subseteq D$ such that $|S_{\vec{x}}| = k$ and

$\forall (x', y') \in D \setminus S_{\vec{x}}$ distance (\vec{x}, \vec{x}') \geq max distance (\vec{x}, \vec{x}'')
complement $S_{\vec{x}}$ $(\vec{x}'', y'') \in S_{\vec{x}}$

นั่นคือ $S_{\vec{x}}$ ภายใน D ที่มีขนาดของ $S_{\vec{x}} = k$

และพิจารณา \vec{x}', y' ใน D ที่ไม่อยู่ใน $S_{\vec{x}}$

- The classifier hc is define as:

$$hc(\vec{x}) = \text{mode}(\{ y' : (\vec{x}, y') \in S_{\vec{x}} \})$$

Distance function

- วัดการต่าง = ความต่างของ x' กับ x โดยปกติแล้ว เพื่อหา km ที่ใกล้ที่สุด
- the most common choice is the "minkowski distance"

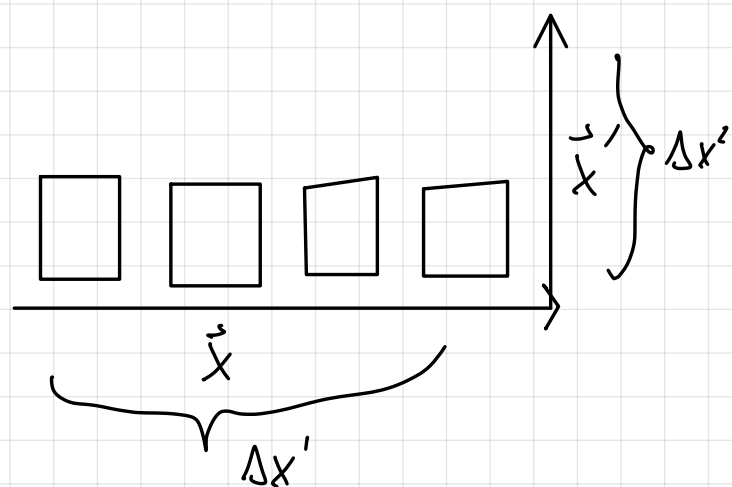
$$\text{distance}(\vec{x}, \vec{x}') = \left(\sum_{i=1}^d |\vec{x}^i - \vec{x}'^i|^p \right)^{1/p}$$

Example

① $p=1$ (Manhattan distance)

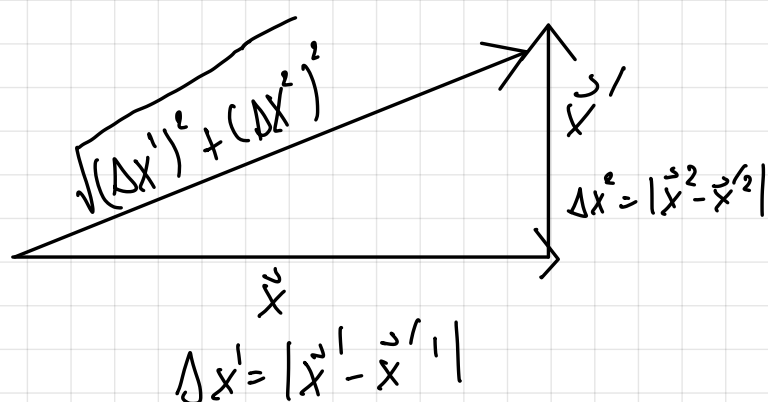
$$\vec{z} = \begin{bmatrix} 2 \\ 5 \end{bmatrix} \quad \vec{z}' = \begin{bmatrix} 5 \\ 3 \end{bmatrix}$$

$$\text{dist}(\vec{z}, \vec{z}') = \left| \begin{bmatrix} 2 \\ 5 \end{bmatrix} - \begin{bmatrix} 5 \\ 3 \end{bmatrix} \right| = [|2-5| + |5-3|] = 5$$



② $p=2$ (Euclidean distance)

$$\begin{aligned} \text{dist}(\vec{z}, \vec{z}') &= \left[|2-5|^2 + |5-3|^2 \right]^{\frac{1}{2}} : \\ &= \sqrt{9+4} = \sqrt{13} \end{aligned}$$



③ note - $p: +\infty$
- $p: -\infty$

$p: \infty$

$$\begin{aligned} &: (|2-5|^\infty + |5-3|^\infty)^{\frac{1}{\infty}} \\ &: (3^\infty + 2^\infty)^{\frac{1}{\infty}} = (3^\infty)^{\frac{1}{\infty}} \approx 3 \end{aligned}$$

$p: -\infty$

$$(3^{-\infty} + 2^{-\infty})^{-\frac{1}{\infty}} \approx 2$$

Implement of k-NN:

- Assume we already store all the training data

Step 1: The algorithm simply "sweeps" through all the training data points, and calculates the distance from the test point to each particular point at the same time. (complexity of step 1 is $O(dn)$)

Step 2: the algorithm sorts the training data points using its distance to the test point as keys, and pick the points whose distance is k-smallest ($O(n \log n)$)

Overall: $O(dn) + O(n \log n) = O(dn)$

Pros and Cons of k-NN

- Given a good and small-dimensional data set ($d \approx \text{small}$), k-NN can give a good Classifier (with high accuracy of prediction)

- However, k-NN has 2 main downsides when used with high-dimensional data set:

- ① The classifier can be very slow when either n or d is large (or both).

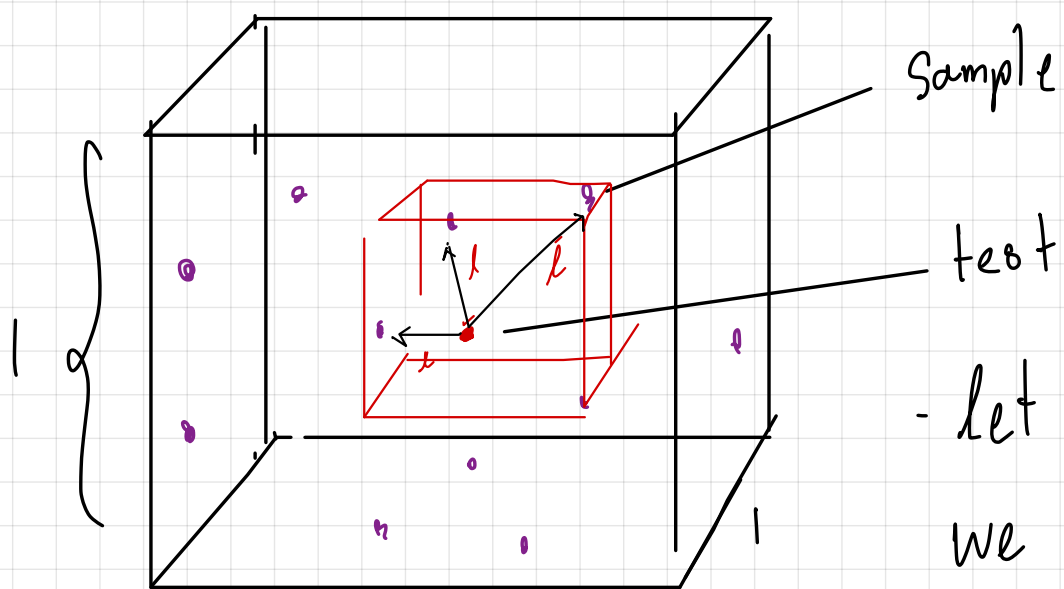
- ② If our data is high-dimensional ($d \gg 0$; e.g. image data),

k-NN will give a bad Classifier (known as "curse of Dimensionality")

Curse of Dimensionality

- The phenomenon that data points in high dimensional space are drawn from a probability distribution, and they tend to be far apart from each other (there is no close points in high dimensional space; this is where the k-NN's assumption can go very wrong)

Illustration: Imagine the unit cube $[0,1]^d$, and all training data points are sampled uniformly at random within the cube.



- We can think of this cube as a box of Vol $1 \times 1 \dots 1 = 1^d$ containing n points

- Let think of $k=10$, Then in expectation we can hope to find $S_{\vec{x}}$ of a test point \vec{x} within a smaller cube whose length is some value l . This means we should imagine that this smaller cube is of volume l^d and it should contain k points out of n points

- $\frac{k}{n} \approx l^d \Rightarrow l \approx (\frac{k}{n})^{\frac{1}{d}}$

Q: How big is l ?

A: $l^d \approx \frac{k}{n} \Rightarrow l \approx (\frac{k}{n})^{\frac{1}{d}}$
 l decreases as n, k

Simulation under $k=10$

n	d	l
1000	2	0.1
1000	10	0.63
1000	100	0.955
1000	1000	0.9934

→ observation: As d becomes larger, all k-NN tend to be far from \vec{x}