

# Ch2 - Vector Space

2020年9月21日 上午 10:46



Ch2\_Vec...

## Ch. 2. Vector spaces, metric and convergence

### § 2.1. Vector spaces

- Definition: A vector space over  $\mathbb{R}$  (the real domain) is a set  $V$  together with two functions:

vector addition:  $+ : V \times V \rightarrow V$  (i.e.,  $x+y$ , where  $x, y \in V$ )

scalar multiplication:  $\cdot : \mathbb{R} \times V \rightarrow V$  (i.e.,  $\alpha \cdot x$ , where  $\alpha \in \mathbb{R}, x \in V$ )  
that satisfying the following.

① Associativity of addition:  $x + (y + z) = (x + y) + z \quad \forall x, y, z \in V$

② Commutativity of addition:  $x + y = y + x \quad \forall x, y \in V$ .

③ Zero vector:  $\exists$  an element, denoted by  $0$ , in  $V$ , s.t.

$$x + 0 = 0 + x = x \quad \forall x \in V.$$

④ Negative vector:  $\forall x \in V$ ,  $\exists$  an element, denoted by  $-x$ , in  $V$ , s.t.

$$x + (-x) = (-x) + x = 0$$

⑤  $\forall x \in V$ ,  $1x = x$ .

⑥  $\forall x \in V$  and  $\alpha, \beta \in \mathbb{R}$ ,  $\alpha(\beta x) = (\alpha\beta)x$ .

⑦  $\forall x \in V$  and  $\alpha, \beta \in \mathbb{R}$ ,  $(\alpha + \beta)x = \alpha x + \beta x$

⑧  $\forall x, y \in V$  and  $\alpha \in \mathbb{R}$ ,  $\alpha(x+y) = \alpha x + \alpha y$   $\square$

Remark: • We can define vector space over  $\mathbb{C}$  (the complex domain) similarly.  
• We will assume vector space over  $\mathbb{R}$  for default. Vector space over  $\mathbb{C}$  is used very rarely.

Example 1:  $\mathbb{R}$  is a vector space, with "+" the standard addition of real numbers and " $\cdot$ " the standard multiplication of real numbers.

Example 2:  $\mathbb{R}^n$  is a vector space, with "+" and " $\cdot$ " defined by:

addition:  $\forall \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$ ,  $\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix}$

scalar multiplication:  $\forall \alpha \in \mathbb{R}$  and  $\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n$ ,  $\alpha \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \alpha x_1 \\ \alpha x_2 \\ \vdots \\ \alpha x_n \end{bmatrix}$ .

Many input data can be modeled by vectors in  $\mathbb{R}^n$ .

- Digital sound signals of length  $n$ .
- Time series of length  $n$ . Stock price in base  $n$
- $n$  different attributes/features of a single thing or object.

Example 3: All real  $m \times n$  matrices is a vector space, with standard matrix addition and standard scalar multiplication.

- This vector space is the same as  $\mathbb{R}^{mn}$ . A long vector
- An  $m \times n$  matrix can be used to represent a black-white digital image of  $m \times n$  pixels.

All real 3-array(tensor) of size  $m \times n \times l$  is a vector space

Example 4: All 3-arrays of size  $m \times n \times l$  is a vector space,

if "+" and ":" is defined by Three dimensional matrix

$$+: \forall x = [x_{ijk}]_{i=1}^n \begin{smallmatrix} m \\ j=1 \\ k=1 \end{smallmatrix} \text{ and } y = [y_{ijk}]_{i=1}^n \begin{smallmatrix} m \\ j=1 \\ k=1 \end{smallmatrix}, 1 \leq i \leq n, 1 \leq j \leq m, 1 \leq k \leq l$$

$$x+y = [x_{ijk} + y_{ijk}]_{i=1}^n \begin{smallmatrix} m \\ j=1 \\ k=1 \end{smallmatrix}$$

$$\cdot : \forall x = [x_{ijk}]_{i=1}^n \begin{smallmatrix} m \\ j=1 \\ k=1 \end{smallmatrix} \text{ and } \alpha \in \mathbb{R},$$

$$\alpha x = [\alpha x_{ijk}]_{i=1}^n \begin{smallmatrix} m \\ j=1 \\ k=1 \end{smallmatrix}$$

- This vector space is the same as  $\mathbb{R}^{mnl}$ .

- An  $n \times m \times 3$  3-array can be used to model a color digital image, where  $x_{ijk}$  means  $(i,j)$ -th pixel in channel  $k$ , and channel 1, 2, 3 means Red, Green, Blue channels of the image.

- An  $n \times m \times l$  3-array can be used to model a <sup>black-white</sup> video, where  $x_{ijk}$  means the  $(i,j)$ -th pixel at  $k$ -th frame.

$$\forall X, Y \in \mathbb{R}^{m \times n}, X+Y = \begin{bmatrix} x_{11} + y_{11} & x_{12} + y_{12} & \dots \\ x_{21} + y_{21} & \dots & \dots \\ \vdots & \dots & \dots \\ x_{m1} + y_{m1} & \dots & \dots \end{bmatrix}$$

$$\forall \alpha \in \mathbb{R}$$

$$\forall X \in \mathbb{R}^{m \times n},$$

$$\alpha X = \begin{bmatrix} \alpha x_{11} & \dots & \alpha x_{1n} \\ \vdots & \ddots & \vdots \end{bmatrix}$$

$$\forall x \in \mathbb{R}^{m \times n}, \quad dX = \begin{bmatrix} dx_{11} & \dots & dx_{1n} \\ \vdots & \ddots & \vdots \\ dx_{m1} & \dots & dx_{mn} \end{bmatrix}$$

Example 5: Consider the set of all strings.

Define the addition by, e.g.,

$$'I' + 'am' = 'I am'$$

and some scalar multiplication.

Not commutative!

Then it doesn't form a vector space.

- Therefore, we cannot use vector space to model text data in this naive way.

How to "vectorize" the texts is a fundamental research topic in text data analysis and natural language processing.

Example 6: The function space  $C[a,b] = \{f \mid f \text{ is continuous on } [a,b]\}$  is a vector space if we define "+" and ":" by:

$$+ : \forall f, g \in C[a,b], \quad (f+g)(t) = f(t)+g(t), \quad \forall t \in [a,b].$$

$$\cdot : \forall f \in C[a,b], \alpha \in \mathbb{R}, \quad (\alpha f)(t) = \alpha f(t).$$

- $C[a,b]$  is referred to as a function space, since any vector in the vector space is a function.

- $C[a,b]$  could be the hypothesis space of a learner with

In unsupervised learning, one input and one output, i.e.,

從 hypothesis space 穩適  $x_i \rightarrow ? \rightarrow y_i$ , with  $x_i \in [a, b]$  and  $y_i \in \mathbb{R}$ .

合嘅 function

Leave a  $f \in C[a,b]$  s.t.  $f(x_i) \approx y_i$  for all  $i$ .

Example 7: The infinite sequence

$$l_\infty = \left\{ \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} \mid \exists \text{ a finite } C \text{ s.t. } |a_i| \leq C \quad \forall i \right\}$$

with: "+" :  $(a+b)_i = a_i + b_i$  for all  $a, b \in l_\infty$  for all  $i$

"." :  $(\alpha a)_i = \alpha a_i$  for all  $\alpha \in \mathbb{R}$ , for all  $i$

It forms a vector space.

$\exists \text{ a finite } C$   
 $\text{ s.t. } |a_i| \leq C$

因為+完同乘完constant之後都會被bounded,所以成立

- This vector space can be used to model, e.g.,  
time series with a very long time and/or a very fine time  
resolution.

A very long time signal or very very fine resolution

## § 1.2. Metric in vector spaces

In order to do calculus on vector spaces, we need to define 'distance, closeness between vectors.'

Let  $V$  be a vector space. Let  $x, y \in V$ . Then,

$$\text{distance}(x, y) = \text{distance}(x-y, y-y) = \text{distance}(x-y, 0)$$

(distance should be shift invariant.)

Shift invariant = 搬  $x$  同  $y$  去其他地方都唔應該影響呢個 distance

Therefore, to define a distance, we only need to define a length for each vector in  $V$ .

Let  $x \in V$ . Let  $\|x\|$  be its length, called **norm**, which should satisfy

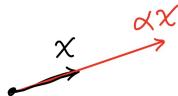
① A length should be nonnegative, i.e.

$$\|x\| \geq 0 \quad \forall x \in V.$$

Moreover, only the zero vector has a zero length, i.e.,  
 $\|x\| = 0 \Leftrightarrow x = 0$ .

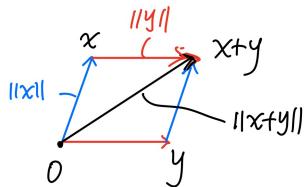
② The length of a multiple of a vector should be the multiple of the length of the vector, i.e.,

$$\forall \alpha \in \mathbb{R}, \quad \|\alpha x\| = |\alpha| \|x\|$$



③ Triangular inequality: the length of the direct path is the smallest

$$\|x+y\| \leq \|x\| + \|y\|$$



Definition: Let  $V$  be a vector space over  $\mathbb{R}$ . A norm on  $V$  is a function  $\|\cdot\| : V \rightarrow \mathbb{R}$  such that:

- ①  $\|x\| \geq 0, \forall x \in V$  and  $\|x\|=0 \Leftrightarrow x=0$ .
- ②  $\|\alpha x\| = |\alpha| \|x\|, \forall x \in V$  and  $\alpha \in \mathbb{R}$ .
- ③  $\|x+y\| \leq \|x\| + \|y\|, \forall x, y \in V$ .

Example 1:  $\mathbb{R}$  is a vector space over  $\mathbb{R}$ .

Let  $\|x\| = |x| \quad \forall x \in \mathbb{R}$ . Then it is a norm on  $\mathbb{R}$ .

(Can you find other norms on  $\mathbb{R}$ ? ) 2 \* absolute 就係

Example 2:  $\mathbb{R}^n$  is a vector space over  $\mathbb{R}$ .

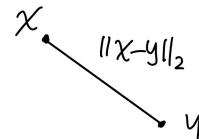
There are many norms on  $\mathbb{R}^n$ .

\* 2-norm: (Euclidean norm)

$$\|x\|_2 = \left( \sum_{i=1}^n x_i^2 \right)^{\frac{1}{2}}$$

The induced distance

$$\|x-y\|_2 = \left( \sum_{i=1}^n (x_i - y_i)^2 \right)^{\frac{1}{2}} \text{ is the Euclidean distance}$$



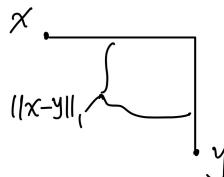
\* 1-norm:

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$

The induced distance

$$\|x-y\|_1 = \sum_{i=1}^n |x_i - y_i|$$

is known as Manhattan distance (You can walk only horizontally and vertically)

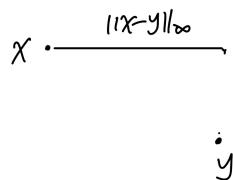


\*  $\infty$ -norm:

$$\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$$

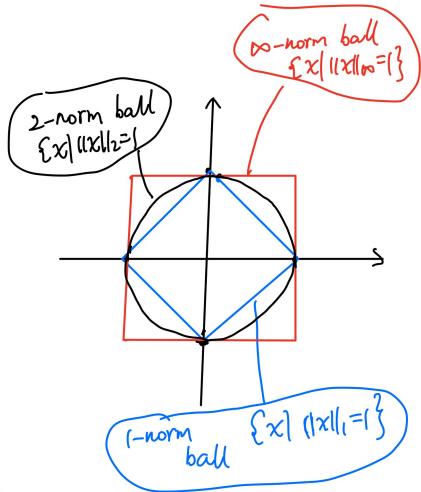
The induced distance is

$$\|x-y\|_\infty = \max_{1 \leq i \leq n} |x_i - y_i|$$



\*  $p$ -norm ( $p \geq 1$ )  
 $\|x\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$

- Comparison of unit balls.



- Note that  $(\mathbb{R}^n, \|\cdot\|_1)$ ,  $(\mathbb{R}^n, \|\cdot\|_2)$ ,  $(\mathbb{R}^n, \|\cdot\|_\infty)$ , ... are all different normed spaces. So, for a given vector space, we can obtain various normed space by choosing different norms.

- Calculate the norms of  $x = \begin{pmatrix} 3 \\ 4 \end{pmatrix}$

$$\|x\|_2 = (3^2 + 4^2)^{1/2} = 5 \quad \|x\|_1 = |3| + |4| = 7$$

$$\|x\|_\infty = \max\{3, 4\} = 4.$$

Example 3:  $\mathbb{R}^{m \times n}$  is a vector space

- Since  $\mathbb{R}^{m \times n}$  can be viewed as  $\mathbb{R}^{mn}$ , we can define vector  $p$ -norms for matrices

$$\|A\|_{p,\text{vector}} = \left( \sum_{j=1}^n \sum_{i=1}^m |a_{ij}|^p \right)^{1/p}$$

- $p=1$ : sum of absolute values of entries of  $A$ .
- $p=2$ : sum of squares and then take square root  
(also known as Frobenius norm)

$$\|A\|_F = \left( \sum_{j=1}^n \sum_{i=1}^m |a_{ij}|^2 \right)^{\frac{1}{2}}$$

$\rightarrow p=\infty$ : max entry of  $A$  in magnitude

$\rightarrow$  calculate the Frobenius norm of  $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \in \mathbb{R}^{2 \times 2}$

$$\|A\|_F = \left( \sum_{j=1}^2 \sum_{i=1}^2 |a_{ij}|^2 \right)^{\frac{1}{2}} = (1^2 + 2^2 + 3^2 + 4^2)^{\frac{1}{2}} = \sqrt{30}$$

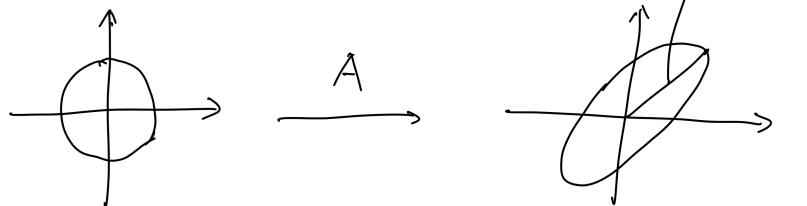
$A \in \mathbb{R}^{m \times n}$  is also a linear transformation from  $\mathbb{R}^n$  to  $\mathbb{R}^m$

Define  $p$ -norms on  $\mathbb{R}^n$  and  $\mathbb{R}^m$  respectively.

We can define matrix  $p$ -norm of a matrix by

$$\begin{aligned} \|A\|_{p \rightarrow p} &= \sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{\|Ax\|_p}{\|x\|_p} \quad (\text{check it is a norm!}) \\ &= \sup_{\substack{x \in \mathbb{R}^n \\ \|x\|_p=1}} \|Ax\|_p \end{aligned}$$

$p=2$



Find the vector 2-norm (Frobenius norm) and the matrix

所以如果A是一個Data Matrix

咁 $\|A\|_F$  makes more sense

如果係linear transformation

咁 $\|A\|_2$  makes more sense

$$2\text{-norm of } A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\text{Frobenius norm : } \|A\|_F = (1^2 + 0^2 + 0^2 + 1^2)^{\frac{1}{2}} = \sqrt{2}$$

$$\text{matrix 2-norm : } \|A\|_{2 \rightarrow 2} = \sup_{\substack{x \in \mathbb{R}^2 \\ x \neq 0}} \frac{\|Ax\|_2}{\|x\|_2} = \sup_{\substack{x \in \mathbb{R}^2 \\ x \neq 0}} \frac{\|x\|_2}{\|x\|_2} = 1$$

Example 4:  $C[a, b]$  is a vector space over  $\mathbb{R}$ .

The 0 vector in  $C[a, b]$  is the function that takes value 0 on  $[a, b]$ .

To measure how large a function  $f$  is,

we can use the following norms

$$\|f\|_\infty = \sup_{x \in [a,b]} |f(x)|$$

Then distance of two functions  $f, g \in C[a,b]$

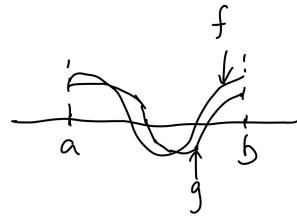
$$\text{is } \|f-g\|_\infty = \sup_{x \in [a,b]} |f(x)-g(x)|.$$

Some other norms of  $C[a,b]$  can be

$$\|f\|_1 = \int_a^b |f(x)| dx$$

$$\|f\|_2 = \left( \int_a^b |f(x)|^2 dx \right)^{\frac{1}{2}}$$

$$\|f\|_p = \left( \int_a^b |f(x)|^p dx \right)^{\frac{1}{p}}$$



Example 5:  $\ell_\infty$  is a vector space over  $\mathbb{R}$

$$\forall a \in \ell_\infty, \text{ define } \|a\|_\infty = \sup_i |a_i|$$

It is a norm on  $\ell_\infty$ , called  $\ell_\infty$ -norm.

- Similarly, for any infinite sequence  $a = (a_i)_{i \in \mathbb{Z}}$

$$\text{we define } \|a\|_p = \left( \sum_{i=1}^{\infty} |a_i|^p \right)^{\frac{1}{p}}, \text{ where } p \geq 1.$$

Consider the space  $\ell_p = \{a \mid \|a\|_p < +\infty\} \subset \ell_\infty$

We can show that  $\ell_p$  is a vector space

and  $\|\cdot\|_p$  is a norm on  $\ell_p$ , call  $\ell_p$ -norm.

- Example:

$$a = \begin{pmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \end{pmatrix}$$

Then  $\|a\|_\infty = \sup_i |a_i| = 1$  呢個係從  $\sin(x)$  Taylor series 推的

$$\|a\|_2 = \left( \sum_{i=1}^{\infty} (a_i)^2 \right)^{\frac{1}{2}} = \left( \frac{\pi^2}{6} \right)^{\frac{1}{2}} = \frac{\pi}{\sqrt{6}}$$

$$\|a\|_1 = \left( \sum_{i=1}^{\infty} \frac{1}{i} \right) = +\infty$$

Therefore,  $a \in l_\infty$ ,  $a \in l_2$ , but  $a \notin l_1$

Remarks: 1. For the same vector space, we can define infinitely many different norms.

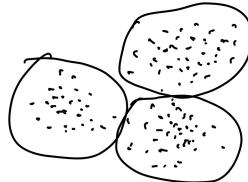
2. A common technique in machine learning to find the solution is to optimize some norm of the unknown. Different norms lead to very different solutions.

## § 1.2 Case study: Clustering, k-means, k-medians

### Clustering

Suppose we are given  $N$  vectors  $x_1, x_2, \dots, x_N \in \mathbb{R}^n$

The goal of clustering is to group or partition the vectors into  $K$  groups or clusters, with the vectors in each group close to each other.



- We use  $\mathbb{R}^n$  because it is simple yet able to model a variety of data sets (e.g., signals, images, videos, attributes of things)
- Actually, the methods can be extended to any normed spaces.
- Applications:
  - Topic discovery. Suppose the  $N$  vectors are word histograms with  $N$  documents respectively, i.e., the  $j$ -th component in  $x_i$  is the counts of the  $j$ -th word in document  $i$ .  
A clustering algorithm partitions the documents into  $K$  groups, which typically can be interpreted as groups of documents with the same topics, genre, or author.
  - Patient Clustering. If  $\{x_i\}_{i=1}^N$  are feature vectors associated with  $N$  patients admitted to a hospital, a clustering algorithm clusters the patient into  $K$  groups of similar patients.
  - Recommendation system. A group of  $N$  people respond to ratings of  $n$  movies. A clustering algorithm can be used to cluster the people into  $K$  groups, each with similar taste.

Then we can recommend new movies liked by someone to people in the same group as him/her.

— Many other applications

Mathematical formulation:

c係每個point屬於咩group

G係每個group有咩point

• Representation:

Let  $C_i \in \{1, 2, \dots, k\}$  be the

group that  $X_i$  belongs to,  
 $i=1, 2, \dots, N$

Let  $G_j \subset \{1, 2, \dots, N\}$  be the data points of group  $j$ .  
 $j=1, 2, \dots, k$ .

then  $\boxed{G_j = \{i \mid C_i=j\}}$  and  $\boxed{C_i = j \text{ for all } i \in G_j}$

Then, group  $j$ , denoted by

Then, to find the grouping, we only need to find  $C_i$ ,  $i=1, \dots, N$

We assign each group a

or  $\boxed{G_j, j=1, \dots, N}$

The representative vectors are not necessarily one of given vectors.

• Evaluation: 白話: 同一group嘅vector應該close啲

First of all, within one specific group  $j$ ,  $G_j$ , all vectors should be close to the representative vector  $z_j$ . More precisely, let

$$J_j = \sum_{i \in G_j} \|X_i - z_j\|_2^2$$

Then,  $J_j$  should be small.

$J_j$  係 total closeness for group  $j$ , 但我地通常唔想用呢舊野因為用2-norm有可能會令呢舊野

undifferentiable(not smooth), 所以我地take square

Secondly, consider all groups, since each  $J_j$  is small,

$$J = J_1 + J_2 + \dots + J_k$$

should be small.

Altogether, we solve the following

姐係咁多group一齊都要最細

$$\min_{\substack{G_1, \dots, G_k \\ z_1, \dots, z_k}} J \iff \min_{\substack{G_1, \dots, G_k \\ z_1, \dots, z_k}} \sum_{j=1}^k J_j \iff \min_{\substack{G_1, \dots, G_k \\ z_1, \dots, z_k}} \sum_{j=1}^k \left( \sum_{i \in G_j} \|X_i - z_j\|_2^2 \right)$$

Optimization.

We may use an alternating minimization to solve the minimization.

Step 1: Fix the representatives  $z_1, \dots, z_k$ , find the best partitions

$G_1, \dots, G_k$ , i.e., solve

$$\min_{\substack{G_1, \dots, G_k \\ z_1, \dots, z_k}} \sum_{j=1}^k \left( \sum_{i \in G_j} \|X_i - z_j\|_2^2 \right) \quad \dots \quad (1)$$

fix咗一個然後計第二個

min擺底嘅野

代表呢嘅野係adjustable嘅

可以當係given (= parameter)

Step 2: Fix the groups  $G_1, \dots, G_K$ , find the best representatives

$\bar{z}_1, \dots, \bar{z}_K$ , i.e., solve

$$\min_{\bar{z}_1, \dots, \bar{z}_K} \sum_{j=1}^K \left( \sum_{i \in G_j} \|x_i - \bar{z}_j\|_2^2 \right) \quad \dots \quad (2)$$

The two steps are repeated until convergence.

Let's find the solutions of the sub-problems (1) and (2) respectively.

For (1):

finding the partition  $G_1, \dots, G_K$  is equivalent to

finding  $c_1, c_2, \dots, c_N$ . So (1) becomes

$$\min_{c_1, c_2, \dots, c_N} \left( \underbrace{\|x_1 - z_{c_1}\|_2^2}_{\text{depends on } c_1 \text{ only}} + \underbrace{\|x_2 - z_{c_2}\|_2^2}_{\text{depends on } c_2 \text{ only}} + \dots + \underbrace{\|x_N - z_{c_N}\|_2^2}_{\text{depends on } c_N \text{ only}} \right)$$

$G_j$ 係 disjoint嘅，淨係出現一次

↓

$$\min_{c_i} \|x_i - z_{c_i}\|_2^2 \quad i = 1, 2, \dots, N.$$

Since  $c_i \in \{1, 2, \dots, K\}$ , to get  $c_i$ , we only need to compare

$$\|x_i - z_1\|_2^2, \|x_i - z_2\|_2^2, \dots, \|x_i - z_K\|_2^2$$

and choose the minimum from it. i.e.,

$$c_i = \arg \min_{j \in \{1, \dots, K\}} \|x_i - z_j\|_2^2, \quad i = 1, 2, \dots, N.$$

In other words,

$x_i$  is assigned to the group whose representative vector is the closest to  $x_i$ .

For (2): It is rewritten as

$$\min_{\bar{z}_1, \dots, \bar{z}_K} \left( \underbrace{\sum_{i \in G_1} \|x_i - \bar{z}_1\|_2^2}_{\text{depends on } \bar{z}_1 \text{ only}} + \underbrace{\sum_{i \in G_2} \|x_i - \bar{z}_2\|_2^2}_{\text{depends on } \bar{z}_2 \text{ only}} + \dots + \underbrace{\sum_{i \in G_K} \|x_i - \bar{z}_K\|_2^2}_{\text{depends on } \bar{z}_K \text{ only}} \right)$$

Obviously, it is equivalent to minimize each term independently,

$$(1) \Leftrightarrow \min_{c_1, \dots, c_N} \sum_{j=1}^K \sum_{i \in G_j} \|x_i - \bar{z}_j\|_2^2$$

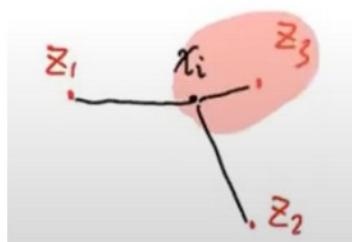
$$\Leftrightarrow \min_{c_1, \dots, c_N} (\|x_1 - \bar{z}_{c_1}\|_2^2 + \|x_2 - \bar{z}_{c_2}\|_2^2 + \dots + \|x_N - \bar{z}_{c_N}\|_2^2)$$

(其實呢到就係kmean個定義黎 $z_{c_1}$ 係指 $x_i$ 呢點所屬嘅group嘅representative, 每個點同佢group相減得出total distance)

$$\Leftrightarrow \min_{C_i \in \{1, \dots, k\}} \|x_i - z_{C_i}\|_2^2, \quad 1 \leq i \leq N$$

$$\Leftrightarrow C_i = \arg \min_{j=1, \dots, k} \|x_i - z_j\|_2^2$$

另一種寫法，咁樣loop唔洗compare all point, 只係compare all  $z_j$



例子如左圖:  $x_i$  穩佢同邊個 $z_j$ 最近即可

Then we set  $G_j = \{i \mid C_i = j\}, j = 1, \dots, k$

Step 2 的推導:

$$(2) \Leftrightarrow \min_{z_1, \dots, z_k} \sum_{j=1}^k \sum_{i \in G_j} \|x_i - z_j\|_2^2 \quad \text{係交換summation嘅次序, 將} j=1 \text{tok sub落去}$$

$$\Leftrightarrow \min_{z_1, \dots, z_k} \left( \sum_{i \in G_1} \|x_i - z_1\|_2^2 + \sum_{i \in G_2} \|x_i - z_2\|_2^2 + \dots + \sum_{i \in G_k} \|x_i - z_k\|_2^2 \right)$$

$$\Leftrightarrow \min_{z_j} \sum_{i \in G_j} \|x_i - z_j\|_2^2, \quad j = 1, \dots, k \quad \begin{array}{l} \text{係觀察到} z_1..z_k \text{都係獨立事件} \\ \text{可以抽出黎sum} \end{array}$$

例子:



已知  
呢三個應該要同一group,  
要移到 $z_j$ , s.t.  $\|x_i - z_j\|_2^2$  is minimized

i.e., solve  $K$  independent problems.

$$\min_{z_j} \left( \sum_{i \in G_j} \|x_i - z_j\|_2^2 \right), \quad j = 1, 2, \dots, k.$$

Note that

$$\begin{aligned} \sum_{i \in G_j} \|x_i - z_j\|_2^2 &= \sum_{i \in G_j} \sum_{l=1}^n (x_{il} - z_{jl})^2 \\ &= \sum_{l=1}^n \left( \sum_{i \in G_j} (x_{il} - z_{jl})^2 \right) \end{aligned}$$

$x_{il}$  are  $l$ -th component  
 $z_{jl}$  of  $x_i$  respectively

即  $l$  从 1 到  $n$  的求和是 independent sum

i.e., solve  $K$  independent problems.

$$\min_{\bar{z}_j} \left( \sum_{i \in G_j} \|x_i - \bar{z}_j\|_2^2 \right), \quad j=1, 2, \dots, K.$$

Note that

$$\begin{aligned} \sum_{i \in G_j} \|x_i - \bar{z}_j\|_2^2 &= \sum_{i \in G_j} \sum_{l=1}^n (x_{il} - \bar{z}_{jl})^2 \\ &= \sum_{l=1}^n \left( \sum_{i \in G_j} (x_{il} - \bar{z}_{jl})^2 \right) \end{aligned}$$

(  $x_{il}$  are  $l$ -th component  
of  $\bar{z}_j$  respectively )

Each term in this summation are independent again.

$$\text{Thus, } \min_{\bar{z}_j} \left( \sum_{i \in G_j} \|x_i - \bar{z}_j\|_2^2 \right) \Leftrightarrow \min_{\bar{z}_{jl}} \sum_{i \in G_j} (x_{il} - \bar{z}_{jl})^2, \quad l=1, 2, \dots, n.$$

↑  
One variable minimization.

Taking derivative w.r.t.  $\bar{z}_{jl}$  and setting it to 0, we obtain that the solution  $\bar{z}_{jl}$  satisfies

$$\begin{aligned} 2 \sum_{i \in G_j} (\bar{z}_{jl} - x_{il}) &= 0 \\ \Rightarrow \bar{z}_{jl} &= \left( \sum_{i \in G_j} x_{il} \right) / |G_j| \quad ( |G_j| \text{ is the number of elements in } G_j ) \\ l &= 1, 2, \dots, n. \end{aligned}$$

In vector form,

$$\begin{pmatrix} \bar{z}_{j1} \\ \bar{z}_{j2} \\ \vdots \\ \bar{z}_{jn} \end{pmatrix} = \frac{1}{|G_j|} \cdot \sum_{i \in G_j} \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{in} \end{pmatrix} \Leftrightarrow \bar{z}_j = \frac{1}{|G_j|} \left( \sum_{i \in G_j} x_i \right)$$

In other words,

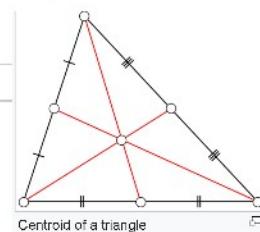
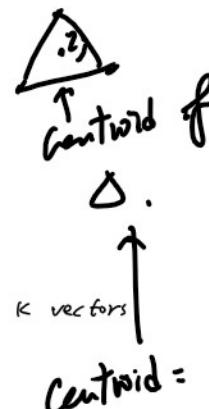
$\bar{z}_j$  is the mean of all vectors in  $G_j$ .

Altogether, we get the following clustering algorithm.

Input:  $x_1, x_2, \dots, x_N \in \mathbb{R}^n$ .

Output:  $C_1, C_2, \dots, C_K$  and  $\bar{z}_j, j=1, \dots, K$ .

Initialization: Initialize  $\bar{z}_1, \bar{z}_2, \dots, \bar{z}_K$  by choosing  $K$  vectors



from  $x_1, x_2, \dots, x_N$  randomly

Step 1: Given  $\bar{z}_1, \bar{z}_2, \dots, \bar{z}_K$ , compute initialization

$$C_i = \arg \min_{j \in \{1, 2, \dots, K\}} \|x_i - \bar{z}_j\|_2^2, \quad i=1, 2, \dots, N.$$

and update

人說：

人說：

$$c_i = \arg \min_{j \in \{1, 2, \dots, k\}} \|x_i - z_j\|_2^2, \quad i=1, 2, \dots, N.$$

and define

$x_i$  is assigned to the group whose representative vector is the closest to  $x_i$ .

$$G_j = \{i \mid c_i = j\}, \quad j=1, 2, \dots, k.$$

Step 2: Given  $G_1, G_2, \dots, G_k$ , compute

$$z_j = \frac{1}{|G_j|} \left( \sum_{i \in G_j} x_i \right)$$

Go back to step 1.

$c_i = j$  such that  $z_j$  is the closest representative to  $x_i$

This algorithm is known as "k-means" algorithm, because it computes K means of vectors at step 2.

### K-medians Algorithm

In k-means, the Euclidean norm is used. We can replace it by  $\ell$ -norm. We solve

$$\min_{\substack{G_1, \dots, G_k \\ z_1, \dots, z_k}} \sum_{j=1}^k \left( \sum_{i \in G_j} \|x_i - z_j\|_1 \right)$$

The numerical solver is

Step 1: Fix  $z_1, \dots, z_k$ , solve

$$\min_{\substack{G_1, \dots, G_k \\ z_1, \dots, z_k}} \sum_{j=1}^k \left( \sum_{i \in G_j} \|x_i - z_j\|_1 \right).$$

Similar to the discussion in k-means, the solution is

$$c_i = \arg \min_{j \in \{1, 2, \dots, k\}} \|x_i - z_j\|_1, \quad i=1, 2, \dots, N.$$

and  $G_j = \{i \mid c_i = j\}$ .

Step 2: Fix  $G_1, G_2, \dots, G_k$ , solve

$$\min_{\substack{z_1, \dots, z_k \\ G_1, \dots, G_k}} \sum_{j=1}^k \left( \sum_{i \in G_j} \|x_i - z_j\|_1 \right)$$