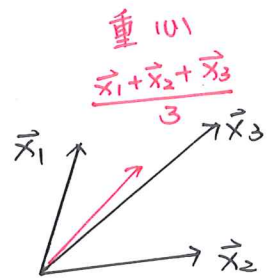


Linear algebra: k-means clustering

$$\vec{x}_1, \vec{x}_2, \dots, \vec{x}_t \in \mathbb{R}^d$$

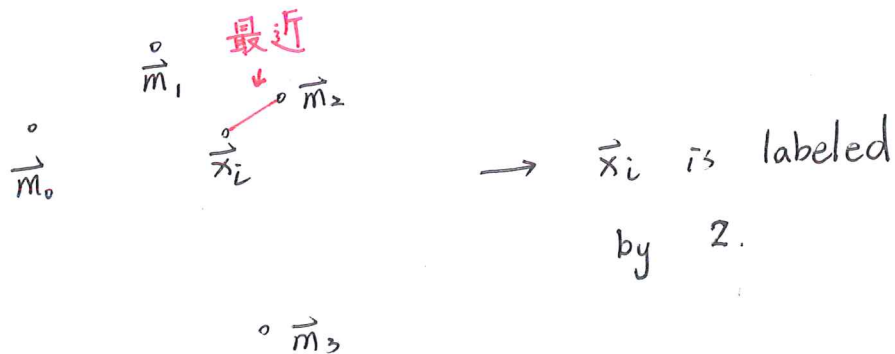
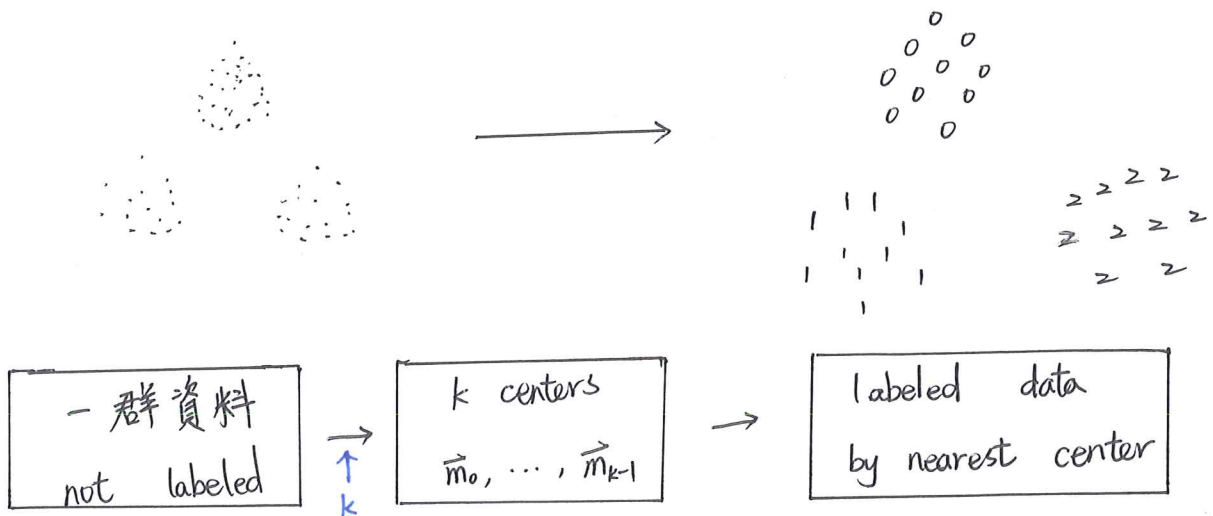
The center is $\frac{1}{t} (\vec{x}_1 + \vec{x}_2 + \dots + \vec{x}_t)$
mean



The distance between \vec{x} and \vec{y}

$$\begin{aligned} \text{is } |\vec{x} - \vec{y}| &= \sqrt{(\vec{x} - \vec{y}) \cdot (\vec{x} - \vec{y})} \\ &= \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_d - y_d)^2} \end{aligned}$$

clustering (unsupervised)

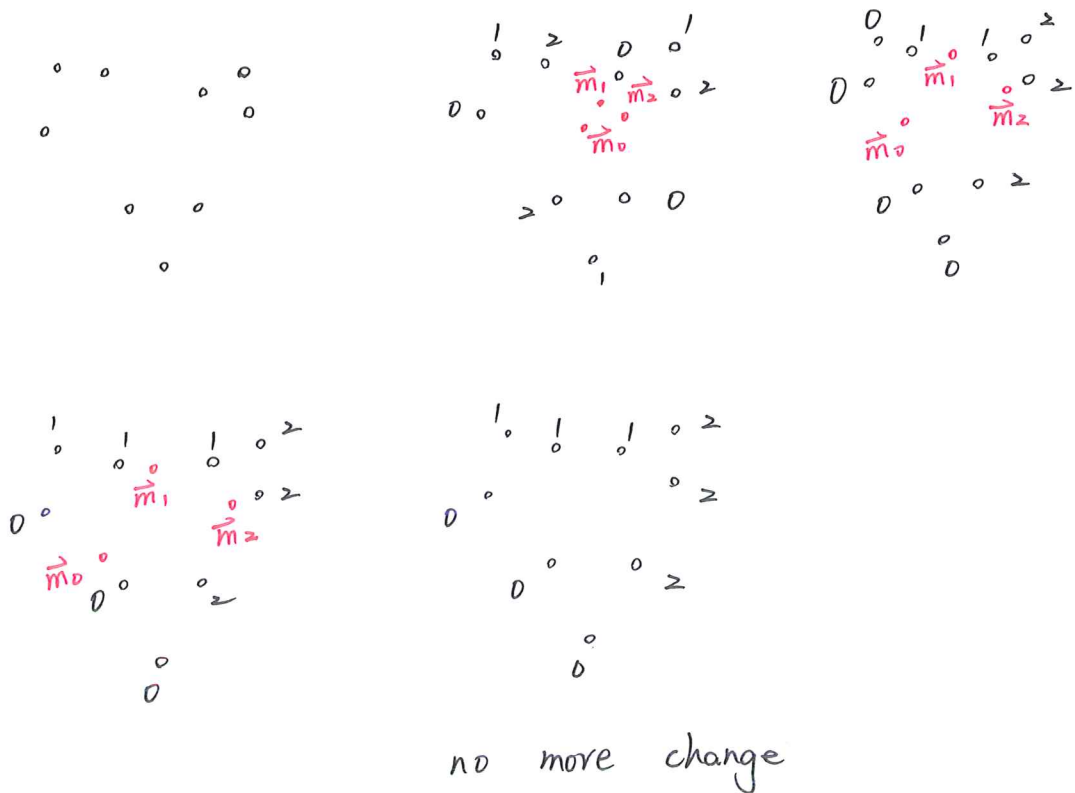


k-means clustering algorithm

Input: unlabeled data $\vec{x}_1, \dots, \vec{x}_N$, and integer k .
(number of groups)

Output: label $\vec{y} = (y_1, \dots, y_N)$, each $y_i \in \{0, 1, \dots, k-1\}$

- ① Assign random \vec{y} ($y_i \in \{0, \dots, k-1\}$)
- ② For each group j $\{\vec{x}_i : y_i = j\}$,
compute the center \vec{m}_j .
- ③ For each \vec{x}_i , if \vec{m}_j is the nearest center
assign $y_i = j$.
- ④ Repeat ②, ③ until \vec{y} does not change.



Thm.

The k-means clustering algorithm will always stop.

pf.

$$\text{Err}(\vec{m}_0, \dots, \vec{m}_{k-1}) = \sum |\vec{x}_i - \vec{m}_{y_i}|^2$$

Claim: Err will strictly decrease.

e.g. If y_i is changed from 0 to 1,



\vec{x}_i

\vec{m}_1 \vec{m}_0

$$\Rightarrow |\vec{x}_i - \vec{m}_1| < |\vec{x}_i - \vec{m}_0|$$

$$\Rightarrow \text{new Err} = \text{Err} - |\vec{x}_i - \vec{m}_0|^2 + |\vec{x}_i - \vec{m}_1|^2 < \text{Err}$$

Claim: There are only k^N different labels.

$$\vec{y} = (y_1, \dots, y_N)$$

$\uparrow \uparrow \uparrow \uparrow$
0 ~ k-1 k^N

The algorithm will reach min Err and stop.