

Class (Sept-10, Thursday)

---

---

---

---

---

---



Application of concepts of distance / similarity.

task: certain collection of objects is given.

Put "similar" objects in one class.

vectors / feature vectors  $\in \mathbb{R}^n$

collection vectors which are "close"

Objects  $\leftrightarrow$  vectors

similarity  $\leftrightarrow$  distance

groups  $\leftrightarrow$  clusters (unsupervised / semi-supervised)

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

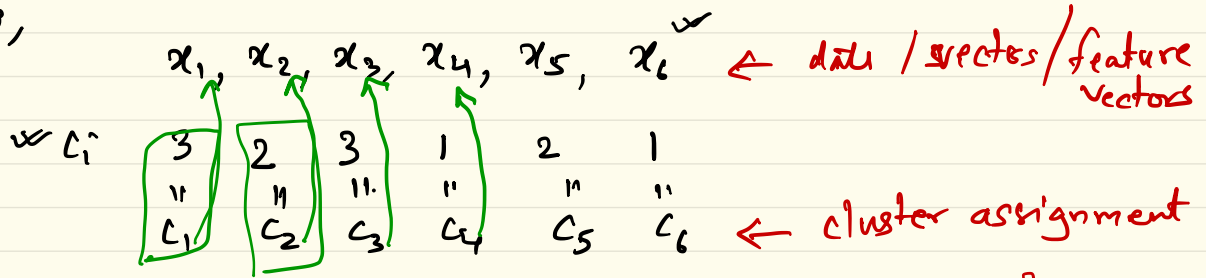
$c_1, c_2, \dots, c_n \leftarrow$  numbers / cluster assignment.

$x_1, x_2, \dots, x_N \in \mathbb{R}^n$  (Data / vectors)

Total - 3 clusters.

$c_i \in \{1, 2, 3\}$  for  $i=1, 2, \dots, N$

$N = 6,$



$G_1 = \{4, 5\}, G_2 = \{x_2, x_5\}, G_3 = \{1, 3\}$

clusters


cluster representative: (Group representative)

$N$  - vectors,  $k$  - clusters

$N$  - data points       $k$  - clusters

Let  $z_1, z_2, \dots, z_k$  : cluster representatives.

Representatives are "close" to the vectors in a cluster.

$\|x_i - z_{c_i}\|$  should be small.  


$x_1, x_2, \dots, x_N$  - vectors.

$c_1, c_2, \dots, c_N$  - cluster assignments \*  
 $c_1, \dots, c_N \in \{1, 2, \dots, k\}$

$z_1, z_2, \dots, z_k$  : cluster representatives.

$z_{c_i} \in \mathbb{R}^n$

$$N=100, k=5$$

$x_1, x_2, \dots, x_{100}$

1 5 1

" " "

$c_1, c_2, \dots$

$c_{100}$

$x_i \longleftrightarrow c_i$

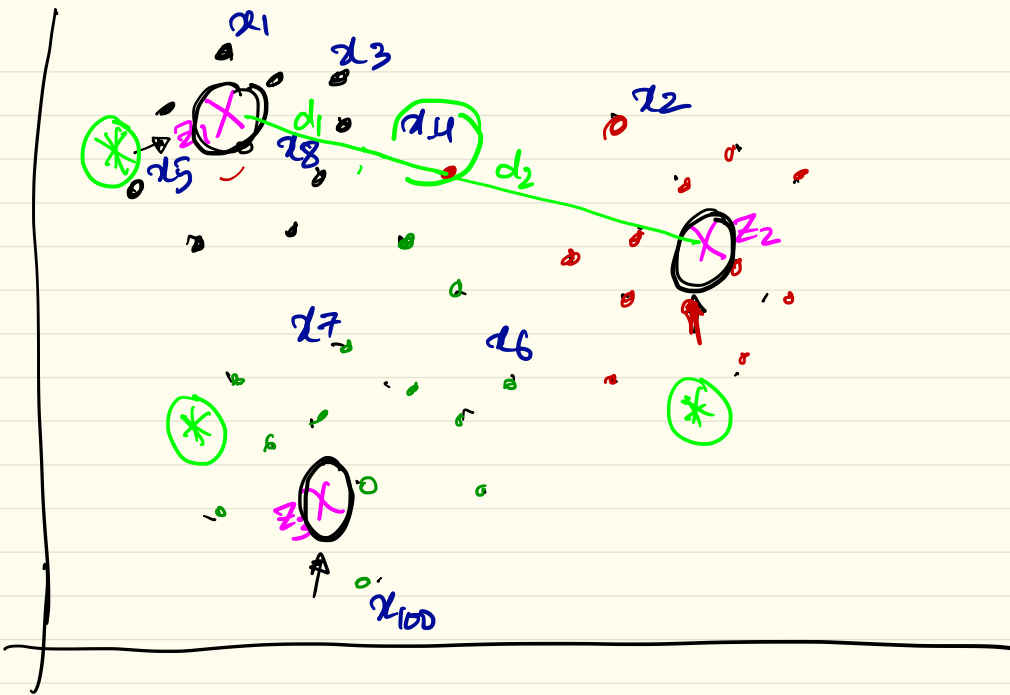
$z_{c_i}$

$c_i=1,$

$z_{c_i} = z_1$

clustering objective:

$$J^{\text{clust}} = \frac{(\|x_1 - z_{c_1}\|^2 + \dots + \|x_N - z_{c_N}\|^2)}{N}$$



$$\underline{k=3}$$

$$N=100$$

$z_1, z_2, z_3$   
: cluster  
representatives.

$$c_1, c_2, \dots, c_{100}$$

$$\in \{1, 2, 3\}$$

$$d_1 = d_2$$

$$x_1$$

$$c_1 = 1$$

$$z_{c_1} = z_1$$

$$\|x_1 - z_{c_1}\|$$

$$x_2$$

$$c_2 = 2$$

$$z_{c_2} = z_2$$

$$\|x_2 - z_{c_2}\|$$

$$x_3$$

$$c_3 = 1$$

$$z_{c_3} = z_1$$

$$\|x_3 - z_{c_3}\|$$

$$J^{\text{clust}} = \frac{\|x_1 - z_{c_1}\|^2 + \|x_2 - z_{c_2}\|^2 + \dots + \|x_N - z_{c_N}\|^2}{N}$$

$J^{\text{clust}}$  = sum of squares  $\geq 0$

= 0 (best case)

clustering objective: minimize  $J^{\text{clust}}$ .

$k$ : How many number of clusters. ( $k \leq N$ ) ← (1)  
 $c_1, \dots, c_N$ : assignment. ← (\*)  
 $z_1, \dots, z_k$ : choice of cluster representatives. ← (2)

Optimal clustering: minimize  $J^{\text{clust}}$  : decision variables ( $z$ )

Suboptimal clustering:

