

AI

AI VIET NAM

@aivietnam.edu.vn

# K-Means

Quang-Vinh Dinh  
Ph.D. in Computer Science

# Objectives

## Introduction



Data processing  
and select K



Select K centroids



Assign each data  
point to a cluster

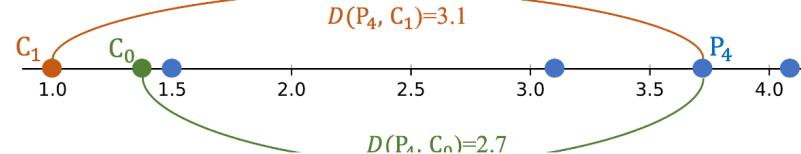


Update centroids



## Tabular Data

	Index	Length	Distance to C0	Distance to C1	Centroid
C <sub>0</sub>	0	1.4	0.0	0.4	0
C <sub>1</sub>	1	1	0.4	0.0	1
	2	1.5	0.1	0.5	-
	3	3.1	1.7	2.1	-
	4	3.7	2.4	2.8	-
	5	4.1	2.7	3.1	-



## Image&Text Data



# Outline

SECTION 1

## Introduction

SECTION 2

## Tabular Data

SECTION 3

## Image&Text Data

Data processing  
and select K

Select K centroids

Assign each data  
point to a cluster

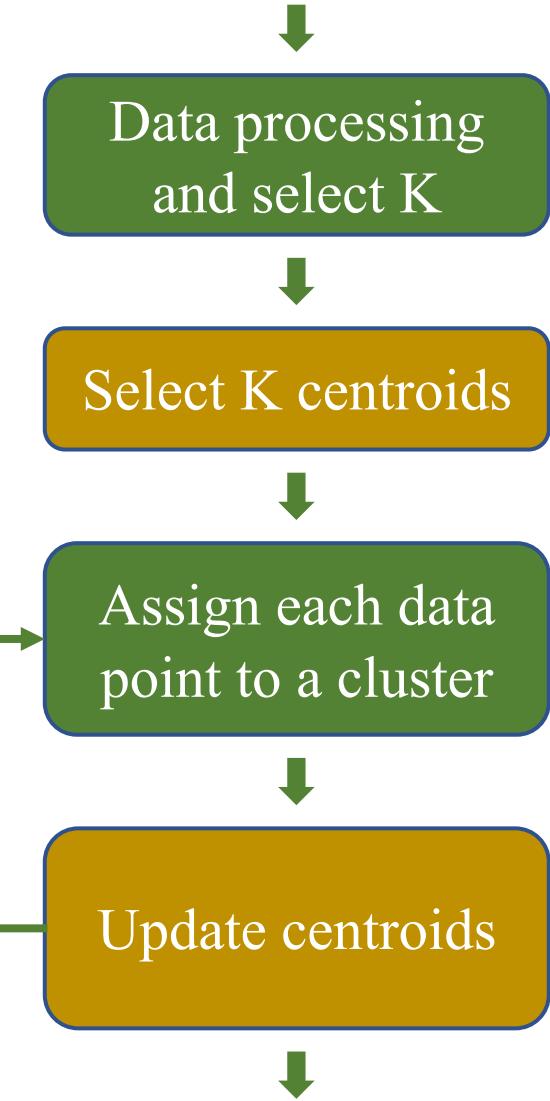
Update centroids



# K-Means

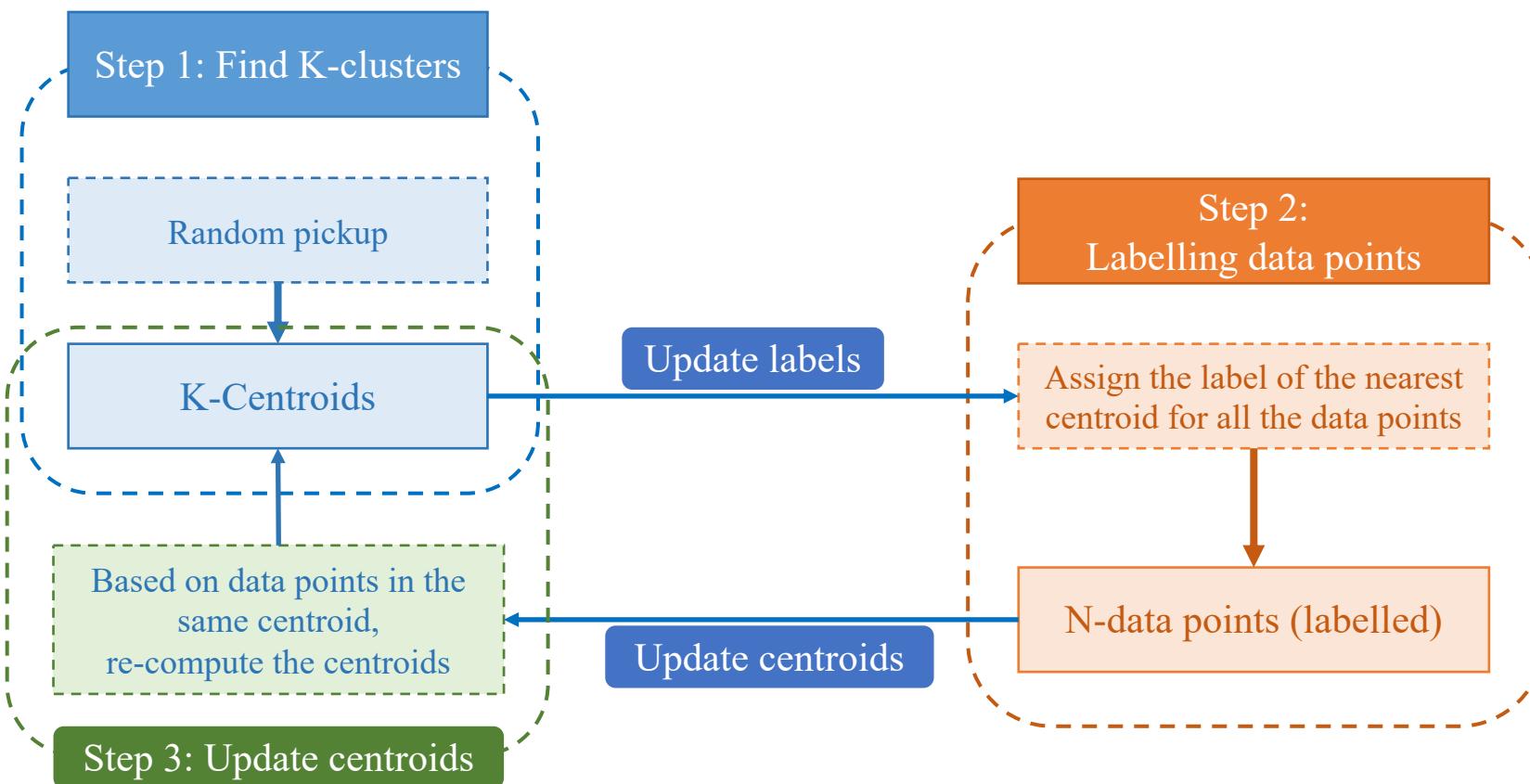
## ❖ Procedure

1. Initialize the value of k
2. Select random K centroids.
3. Assign each data point to the closest centroid
4. Calculate a new centroid of each cluster
5. If not convergent, go to step 3. If yes, stop.



# K-Means

## ❖ K-means algorithm



# K-Means

## ❖ Formulae

Given  $(x_1, \dots, x_n)$  where  $x_i \in \mathbb{R}^d$

Given an initial set of k centroids  $c_0^{(1)}, c_1^{(1)}, \dots$

Let  $S = (S_0, S_1, \dots)$  be k clusters

$$\operatorname{argmin}_S \sum_{i=0}^{k-1} \sum_{x \in S_i} \|x - c_i\|^2$$

Norm

$$\|\vec{u}\| = \sqrt{u_1^2 + \dots + u_n^2}$$

Squared norm

$$\|\vec{u}\|^2 = u_1^2 + \dots + u_n^2$$

# K-Means

## ❖ Formulae

Given  $(x_1, \dots, x_n)$  where  $x_i \in \mathbb{R}^d$

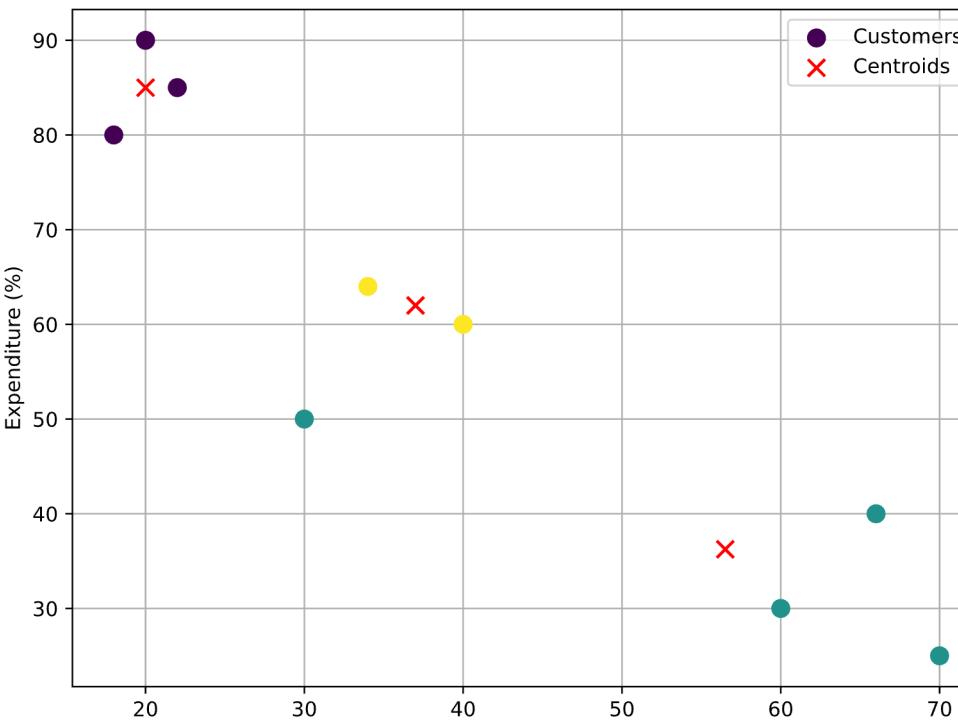
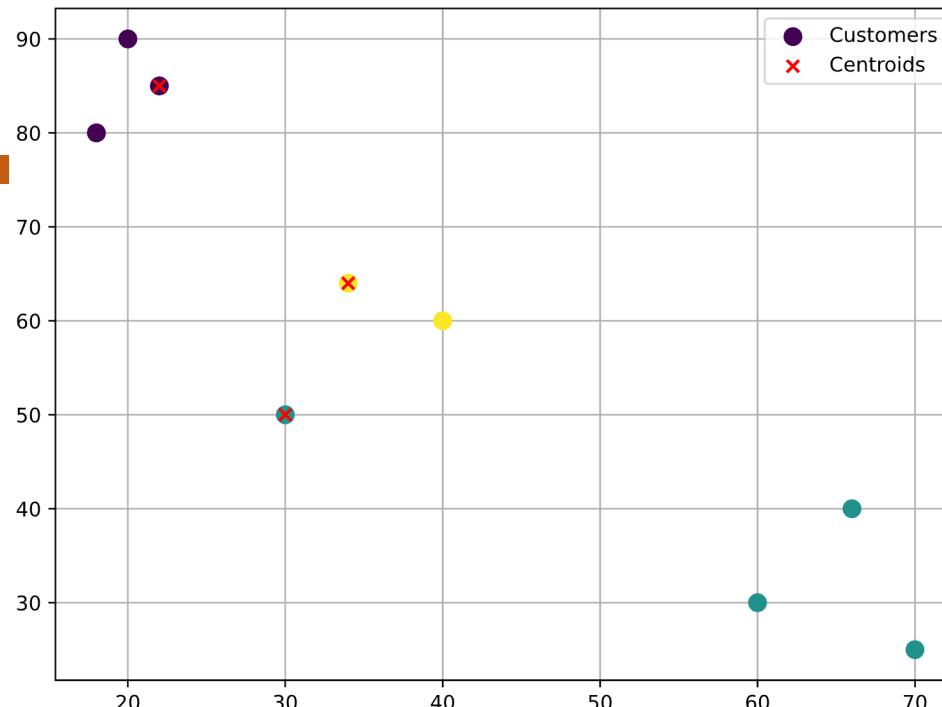
Given an initial set of k centroids  $c_0^{(1)}, c_1^{(1)}, \dots$

## Assignment step

$$S_i^{(t)} = \left\{ x_p : \|x_p - c_i^{(t)}\|^2 \leq \|x_p - c_j^{(t)}\|^2 \forall j, 0 \leq j < k \right\}$$

## Update step

$$c_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x \in S_i^{(t)}} x$$



# Outline

SECTION 1

## Introduction

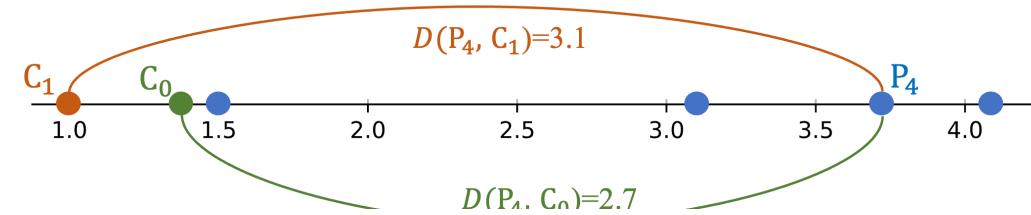
SECTION 2

## Tabular Data

SECTION 3

## Image&Text Data

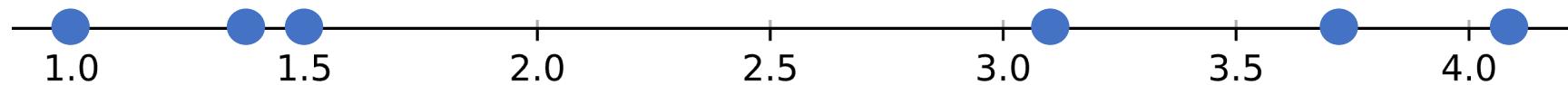
Index	Length	Distance to C0	Distance to C1	Centroid
C <sub>0</sub>	0	1.4	0.0	0.4
C <sub>1</sub>	1	1	0.4	0.0
	2	1.5	0.1	0.5
	3	3.1	1.7	2.1
	4	3.7	2.4	2.8
	5	4.1	2.7	3.1



# K-Means

## ❖ Simple Example

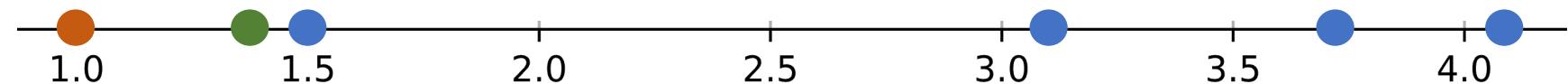
Index	Length
0	1.4
1	1.0
2	1.5
3	3.1
4	3.7
5	4.1



$k = 2$   
 $C_0 = 1.4$   
 $C_1 = 1.0$

Index	Length	
1	1.4	centroid 0
2	1	centroid 1
3	1.5	
4	3.1	
5	3.7	
6	4.1	

1. Initialize k centroids (randomly)

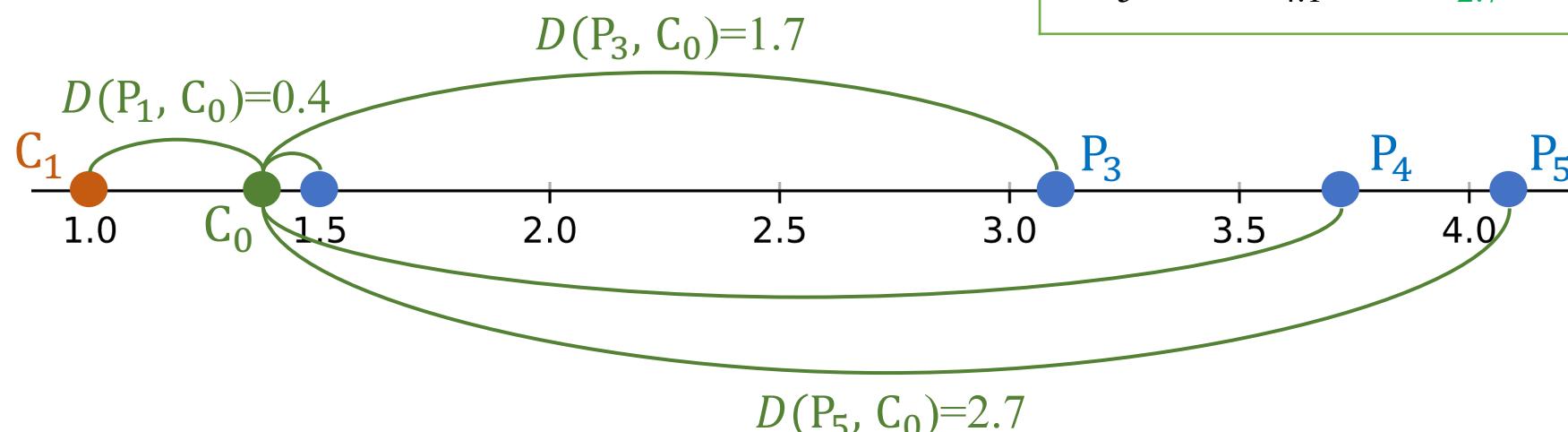


# K-Means

## ❖ Simple Example

1. Initialize k centroids (randomly)
2. Compute distances to each centroids

$$\boxed{k = 2 \quad C_0 = 1.4 \quad C_1 = 1.0}$$



Index	Length	Distance to C0	Distance to C1	Centroid
$C_0$	0	1.4	0.0	0.4
	1	1	0.4	0.0
2	1.5	0.1	0.5	-
3	3.1	1.7	2.1	-
4	3.7	2.4	2.8	-
5	4.1	2.7	3.1	-

# K-Means

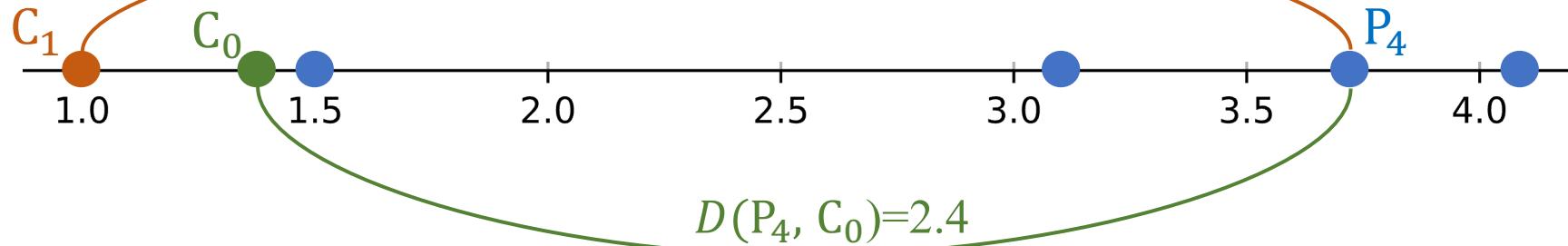
## ❖ Simple Example

1. Initialize k centroids (randomly)
2. Compute distances to each centroids
3. Assign a label to each data point

$k = 2$        $C_0 = 1.4$        $C_1 = 1.0$

$$D(P_4, C_1) = 3.1$$

$$D(P_4, C_0) = 2.7$$



Index	Length	Distance to C0	Distance to C1	Centroid	Label
0	1.4	0.0	0.4	0	0
1	1	0.4	0.0	1	1
2	1.5	0.1	0.5	-	-
3	3.1	1.7	2.1	-	-
4	3.7	2.4	2.8	-	-
5	4.1	2.7	3.1	-	-

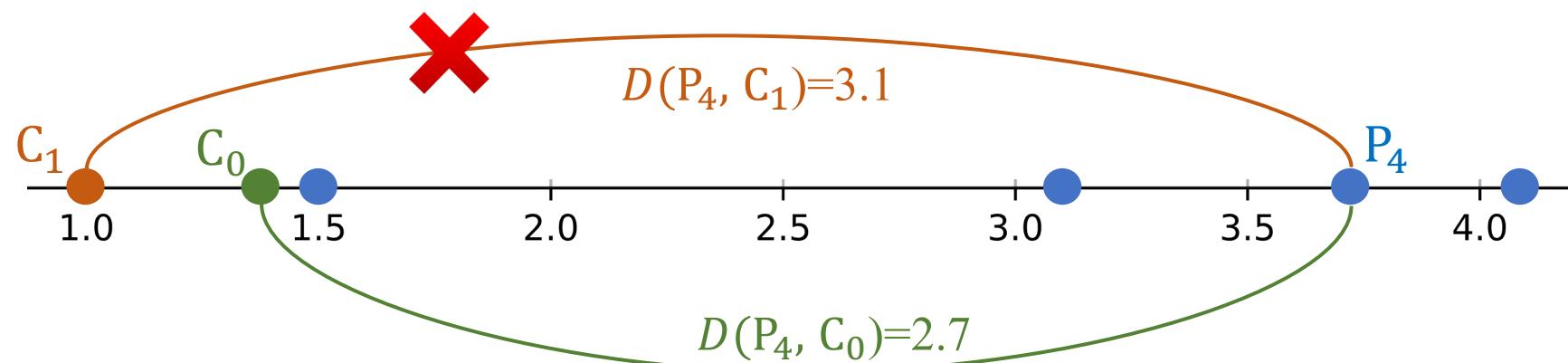
# K-Means

## ❖ Simple Example

1. Initialize k centroids (randomly)
2. Compute distances to each centroids
3. Assign a label to each data point

$$\boxed{k = 2 \quad C_0 = 1.4 \quad C_1 = 1.0}$$

$$D(P_4, C_1) > D(P_4, C_0)$$



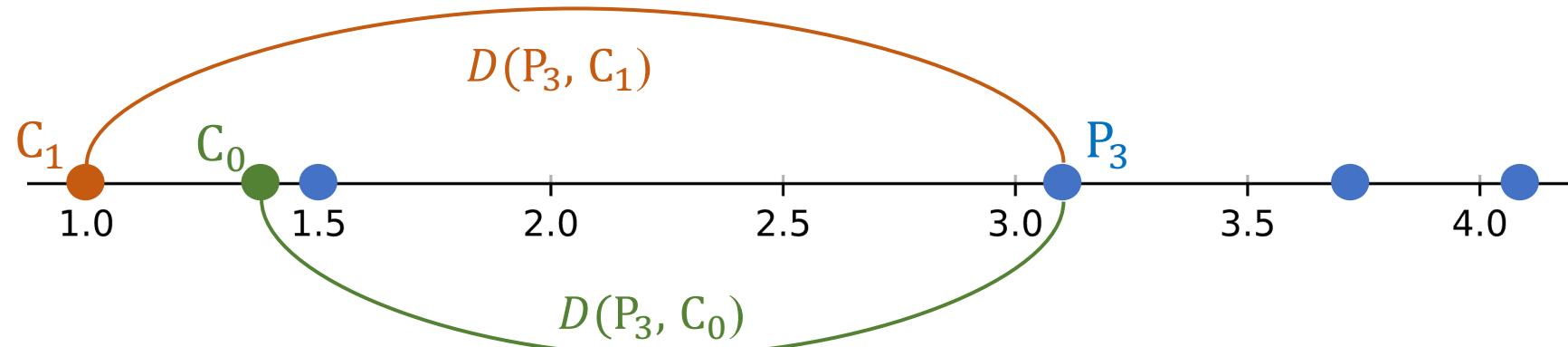
Index	Length	Distance to C0	Distance to C1	Centroid	Label
$C_0$	0	1.4	0.0	0.4	0
$C_1$	1	1	0.4	0.0	1
	2	1.5	0.1	0.5	-
	3	3.1	1.7	2.1	-
	4	3.7	2.4	2.8	-
	5	4.1	2.7	3.1	-

## ❖ Simple Example: Quiz

1. Initialize k centroids (randomly)
2. Compute distances to each centroids
3. Assign a label to each data point

$$k = 2 \quad C_0 = 1.4 \quad C_1 = 1.0$$

Index	Length	Distance to C0	Distance to C1	Centroid	Label
C <sub>0</sub>	0	1.4	0.0	0.4	0
C <sub>1</sub>	1	1	0.4	0.0	1
	2	1.5	0.1	0.5	-
	3	3.1	1.7	2.1	?
	4	3.7	2.4	2.8	0
	5	4.1	2.7	3.1	-

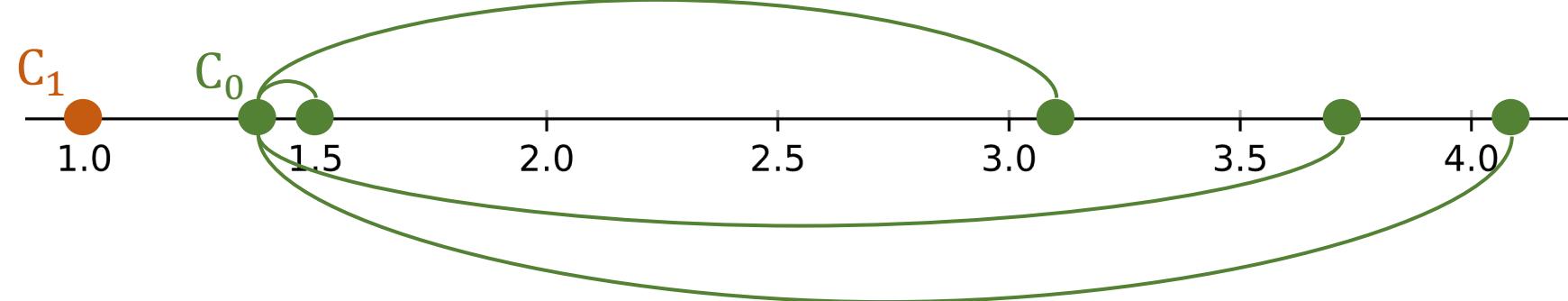


# K-Means

## ❖ Simple Example

1. Initialize k centroids (randomly)
2. Compute distances to each centroids
3. Assign a label to each data point

$k = 2$        $C_0 = 1.4$        $C_1 = 1.0$



Index	Length	Distance to C0	Distance to C1	Centroid	Label
$C_0$	0	1.4	0.0	0	0
	1	1	0.4	0.0	1
2	1.5	0.1	0.5	-	0
3	3.1	1.7	2.1	-	0
4	3.7	2.4	2.8	-	0
5	4.1	2.7	3.1	-	0

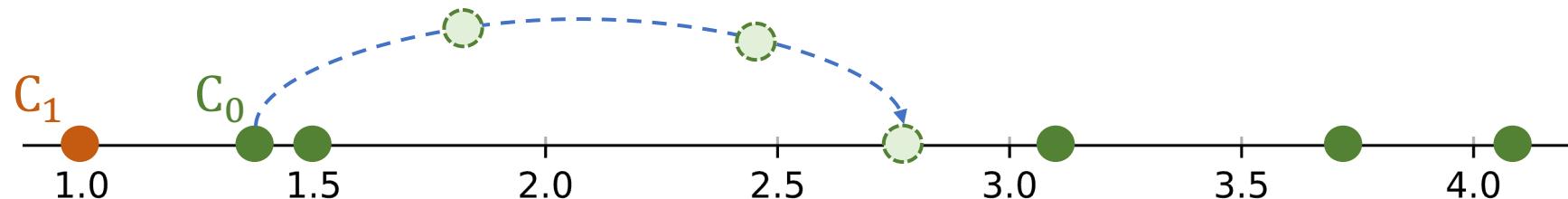
# K-Means

## ❖ Simple Example

1. Initialize k centroids (randomly)
2. Compute distances to each centroids
3. Assign a label to each data point
4. Update centroids

$$C_0 = \frac{1.4 + 1.5 + 3.1 + 3.7 + 4.1}{5} = 2.78$$

$$C_1 = 1.0$$



Index	Length	Distance to C0	Distance to C1	Centroid	Label
$C_0$	0	1.4	0.0	0	0
	1	1	0.4	0.0	1
2	1.5	0.1	0.5	-	0
3	3.1	1.7	2.1	-	0
4	3.7	2.4	2.8	-	0
5	4.1	2.7	3.1	-	0

# K-Means

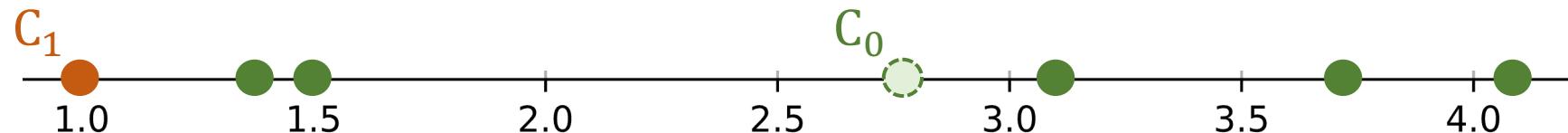
## ❖ Simple Example

1. Initialize k centroids (randomly)
2. Compute distances to each centroids
3. Assign a label to each data point
4. Update centroids

$$C_0 = 2.78$$

$$C_1 = 1.0$$

Index	Length	Distance to C0	Distance to C1	Label
0	1.4	-	-	0
1	1	-	-	1
2	1.5	-	-	0
3	3.1	-	-	0
4	3.7	-	-	0
5	4.1	-	-	0



# K-Means

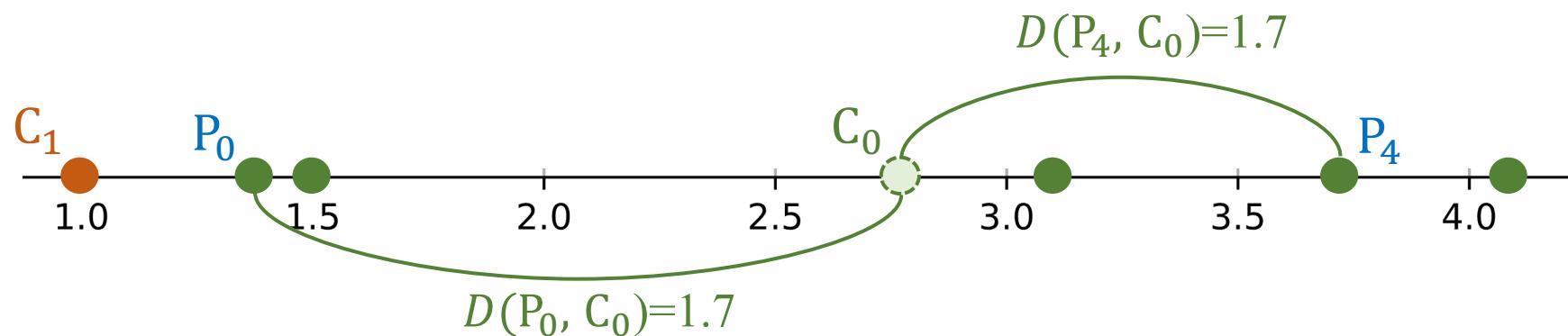
## ❖ Simple Example

1. Initialize k centroids (randomly)
- 2. Compute distances to each centroids
3. Assign a label to each data point
4. Update centroids

$$C_0 = 2.78$$

$$C_1 = 1.0$$

Index	Length	Distance to C0	Distance to C1	Label
0	1.4	1.38	0.4	-
1	1	1.78	0.0	-
2	1.5	1.28	0.5	-
3	3.1	0.32	2.1	-
4	3.7	1.02	2.8	-
5	4.1	1.32	3.1	-



# K-Means

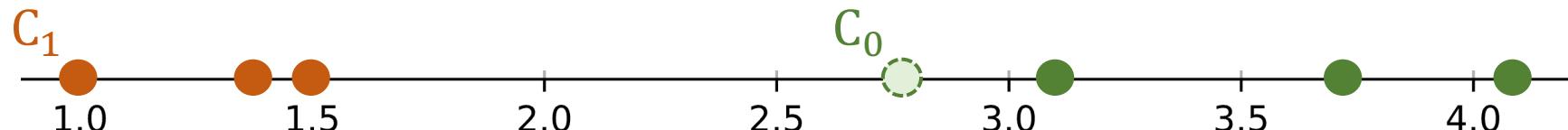
## ❖ Simple Example

1. Initialize k centroids (randomly)
2. Compute distances to each centroids
- 3. Assign a label to each data point
4. Update centroids

$$C_0 = 2.78$$

$$C_1 = 1.0$$

Index	Length	Distance to C0	Distance to C1	Label
0	1.4	1.38	0.4	1
1	1	1.78	0.0	1
2	1.5	1.28	0.5	1
3	3.1	0.32	2.1	0
4	3.7	1.02	2.8	0
5	4.1	1.32	3.1	0



# K-Means

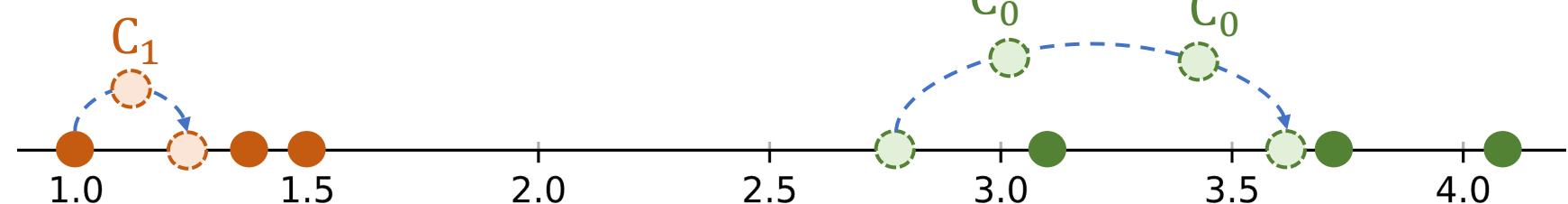
## ❖ Simple Example

1. Initialize k centroids (randomly)
2. Compute distances to each centroids
3. Assign a label to each data point
4. Update centroids

$$C_0 = \frac{3.1 + 3.7 + 4.1}{3} = 3.67$$

$$C_1 = \frac{1.4 + 1.0 + 1.5}{3} = 1.3$$

Index	Length	Distance to C0	Distance to C1	Label
0	1.4	1.38	0.4	1
1	1	1.78	0.0	1
2	1.5	1.28	0.5	1
3	3.1	0.32	2.1	0
4	3.7	1.02	2.8	0
5	4.1	1.32	3.1	0



# K-Means

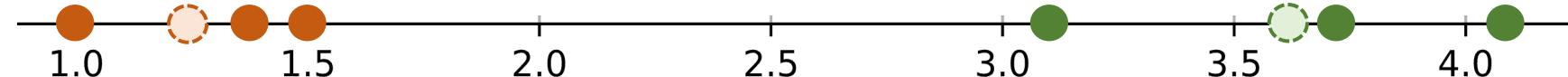
## ❖ Simple Example

1. Initialize k centroids (randomly)
2. Compute distances to each centroids
3. Assign a label to each data point
- 4. Update centroids

$$C_0 = \frac{3.1 + 3.7 + 4.1}{3} = 3.67$$

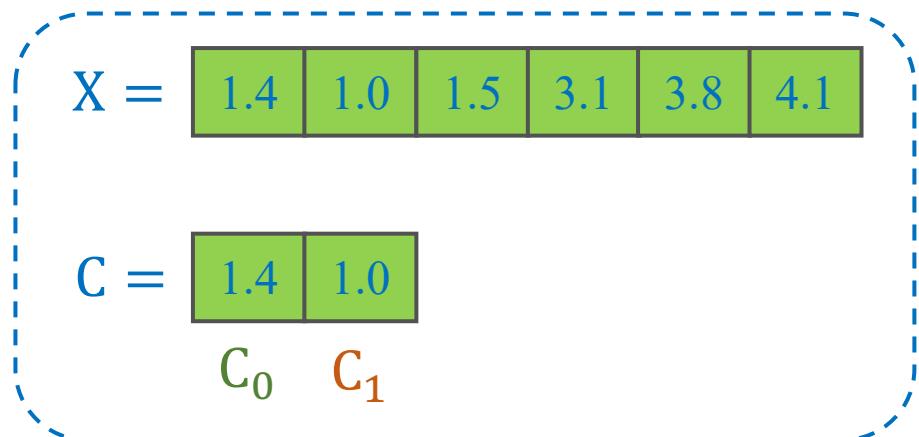
$$C_1 = \frac{1.4 + 1.0 + 1.5}{3} = 1.3$$

Index	Length	Distance to C0	Distance to C1	Label
0	1.4			1
1	1			1
2	1.5			1
3	3.1			0
4	3.7			0
5	4.1			0



# K-Means

## ❖ Implementation using Numpy



Distance to $C_0$	Distance to $C_1$
0.0	0.4
0.4	0.0
0.1	0.5
1.7	2.1
2.4	2.8
2.7	3.1

$\text{np.abs}([1.4, 1.0, 1.5, 3.1, 3.8, 4.1] - [1.4])$

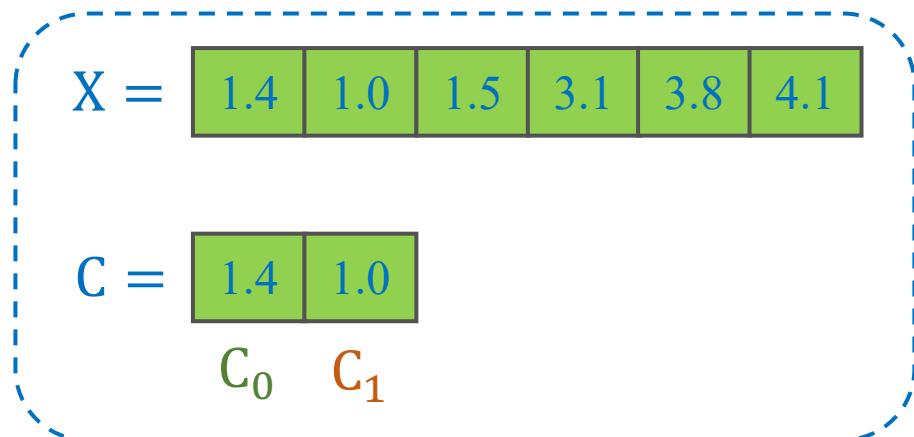
$\text{np.abs}([1.4, 1.0, 1.5, 3.1, 3.8, 4.1] - [1.0])$

```
np.abs(X-C[0])
✓ 0.0s
array([0. , 0.4, 0.1, 1.7, 2.4, 2.7])

np.abs(X-C[1])
✓ 0.0s
array([0.4, 0. , 0.5, 2.1, 2.8, 3.1])
```

# K-Means

## ❖ Implementation using Numpy



	Distance to C0	Distance to C1
	0.0	0.4
	0.4	0.0
	0.1	0.5
	1.7	2.1
	2.4	2.8
	2.7	3.1

```
X = np.array([1.4])
print(X.shape)
print(X)

✓ 0.0s
(1,)
[1.4]
```

```
print(C.shape)

✓ 0.0s
(2,)

np.abs(X-C)

✓ 0.0s
array([0. , 0.4])
```

```
print(X.shape)
print(X)

✓ 0.0s
(6,)
[1.4 1. 1.5 3.1 3.8 4.1]

print(C.shape)

✓ 0.0s
(2,)

np.abs(X-C)

✓ 0.0s
ValueError
Cell In[21], line 1
----> 1 np.abs(X-C)

ValueError: operands could not
```

# K-Means

## ❖ Implementation using Numpy

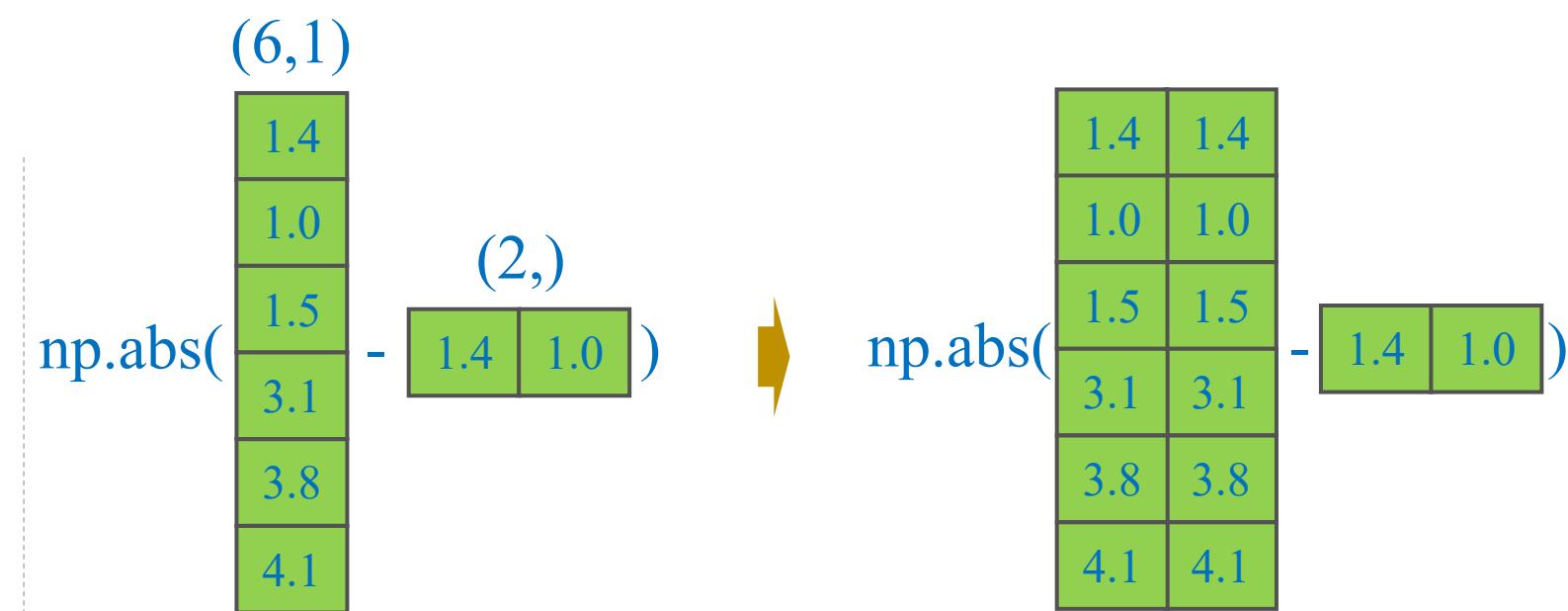
$$X = \begin{bmatrix} 1.4 & 1.0 & 1.5 & 3.1 & 3.8 & 4.1 \end{bmatrix}$$

$$C = \begin{bmatrix} 1.4 & 1.0 \end{bmatrix}$$

$C_0$     $C_1$



	Distance to $C_0$	Distance to $C_1$
	0.0	0.4
	0.4	0.0
	0.1	0.5
	1.7	2.1
	2.4	2.8
	2.7	3.1



```
print(C.shape)
print(X.shape)
✓ 0.0s
(2,)
(6,)
```

```
np.abs(X.reshape(-1, 1) - C)
✓ 0.0s
array([[0. , 0.4],
       [0.4, 0. ],
       [0.1, 0.5],
       [1.7, 2.1],
       [2.4, 2.8],
       [2.7, 3.1]])
```

# K-Means

## ❖ Implementation using Numpy

X =	1.4	1.0	1.5	3.1	3.8	4.1
-----	-----	-----	-----	-----	-----	-----

C =	1.4	1.0
	$C_0$	$C_1$

Distance to C0	Distance to C1
0.0	0.4
0.4	0.0
0.1	0.5
1.7	2.1
2.4	2.8
2.7	3.1

Label
0
1
0
0
0
0



(1) np.argmin(D, axis=0)

(2) np.argmin(D, axis=1)

0	1	0	0	0	0
---	---	---	---	---	---

(a)

0	1
---	---

(b)

```
D
✓ 0.0s
array([[0., 0.4],
       [0.4, 0. ],
       [0.1, 0.5],
       [1.7, 2.1],
       [2.4, 2.8],
       [2.7, 3.1]])
```

np.argmin(D, axis=0)

```
✓ 0.0s
array([0, 1])
```

np.argmin(D, axis=1)

```
✓ 0.0s
array([0, 1, 0, 0, 0, 0])
```

# K-Means

## ❖ Implementation using Numpy

$$C = \begin{bmatrix} 1.4 & 1.0 \end{bmatrix}$$

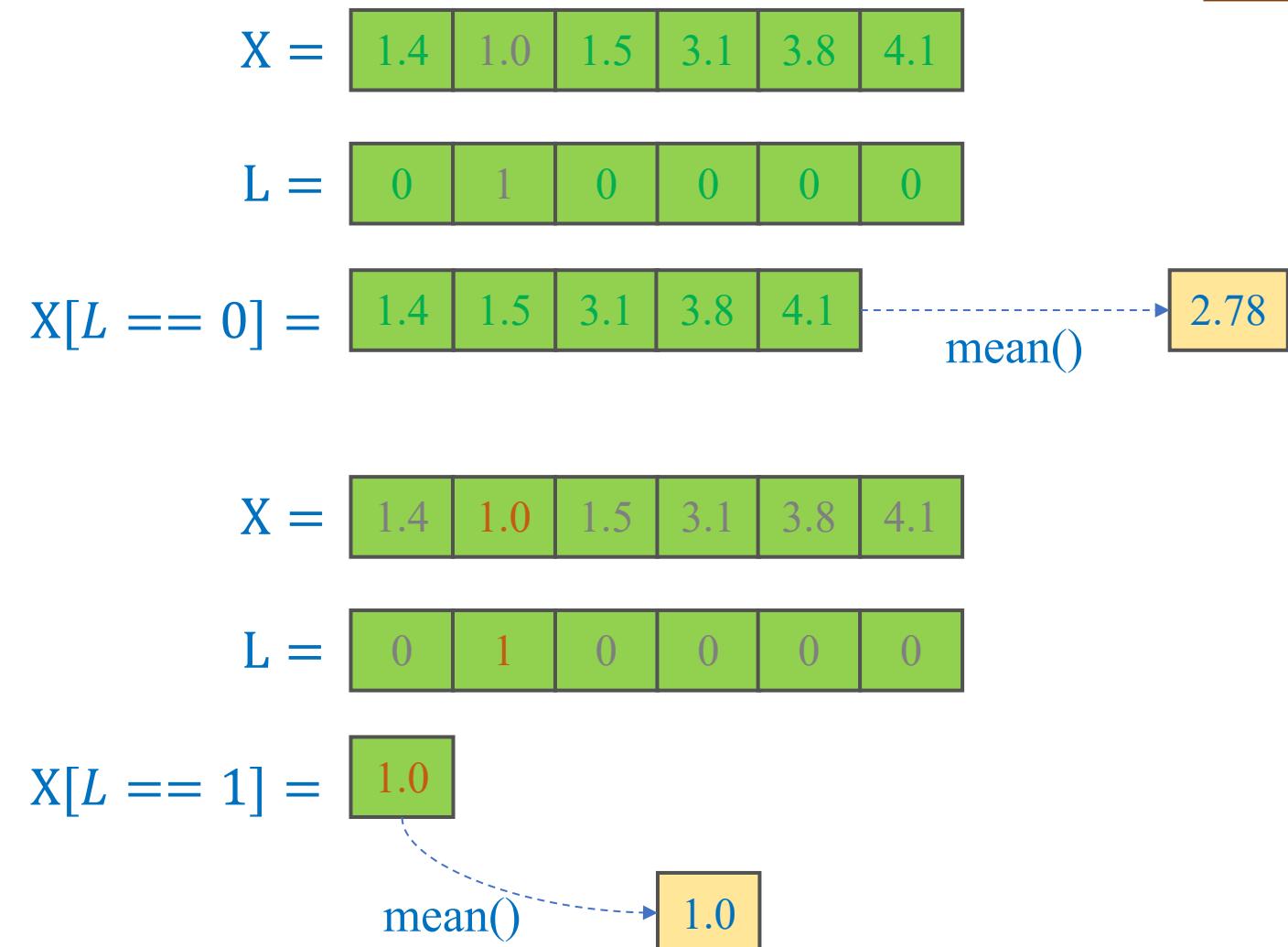
$C_0 \quad C_1$

$$X = \begin{bmatrix} 1.4 & 1.0 & 1.5 & 3.1 & 3.8 & 4.1 \end{bmatrix}$$

$$L = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$C_{\text{new}} = \begin{bmatrix} 2.78 & 1.0 \end{bmatrix}$$

$C_0 \quad C_1$



`np.array([X[L == i].mean() for i in range(k)])`

# K-Means

## ❖ Implementation using Numpy

$$C = \begin{bmatrix} 1.4 & 1.0 \\ C_0 & C_1 \end{bmatrix}$$

$$X = \begin{bmatrix} 1.4 & 1.0 & 1.5 & 3.1 & 3.8 & 4.1 \end{bmatrix}$$

$$L = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

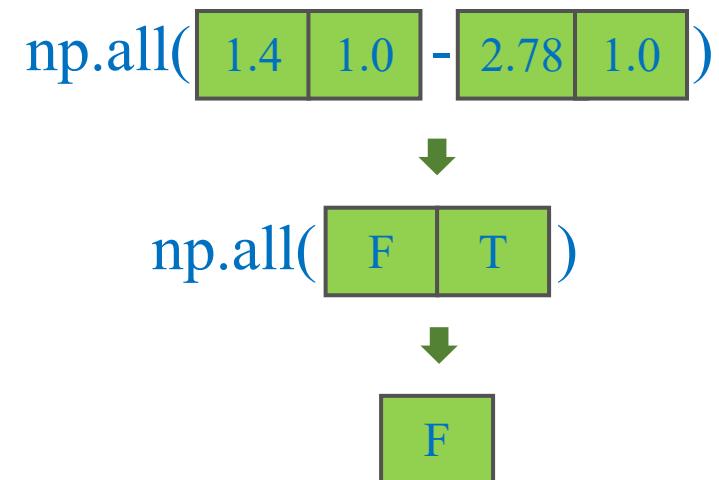
$$C_{\text{new}} = \begin{bmatrix} 2.78 & 1.0 \\ C_0 & C_1 \end{bmatrix}$$

$$(C == C_{\text{new}}) = \begin{bmatrix} F & T \end{bmatrix}$$

```
C = np.array([1.4, 1.0])
C_new = np.array([2.78, 1.0])
print(C==C_new)
✓ 0.0s
[False True]
```

```
if np.all(C == C_new):
    break
```

```
print( np.all([True, True]) )
print( np.all([False, True]) )
print( np.all([True, False]) )
print( np.all([False, False]) )
✓ 0.0s
True
False
False
False
```



# K-Means

$$\text{Norm } \|\vec{u}\| = \sqrt{u_1^2 + \cdots + u_n^2}$$

$$\text{Squared norm } \|\vec{u}\|^2 = u_1^2 + \cdots + u_n^2$$

## ❖ Implementation using Numpy

Given  $(x_1, \dots, x_n)$  where  $x_i \in \mathbb{R}^d$

$$S = (S_0, S_1, \dots)$$

$$\text{WCSS} = \sum_{i=0}^{k-1} \sum_{x \in S_i} \|x - c_i\|^2$$

$$X = \begin{bmatrix} 1.4 & 1.0 & 1.5 & 3.1 & 3.8 & 4.1 \end{bmatrix}$$

$$L = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \end{bmatrix}$$

$$C = \begin{bmatrix} 3.67 & 1.3 \end{bmatrix}$$

$$C_0 \quad C_1$$

$$\text{WCSS} = \sum_{x \in S_0} \|x - c_0\|^2 + \sum_{x \in S_1} \|x - c_1\|^2$$

$$\text{square}(\begin{bmatrix} 1.4 & 1.0 & 1.5 & 1.3 \end{bmatrix})$$

$$\downarrow$$

$$\begin{bmatrix} 0.3249 & 0.0169 & 0.1849 \end{bmatrix}$$

$$\sum$$

$$0.5267$$

$$\text{square}(\begin{bmatrix} 3.1 & 3.8 & 4.1 & 3.67 \end{bmatrix})$$

$$\downarrow$$

$$\begin{bmatrix} 0.01 & 0.09 & 0.04 \end{bmatrix}$$

$$\sum$$

$$0.1399$$

$$\sum$$

$$0.667$$

```
np.sum([np.sum(np.square(X[L == i] - C[i])) for i in range(k)])
```



Demo

#Features=2

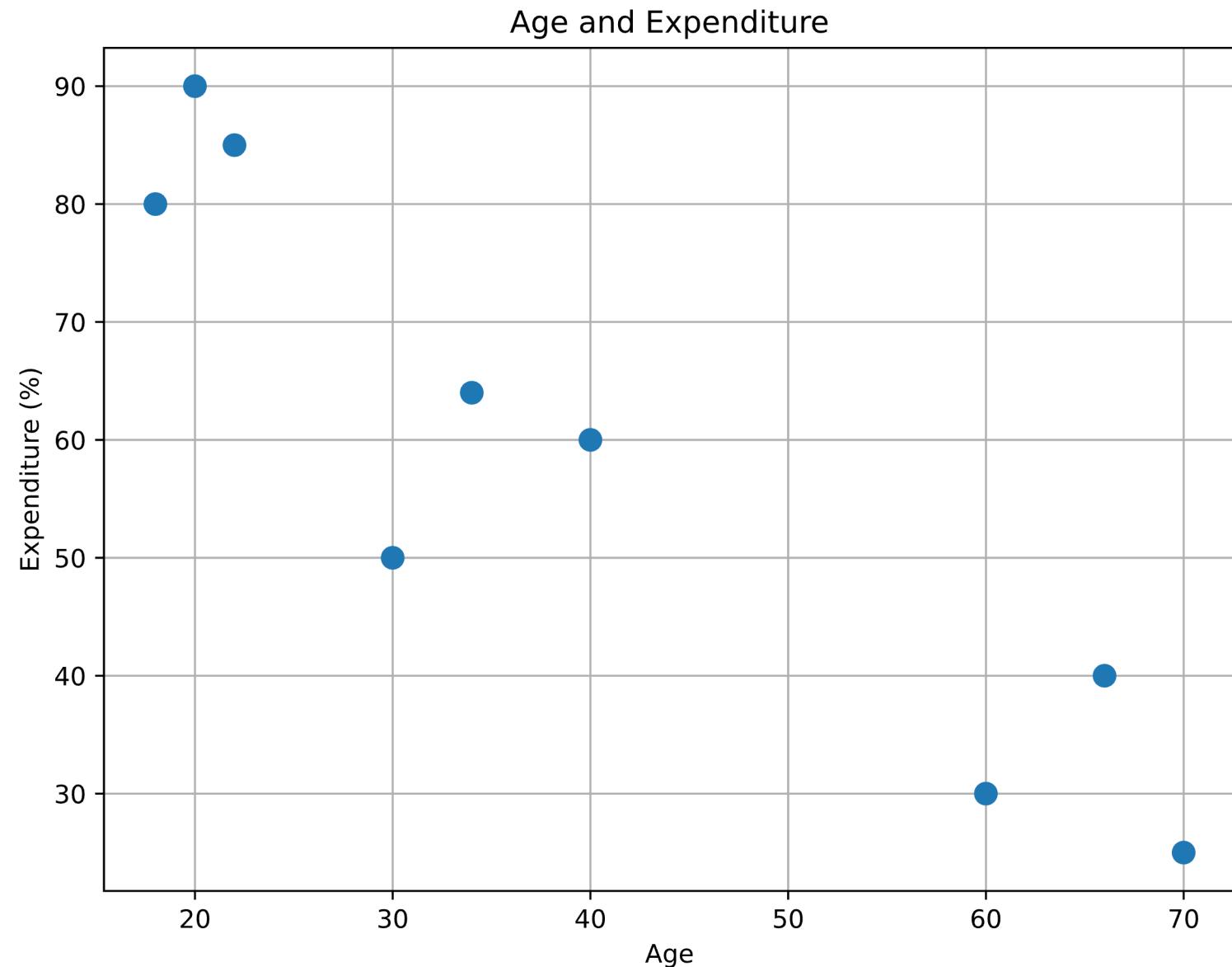
#clusters=3

# K-Means

## ❖ Dataset

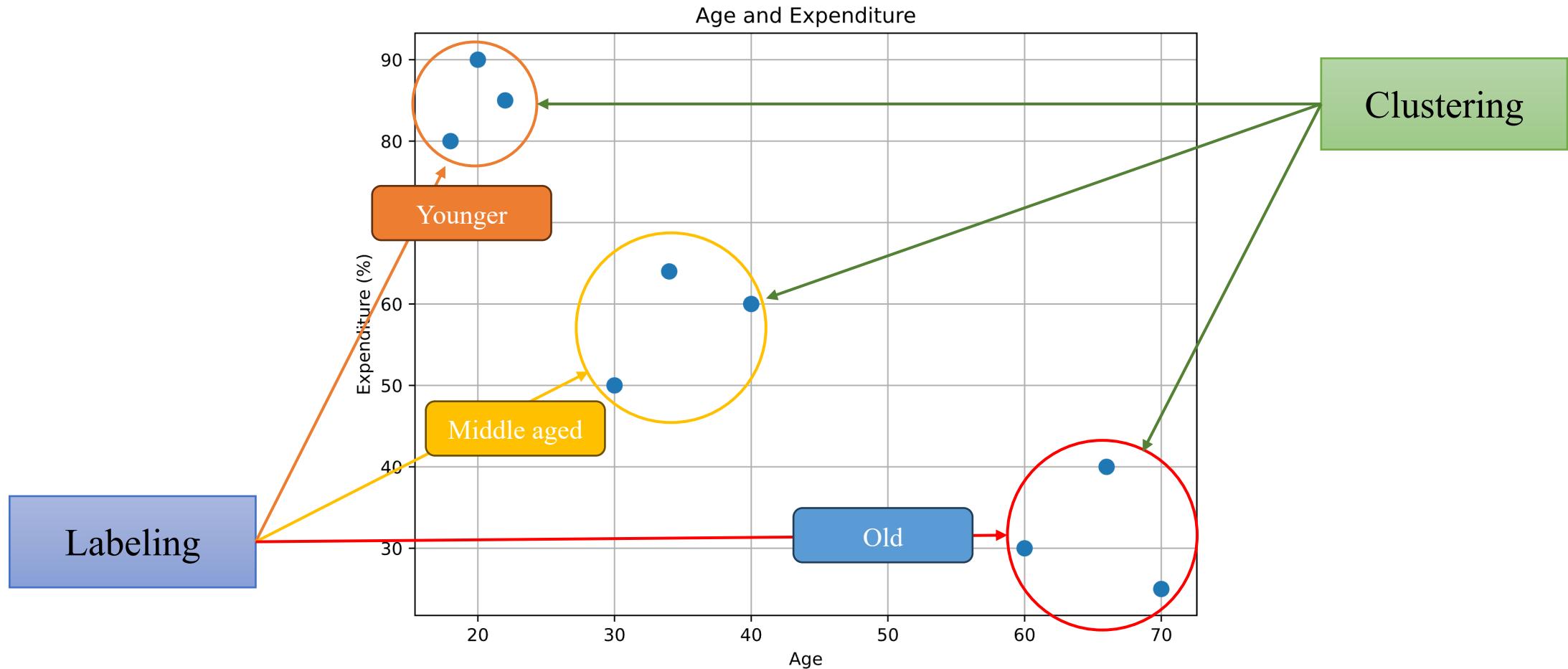
The dataset describes the relationship of monthly expenditure to age.

Index	Age	Expenditure
1	18	80
2	20	90
3	22	85
4	30	50
5	34	64
6	40	60
7	60	30
8	66	40
9	70	25



# K-Means

## ❖ Targets



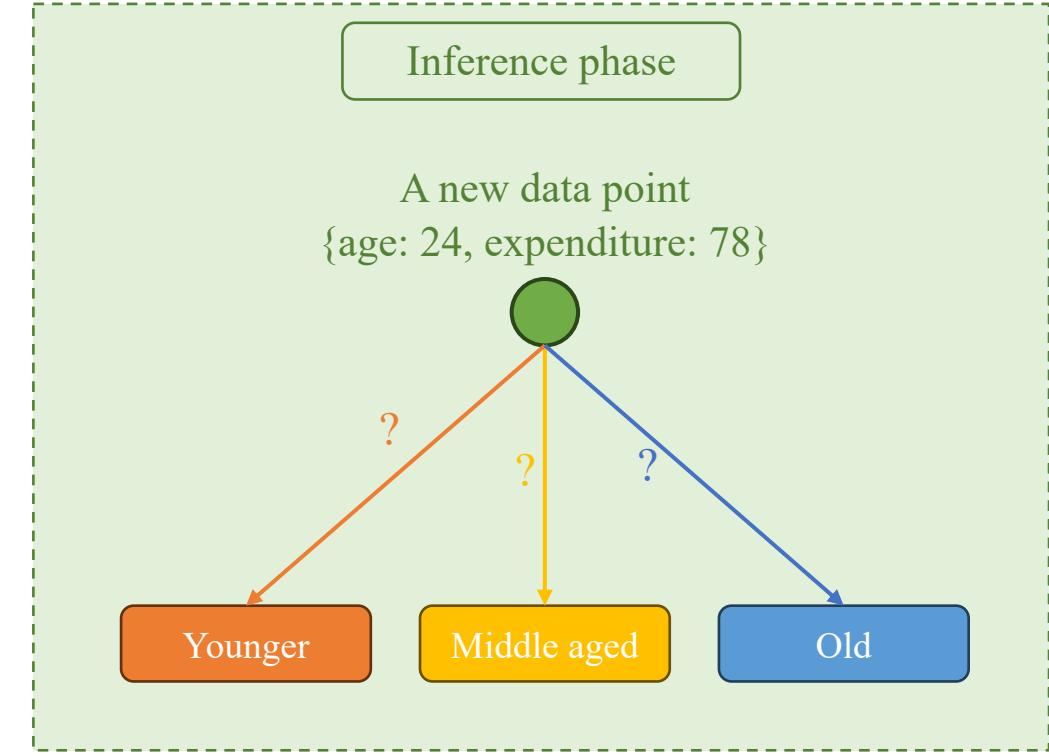
# K-Means

## ❖ Targets

Index	Age	Expenditure
1	18	80
2	20	90
3	22	85
4	30	50
5	34	64
6	40	60
7	60	30
8	66	40
9	70	25

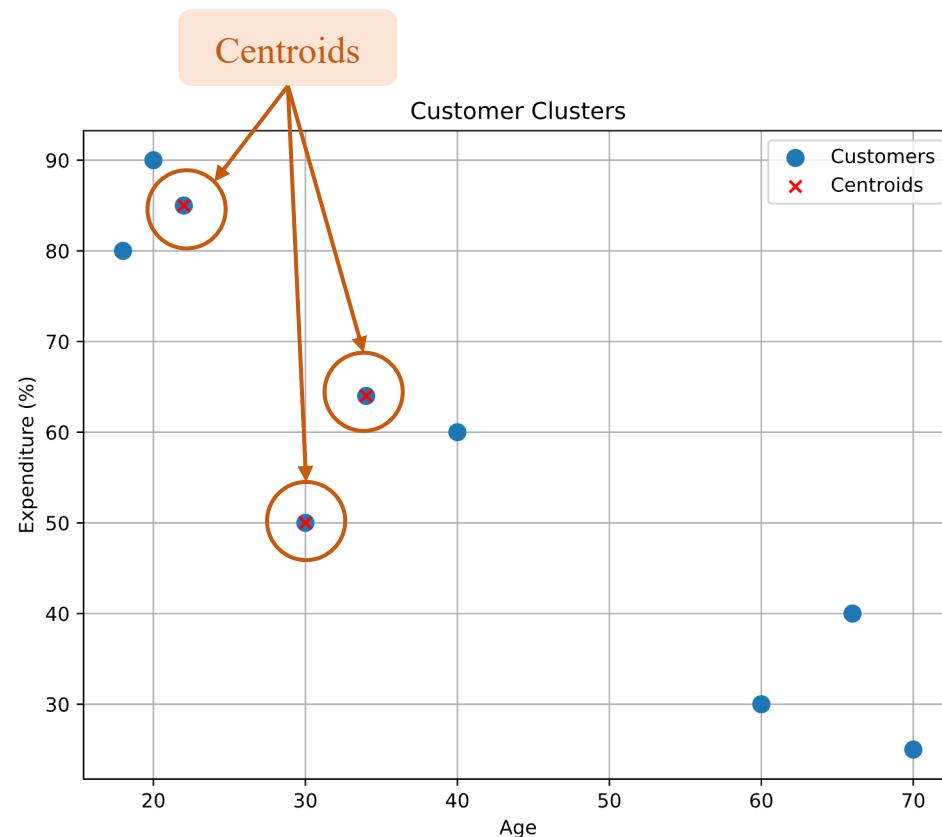
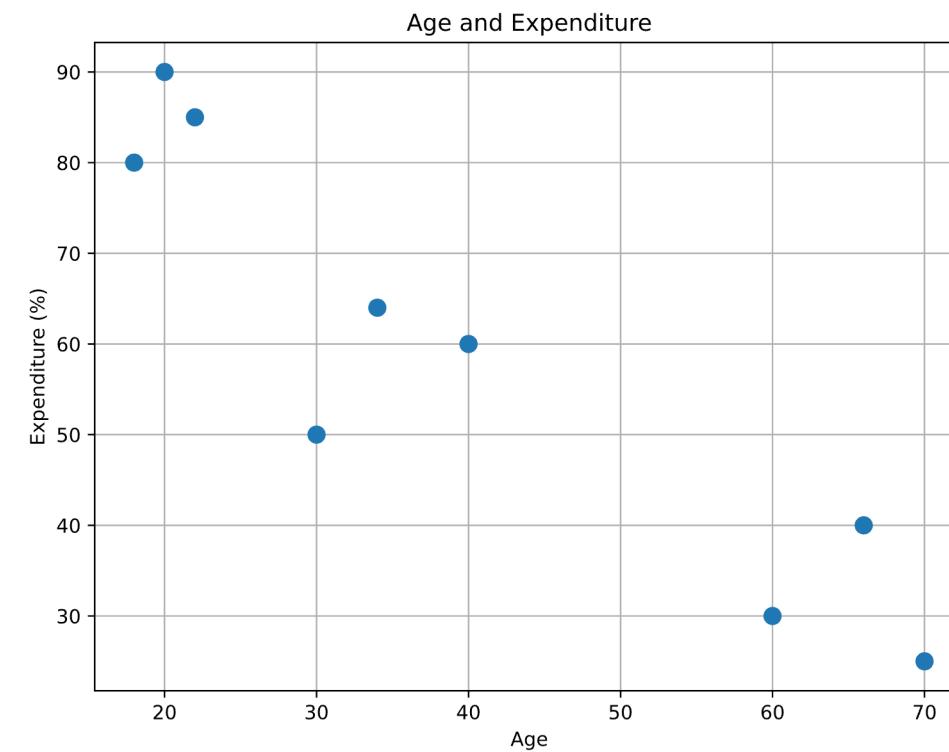
Training phase

Cluster	Label
1	younger
1	younger
1	younger
2	middle-aged
2	middle-aged
2	middle-aged
3	old
3	old
3	old



# K-Means

## ❖ Step 1: Init Centroids by random pickup

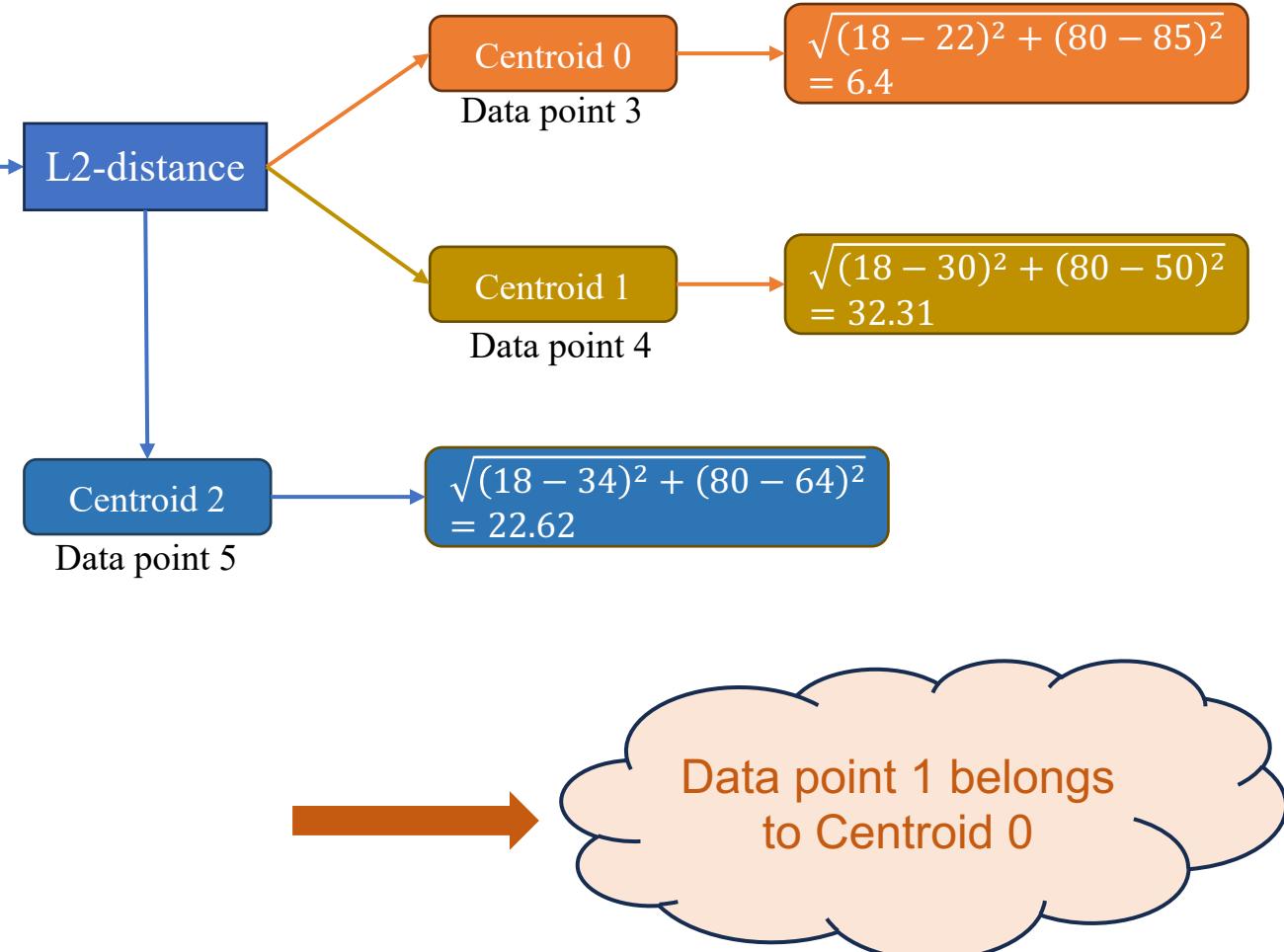


Index	Age	Expenditure
1	18	80
2	20	90
3	22	85
4	30	50
5	34	64
6	40	60
7	60	30
8	66	40
9	70	25

# K-Means

## ❖ Step 2: Assign labels

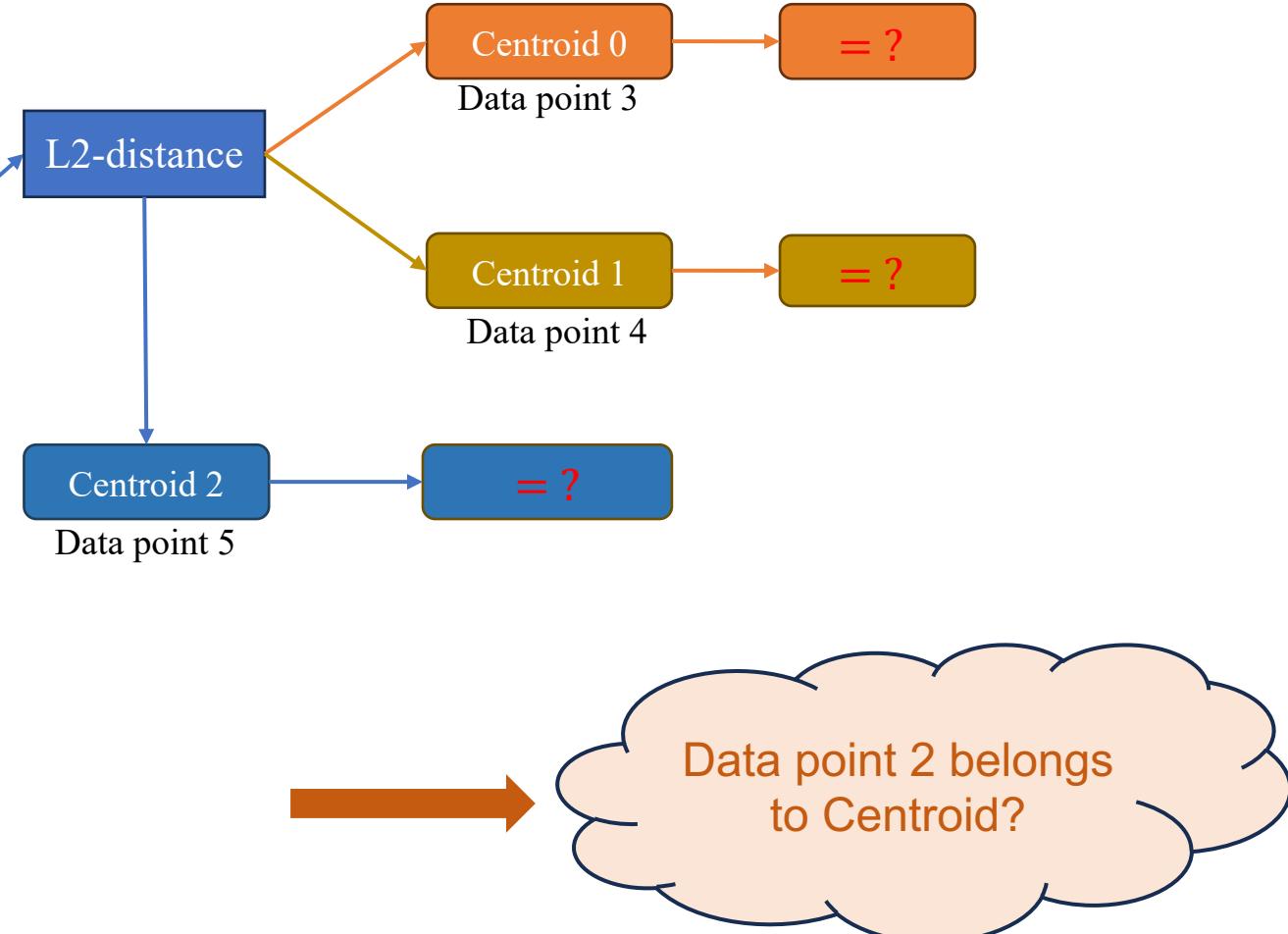
Index	Age	Expenditure	Centroid	Label
1	18	80	-	-
2	20	90	-	-
3	22	85	0	0
4	30	50	1	1
5	34	64	2	2
6	40	60	-	-
7	60	30	-	-
8	66	40	-	-
9	70	25	-	-



# K-Means

## ❖ Step 2: Assign labels

Index	Age	Expenditure	Centroid	Label
1	18	80	-	-
2	20	90	-	-
3	22	85	0	0
4	30	50	1	1
5	34	64	2	2
6	40	60	-	-
7	60	30	-	-
8	66	40	-	-
9	70	25	-	-



# K-Means

## ❖ Step 2: Assign labels

Index	Age	Expenditure	Centroid	Label
1	18	80	-	-
2	20	90	-	-
3	22	85	0	0
4	30	50	1	1
5	34	64	2	2
6	40	60	-	-
7	60	30	-	-
8	66	40	-	-
9	70	25	-	-

L2-distance

L2-distance

L2-distance

L2-distance

L2-distance

L2-distance

L2-distance

L2-distance

L2-distance

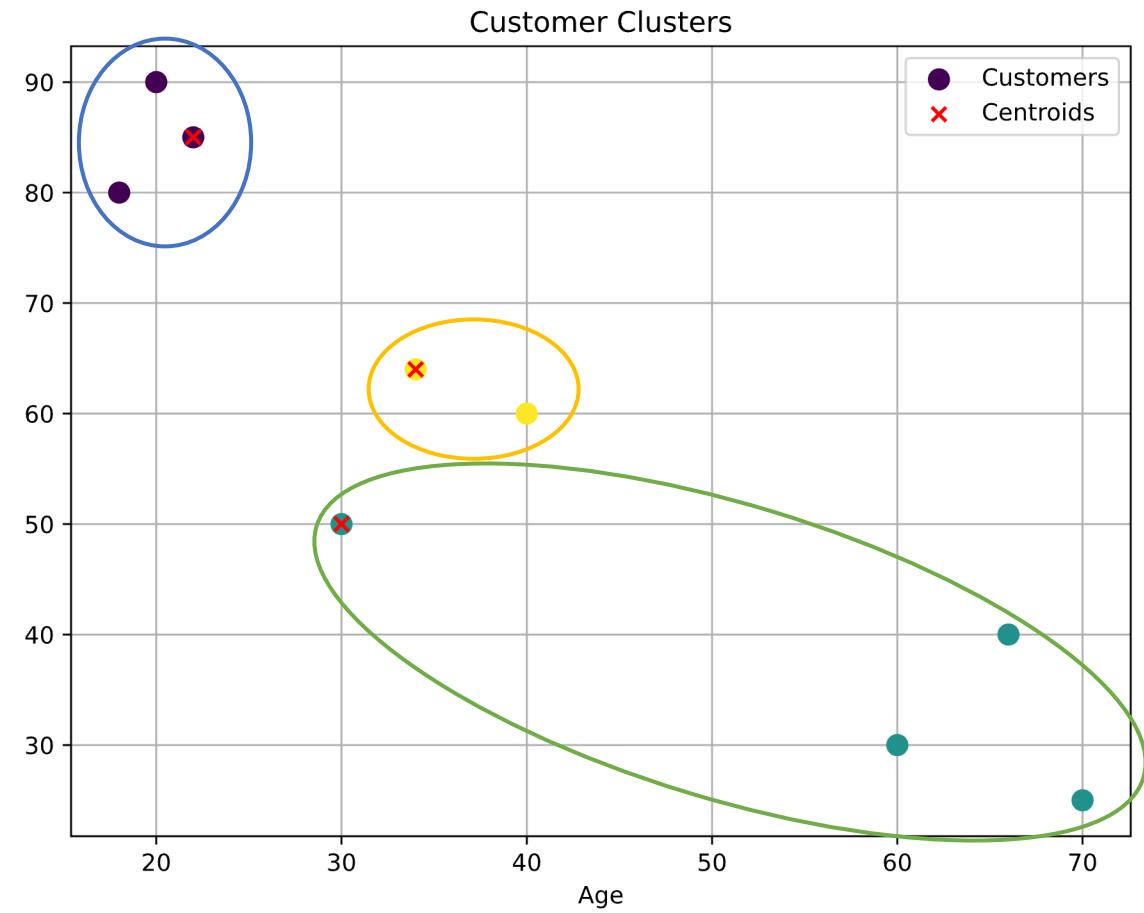


Index	Age	Expenditure	Centroid	Label
1	18	80	-	0
2	20	90	-	0
3	22	85	0	0
4	30	50	1	1
5	34	64	2	2
6	40	60	-	2
7	60	30	-	1
8	66	40	-	1
9	70	25	-	1

# K-Means

## ❖ Step 2: Assign labels

Index	Age	Expenditure	Centroid	Label
1	18	80	-	0
2	20	90	-	0
3	22	85	0	0
4	30	50	1	1
5	34	64	2	2
6	40	60	-	2
7	60	30	-	1
8	66	40	-	1
9	70	25	-	1



# K-Means

## ❖ Step 3: Update Centroids

Index	Age	Expenditure	Centroid	Label
1	18	80	-	0
2	20	90	-	0
3	22	85	0	0
4	30	50	1	1
5	34	64	2	2
6	40	60	-	2
7	60	30	-	1
8	66	40	-	1
9	70	25	-	1

New Centroid 0

$$\left( \frac{18 + 20 + 22}{3}, \frac{80 + 90 + 85}{3} \right) = (20, 85)$$

New Centroid 2

$$\left( \frac{34 + 40}{2}, \frac{64 + 60}{2} \right) = (37, 62)$$

New Centroid 1

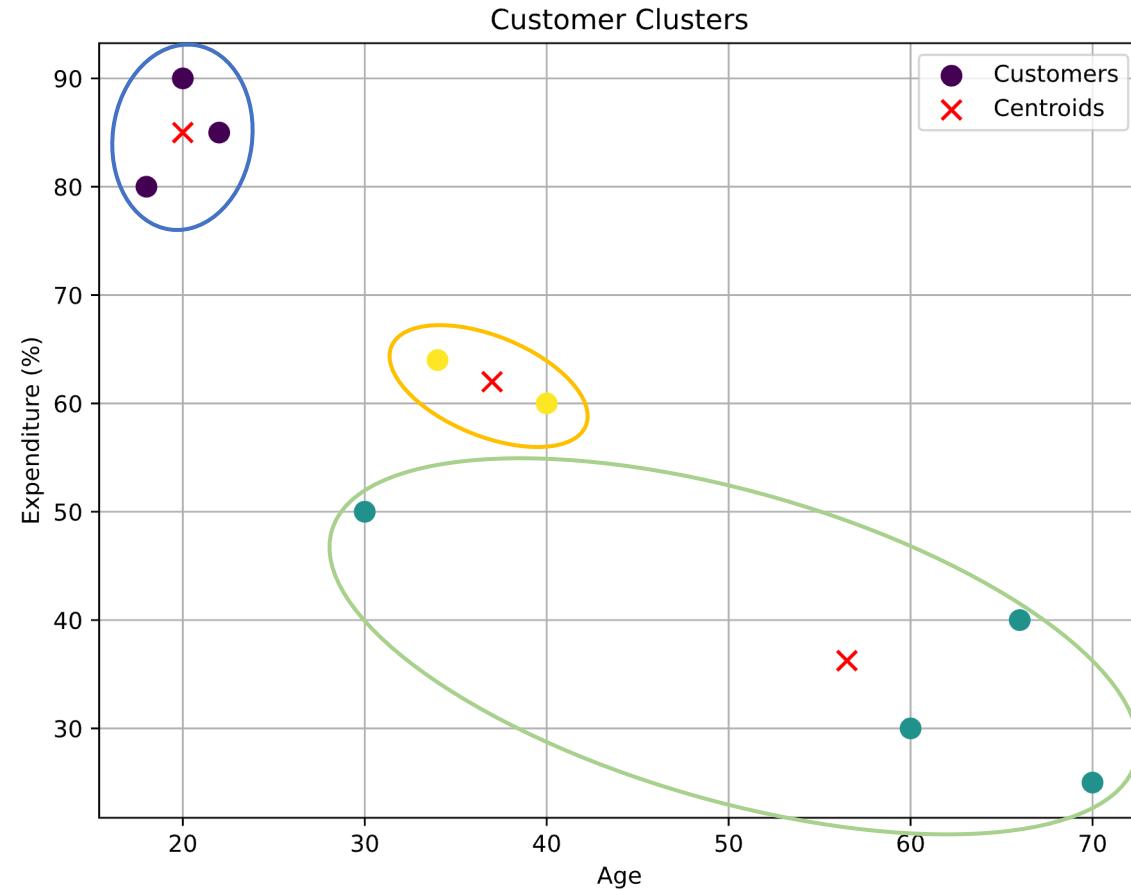
$$\left( \frac{30 + 60 + 66 + 70}{4}, \frac{50 + 30 + 40 + 25}{4} \right) = (56.5, 36.25)$$

Index	Age	Expenditure	Centroid	Label
1	18	80	(20, 85)	0
2	20	90	(20, 85)	0
3	22	85	(20, 85)	0
4	30	50	(56.5,36.25)	1
5	34	64	(37,62)	2
6	40	60	(37,62)	2
7	60	30	(56.5,36.25)	1
8	66	40	(56.5,36.25)	1
9	70	25	(56.5,36.25)	1

# K-Means

## ❖ Step 3: Update Centroids

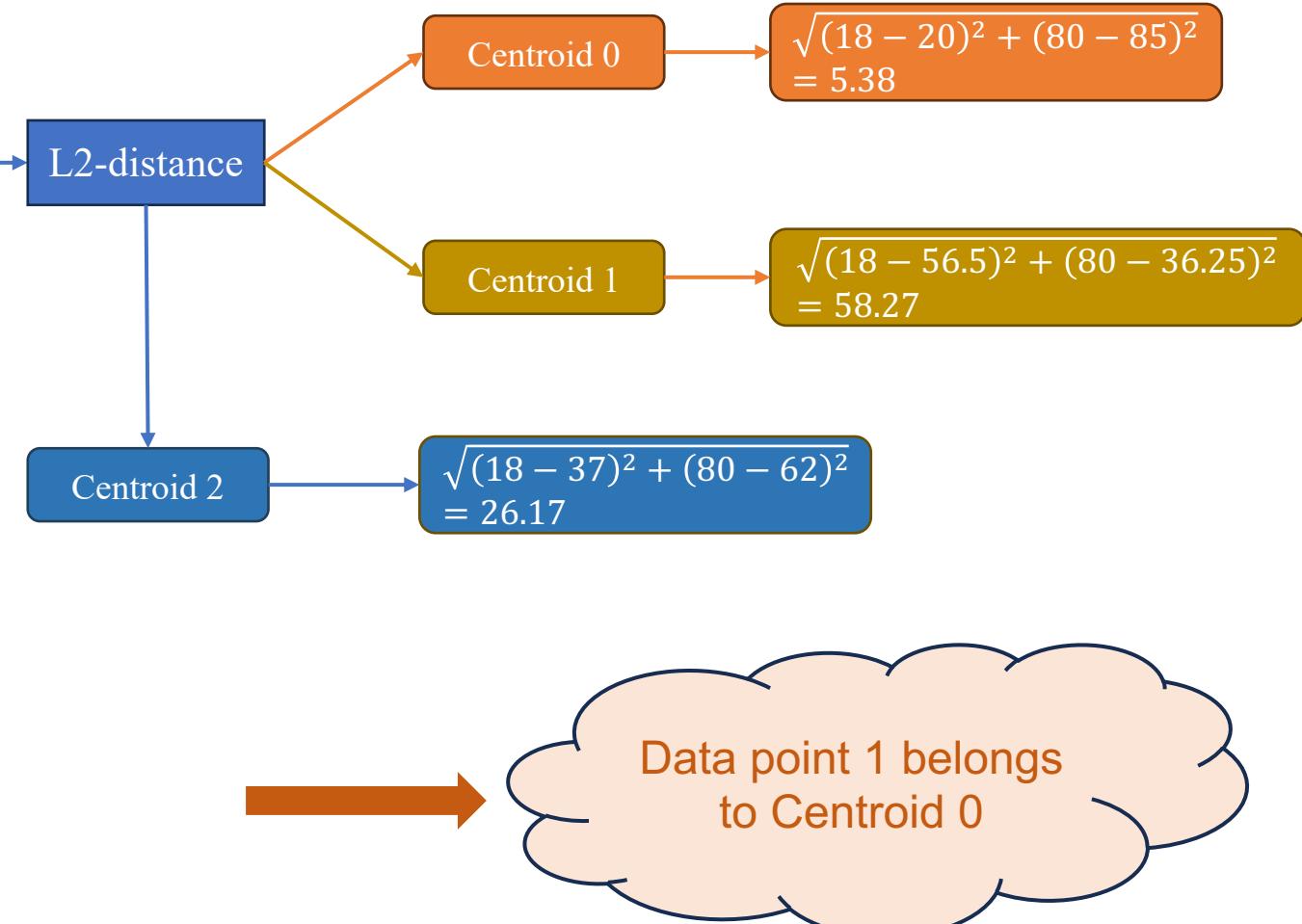
Index	Age	Expenditure	Centroid	Label
1	18	80	(20, 85)	0
2	20	90	(20, 85)	0
3	22	85	(20, 85)	0
4	30	50	(56.5,36.25)	1
5	34	64	(37,62)	2
6	40	60	(37,62)	2
7	60	30	(56.5,36.25)	1
8	66	40	(56.5,36.25)	1
9	70	25	(56.5,36.25)	1



# K-Means

## ❖ Repeat step 2: Calculate distance

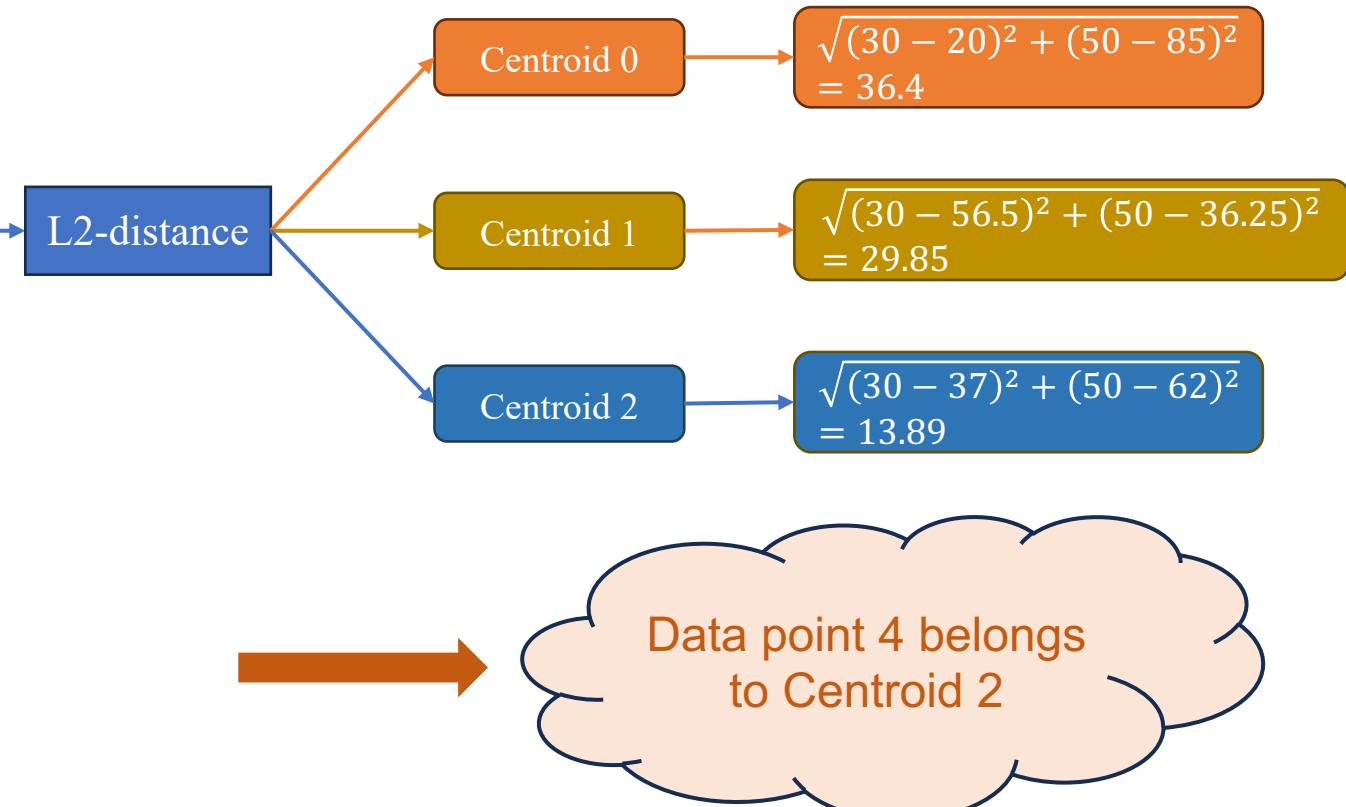
Index	Age	Expenditure	Centroid	Label
1	18	80	(20, 85)	0
2	20	90	(20, 85)	0
3	22	85	(20, 85)	0
4	30	50	(56.5,36.25)	1
5	34	64	(37,62)	2
6	40	60	(37,62)	2
7	60	30	(56.5,36.25)	1
8	66	40	(56.5,36.25)	1
9	70	25	(56.5,36.25)	1



# K-Means

## ❖ Repeat step 2: Calculate distance

Index	Age	Expenditure	Centroid	Label
1	18	80	(20, 85)	0
2	20	90	(20, 85)	0
3	22	85	(20, 85)	0
4	30	50	(56.5,36.25)	1
5	34	64	(37,62)	2
6	40	60	(37,62)	2
7	60	30	(56.5,36.25)	1
8	66	40	(56.5,36.25)	1
9	70	25	(56.5,36.25)	1



# K-Means

## ❖ Repeat step 3: Update Centroids

Index	Age	Expenditure	Centroid	Label
1	18	80	(20, 85)	0
2	20	90	(20, 85)	0
3	22	85	(20, 85)	0
4	30	50	(37,62)	2
5	34	64	(37,62)	2
6	40	60	(37,62)	2
7	60	30	(56.5,36.25)	1
8	66	40	(56.5,36.25)	1
9	70	25	(56.5,36.25)	1

New Centroid 0

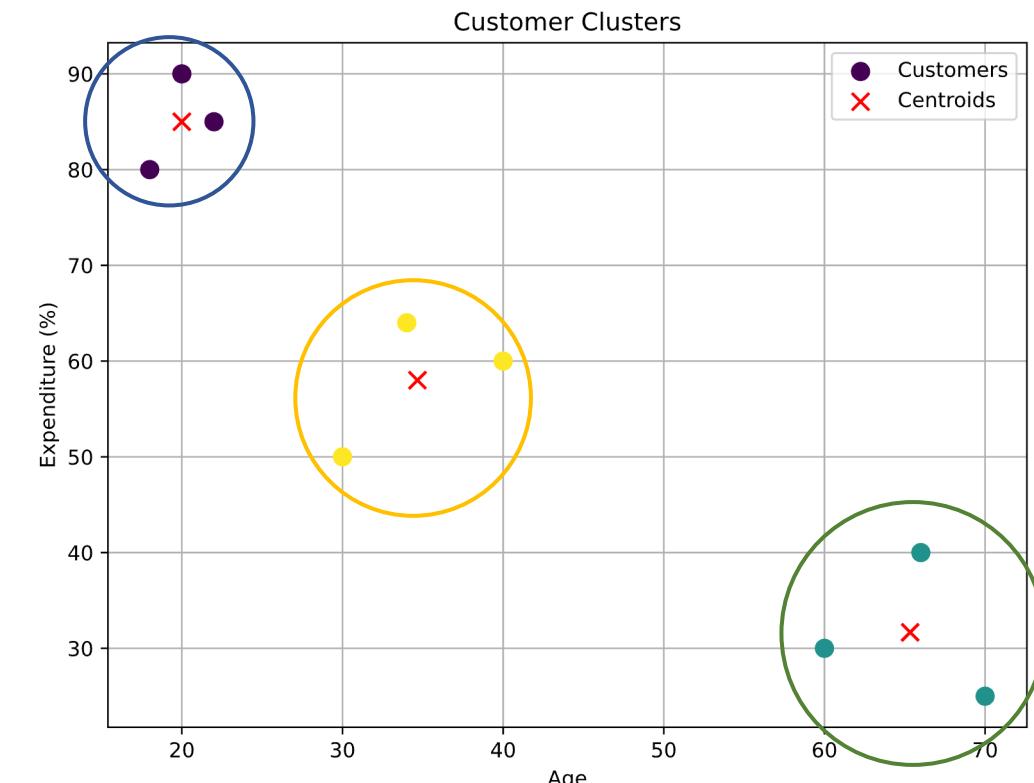
$$\left( \frac{18 + 20 + 22}{3}, \frac{80 + 90 + 85}{3} \right) = (20, 85)$$

New Centroid 2

$$\left( \frac{30 + 34 + 40}{3}, \frac{50 + 64 + 60}{3} \right) = (34.66, 58)$$

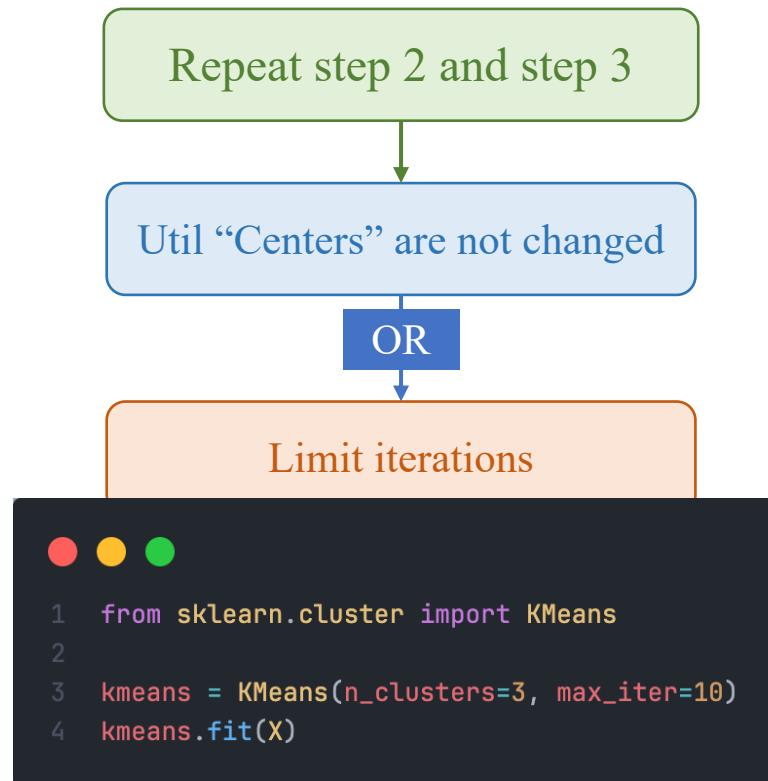
New Centroid 1

$$\left( \frac{60 + 66 + 70}{3}, \frac{30 + 40 + 25}{3} \right) = (65.33, 31.66)$$



# K-Means

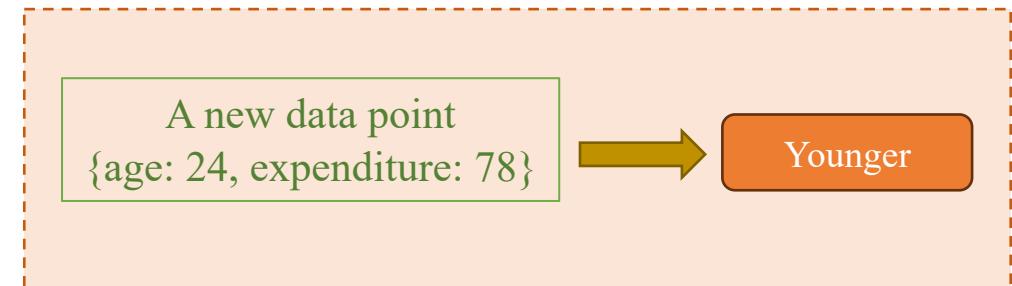
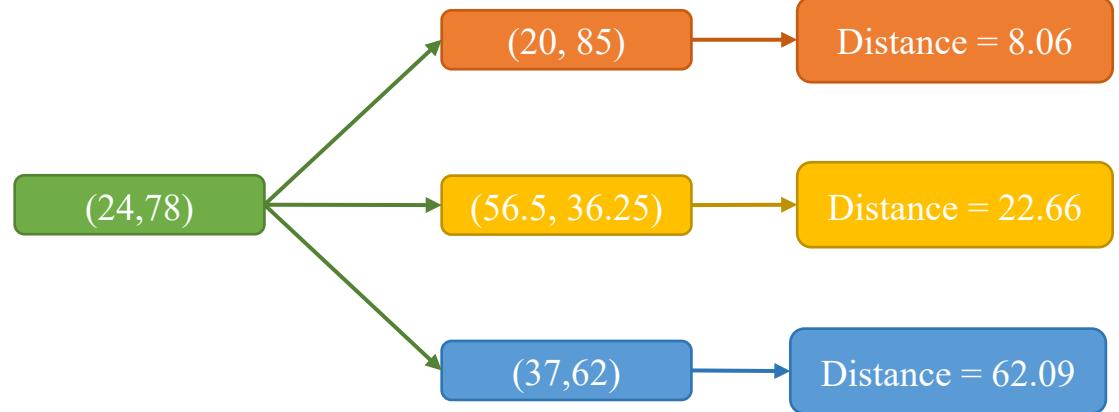
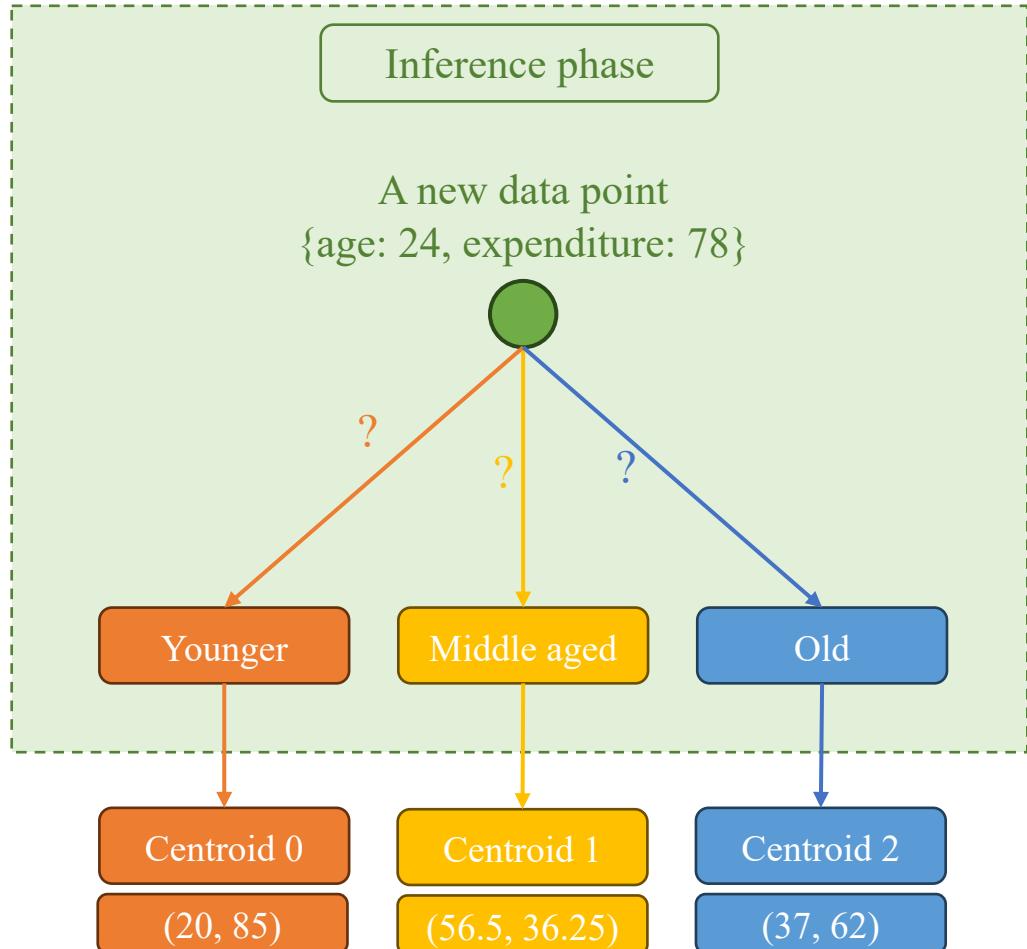
## ❖ Algorithm convergence



Index	Age	Expenditure	Centroid	Label
1	18	80	(20, 85)	0
2	20	90	(20, 85)	0
3	22	85	(20, 85)	0
4	30	50	(37,62)	2
5	34	64	(37,62)	2
6	40	60	(37,62)	2
7	60	30	(56.5,36.25)	1
8	66	40	(56.5,36.25)	1
9	70	25	(56.5,36.25)	1

# K-Means

## Inference





# Question 1

❖ Kết quả chương trình sau là gì?

```
data = np.random.rand(64, 2)
print(data.shape)

centroids = np.random.rand(5, 2)
print(centroids.shape)

distances = (data - centroids)**2
print(distances.shape)
```

✓ 0.0s

(64, 2)  
(5, 2)

(a) (64,2)

(b) (64,5)

(c) (5,2)

(d) Khác

# Question 2

❖ Kết quả chương trình sau là gì?

```
data = np.random.rand(64, 2)

data1 = data[:, np.newaxis, :]
data2 = data.reshape(1, 1, -1)

print(data1.shape)
print(data2.shape)
```

- (a) (64, 1, 2) và (1, 128)
- (b) (64, 2, 1) và (1, 1, 128)
- (c) (64, 1, 2) và (1, 1, 128)
- (d) Khác

# Question 3

❖ Kết quả chương trình sau là gì?

```
data = np.random.rand(32, 4)
print(data.shape)

centroids = np.random.rand(3, 4)
print(centroids.shape)

distances = np.sqrt( (data[:, np.newaxis, :] - centroids)**2 )
print(distances.shape)
```

✓ 0.0s

(32, 4)  
(3, 4)

(a) (32, 3)

(b) (32, 3, 4)

(c) (32, 4, 3)

(d) Khác

# Question 4

- ❖ Cho array chứa **squared norm** giữa các feature của các điểm và của các centroid
- ❖ Ý nghĩa của các tham số như hình bên dưới.
- ❖ Tính distance giữa các point và các centroid (giải thích thêm)?

```
squared_norm = np.random.rand(32, 3, 4)
# 32 points
# 3 dimensions (3 features)
# 4 centroids
```

(a) `np.sqrt( squared_norm.sum(axis=0) )`

(b) `np.sqrt( squared_norm.sum(axis=1) )`

(c) `np.sqrt( squared_norm.sum(axis=2) )`

(d) Khác

# Question 5

- ❖ Cho array chứa thông tin (**features-centroids**)
- ❖ Ý nghĩa của các tham số như hình bên dưới.
- ❖ Tính distance giữa các point và các centroid (giải thích thêm)?

```
distances_detail = np.random.rand(32, 3, 2)
# 32 points
# 3 centroids
# 2 dimensions (2 features)
```

(a) np.linalg.norm(distances\_detail, axis=0)

(b) np.linalg.norm(distances\_detail, axis=1)

(c) np.linalg.norm(distances\_detail, axis=2)

(d) Khác

# Question 6

❖ Kết quả chương trình sau là gì?

```
features = np.random.rand(32, 2)
# 32 points
# each point is of 2 dimensions (2 features)

# k=3
# labels contains numbers from 0 to 2
```

- (a) `np.array([data[labels == i].mean(axis=0) for i in range(k)])`    (b) `np.array([data[labels == i].mean(axis=1) for i in range(k)])`
- (c) `np.array([data[labels == i].mean(axis=2) for i in range(k)])`    (d) `np.array([data[labels == i] for i in range(k)])`



# K-Means

## ❖ Implementation using Numpy

18	80
20	90
22	85
30	50
34	64
40	60
60	30
66	40
70	25

22	85
30	50
34	64

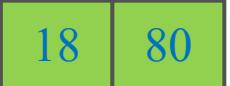
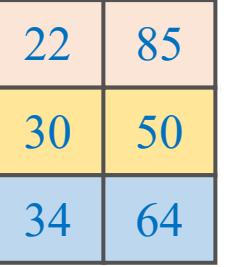
X =

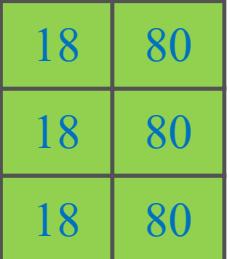
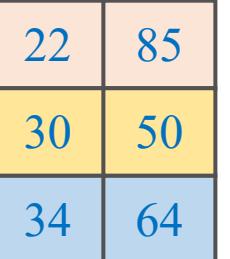
Index	Age	Expenditure	Distance to C0	Distance to C1	Distance to C2
...	...	...	...	...	...

```
np.linalg.norm(X1 - centroids, axis=1)
✓ 0.0s
array([ 6.40312424, 32.31098884, 22.627417 ])
```

```
( (X1 - centroids)**2 ).sum(axis=1)**0.5
✓ 0.0s
array([ 6.40312424, 32.31098884, 22.627417 ])
```

norm( - , axis=1)

norm( - , axis=1)

# K-Means

## ❖ Implementation using Numpy

18	80
20	90
22	85
30	50
34	64
40	60
60	30
66	40
70	25

X =

22	85
30	50
34	64

C =

Index	Age	Expenditure	Distance to C0	Distance to C1	Distance to C2
...	...	...	...	...	...

norm(X[:, np.newaxis, :] - centroids, axis=2)

18	80
20	90
22	85
30	50
34	64
40	60
60	30
66	40
70	25

22	85
30	50
34	64

(3,2)

(9,1,2)

# K-Means

## ❖ Implementation using Numpy

18	80
20	90
22	85
30	50
34	64
40	60
60	30
66	40
70	25

22	85
30	50
34	64

X =

Index	Age	Expenditure	Distance to C0	Distance to C1	Distance to C2
...	...	...	...	...	...

X[:, np.newaxis, :] - centroids

broadcasting

(9,1,2)

18	80
20	90
22	85
30	50
34	64
40	60
60	30
66	40
70	25

(3,2)

18	80
18	80
18	80

(3,2)

22	85
30	50
34	64

...	...
...	...
...	...

(3,2)

(9,3,2)

# K-Means

shape=(9,3)

## ❖ Implementation using Numpy

18	80
20	90
22	85
30	50
34	64
40	60
60	30
66	40
70	25

22	85
30	50
34	64

X =

Index	Age	Expenditure	Distance to C0	Distance to C1	Distance to C2
...	...	...	...	...	...

norm(X[:, np.newaxis, :] - centroids, axis=2)

broadcasting

(9,1,2)

18	80
20	90
22	85
30	50
34	64
40	60
60	30
66	40
70	25

(3,2)

18	80
18	80
18	80

(3,2)

22	85
30	50
34	64

...	...
...	...
...	...

(3,2)

(9,3,2)

# K-Means

## ❖ Implementation using Numpy

18	80
20	90
22	85
30	50
34	64
40	60
60	30
66	40
70	25

22	85
30	50
34	64

$L = \text{np.argmin}(D, \text{axis}=1)$

Distance to C0	Distance to C1	Distance to C2
...	...	...
...	...	...
...	...	...
...	...	...
...	...	...
...	...	...
...	...	...
...	...	...

Label
...
...
...
...
...
...
...
...
...

(9,3)

(9,)



Demo

# Outline

SECTION 1

## Introduction

SECTION 2

## Tabular Data

SECTION 3

## Image&Text Data





# K-Means

## ❖ Image data

■ Training data

■ Cats

■ Dogs



# K-Means

## ❖ Image data

■ Training data



From Cat-Dog dataset

# K-Means

## ❖ Image data



(224,224,3)

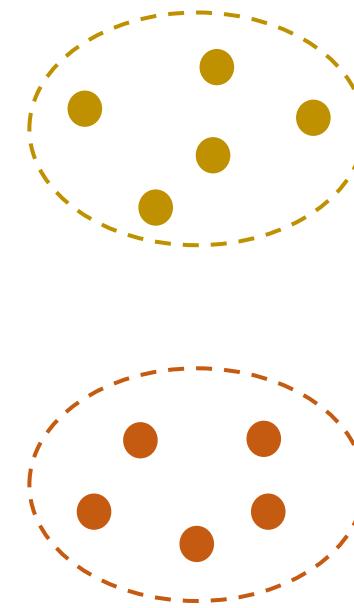


→

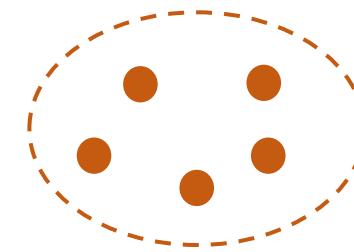
K-Means



(224\*224\*3,)



cluster 0



cluster 1

# K-Means

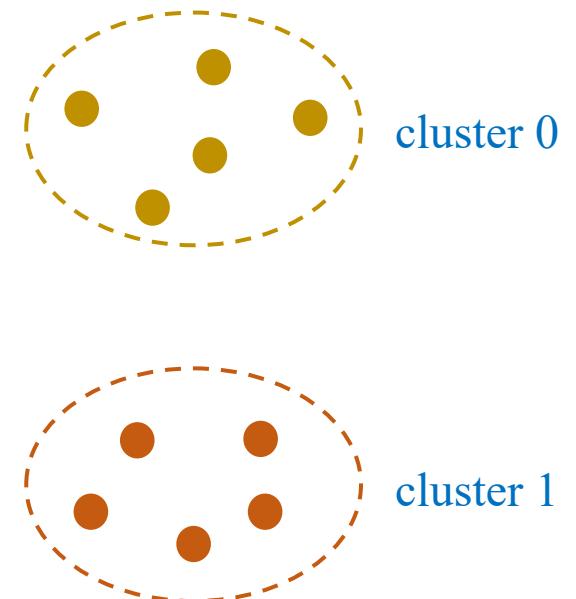
## ❖ Image data



(224,224,3)

**Magic  
Function****K-Means**

(512,)

A vertical black line separating the processing steps from the output results.

cluster 0

cluster 1



Demo

# Text Data

## ❖ Bag of words

### Corpus

**doc1** = “deep learning book”

**doc2** = “machine learning algorithm”

**doc3** = “learning ai from scratch”

**doc4** = “ai vietnam”

### Tokenization



[‘deep’, ‘learning’, ‘book’]

[‘machine’, ‘learning’, ‘algorithm’]

[‘learning’, ‘ai’, ‘from’, ‘scratch’]

[‘ai’, ‘vietnam’]

**Vocabulary** =

deep	learning	book	machine	algorithm	ai	from	scratch	vietnam
------	----------	------	---------	-----------	----	------	---------	---------

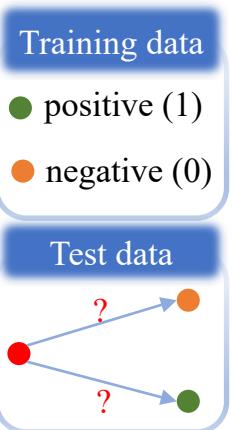
👉 Given a string = “vietnam machine learning deep learning book”

deep	learning	book	machine	algorithm	ai	from	scratch	vietnam
1	2	1	1	0	0	0	0	1
1	1	1	1	0	0	0	0	1

BoW

Binary BoW

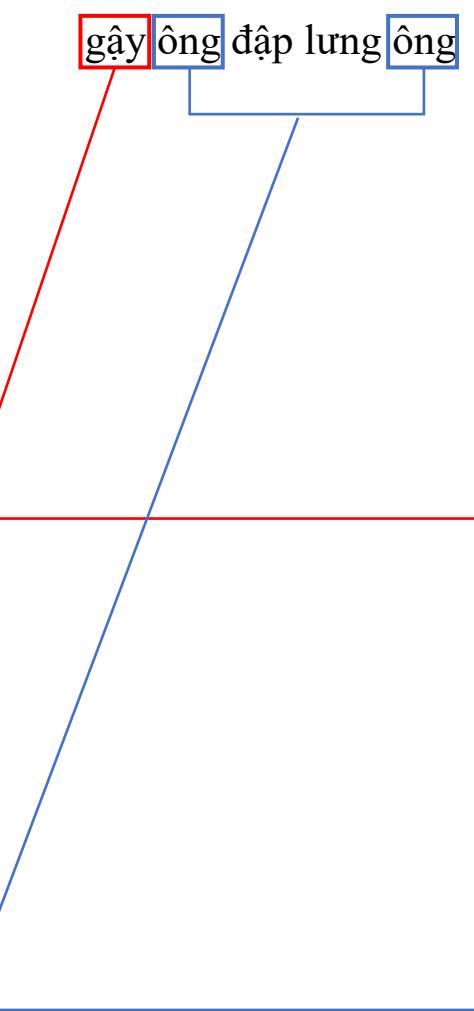
Doc	Label
góp gió gặt bão	1
có làm mới có ăn	1
đất lành chim đậu	1
ăn cháo đá bát	0
gây ông đập lưng ông	0
qua cầu rút ván	0



Vocabulary
bát
bão
chim
cháo
có
cầu
gió
góp
gây
gặt
làm
lành
lưng
mới
qua
rút
ván
ông
ăn
đá
đất
đập
đậu

gây ông đập lưng ông

Tokenization



	doc_0	doc_1	doc_2	doc_3	doc_4	doc_5
bát	0	0	0	1	0	0
bão	1	0	0	0	0	0
chim	0	0	1	0	0	0
cháo	0	0	0	1	0	0
có	0	2	0	0	0	0
cầu	0	0	0	0	0	1
gió	1	0	0	0	0	0
góp	1	0	0	0	0	0
gây	0	0	0	0	1	0
gặt	1	0	0	0	0	0
làm	0	1	0	0	0	0
lành	0	0	1	0	0	0
lưng	0	0	0	0	1	0
mới	0	1	0	0	0	0
qua	0	0	0	0	0	1
rút	0	0	0	0	0	1
ván	0	0	0	0	0	1
ông	0	0	0	0	2	0
ăn	0	1	0	1	0	0
đá	0	0	0	1	0	0
đất	0	0	1	0	0	0
đập	0	0	0	0	1	0
đậu	0	0	1	0	0	0

BoW vectors



Demo

# Summaary

## Introduction

Data processing  
and select K

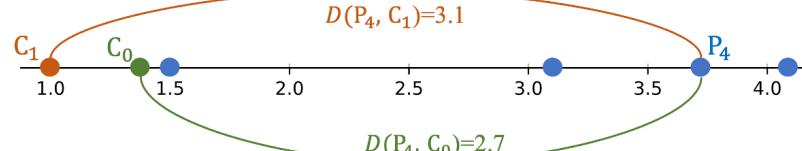
Select K centroids

assign each data  
point to a cluster

update centroids

## Tabular Data

	Index	Length	Distance to C0	Distance to C1	Centroid
C <sub>0</sub>	0	1.4	0.0	0.4	0
C <sub>1</sub>	1	1	0.4	0.0	1
	2	1.5	0.1	0.5	-
	3	3.1	1.7	2.1	-
	4	3.7	2.4	2.8	-
	5	4.1	2.7	3.1	-



## Image Data



