

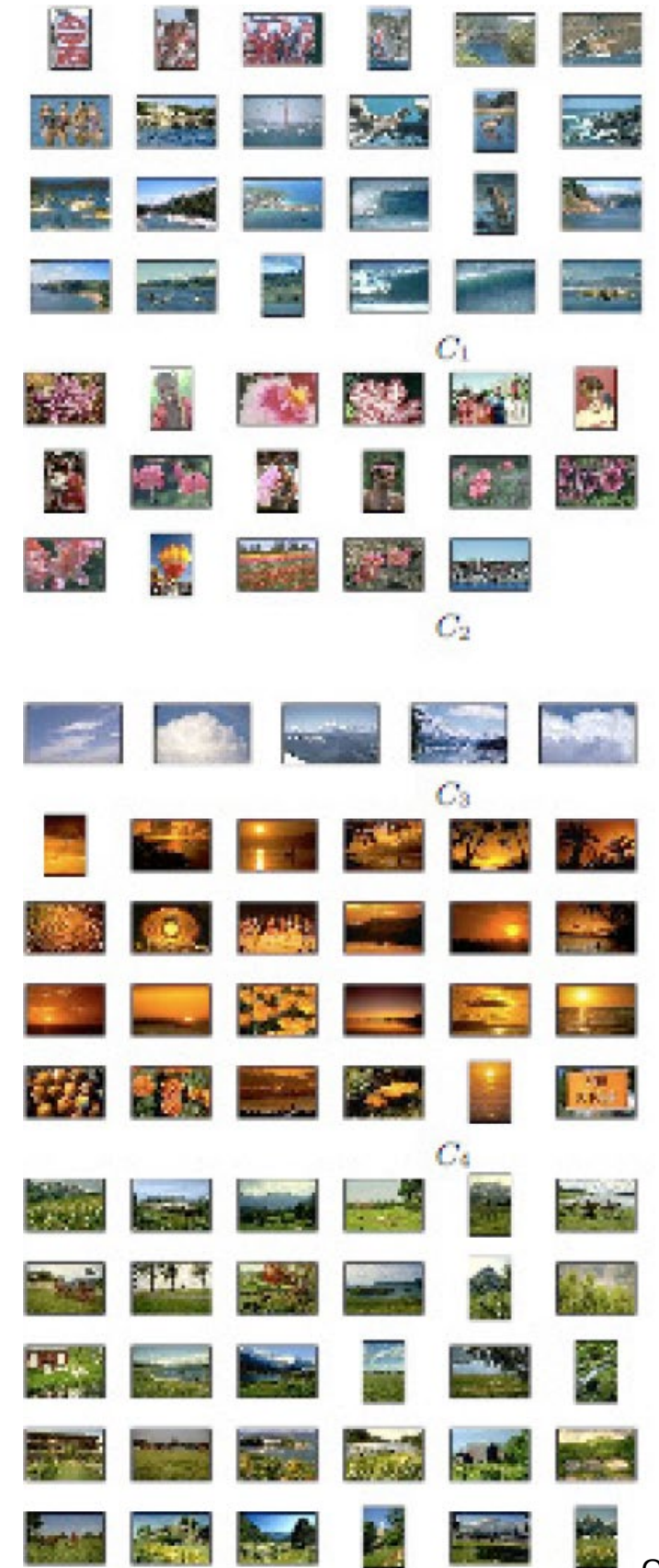
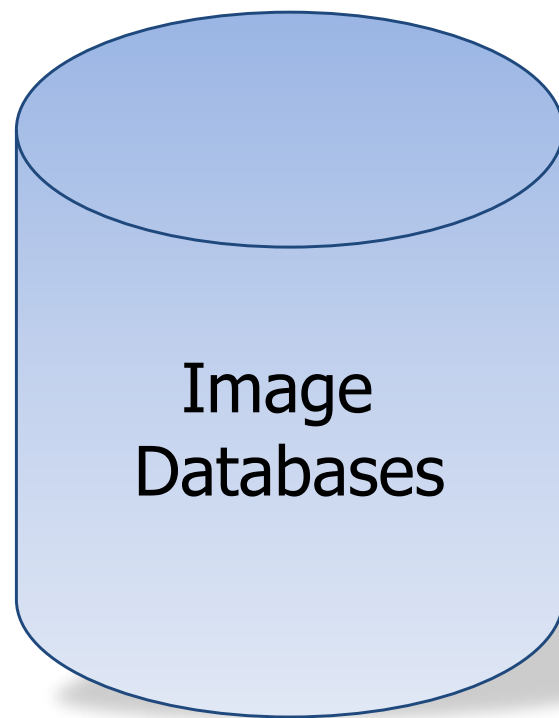
CSE 6740: Computational Data Analysis

Spring 2026

Clustering

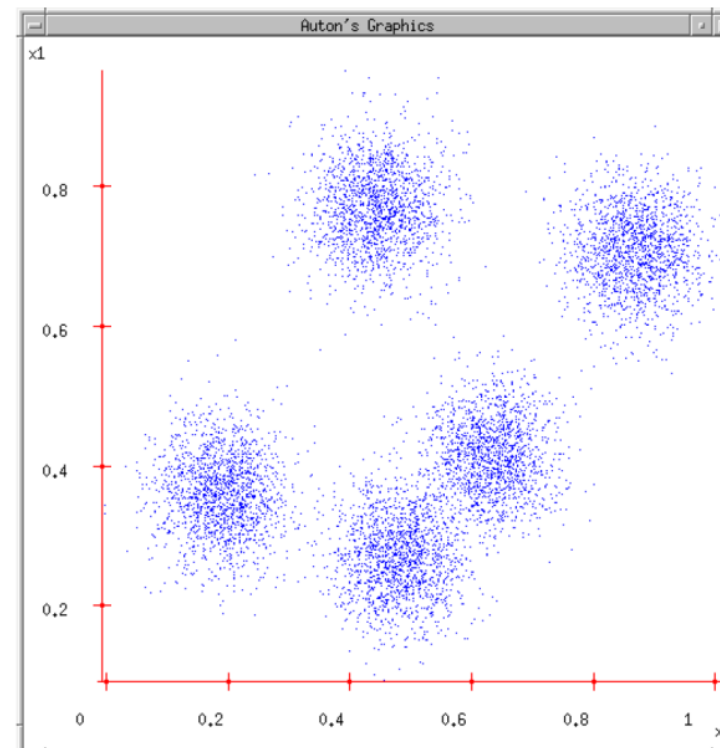
Anqi Wu
01/15

Clustering images



Goal of clustering:

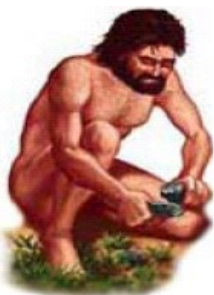
Divide object into groups,
and objects within a group
are more similar than
those outside the group



Cluster other things ...



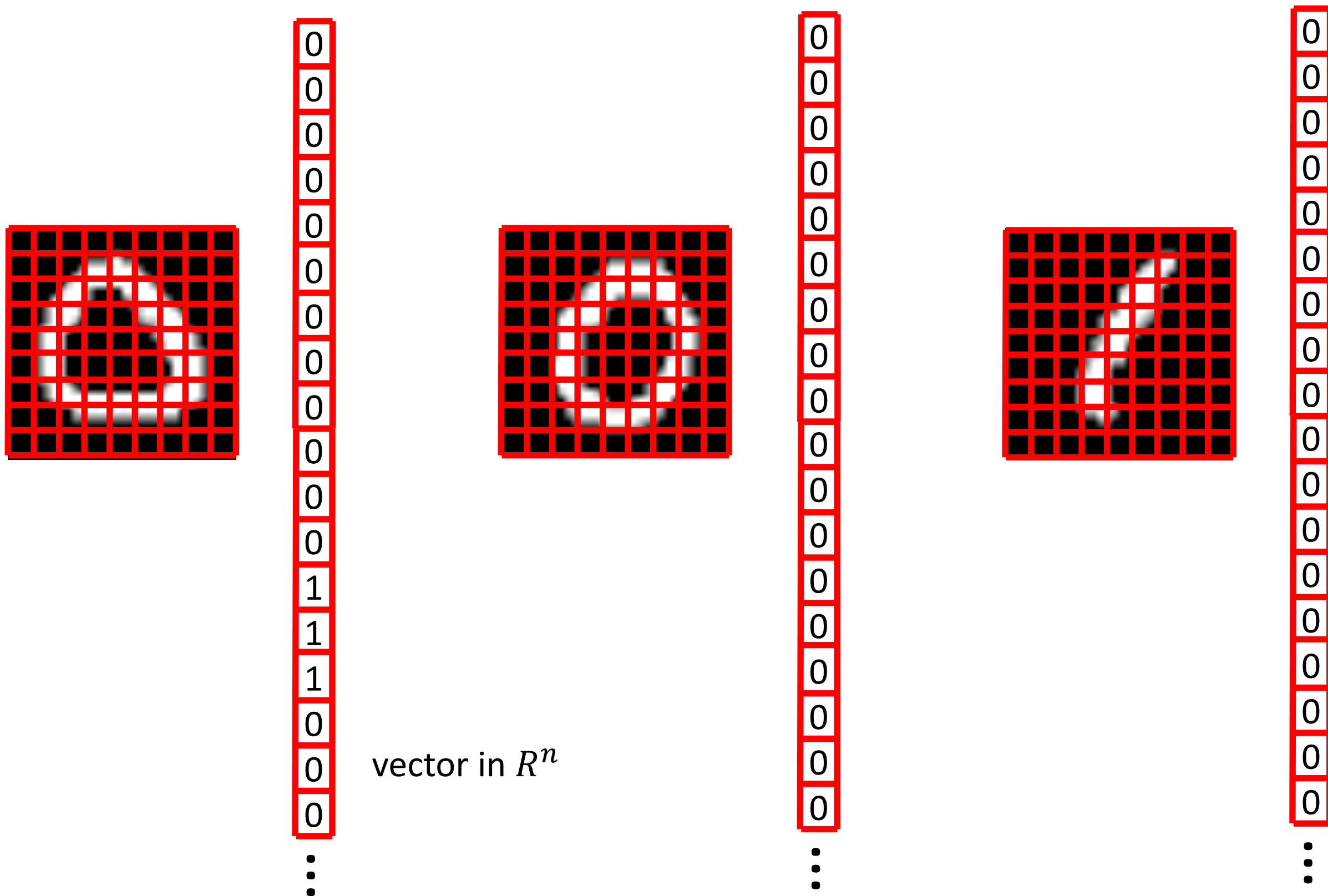
Piotr *Pyotr* *Petros* *Pietro* *Pedro* *Pierre* *Piero* *Peter* *Peder* *Peka* *Peadar*



Cluster handwritten digits

7 2 1 0 4 1 4 9 5 9 0 6 9 0 1 5 9 7 8 4 9 6 6 5 4 0 7 4 0 1 3 1
3 4 7 2 7 1 2 1 1 7 4 2 3 5 1 2 4 4 6 3 5 5 6 0 4 1 9 5 7 8 5 3
7 4 6 4 3 0 7 0 2 9 1 7 3 2 9 7 7 6 2 7 8 4 7 3 6 1 3 6 8 3 1 4
1 7 6 9 6 0 5 4 9 9 2 1 9 4 8 7 3 9 7 4 4 4 9 2 5 4 7 6 7 9 0 5
8 5 6 6 5 7 8 1 0 1 6 4 6 7 3 1 7 1 8 2 0 2 9 9 5 5 1 5 6 0 3 4
4 6 5 4 6 5 4 5 1 4 4 7 2 3 2 7 1 8 1 8 1 8 5 0 8 4 2 5 0 1 1 1
0 9 0 3 1 6 4 2 3 6 1 1 1 3 9 5 2 9 4 5 9 3 9 0 3 6 5 5 7 2 2 7
1 2 8 4 1 7 3 3 8 8 7 9 2 2 4 1 5 9 8 7 2 3 0 4 4 2 4 1 9 5 7 7
2 8 2 6 8 5 7 7 9 1 8 1 8 0 3 0 1 9 9 4 1 8 2 1 2 9 7 5 9 2 6 4
1 5 4 2 9 2 0 4 0 0 2 8 4 7 1 2 4 0 2 7 4 3 3 0 0 3 1 9 6 5 2 5
9 2 9 3 1 4 2 0 7 1 1 2 1 5 3 3 9 7 8 6 5 6 1 3 8 1 0 5 1 3 1 5
5 6 1 8 5 1 1 9 4 6 2 2 5 0 6 5 6 3 7 2 0 8 8 5 4 1 1 4 0 3 3 7
6 1 6 2 1 9 2 8 6 1 9 5 2 5 4 4 2 8 3 8 2 4 5 0 3 1 7 7 5 7 9 7
1 9 2 1 4 2 9 2 0 4 9 1 4 8 1 8 4 5 9 8 8 3 7 6 0 0 3 0 2 6 6 4
9 3 3 3 2 3 9 1 2 6 8 0 5 6 6 6 3 8 8 2 7 5 8 9 6 1 8 4 1 2 5 9
1 9 7 5 4 0 8 9 9 1 0 5 2 3 7 8 9 4 0 6 3 9 5 2 1 3 1 3 6 5 7 8
2 2 6 3 2 6 5 4 8 9 7 1 3 0 3 8 3 1 9 3 4 4 6 4 2 1 8 2 5 4 8 8
4 0 0 2 3 2 7 7 0 8 7 4 4 7 9 6 9 0 9 8 0 4 6 0 6 3 5 4 8 3 3 9
3 3 3 7 8 0 8 7 1 7 0 6 5 4 3 8 0 9 6 3 8 0 9 9 6 8 6 8 5 7 8 6
0 2 4 0 2 2 3 1 9 7 5 1 0 8 4 6 2 4 7 9 3 2 9 8 2 2 9 2 7 3 5 9
1 8 0 2 0 5 1 1 3 7 6 7 1 2 5 8 0 3 7 1 4 0 9 1 8 6 7 7 4 3 4 9
1 9 3 1 7 3 9 7 6 9 1 3 7 8 3 3 6 7 2 9 5 8 5 1 1 4 4 3 1 0 7 7
0 7 9 4 4 8 5 5 4 0 8 2 1 0 8 4 5 0 4 0 6 1 3 3 2 6 7 2 6 9 3 1
4 6 2 5 4 2 0 6 2 1 7 3 4 1 0 5 4 3 1 1 7 4 9 9 4 8 4 0 2 4 5 1
1 6 4 7 1 9 4 2 4 1 5 5 3 8 3 1 4 5 6 8 9 4 1 5 3 8 0 3 2 5 1 2
8 3 4 4 0 8 8 3 3 1 7 3 5 9 6 3 2 6 1 3 6 0 7 2 1 7 1 4 2 4 2 1
7 9 6 1 1 2 4 8 1 7 7 4 8 0 2 3 1 3 1 0 7 7 0 3 5 5 2 7 6 6 9 2
8 3 5 2 2 5 6 0 8 2 9 2 8 8 8 8 7 4 9 3 0 6 6 3 2 1 3 2 2 9 3 0
0 5 7 8 1 4 4 6 0 2 9 1 4 7 4 7 3 9 8 8 4 7 1 2 1 2 2 3 2 3 2 3
9 1 7 4 0 3 5 5 8 6 5 2 6 7 6 6 3 2 7 9 1 1 7 5 6 4 9 5 1 3 3 4
7 8 9 1 1 6 9 1 4 4 5 4 0 6 2 2 3 1 5 1 2 0 3 8 1 2 6 7 1 6 2 3
9 0 1 2 2 0 8 9

How to represent objects?



See demo `kmeans_digit.m`

Formal statement of clustering problem

- Given m data points, $\{x^1, x^2, \dots, x^m\} \in R^n$
- Find k cluster centers, $\{c^1, c^2, \dots, c^k\} \in R^n$
- And assign each data point i to one cluster, $\pi(i) \in \{1, \dots, k\}$
- Such that the averaged square distances from each data point to its respective cluster center is small

$$\min_{c, \pi} \frac{1}{m} \sum_{i=1}^m \|x^i - c^{\pi(i)}\|^2$$

K-means algorithm

- Initialize k cluster centers, $\{c^1, c^2, \dots, c^k\}$, randomly
- Do
 - Decide the cluster memberships of each data point, x^i , by assigning it to the nearest cluster center (**cluster assignment**)

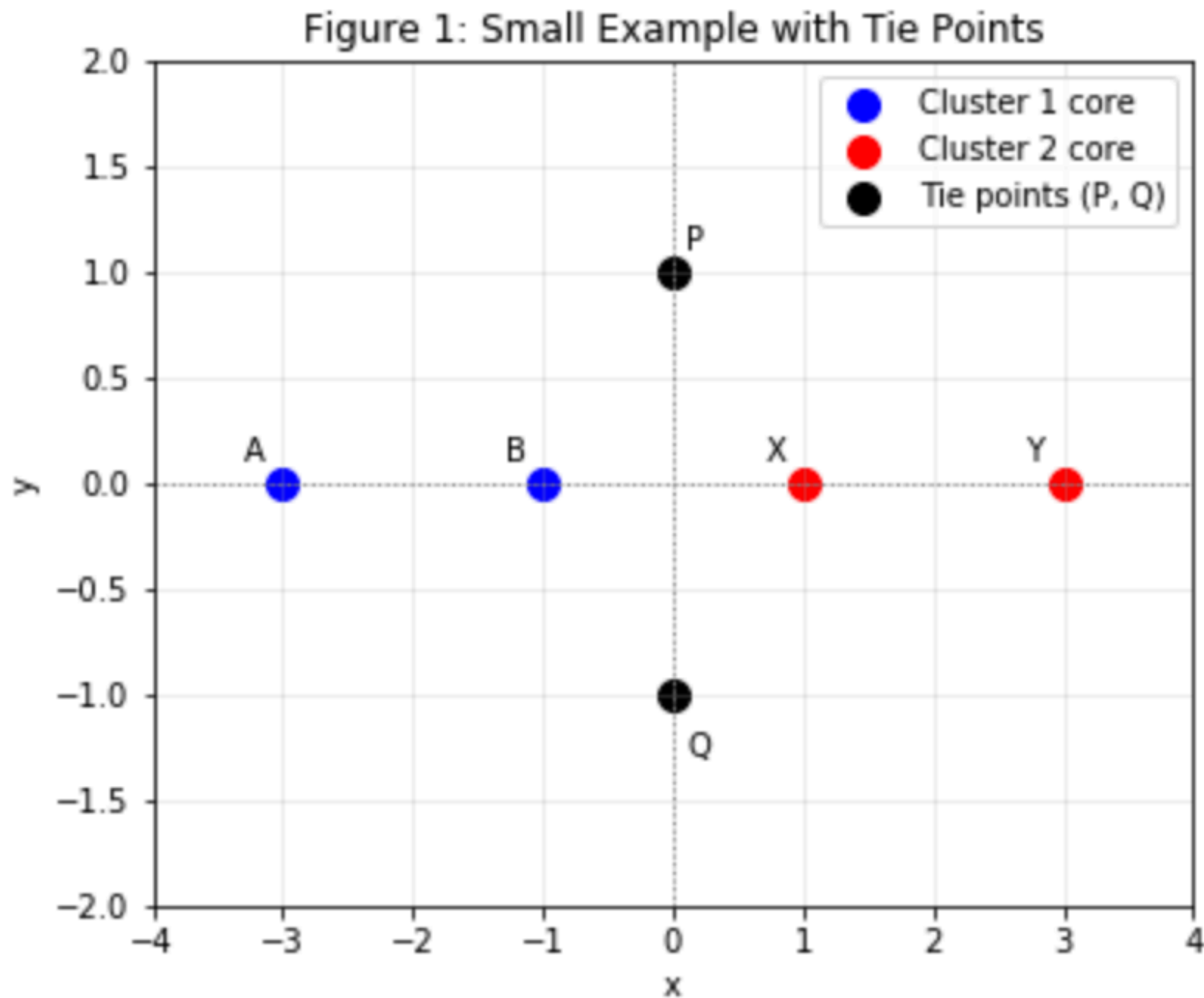
$$\pi(i) = \underset{j=1, \dots, k}{\operatorname{argmin}} \|x^i - c^j\|^2$$

- Adjust the cluster centers (**center adjustment**)

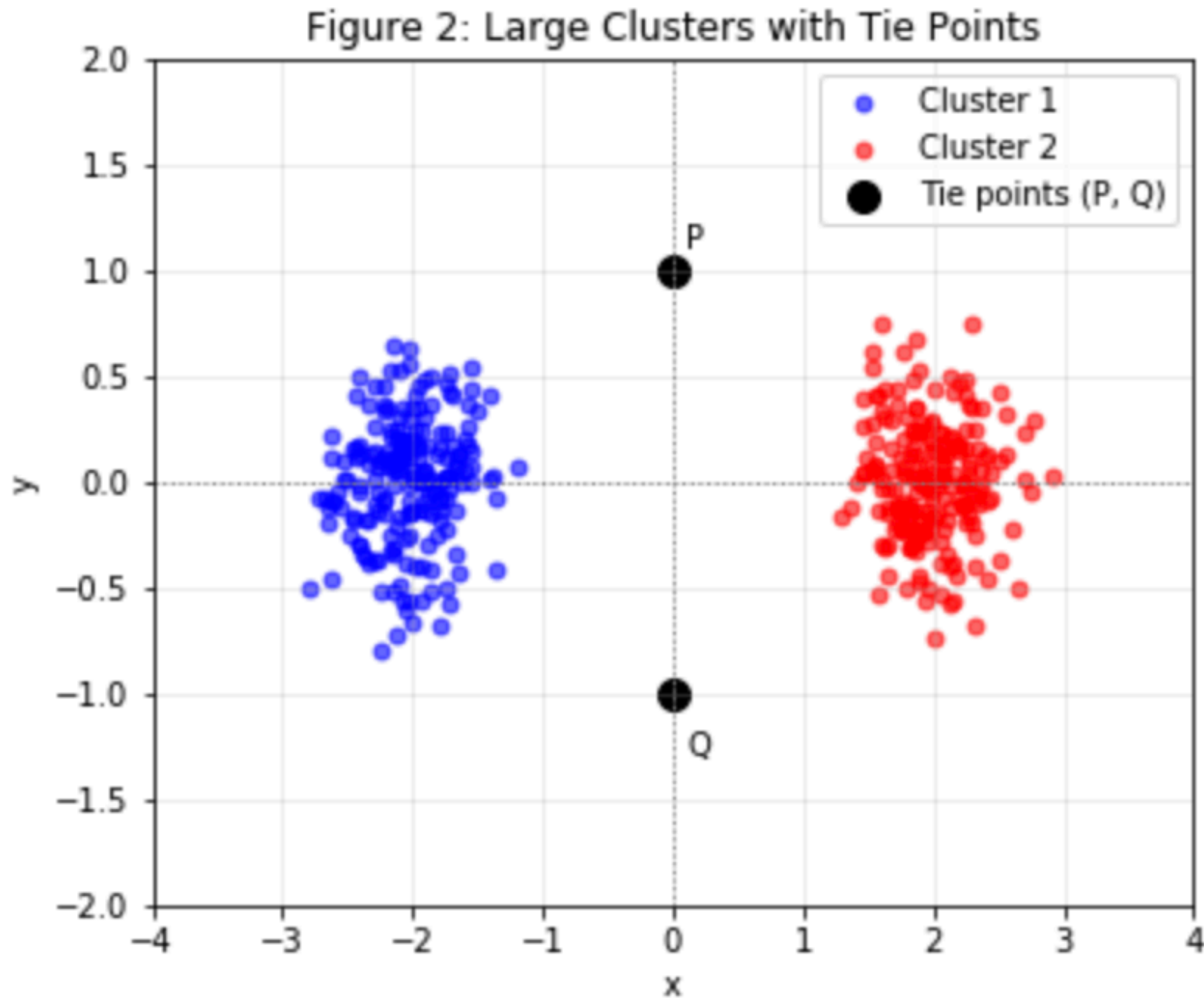
$$c^j = \frac{1}{|\{i: \pi(i) = j\}|} \sum_{i: \pi(i)=j} x^i$$

- While any cluster center has been changed
- Note: Tie points on cluster boundaries might flip between clusters without shifting centers.

K-means algorithm

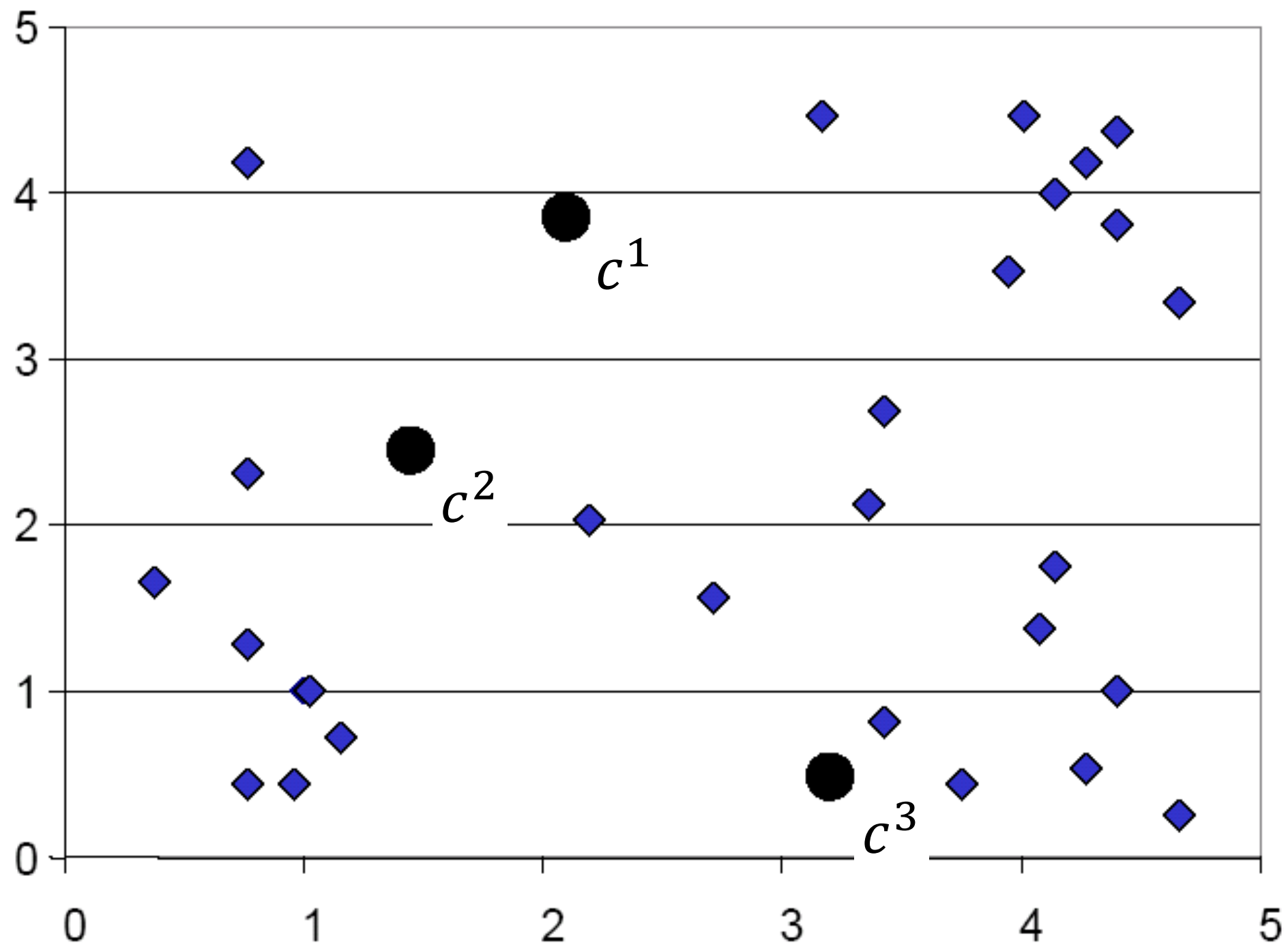


K-means algorithm

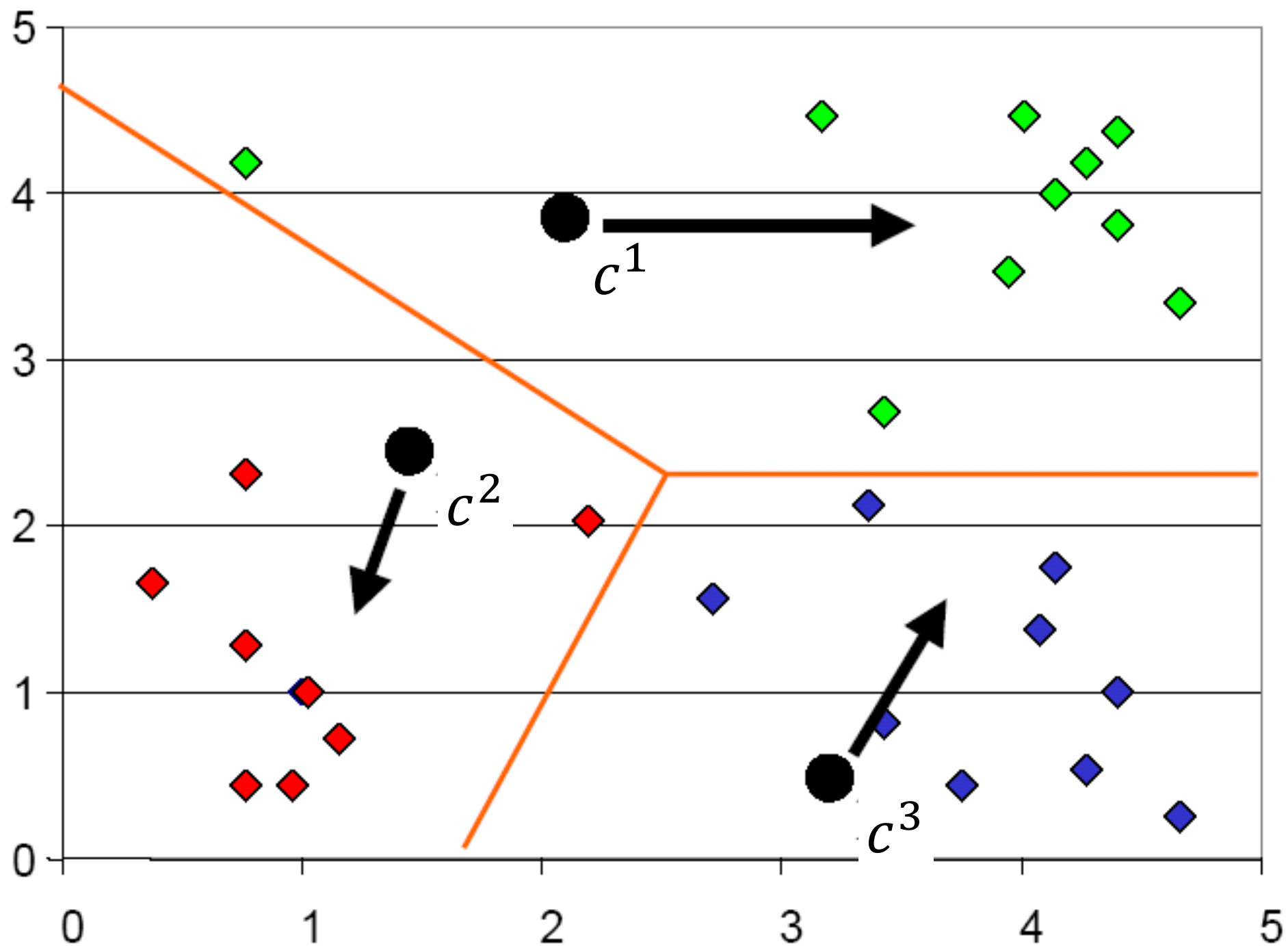


See demo `kmeans_animation.m`

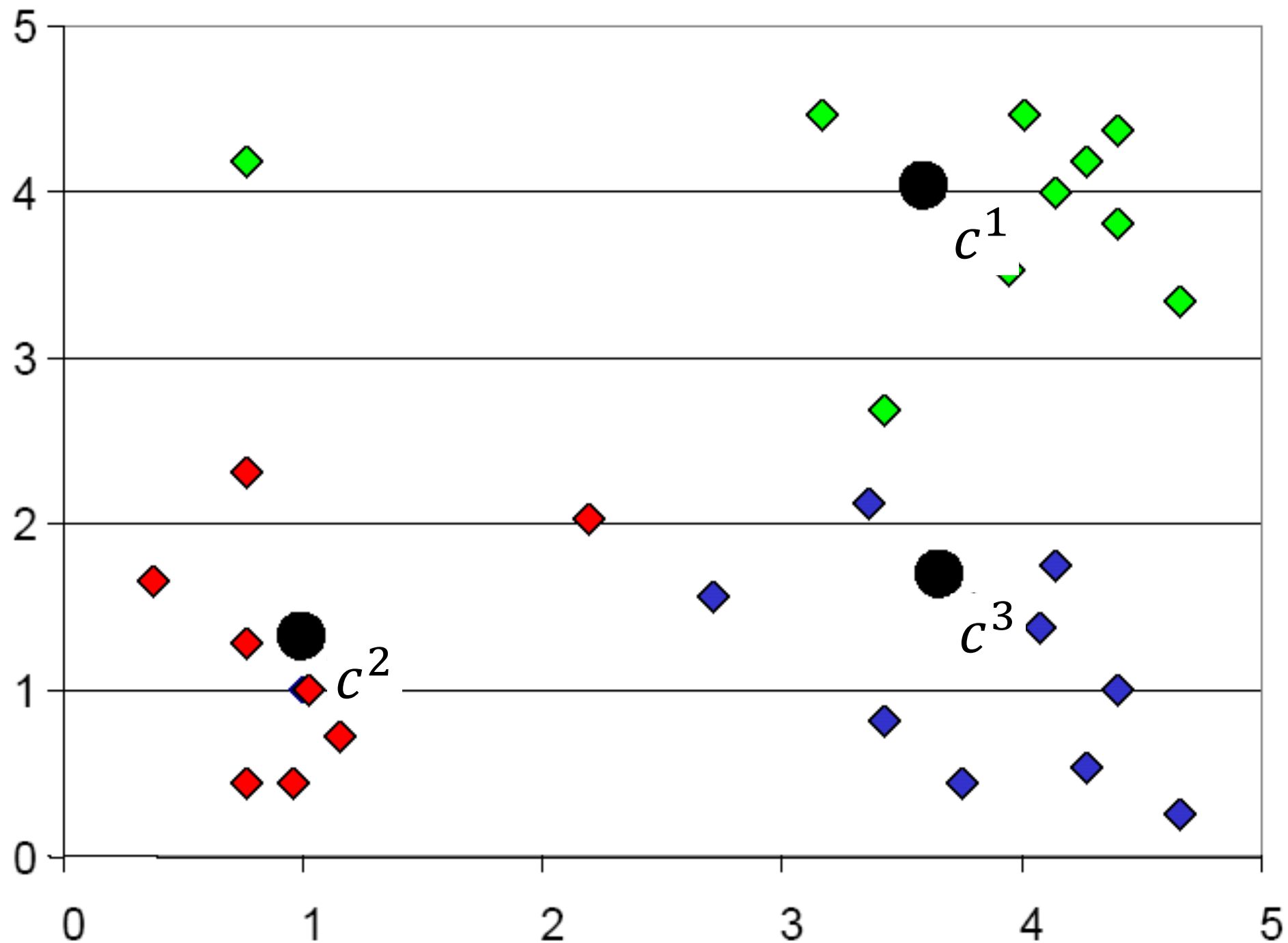
K-means: step 1



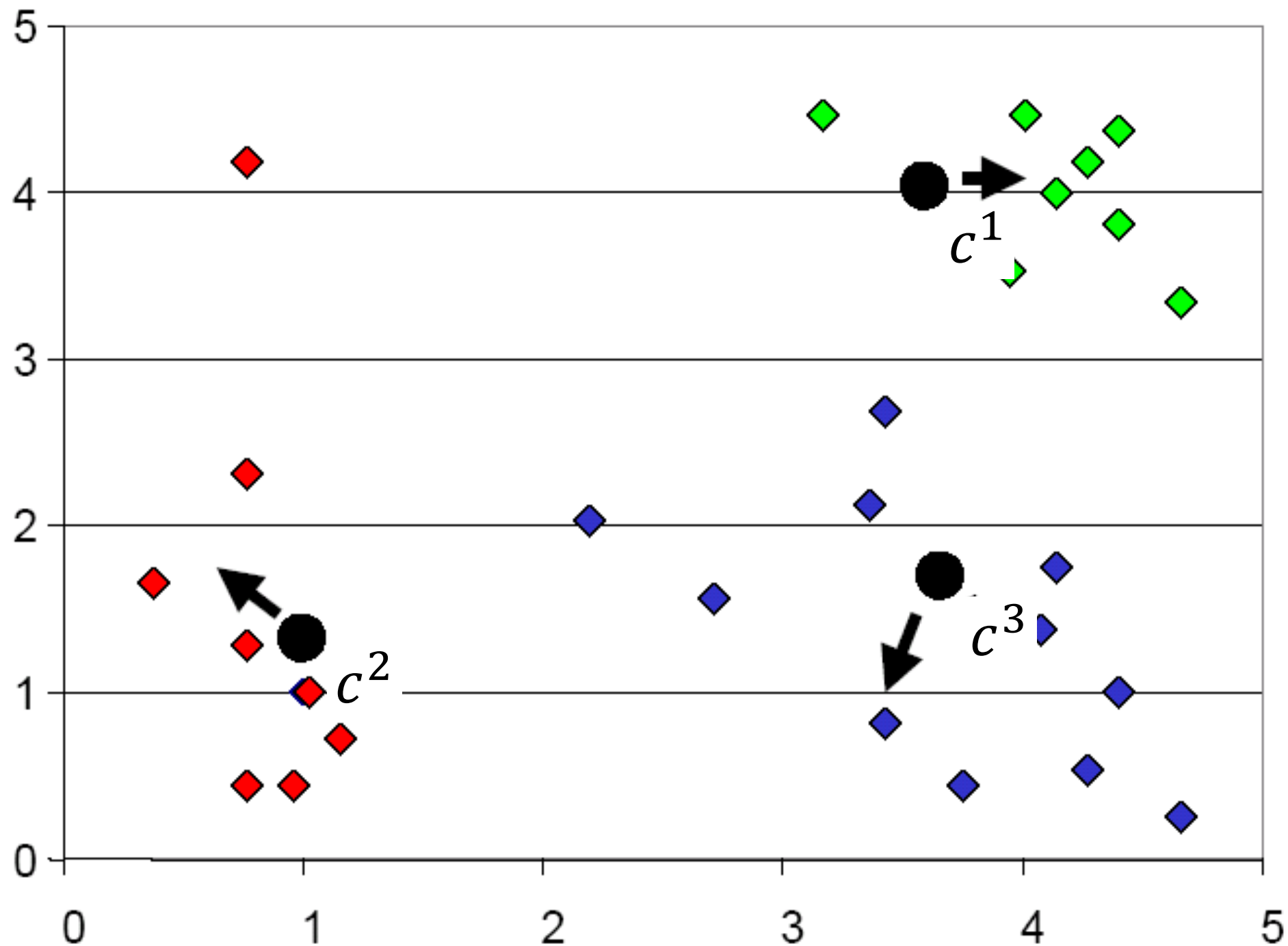
K-means: step 2



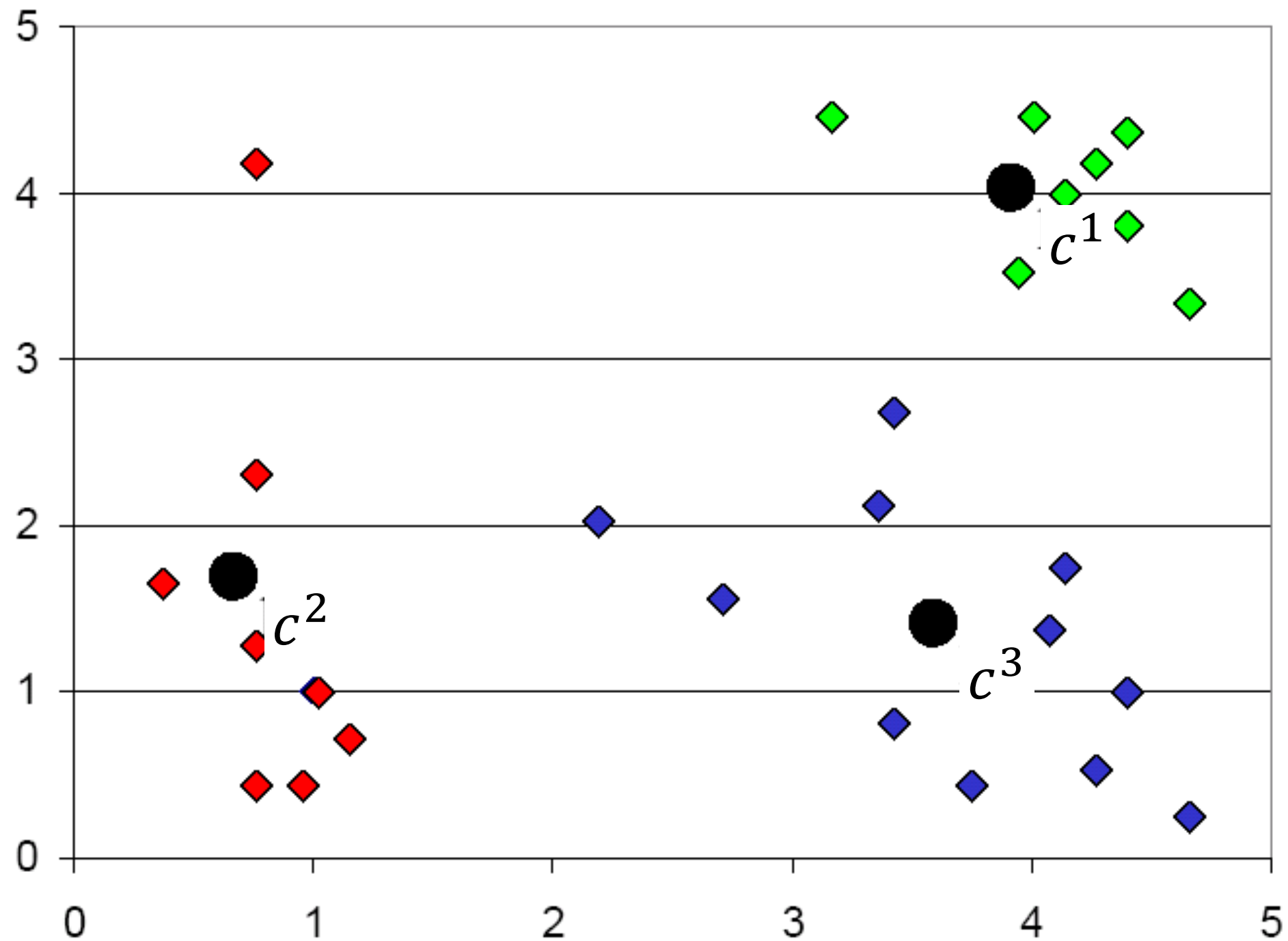
K-means: step 3



K-means: step 4

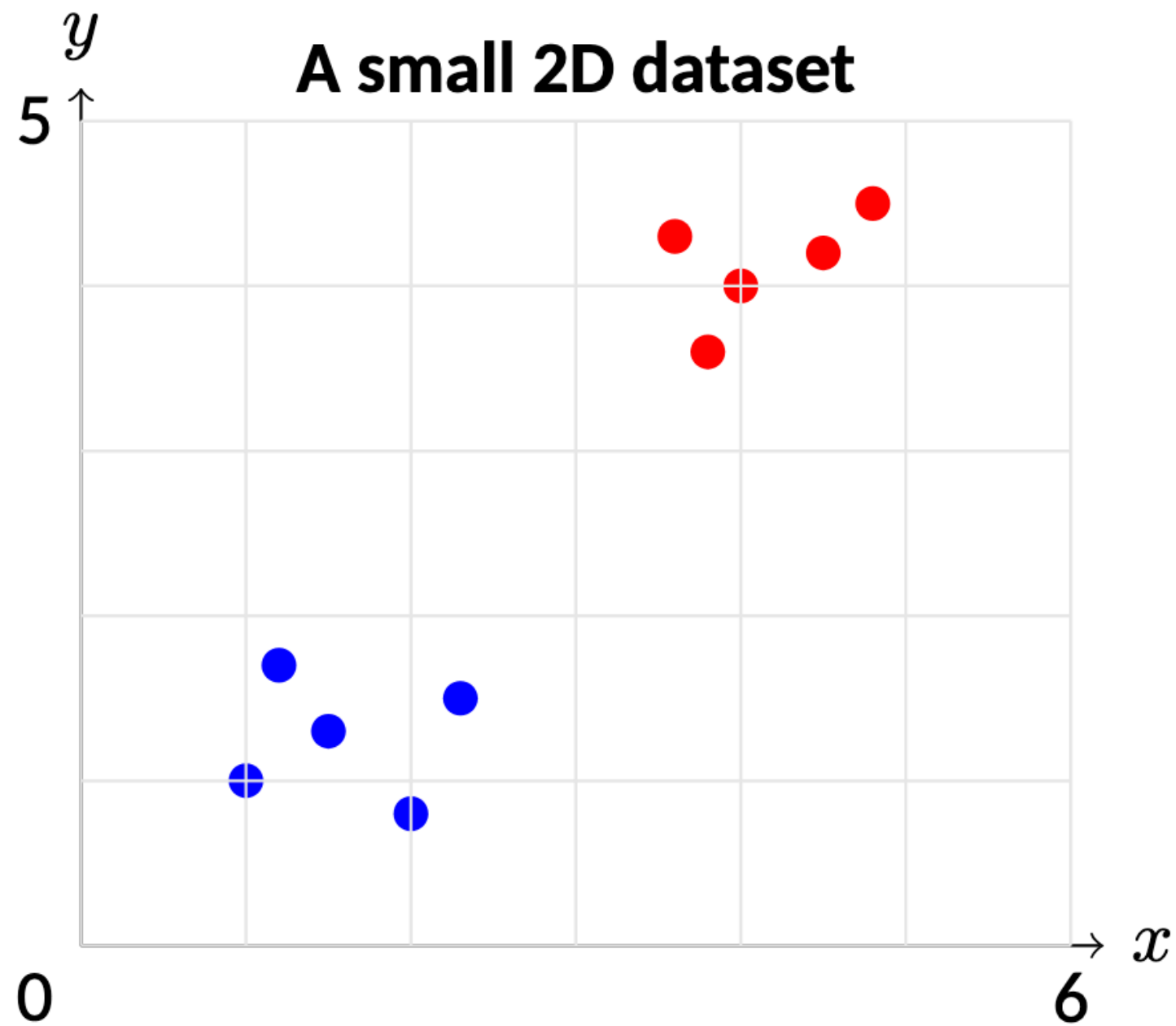


K-means: step 5

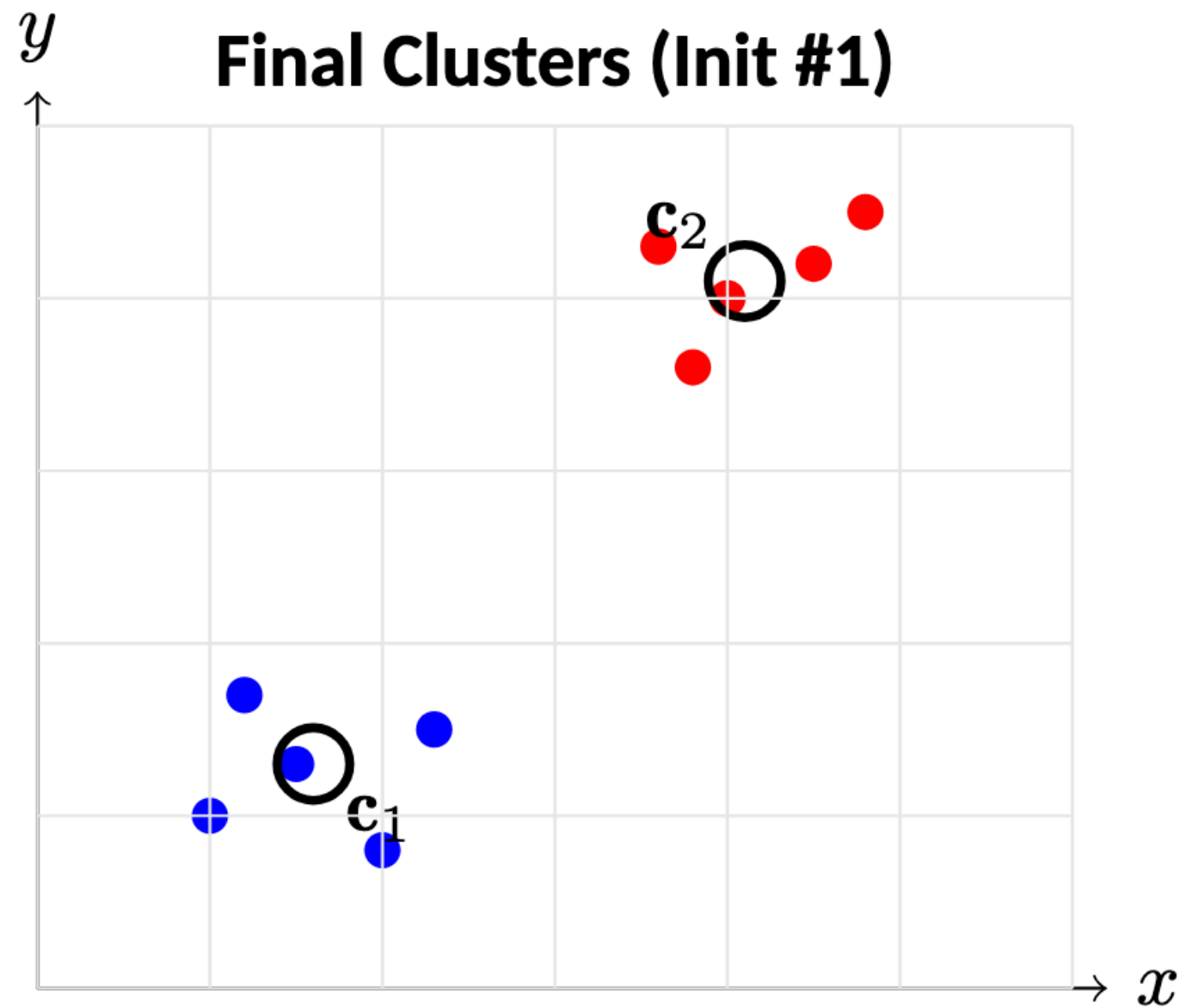


Questions

- Will different initialization lead to different results?
 - Yes
 - No
 - Sometimes
- Will the algorithm always stop after some iteration?
 - Yes
 - No (we have to set a maximum number of iterations)
 - Sometimes

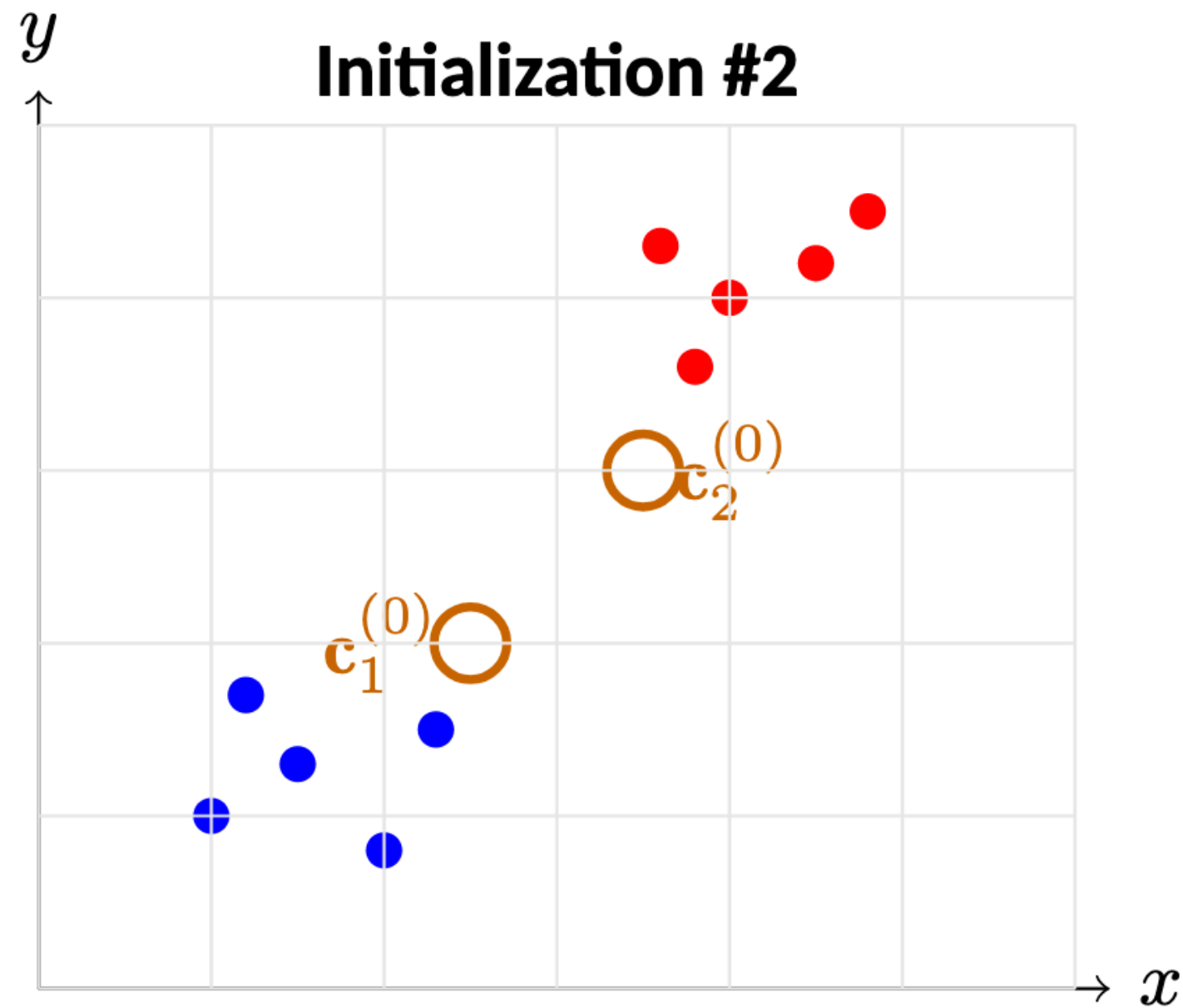


Questions



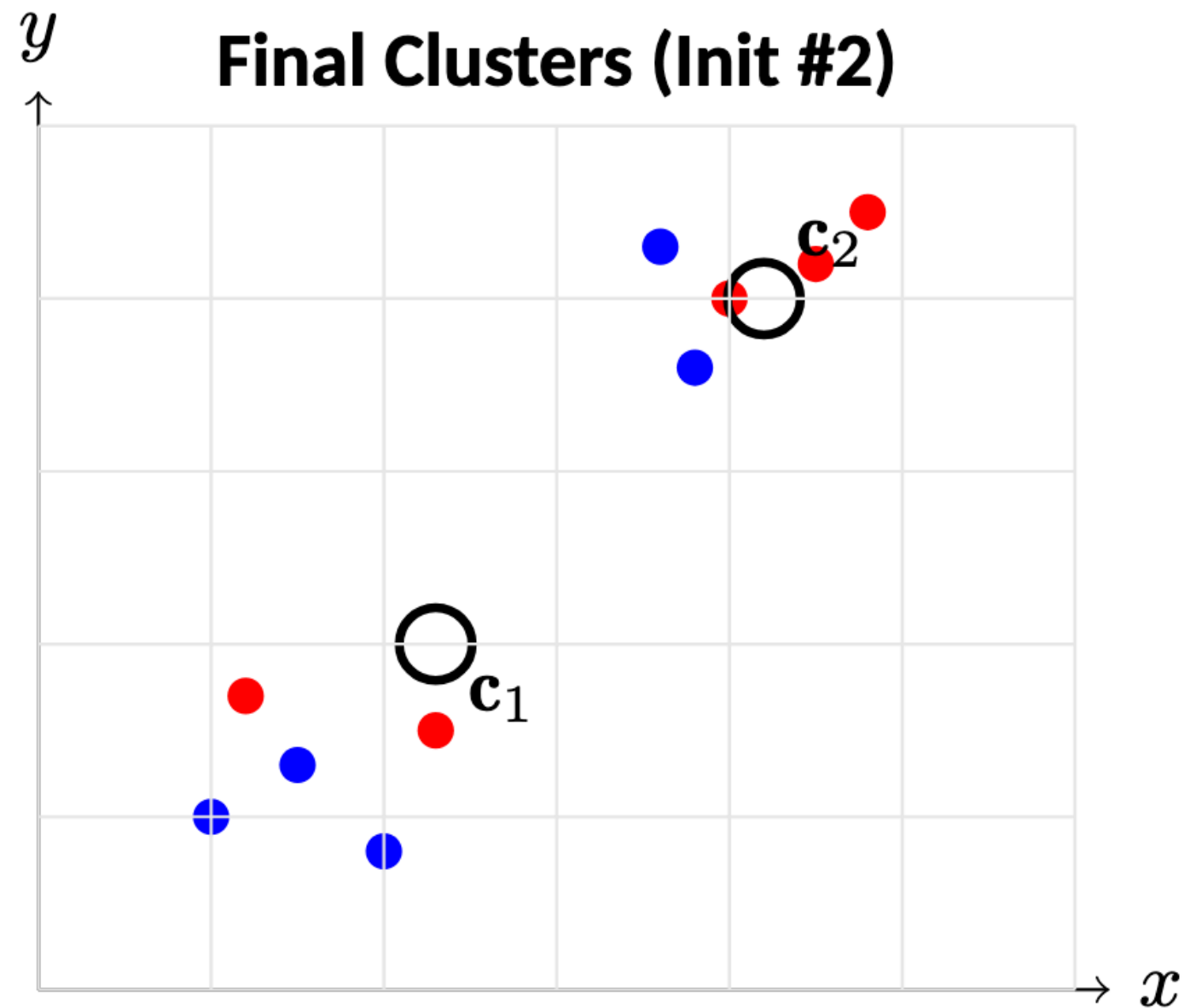
K-Means converges to one partition of the data.

Questions



A different, more central initial placement for centroids.

Questions



A different final grouping emerges, reflecting a different local optimum.

Questions

- Will different initialization lead to different results?
 - Yes
 - No
 - Sometimes
- Will the algorithm always stop after some iteration?
 - Yes
 - No (we have to set a maximum number of iterations)
 - Sometimes

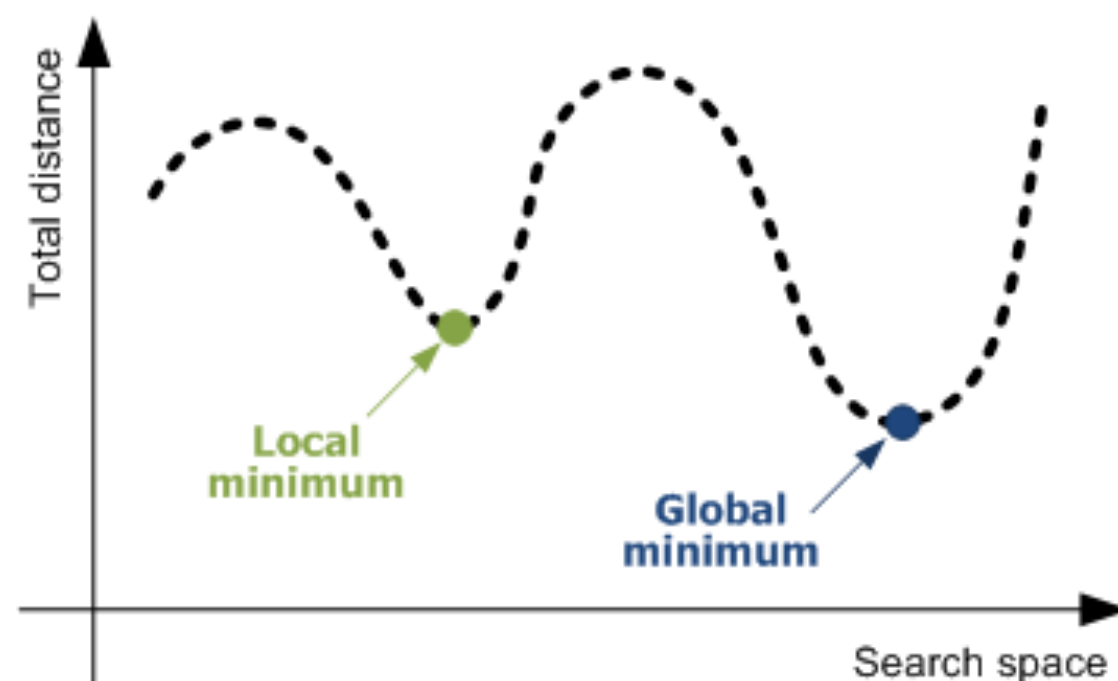
Clustering is NP-hard in general

- Find k cluster centers, $\{c^1, c^2, \dots, c^k\} \in R^n$, and assign each data point i to one cluster, $\pi(i) \in \{1, \dots, k\}$, to minimize

$$\min_{c, \pi} \frac{1}{m} \sum_{i=1}^m \|x^i - c^{\pi(i)}\|^2$$

NP-hard!

- A search problem over the space of discrete assignments
 - For all m data point together, there are k^m possibility
 - The cluster assignment determines cluster centers, and vice versa



Convergence of kmeans algorithm

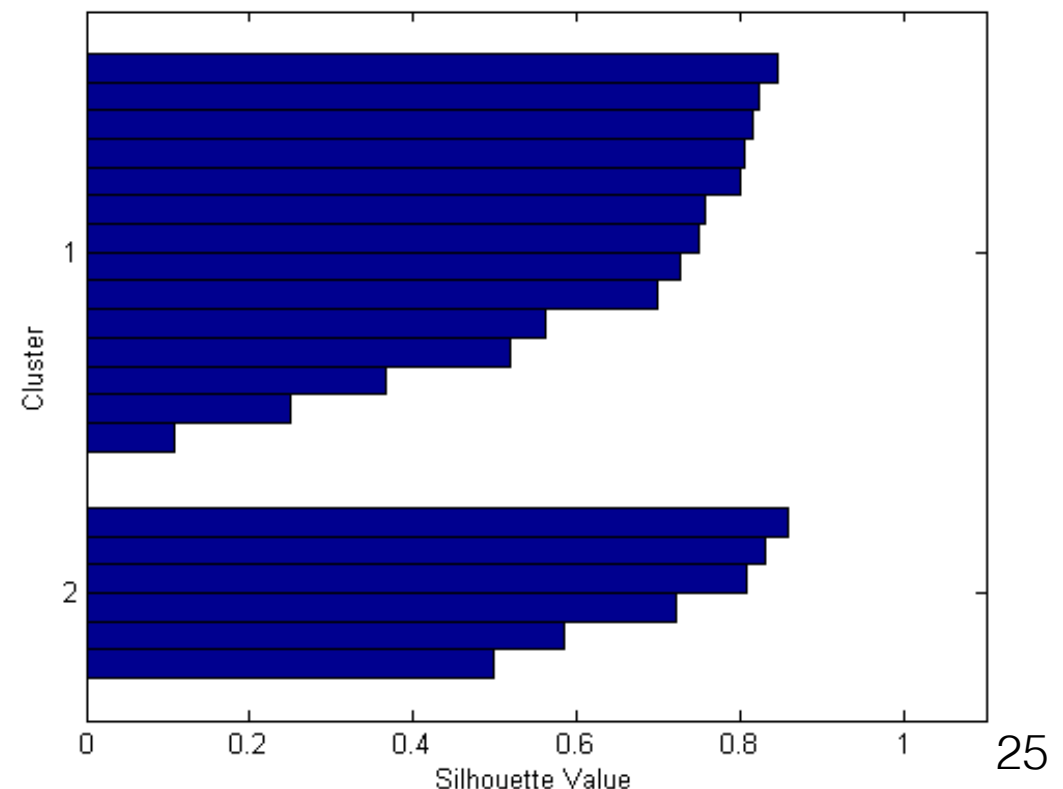
- Will kmeans objective oscillate?

$$\frac{1}{m} \sum_{i=1}^m \|x^i - c^{\pi(i)}\|^2$$

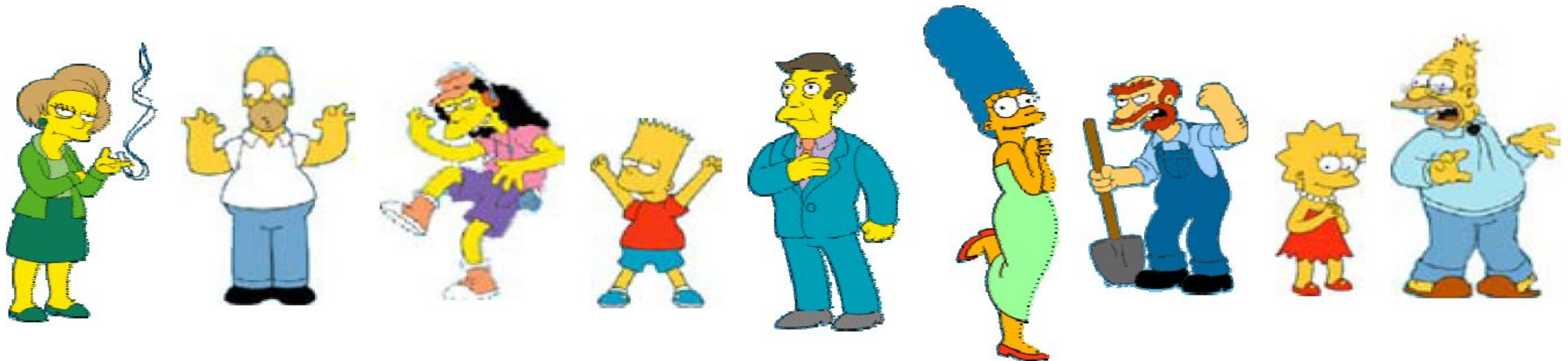
- The minimum value of the objective is finite
- Each iteration of kmeans algorithm decrease the objective
 - Cluster assignment step decreases objective
 - $\pi(i) = \operatorname{argmin}_{j=1,\dots,k} \|x^i - c^j\|^2$ for each data point i
 - Center adjustment step decreases objective
 - $c^j = \frac{1}{|\{i:\pi(i)=j\}|} \sum_{i:\pi(i)=j} x^i = \operatorname{argmin}_c \sum_{i:\pi(i)=j} \|x^i - c\|^2$

How many clusters?

- Fixed a-priori? Data-driven approach?
- Silhouette value: $S_i = \frac{b_i - a_i}{\max(a_i, b_i)}$ (one heuristic)
 - a_i : the average distance from the i th point to the other points in the same cluster as i ,
 - b_i : the minimum average distance from the i th point to points in a different cluster, minimized over clusters.
- No gold standard method
 - Often determine by trial-and-error

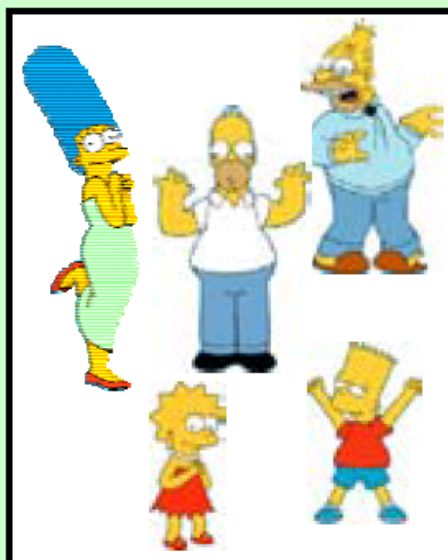


Are these everything about clustering?



What is consider similar/dissimilar?

Clustering is subjective



Simpson's Family



School Employees

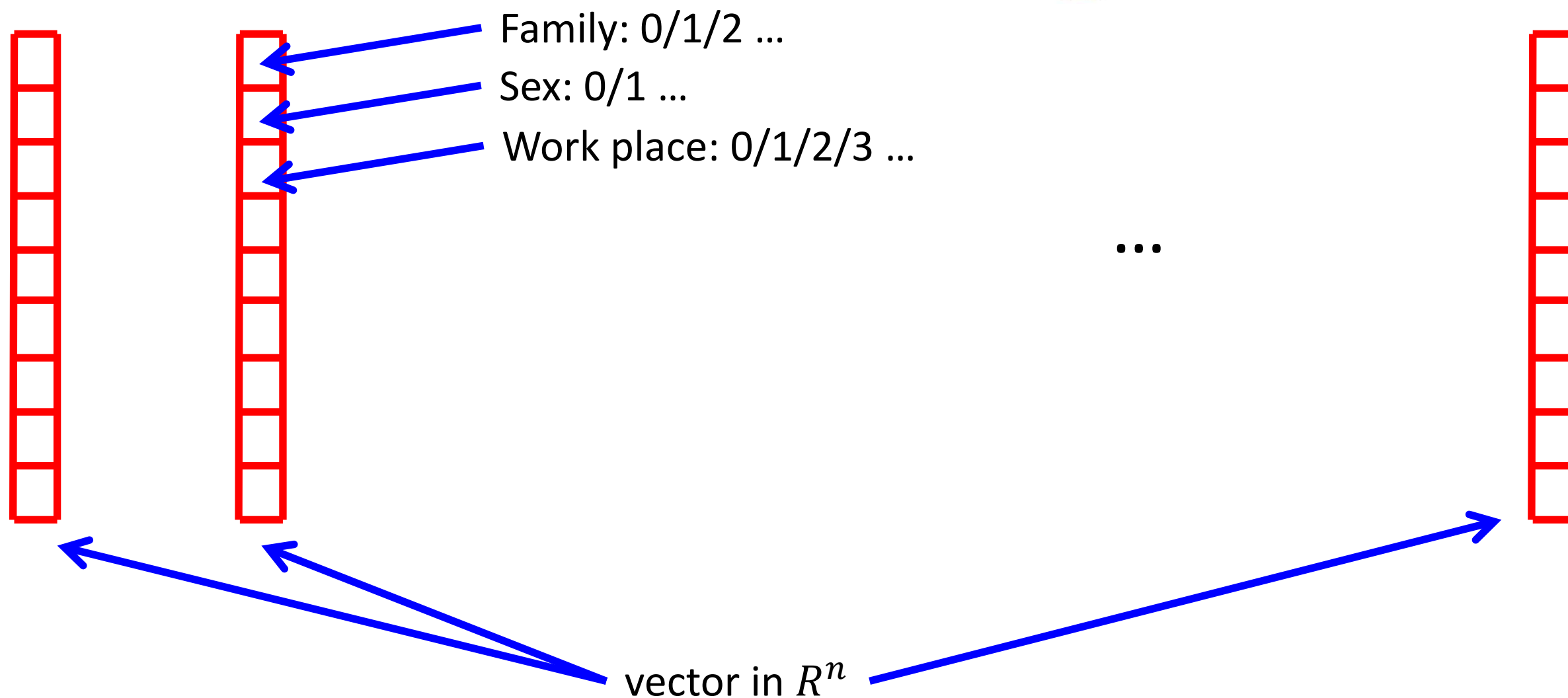
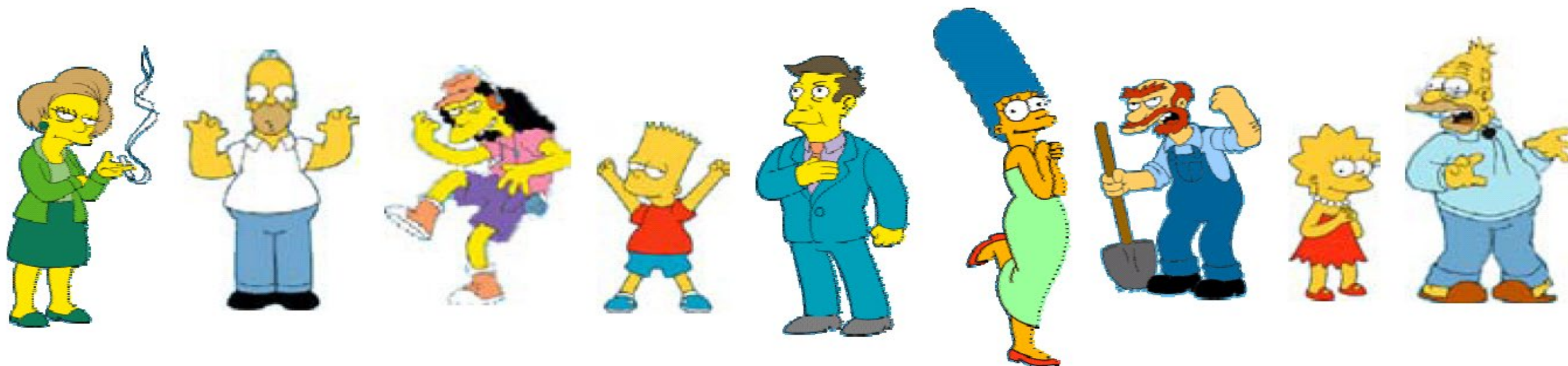


Females

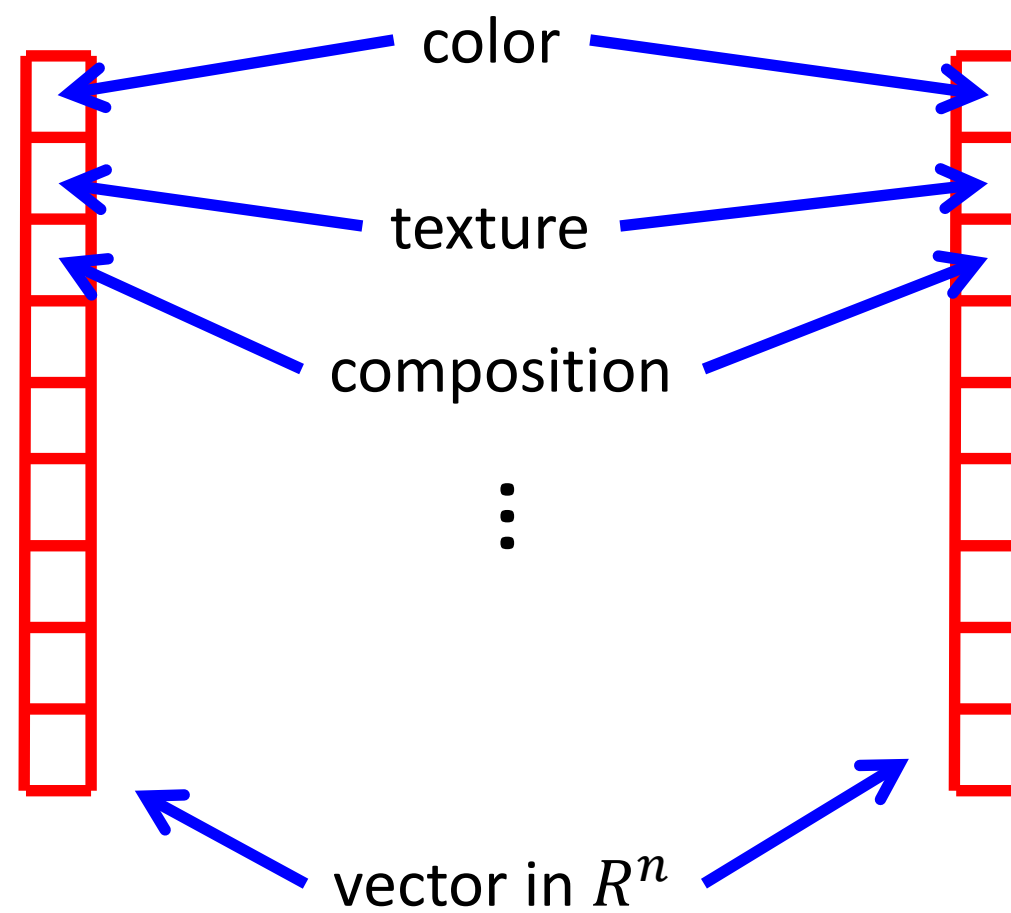


Males

Objects in real life



Images of different sizes



You pick your similarity/dissimilarity



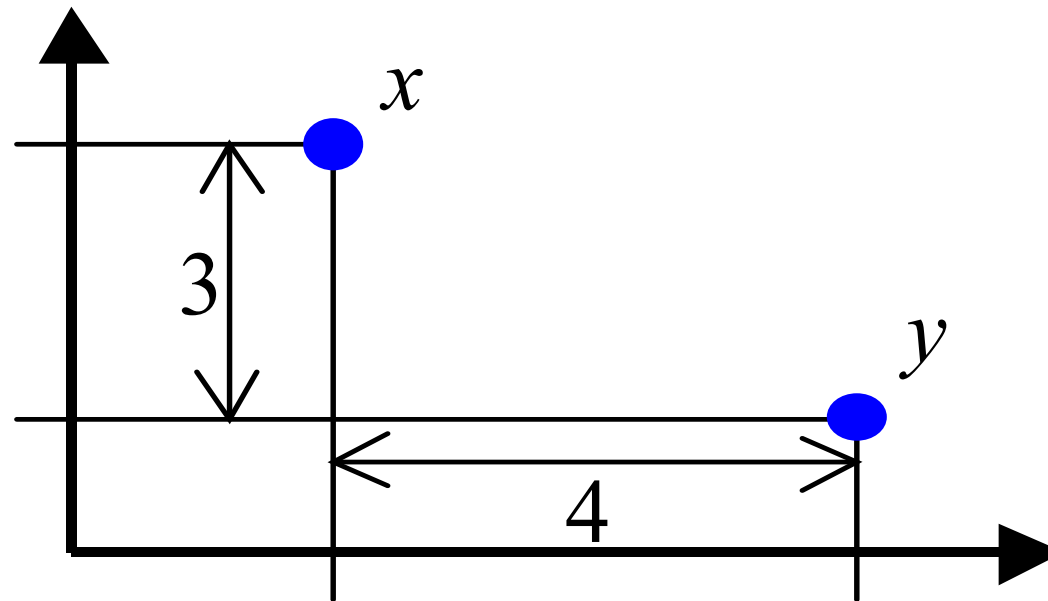
What similarity/dissimilarity function?

- Desired properties of dissimilarity function
 - Symmetry: $d(x, y) = d(y, x)$
 - *Otherwise you could claim "Alex looks like Bob, but Bob looks nothing like Alex"*
 - Positive separability: $d(x, y) = 0$, if and only if $x = y$
 - *Otherwise there are objects that are different, but you cannot tell apart*
 - Triangular inequality: $d(x, y) \leq d(x, z) + d(z, y)$
 - *Otherwise you could claim "Alex is very like Bob, and Alex is very like Carl, but Bob is very unlike Carl"*

Distance functions for vectors

- Suppose two data points, both in R^n
 - $x = (x_1, x_2, \dots, x_n)^\top$
 - $y = (y_1, y_2, \dots, y_n)^\top$
- Euclidian distance: $d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$
- Minkowski distance: $d(x, y) = \sqrt[p]{\sum_{i=1}^n (x_i - y_i)^p}$
 - Euclidian distance: $p = 2$
 - Manhattan distance: $p = 1, d(x, y) = \sum_{i=1}^n |x_i - y_i|$
 - “inf”-distance: $p = \infty, d(x, y) = \max_{i=1}^n |x_i - y_i|$

Distance example



- Euclidian distance: $\sqrt{4^2 + 3^2} = 5$
- Manhattan distance: $4 + 3 = 7$
- “inf”-distance: $\max\{4, 3\} = 4$

Hamming distance

- Manhattan distance is also called *Hamming distance* when all features are binary
- Count the number of difference between two binary vectors
- Example, $x, y \in \{0,1\}^{17}$

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
x	0	1	1	0	0	1	0	0	1	0	0	1	1	1	0	0	1
y	0	1	1	1	0	0	0	0	1	1	1	1	1	1	0	1	1

$$d(x, y) = 5$$

Edit distance

- Transform one of the objects into the other, and measure how much effort it takes

x	I	N	T	E	*	N	T	I	O	N
y	*	E	X	E	C	U	T	I	O	N
	d	s	s		i	s				

d: deletion (cost 5)

s: substitution (cost 1)

i: insertion (cost 2)

$$d(x, y) = 5 \times 1 + 3 \times 1 + 1 \times 2 = 10$$

Generalized K-means algorithm

- Initialize k cluster centers, $\{c^1, c^2, \dots, c^k\}$, randomly
- Do
 - Decide the cluster memberships of each data point, x^i , by assigning it to the nearest cluster center

$$\pi(i) = \operatorname{argmin}_{j=1, \dots, k} d(x^i, c^j)$$

- Adjust the cluster centers

$$c^j = \operatorname{argmin}_{v \in R^n} \sum_{i: \pi(i)=j} d(x^i, v)$$

squared Eclidian distance:

$$c^j = \frac{1}{\#\{\pi(i) = j\}} \sum_{i: \pi(i)=j} x^i$$

- While any cluster center has been changed

-
- A phylogenetic tree illustrating the relationships between various names. The tree is rooted at the top and branches downwards. The names are: Piotr, Pyotr, Petros, Pietro, Pedro, Pierre, Piero, Peter, Peder, Peka, and Peadar. The tree is color-coded: blue for the main branches, red for the branch leading to Piotr and Pyotr, green for the branch leading to Petros, Pietro, Pedro, Pierre, and Piero, pink for the branch leading to Peter and Peder, and black for the branch leading to Peka and Peadar. A red arc highlights the relationship between the red and pink branches.



Bottom up hierarchical clustering

- Assign each data point to its own cluster, $g_1 = \{x_1\}$, $g_2 = \{x_2\}$, ..., $g_m = \{x_m\}$, and let $G = \{g_1, g_2, \dots, g_m\}$

$$D(g_i, g_j) = \min_{x \in g_i, y \in g_j} d(x, y)$$

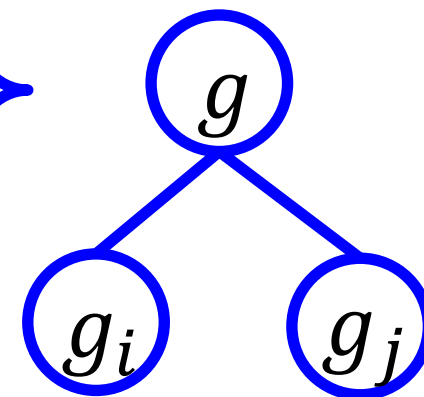
- Do

- Find two clusters to merge: $i, j = \operatorname{argmin}_{1 \leq i, j \leq |G|} D(g_i, g_j)$

- Merge the two clusters to a new cluster: $g \leftarrow g_i \cup g_j$
- keep track of relations

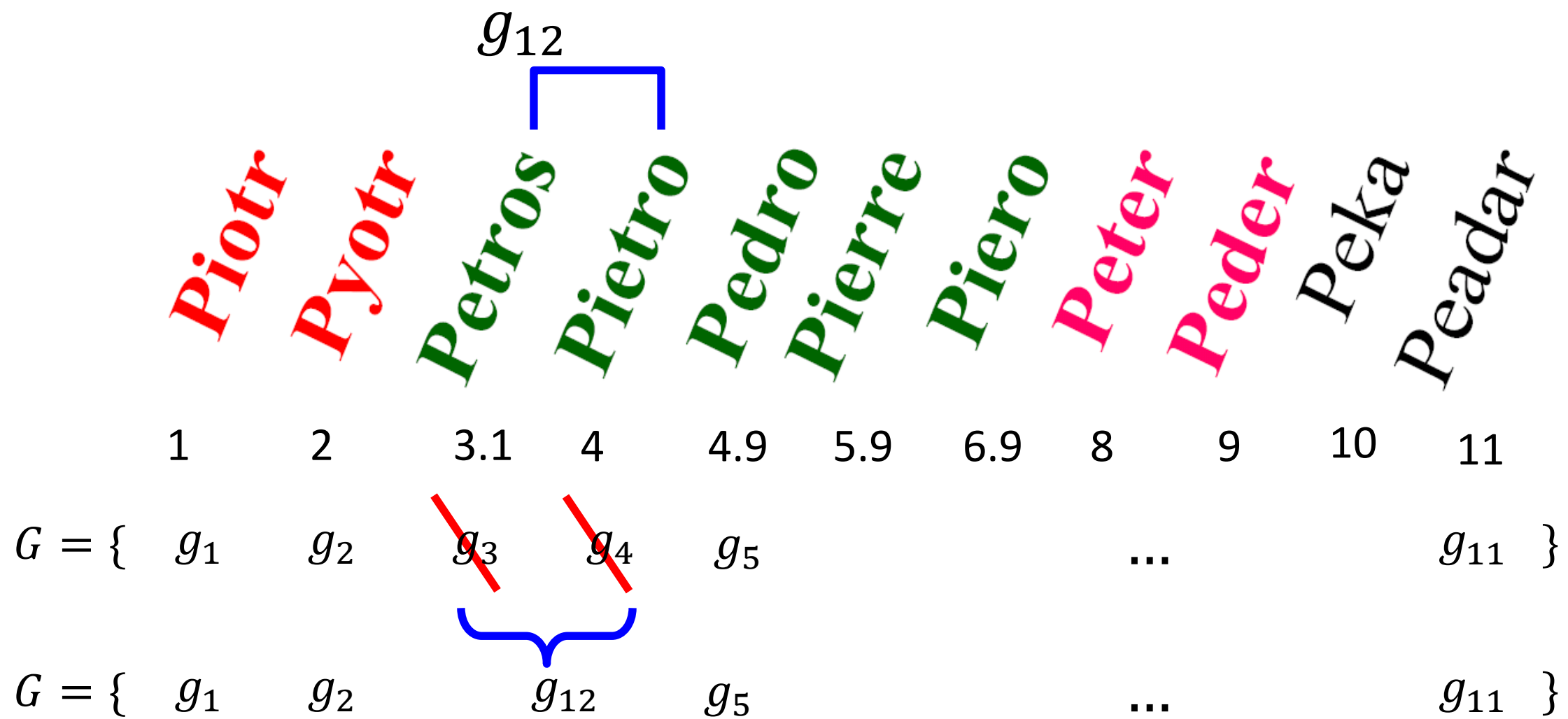
- Remove the merged clusters: $G \leftarrow G \setminus g_i, G \leftarrow G \setminus g_j$

- Add the new cluster: $G \leftarrow G \cup \{g\}$

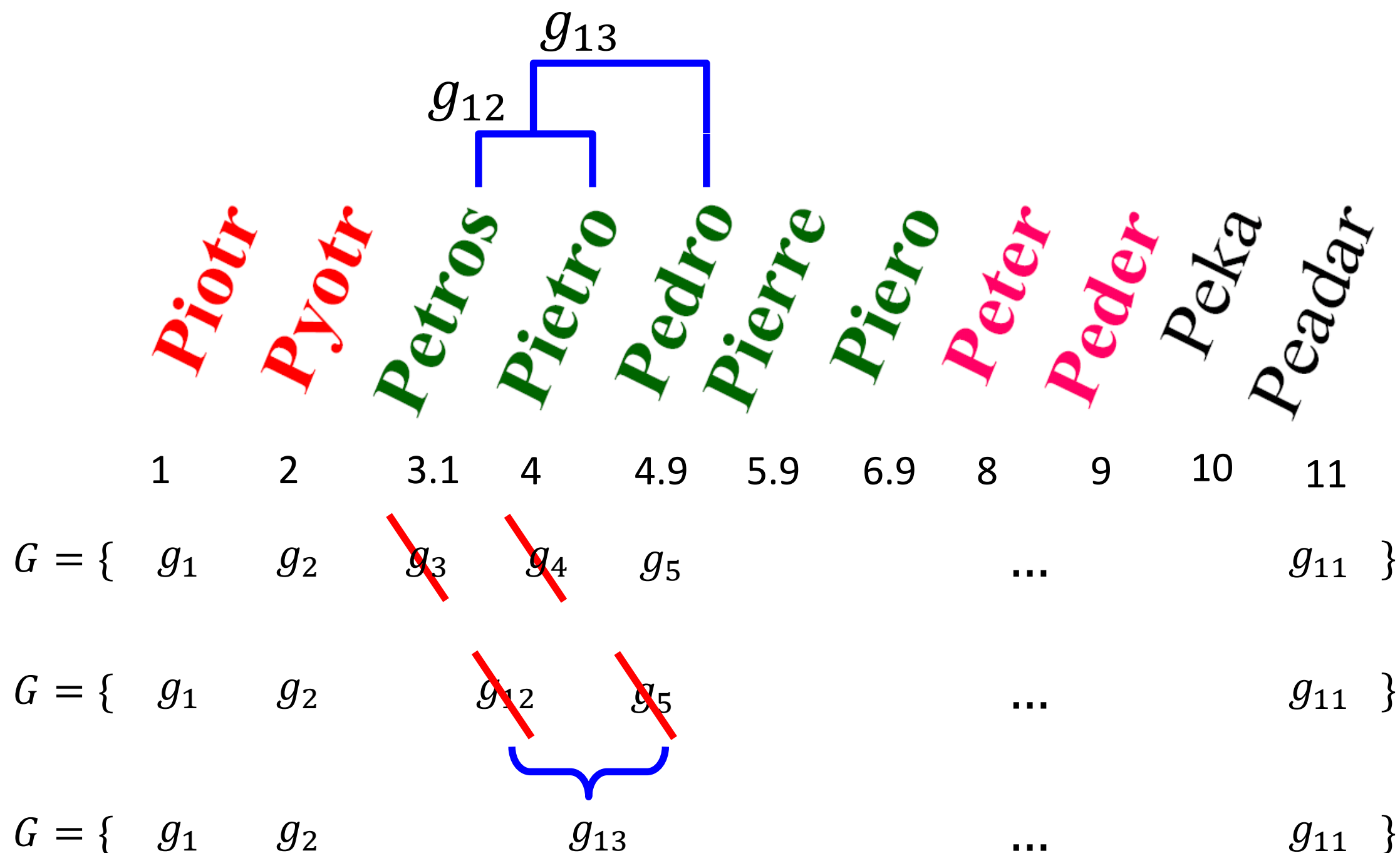


- While $|G| > 1$

Hierarchical clustering: step-2



Hierarchical clustering: step-3



Kmeans vs Hierarchal clustering

	Kmeans	HC
Number of clusters	If there is a specific number of clusters in the dataset, but the group they belong to is unknown	It is easier to determine the number of clusters by hierarchical clustering's dendrogram
Cluster result	unstructured	more interpretable and informative
Time complexity	$O(n \times k \times t)$	$O(n^3)$
Space complexity	$O(n(d + k))$	$O(n^2)$
With a large number of data points	compute faster	compute slower