

Hard Clustering

Alex(ander) Jung

Assistant Professor for Machine Learning

Department of Computer Science

Aalto University

Reading.

Sec. 8.1 of <https://mlbook.cs.aalto.fi>



<https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>

The screenshot shows the scikit-learn website's documentation page for KMeans. The top navigation bar includes links for 'Install', 'User Guide', 'API', 'Examples', 'Community', and 'More'. On the left, there are buttons for 'Prev', 'Up', and 'Next', and a section for 'scikit-learn 1.1.2' with a link to 'Other versions'. A yellow box contains the text 'Please cite us if you use the software.' The main content area has a light blue header with 'sklearn.cluster.KMeans'. Below it, the class definition is shown: `class sklearn.cluster.KMeans(n_clusters=8, *, init='k-means++', n_init=10, max_iter=300, tol=0.0001, verbose=0, random_state=None, copy_x=True, algorithm='lloyd')`. A '[source]' link is at the end of the code block. Below the code, the text 'K-Means clustering' is partially visible.

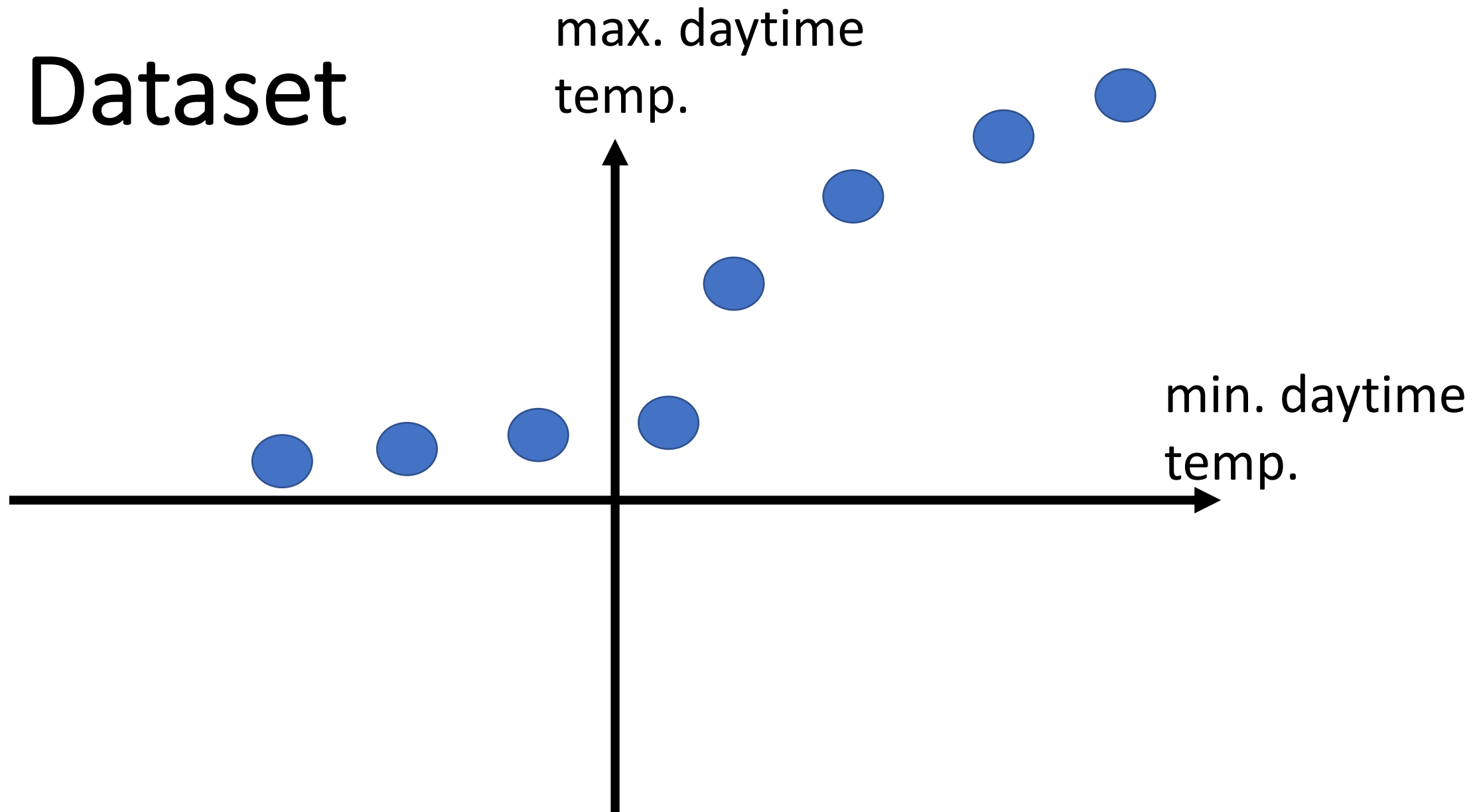
What I want to teach you today:

- basic idea of hard clustering
- k-means method for hard clustering
- optimization problem underlying k-means
- how to choose number of clusters

First things First

What are three main components of Machine Learning ?

A Dataset



What is a Cluster?

Noun [[edit](#)]

cluster (*plural* **clusters**)

1. A **group** or **bunch** of several discrete items that are **close** to each other. [[quotations ▼](#)]

*a **cluster** of islands*

*A **cluster** of flowers grew in the pot.*

*A **leukemia** cluster has developed in the town.*

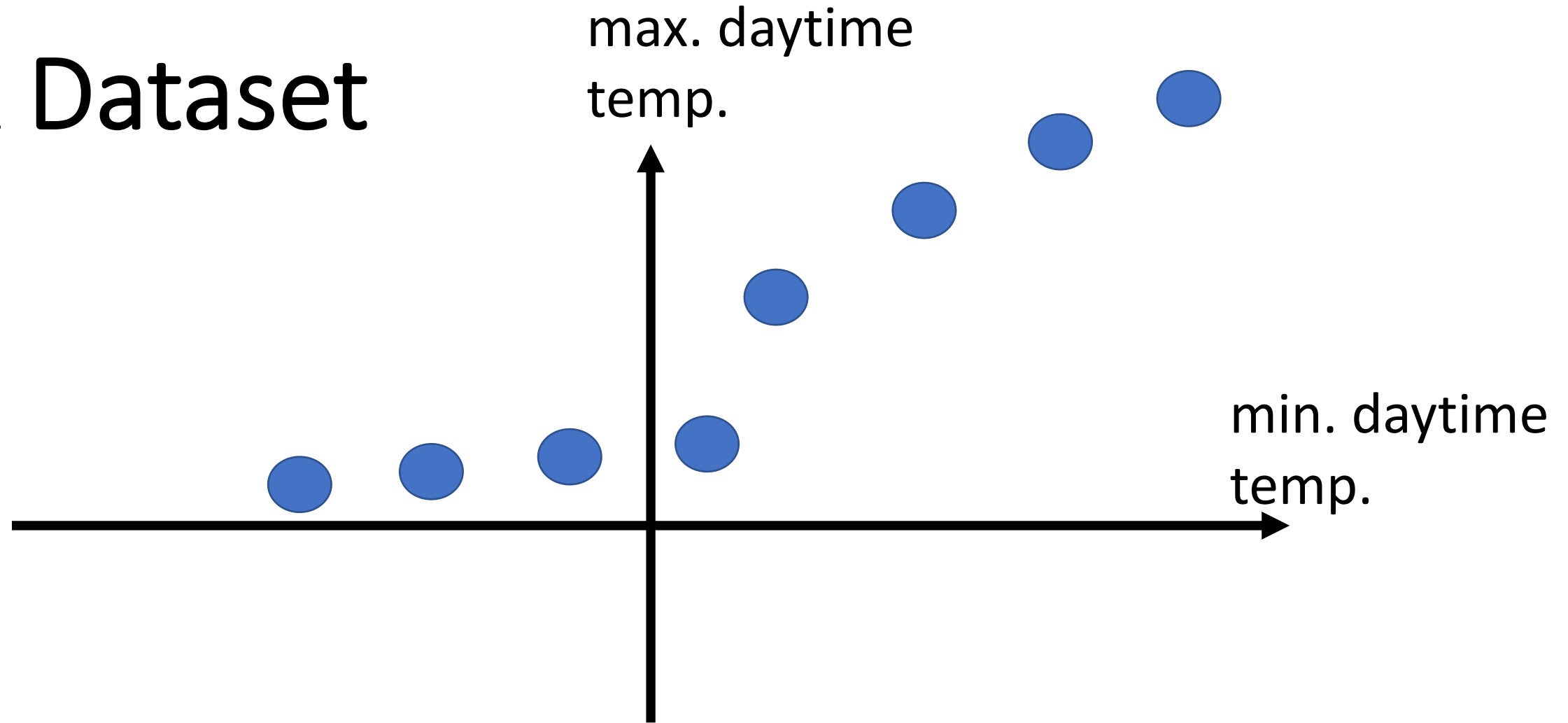
<https://en.wiktionary.org/wiki/cluster>

Informal Definition

a cluster corresponds to a subset of datapoints that are in some sense homogeneous or similar

plethora of different definitions for “homogeneous” and “similar”

A Dataset

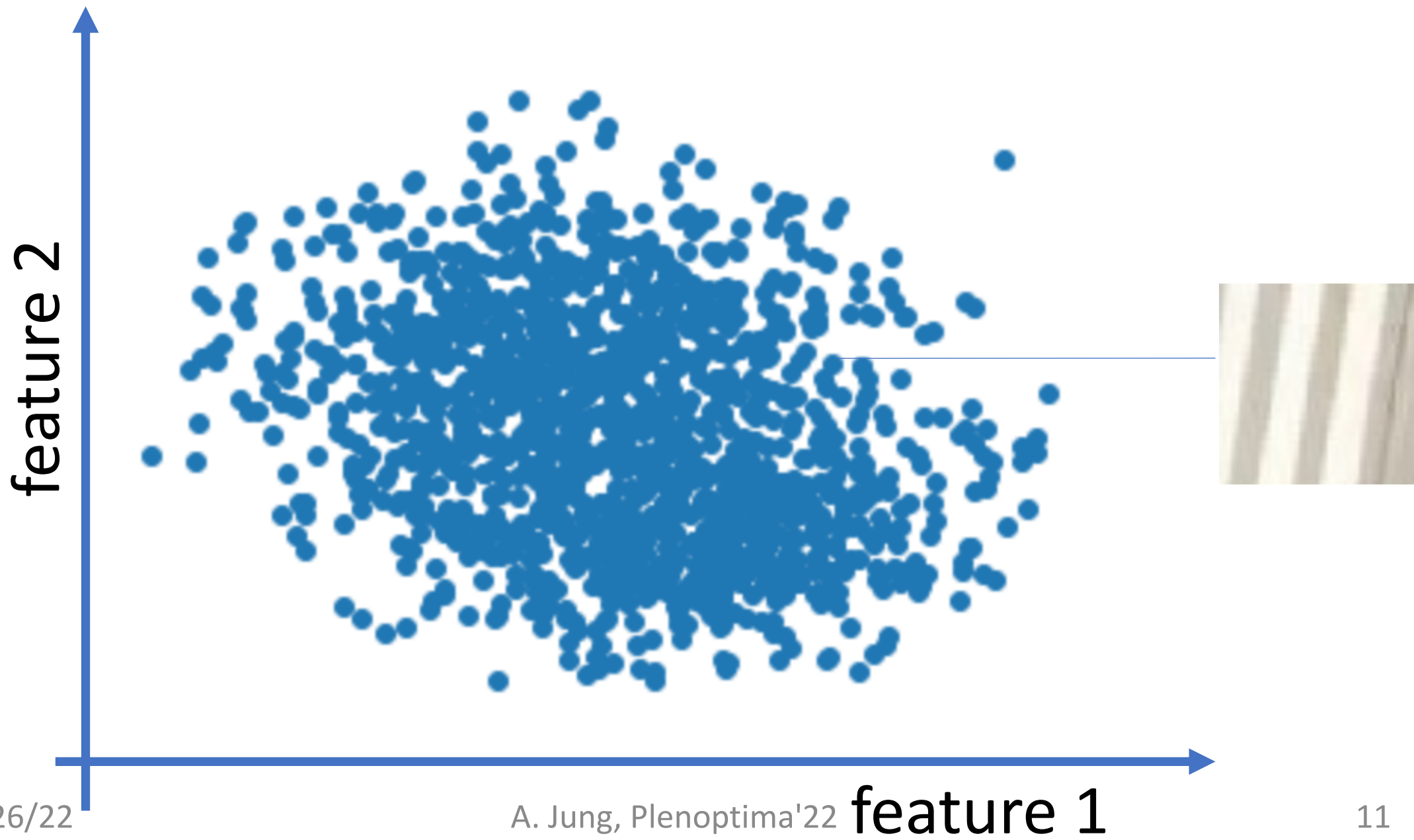


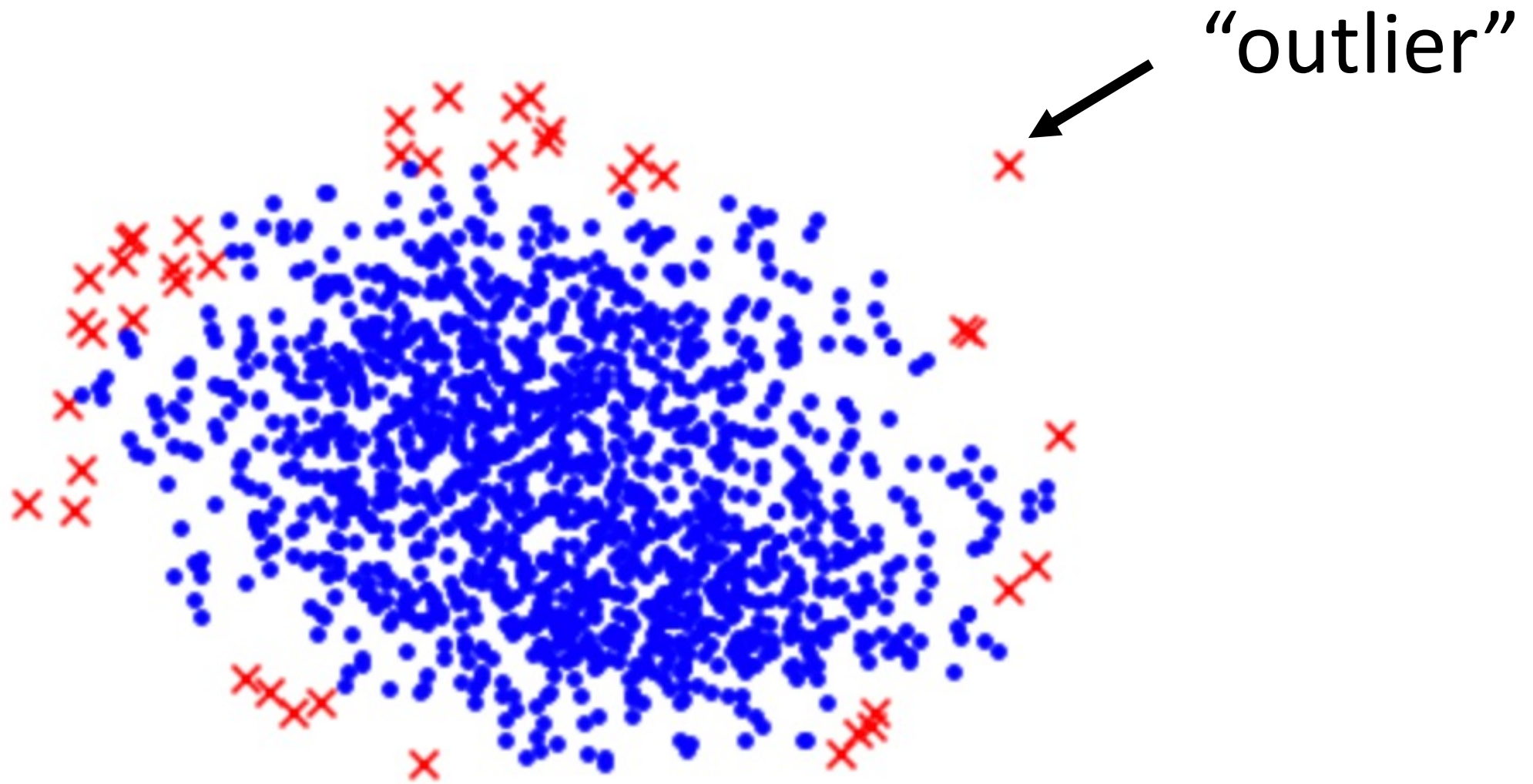
How many clusters do you see ?

Clustering for Outlier/Anomaly Detection

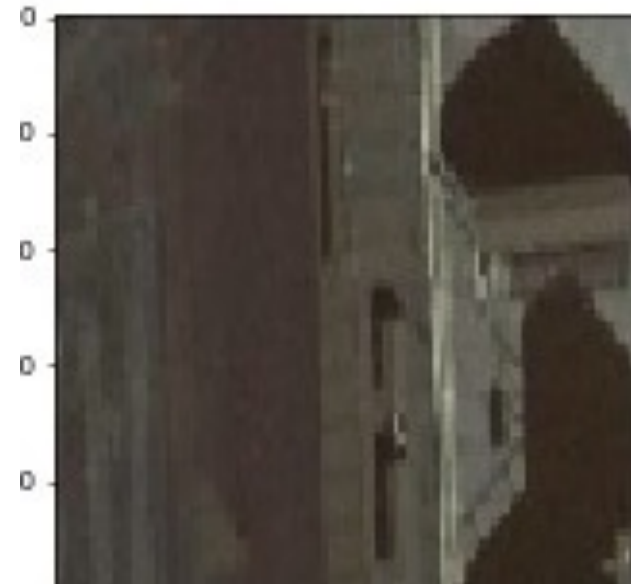
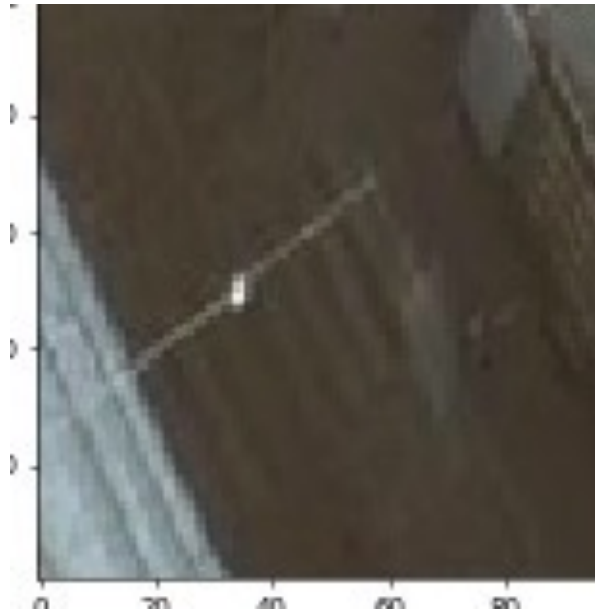
Dataset = “Bunch of Images”





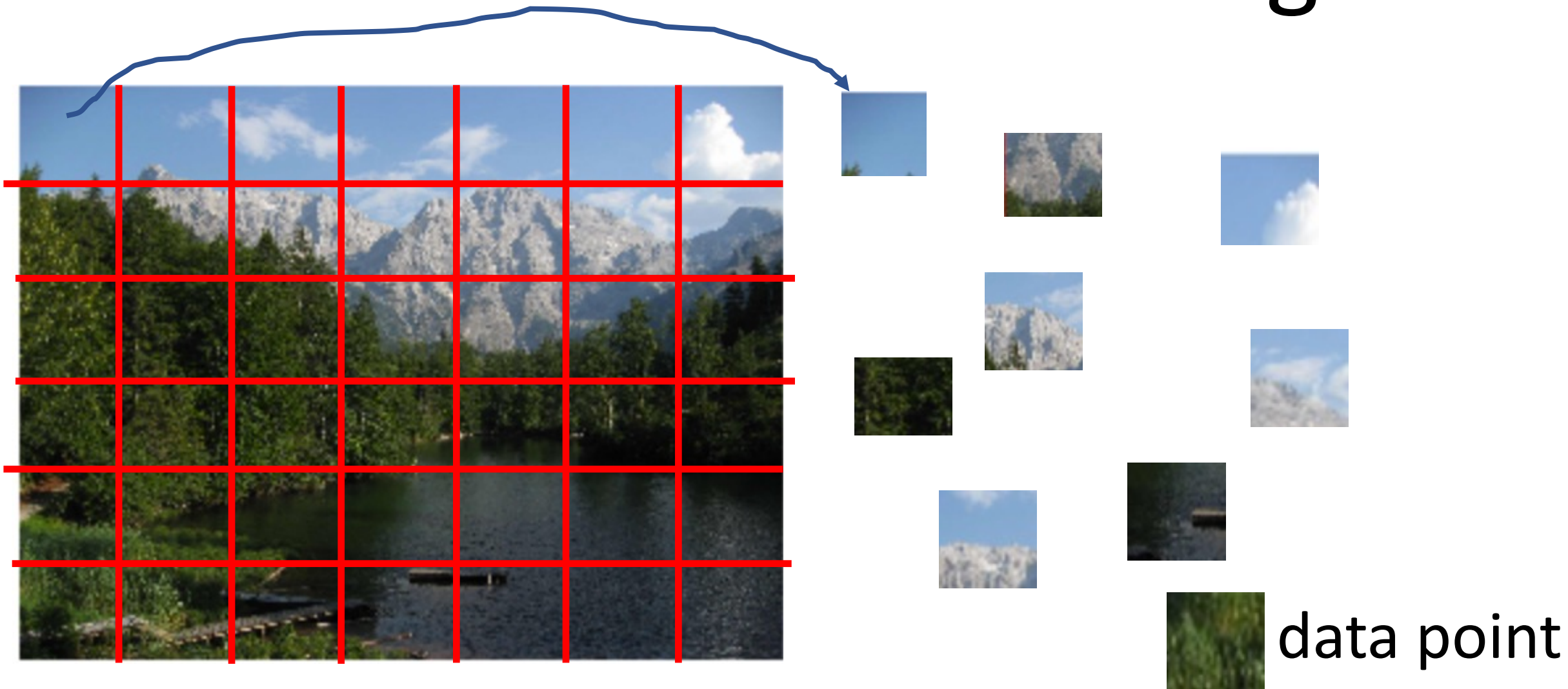


some outliers

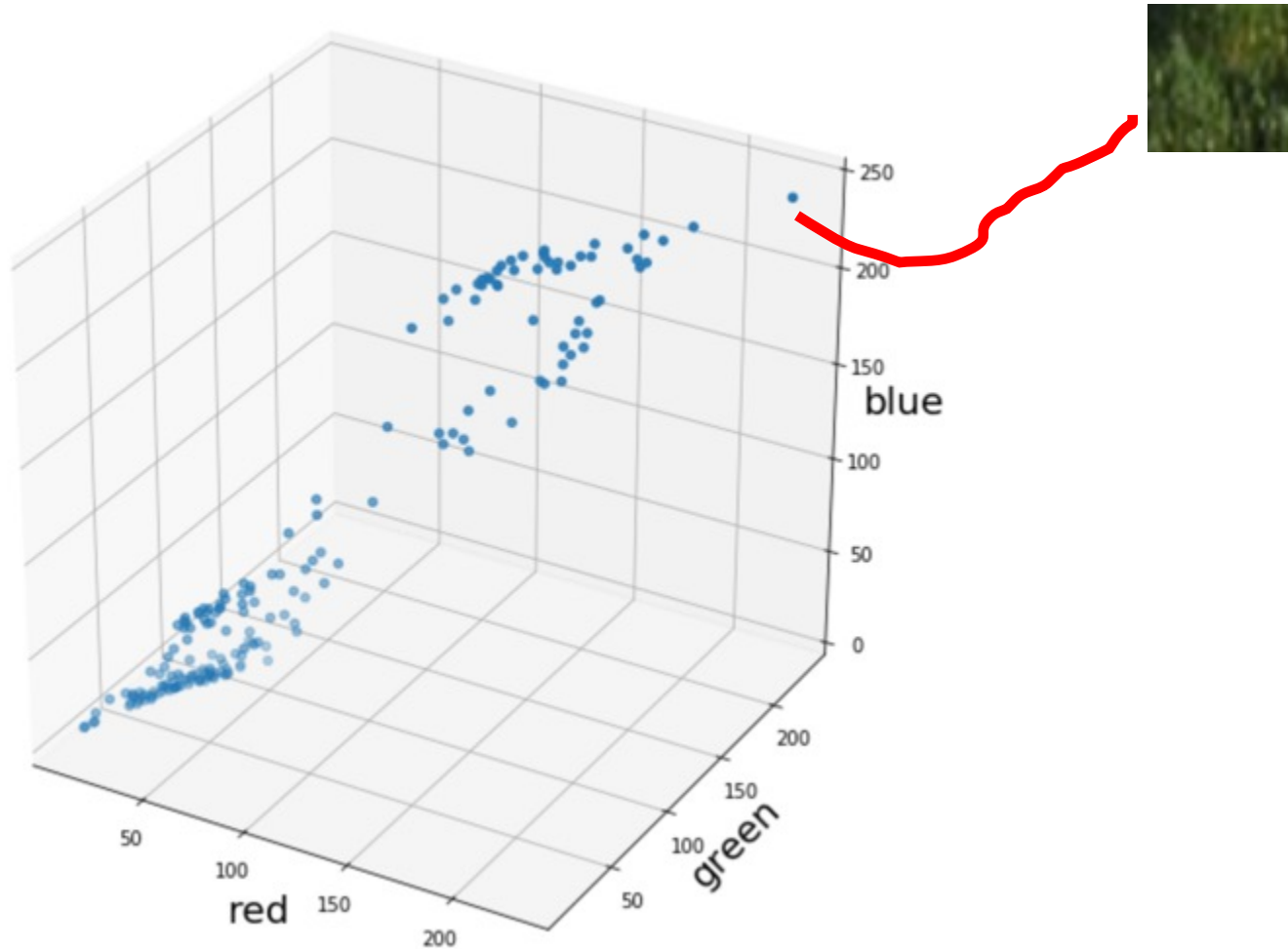


Clustering for Image Segmentation

Dataset = Patches of Image

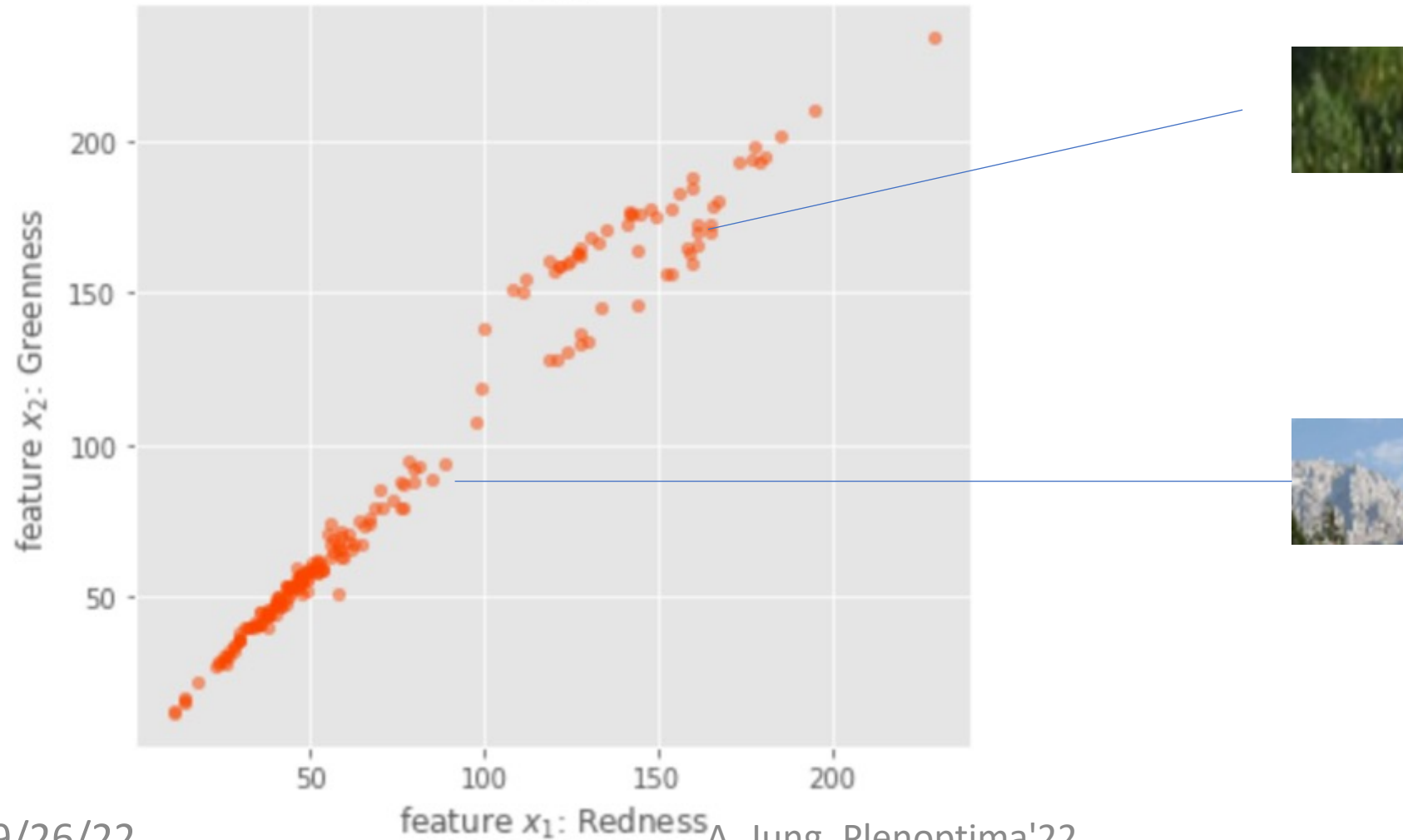


Using Three Features



three features:
average red, green and blue
component

Using Two Features (Red+Green)

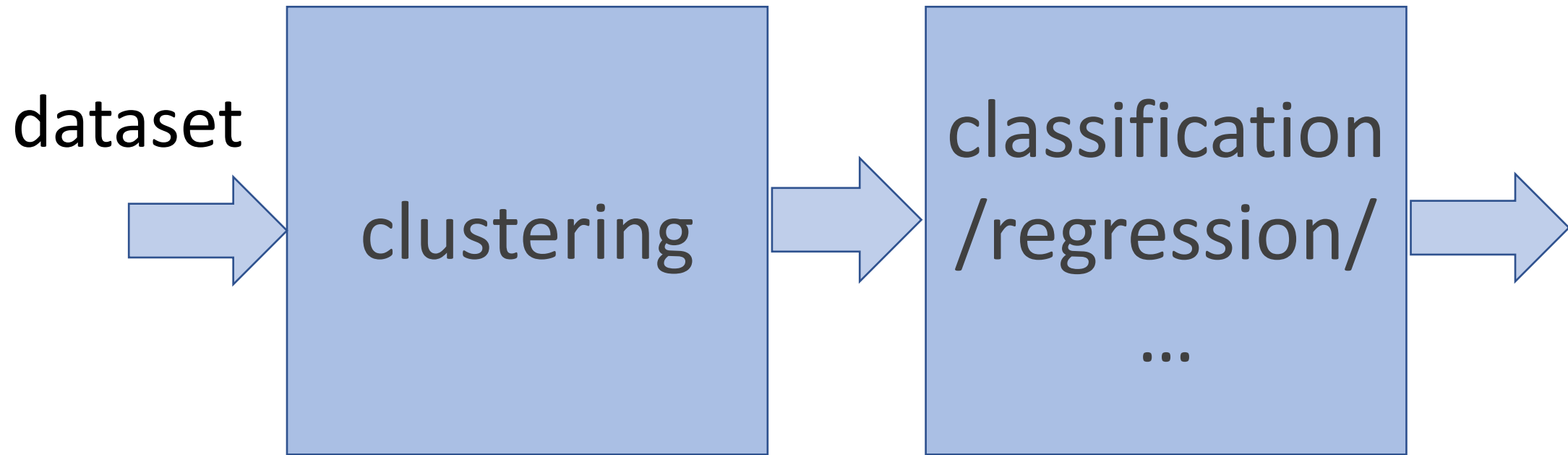


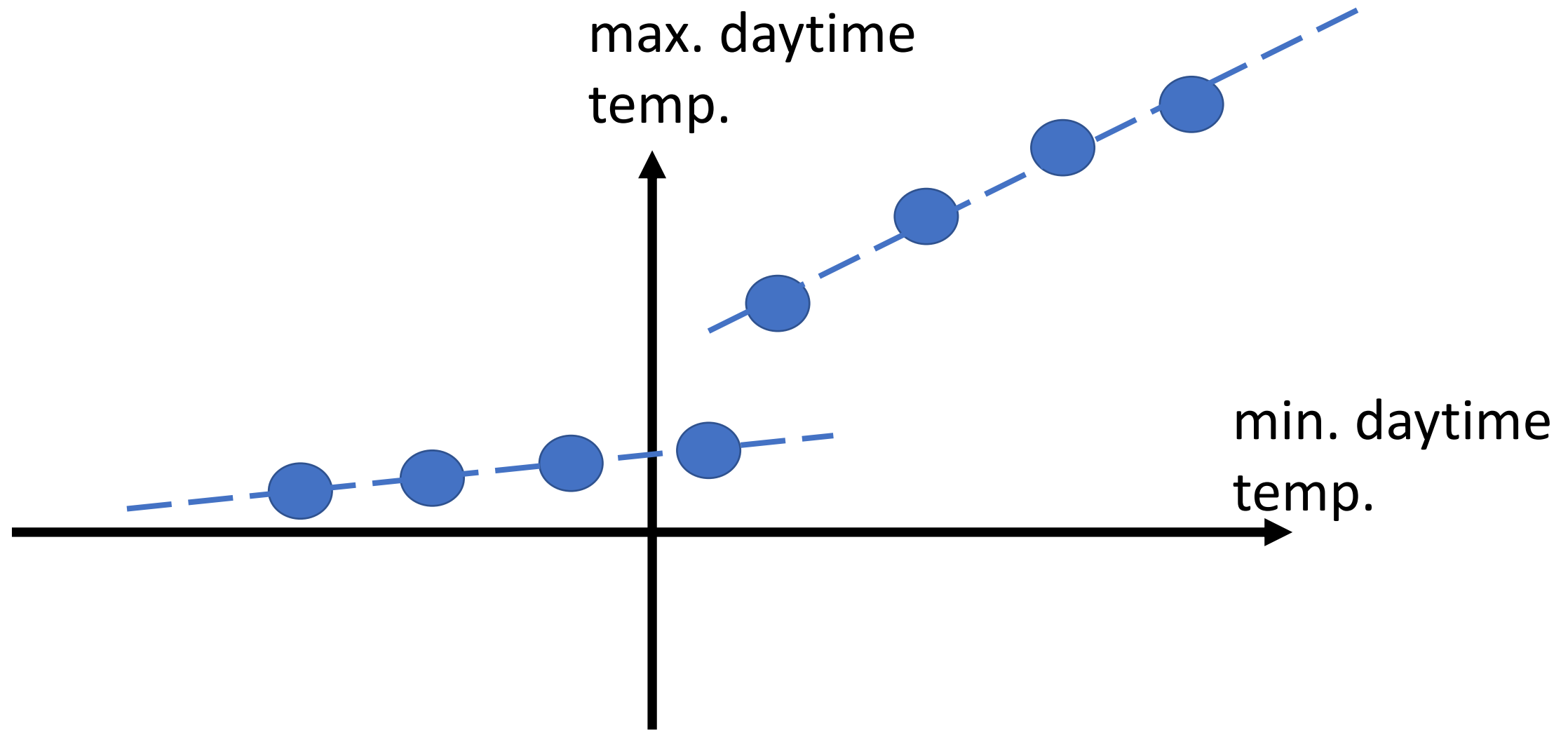
Use Clustering For Image Segmentation



Pre-Processing

Clustering as Pre-Processing





first partition into two clusters. then apply linear regression separately to each cluster

Hard Clustering

- datapoints $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})$
- i-th datapoint characterized by n features

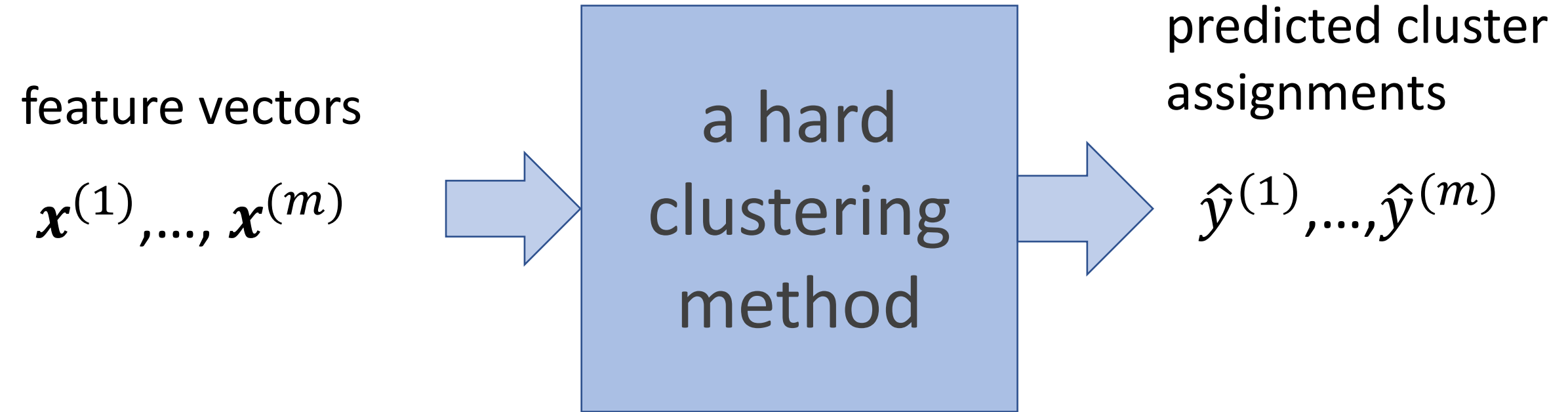
$$\mathbf{x}^{(i)} = (x_1^{(i)}, \dots, x_n^{(i)})$$

- i-th datapoint belongs to one of k clusters
- cluster index of i-th datapoint is $y^{(i)} \in \{1, \dots, k\}$

Hard Clustering Methods

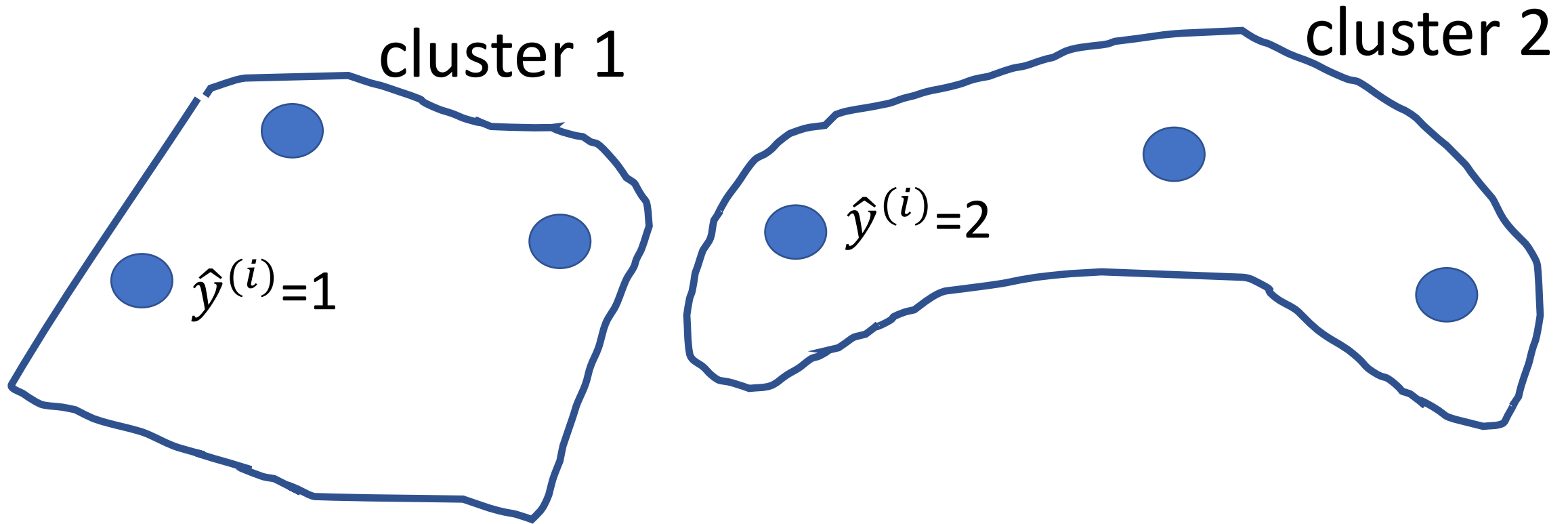
- datapoints $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})$
- cluster index of i -th datapoint is $y^{(i)} \in \{1, \dots, k\}$
- hard clustering methods compute predicted cluster indices $\hat{y}^{(i)}$ based solely on features
- does not require true cluster index $y^{(i)}$ of any datapoint

Hard Clustering Methods



Hard Clustering with k-Means

Representing a Cluster by a Mean

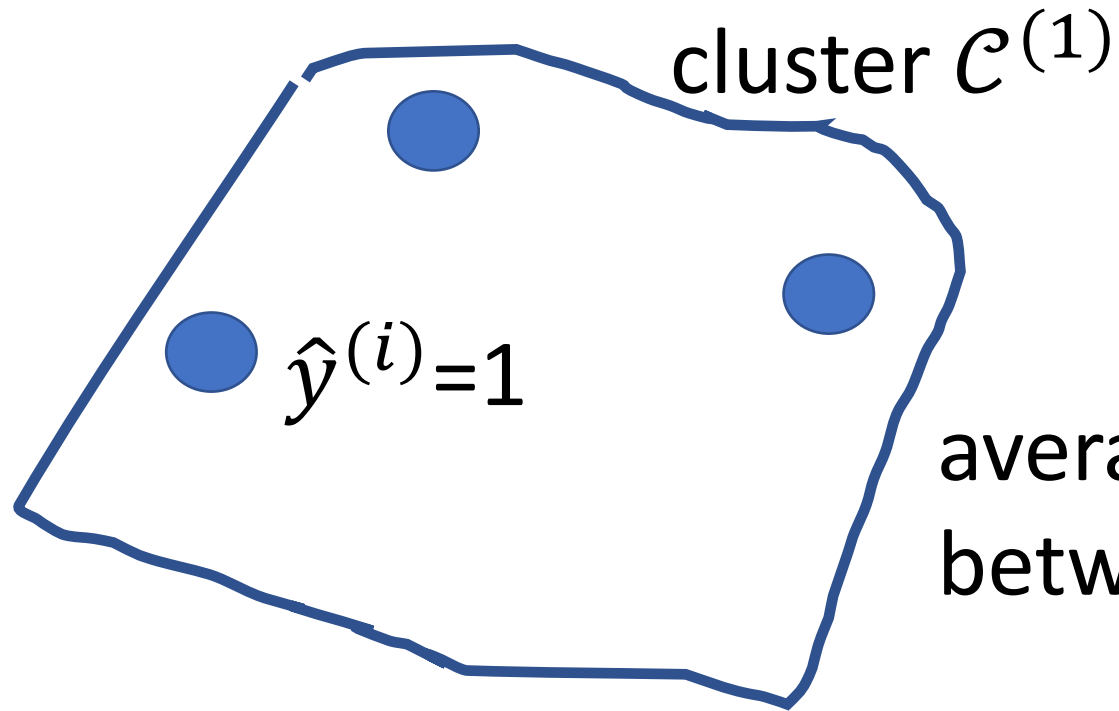


“cluster mean” 1



cluster mean 2

Cluster Spread



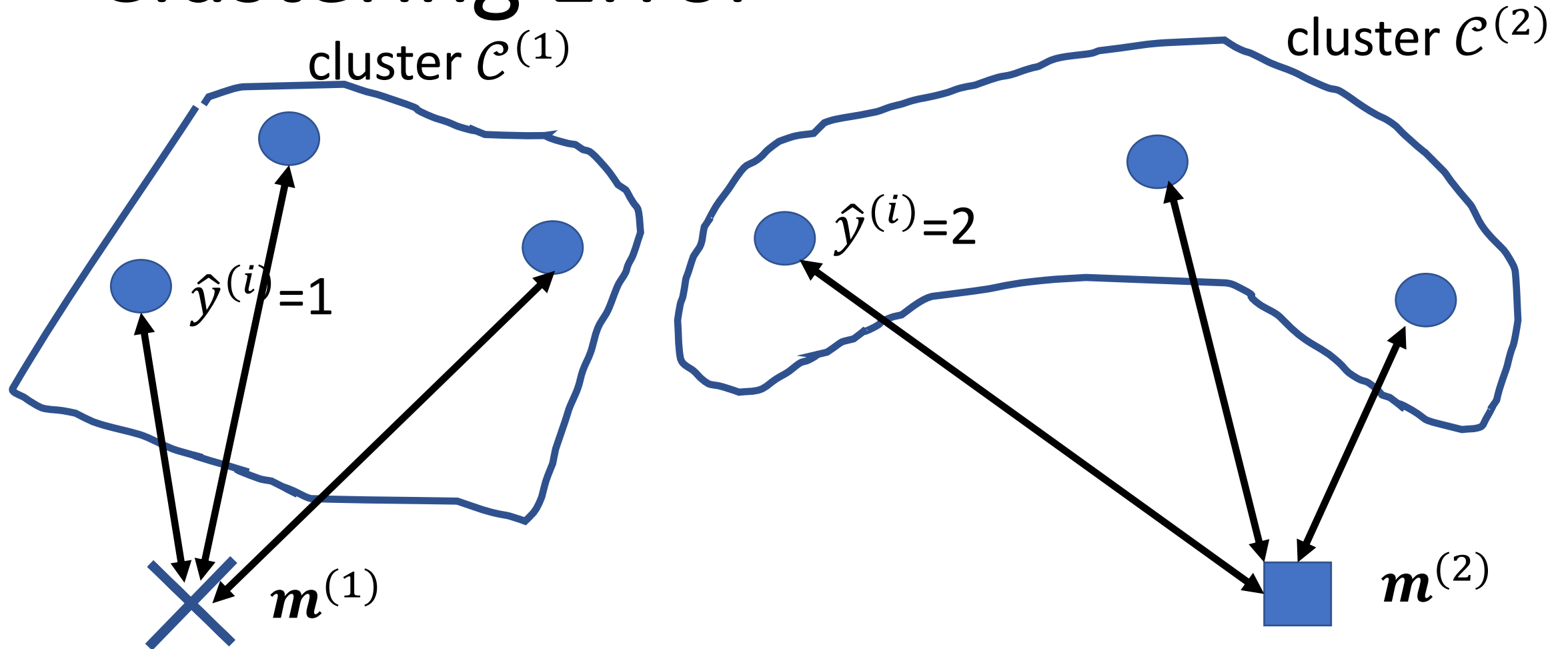
average squared Euclidean distance
between points and mean of cluster

$\times \mathbf{m}^{(1)}$

$$(1/|\mathcal{C}^{(1)}|) \sum_{i \in \mathcal{C}^{(1)}} \|\mathbf{m}^{(1)} - \mathbf{x}^{(i)}\|^2$$

mean for $\mathcal{C}^{(1)}$

Clustering Error



$$(1/m) \sum_{c=1}^2 \sum_{i \in \mathcal{C}^{(c)}} \|m^{(c)} - x^{(i)}\|^2$$

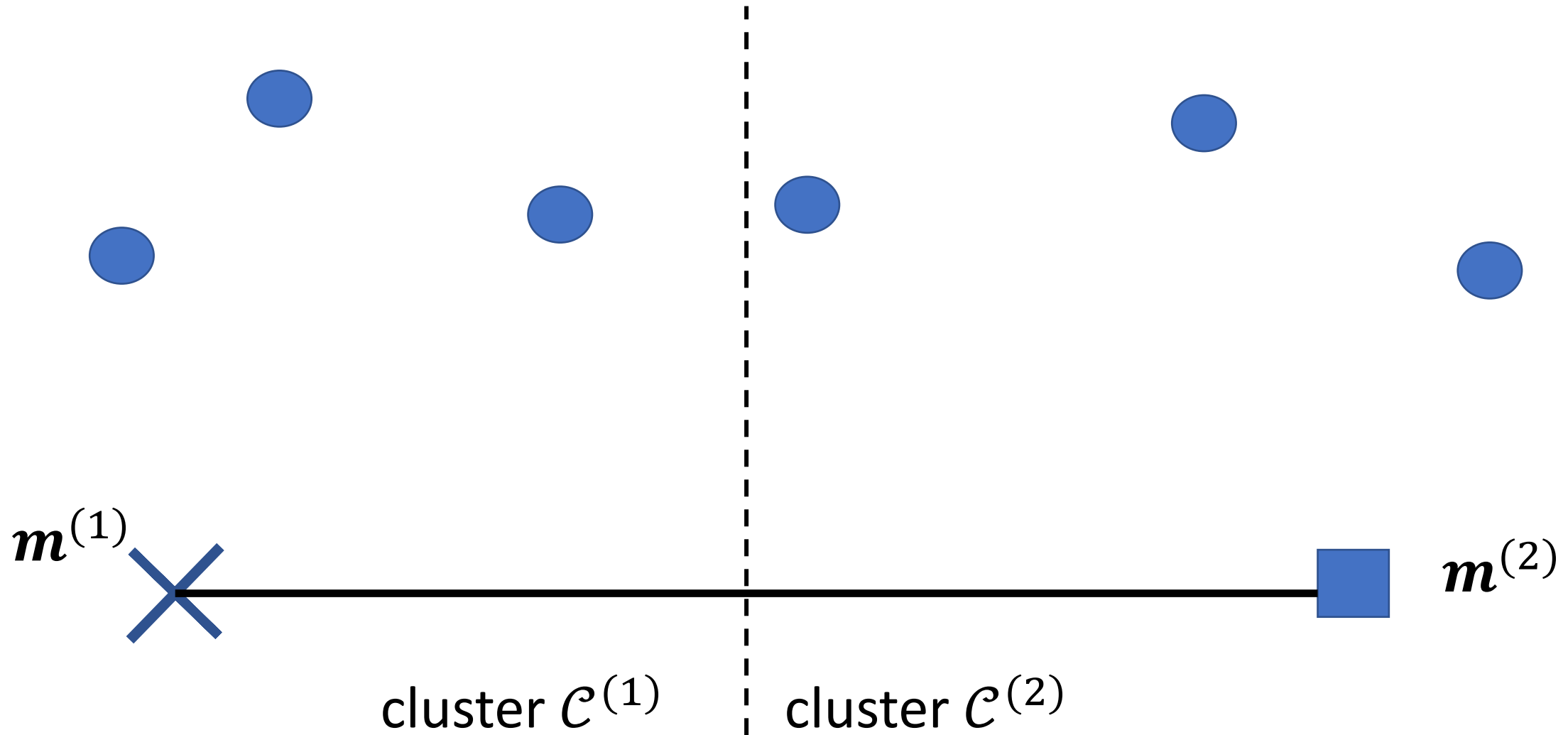
Update Cluster Assignments

for given cluster means, clustering error is minimized by assigning i-th datapoint to cluster with nearest cluster mean

$$\hat{y}^{(i)} := c$$

with $\|\mathbf{m}^{(c)} - \mathbf{x}^{(i)}\|^2 = \min_{c'=1,\dots,k} \|\mathbf{m}^{(c')} - \mathbf{x}^{(i)}\|^2$

Update Cluster Assignment



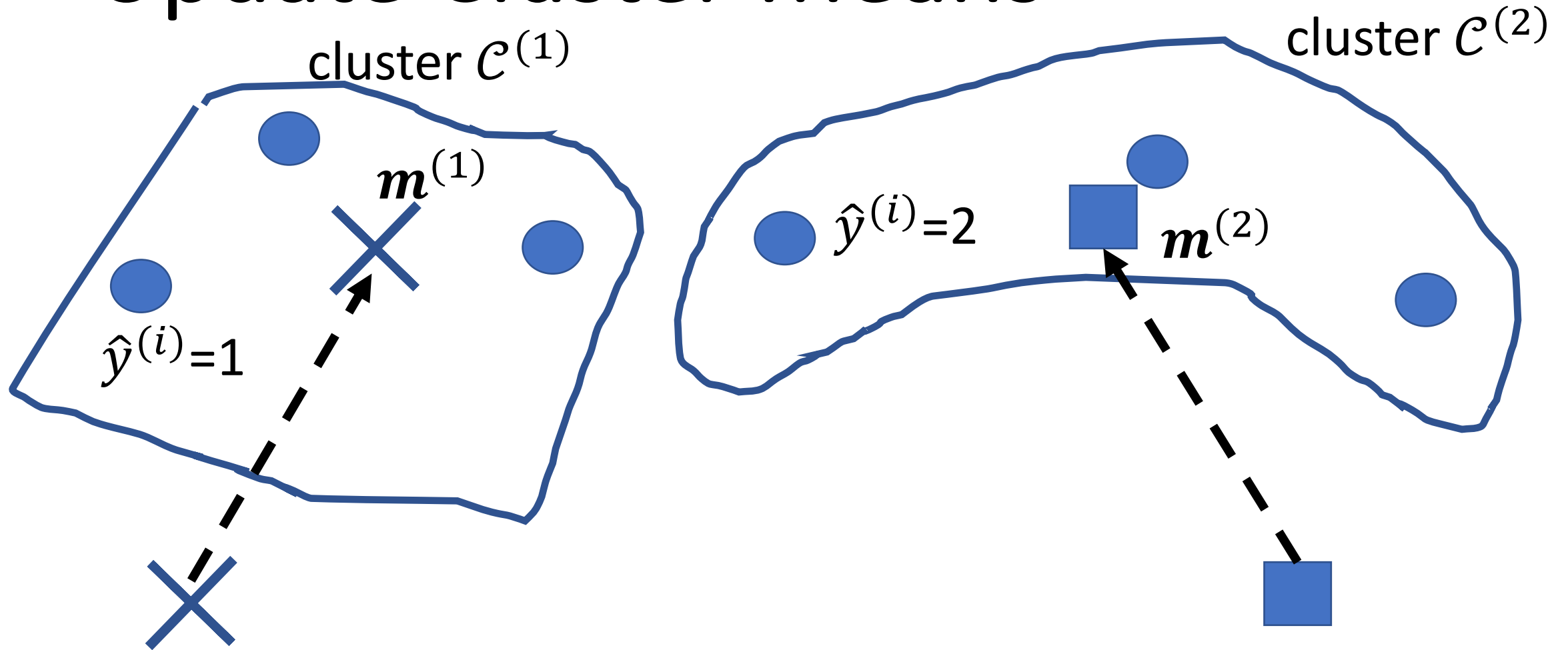
Update Cluster Means

for given cluster assignments, clustering error is minimized by representing c-th cluster by the cluster mean

$$m^{(c)} := \frac{1}{|\mathcal{C}^{(c)}|} \sum_{i \in \mathcal{C}^{(c)}} \mathbf{x}^{(i)}$$

with cluster $\mathcal{C}^{(c)} = \{i: \hat{y}^{(i)} = c\}$

Update Cluster Means



Minimizing the Clustering Error

clustering error

$$\mathcal{E}(\{\mathbf{m}^{(c)}\}, \{\hat{\mathbf{y}}^{(i)}\}) := \frac{1}{m} \sum_{i=1}^m \left\| \mathbf{m}^{(\hat{\mathbf{y}}^{(i)})} - \mathbf{x}^{(i)} \right\|^2$$

simultaneously finding cluster means $\mathbf{m}^{(c)}$
and assignments $\hat{\mathbf{y}}^{(i)}$ that minimize clustering
error is difficult (“NP-hard”)

https://cseweb.ucsd.edu/~avattani/papers/kmeans_hardness.pdf

Alternating Minimization

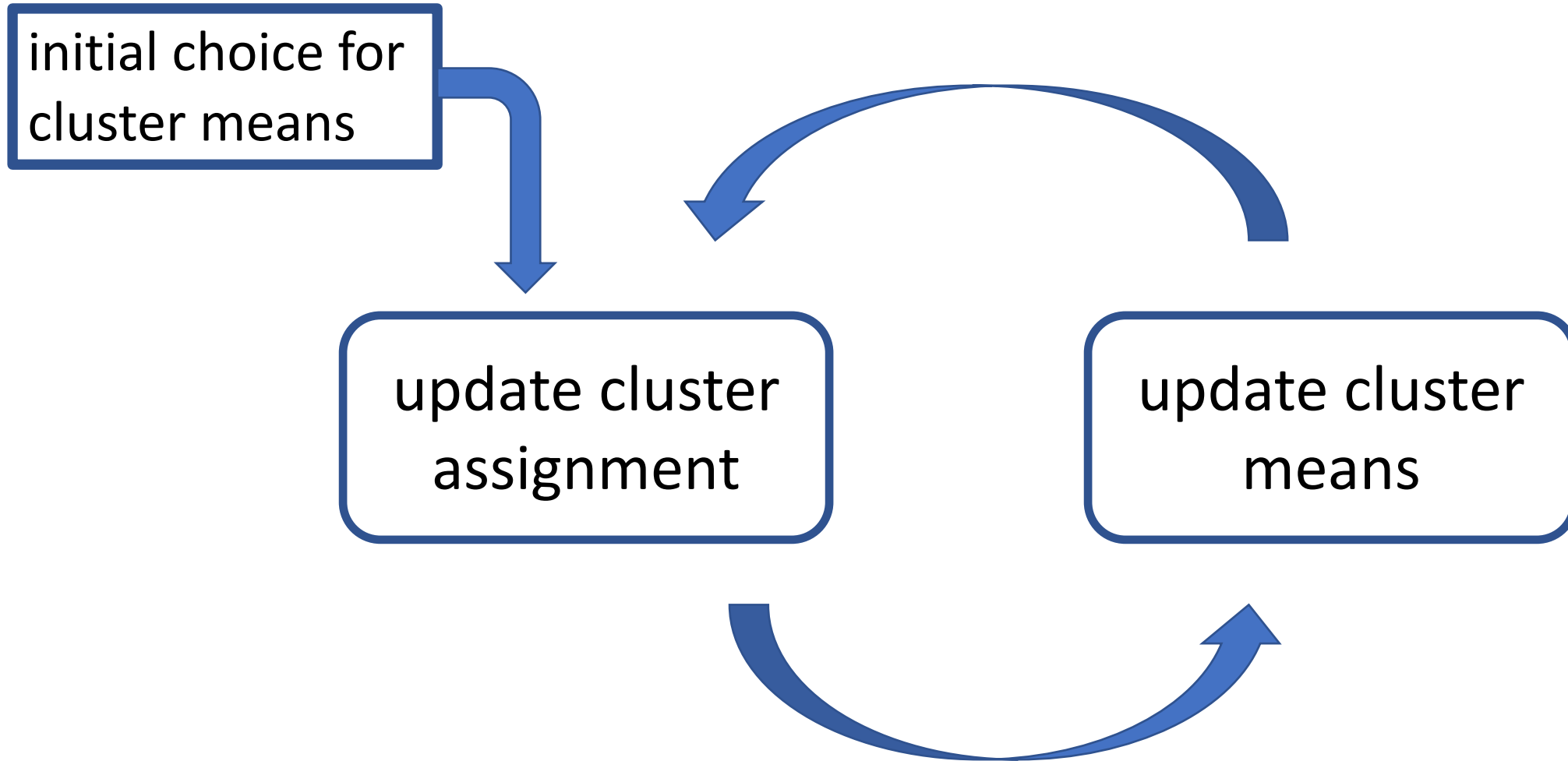
clustering error

$$\varepsilon(\{\mathbf{m}^{(c)}\}, \{\hat{\mathbf{y}}^{(i)}\}) := \frac{1}{m} \sum_{i=1}^m \left\| \mathbf{m}^{(\hat{\mathbf{y}}^{(i)})} - \mathbf{x}^{(i)} \right\|^2$$

for **given assignments** $\hat{\mathbf{y}}^{(i)}$, finding cluster means $\mathbf{m}^{(c)}$
that **minimize clustering error is easy**

for **given cluster means** $\mathbf{m}^{(c)}$, finding assignments $\hat{\mathbf{y}}^{(i)}$
that **minimize clustering error is easy**

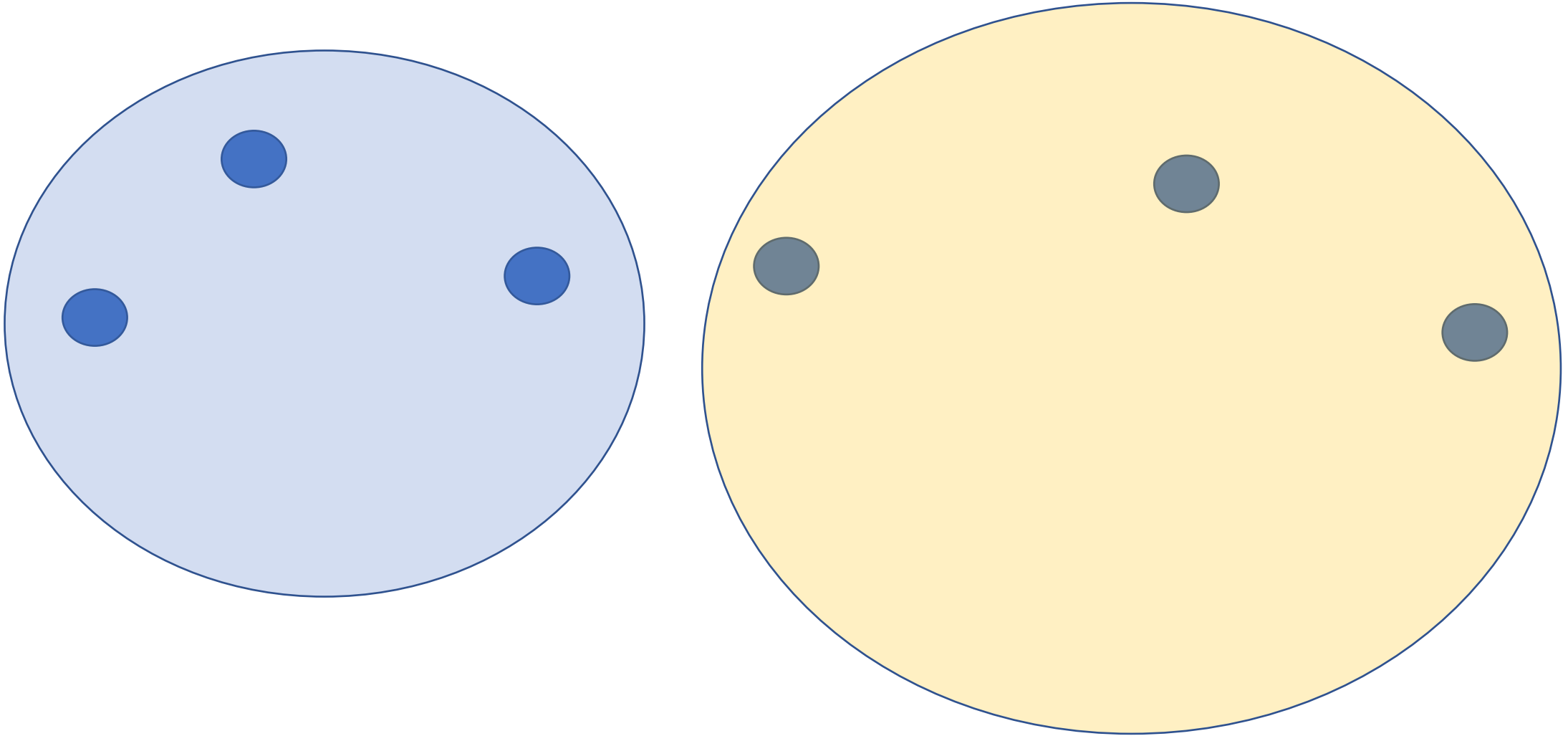
“k-Means”



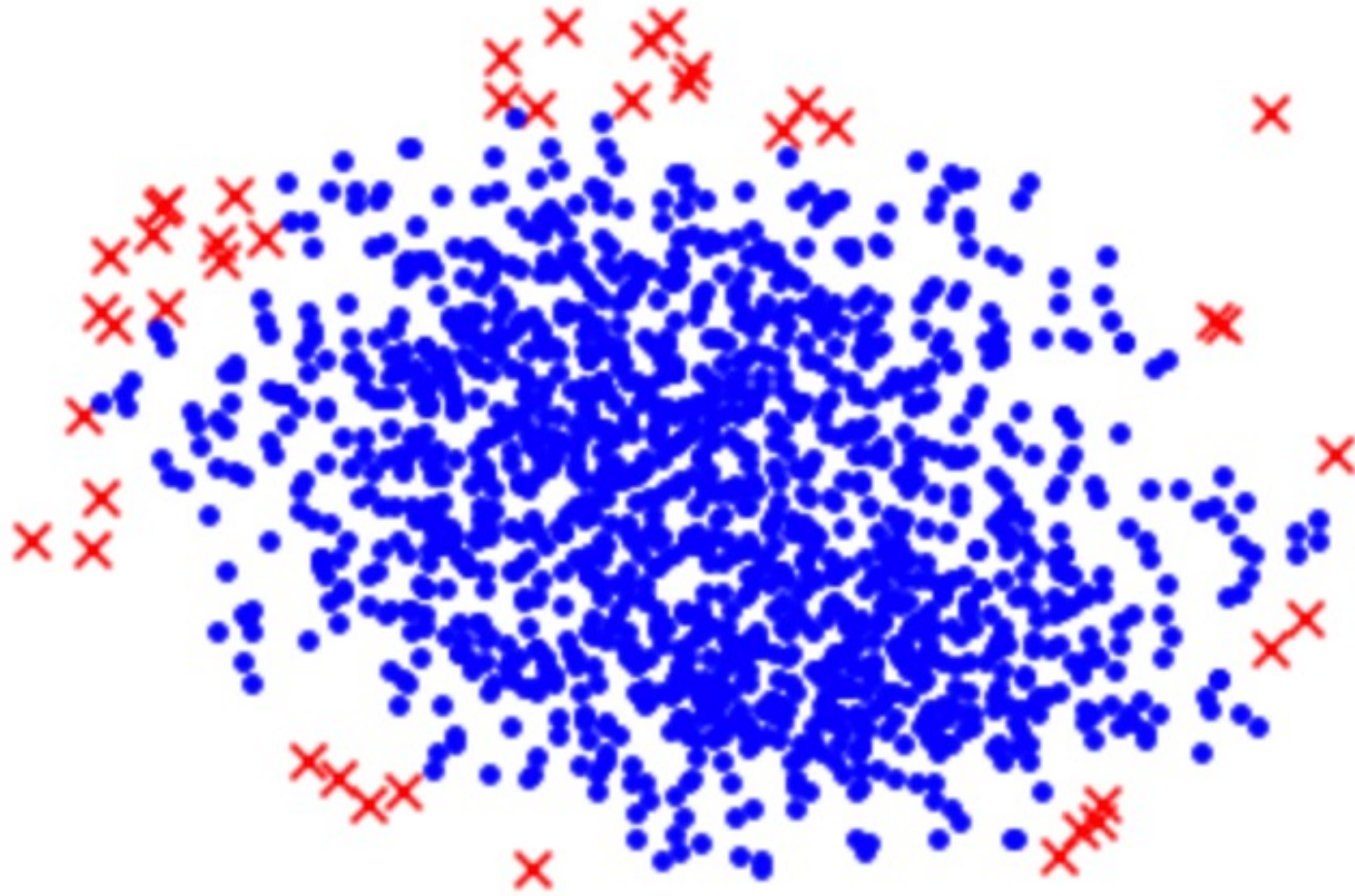
“k-Means” (Algorithm 8 mlbook.cs.aalto.fi)

- **Input:** number k of clusters, initial cluster means
- Step 1: update cluster assignments
- Step 2: update cluster means
- Go to Step 1 unless “Finished”
- **Output:** final cluster means

Cluster Shape of k-means Result



Clustering by k-means?



k-Means never increases Clustering Error !

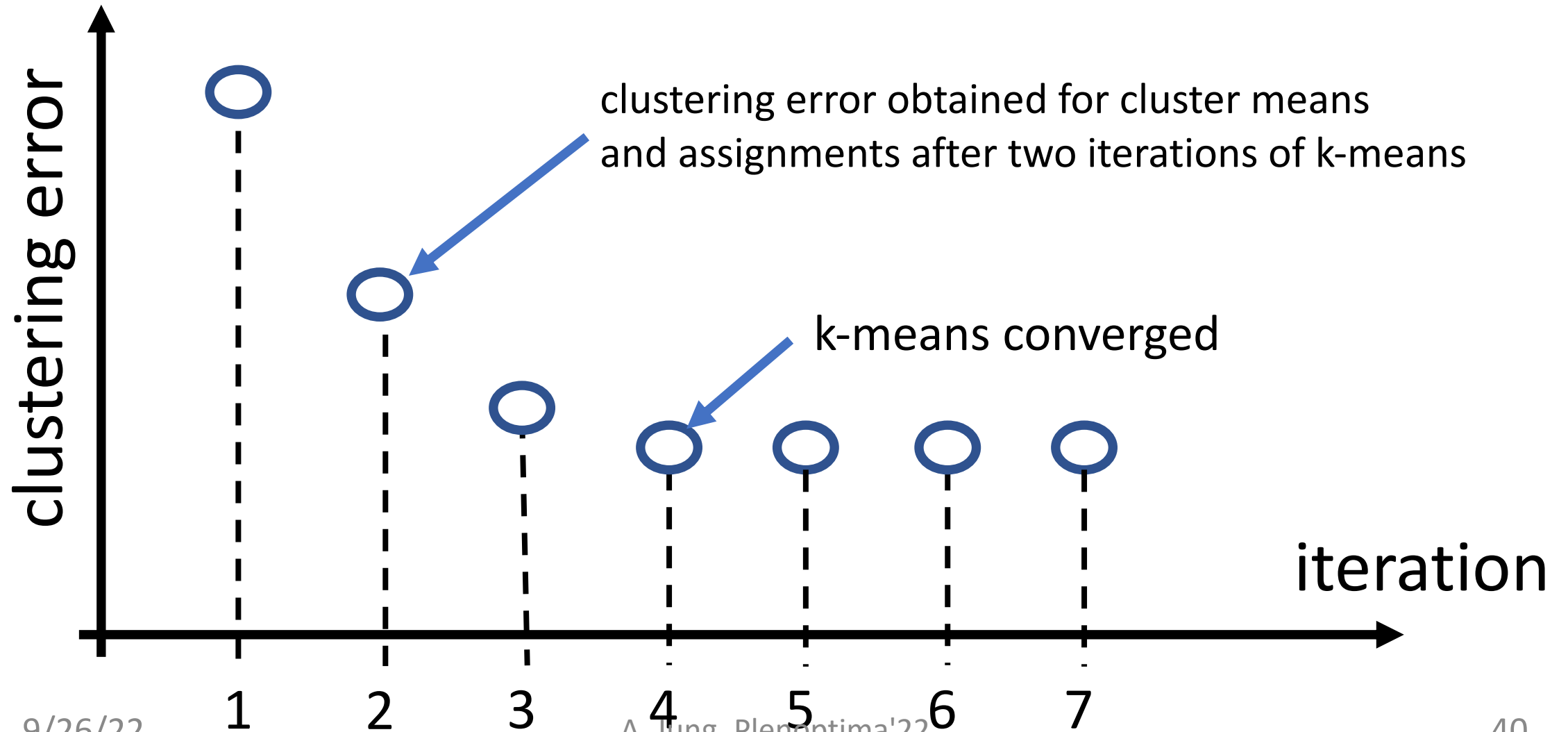
consider cluster means $m^{(c)}$ and assignments $\hat{y}^{(i)}$

run one iteration of k-means

results in new cluster means $\tilde{m}^{(c)}$ and assignments $\tilde{y}^{(i)}$

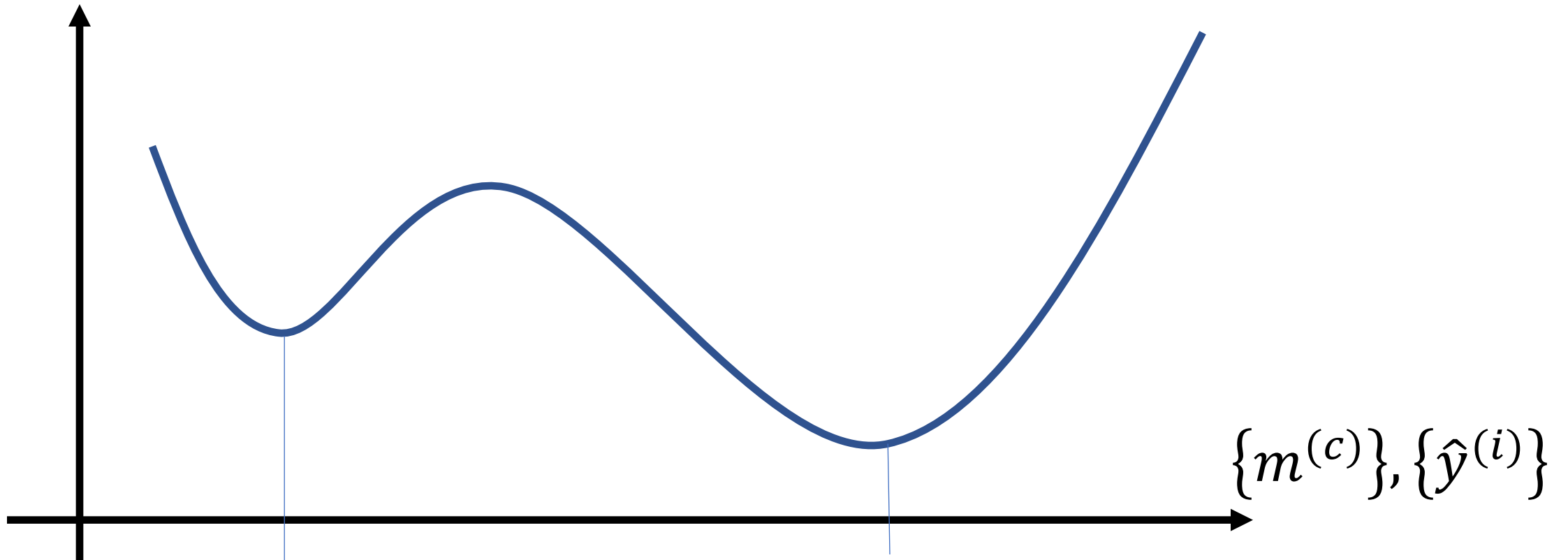
$$\mathcal{E}(\{\tilde{m}^{(c)}\}, \{\tilde{y}^{(i)}\}) \leq \mathcal{E}(\{m^{(c)}\}, \{\hat{y}^{(i)}\})$$

k-Means as Iterative Optimization Method



Non-Convexity of Clustering Error

$$\mathcal{E}(\{m^{(c)}\}, \{\hat{y}^{(i)}\})$$



local optimum

optimal clustering

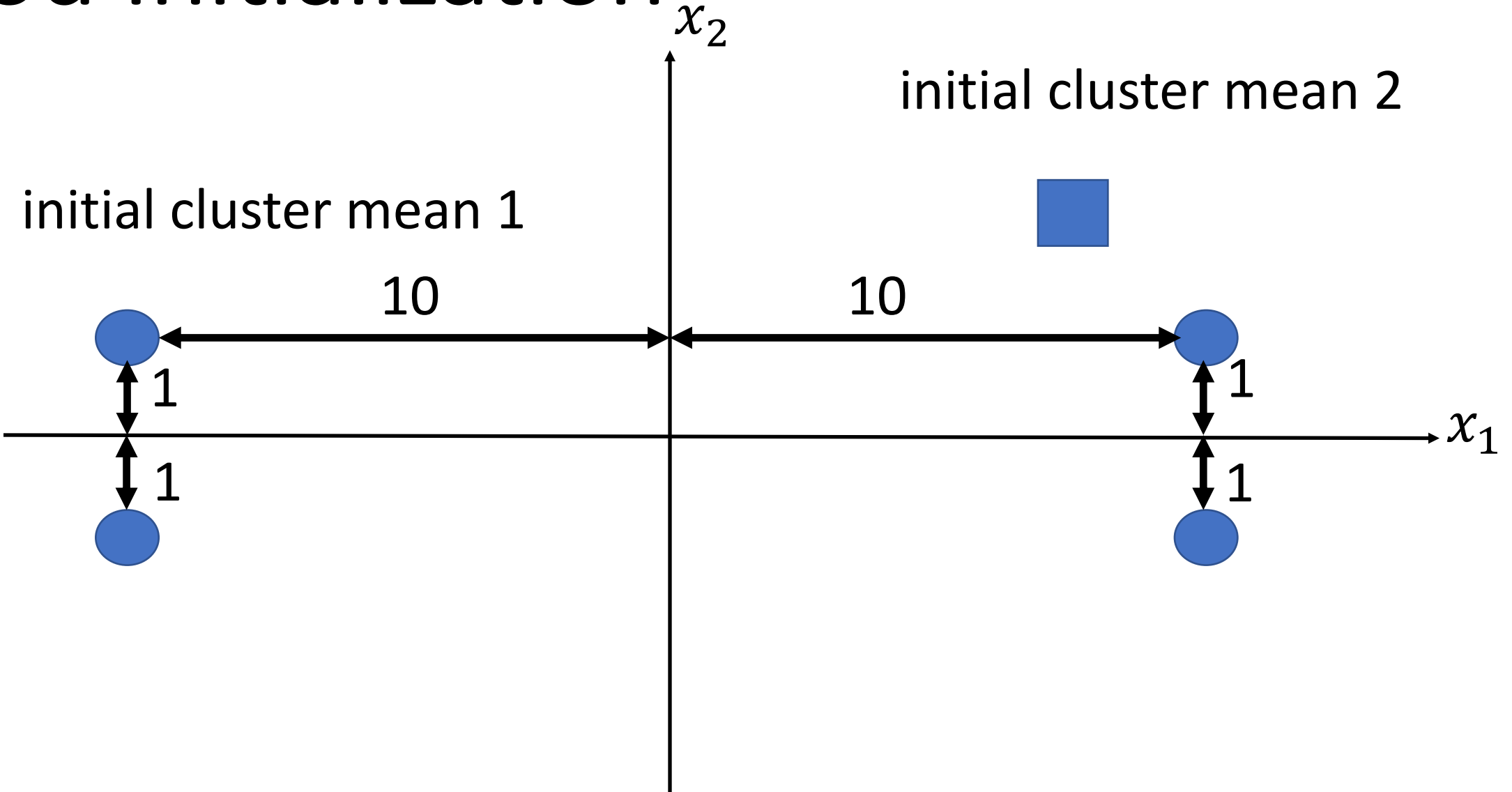
Initialization is Crucial

- k-means requires initial cluster means as inputs
- k-means result depends crucially on init. means
- repeat k-means several times with different init.

Good Initialization



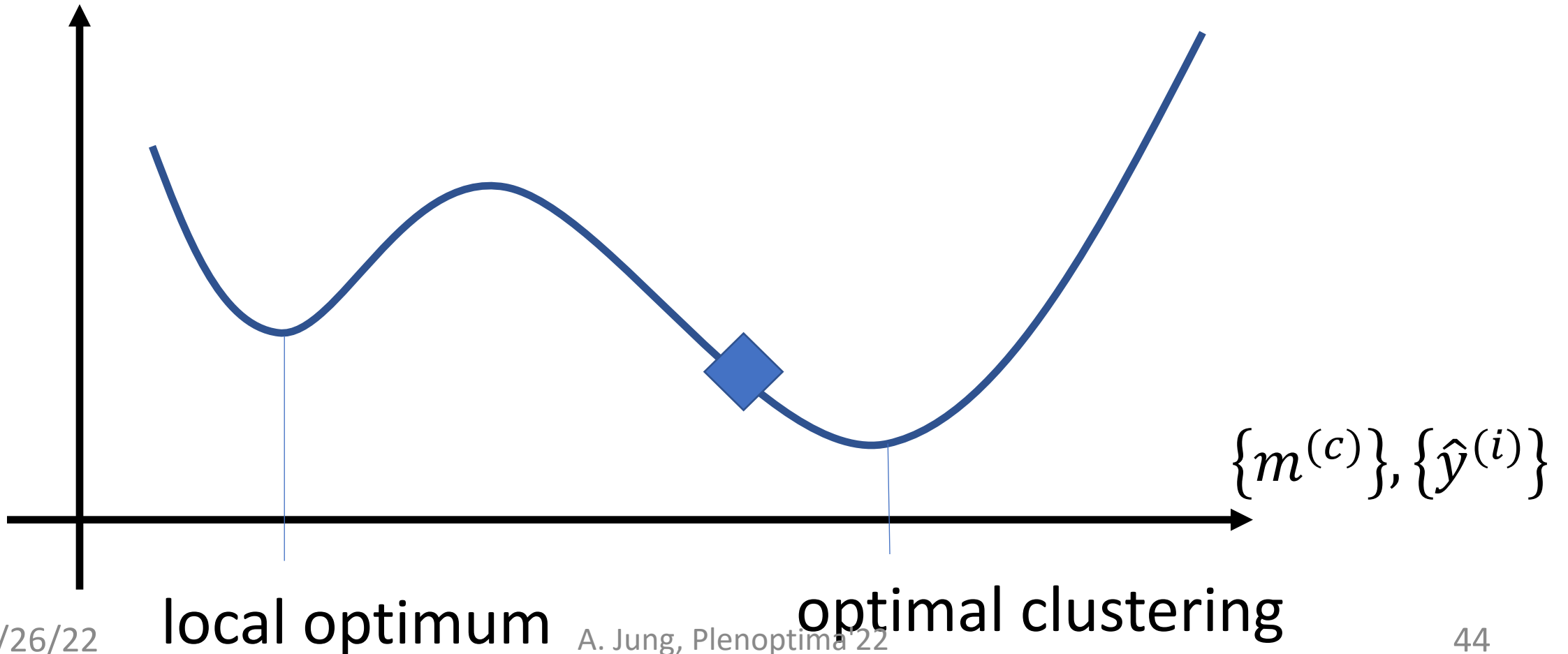
initial cluster mean 1



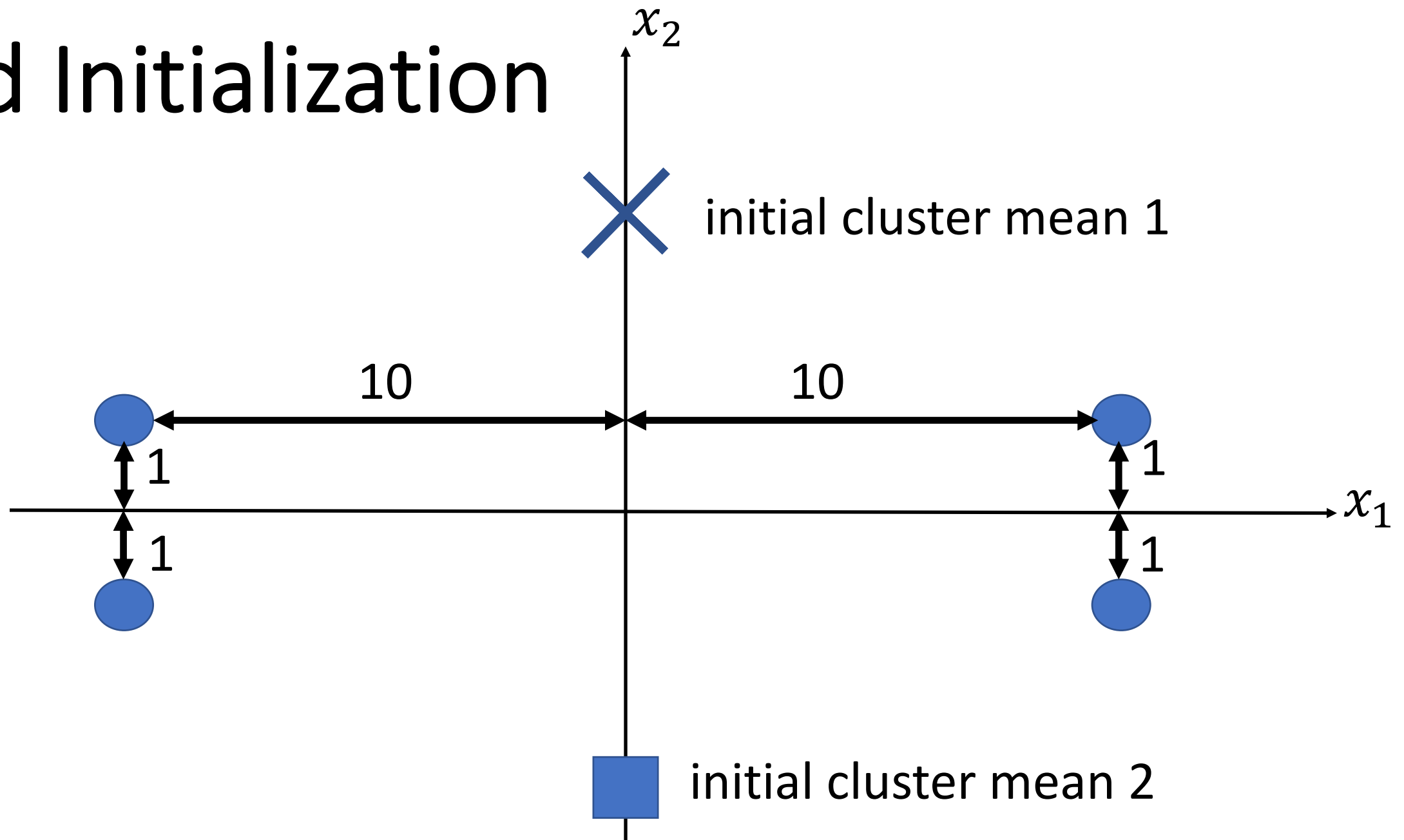
Good Initialization

◆ initial cluster means

$$\mathcal{E}(\{m^{(c)}\}, \{\hat{y}^{(i)}\})$$



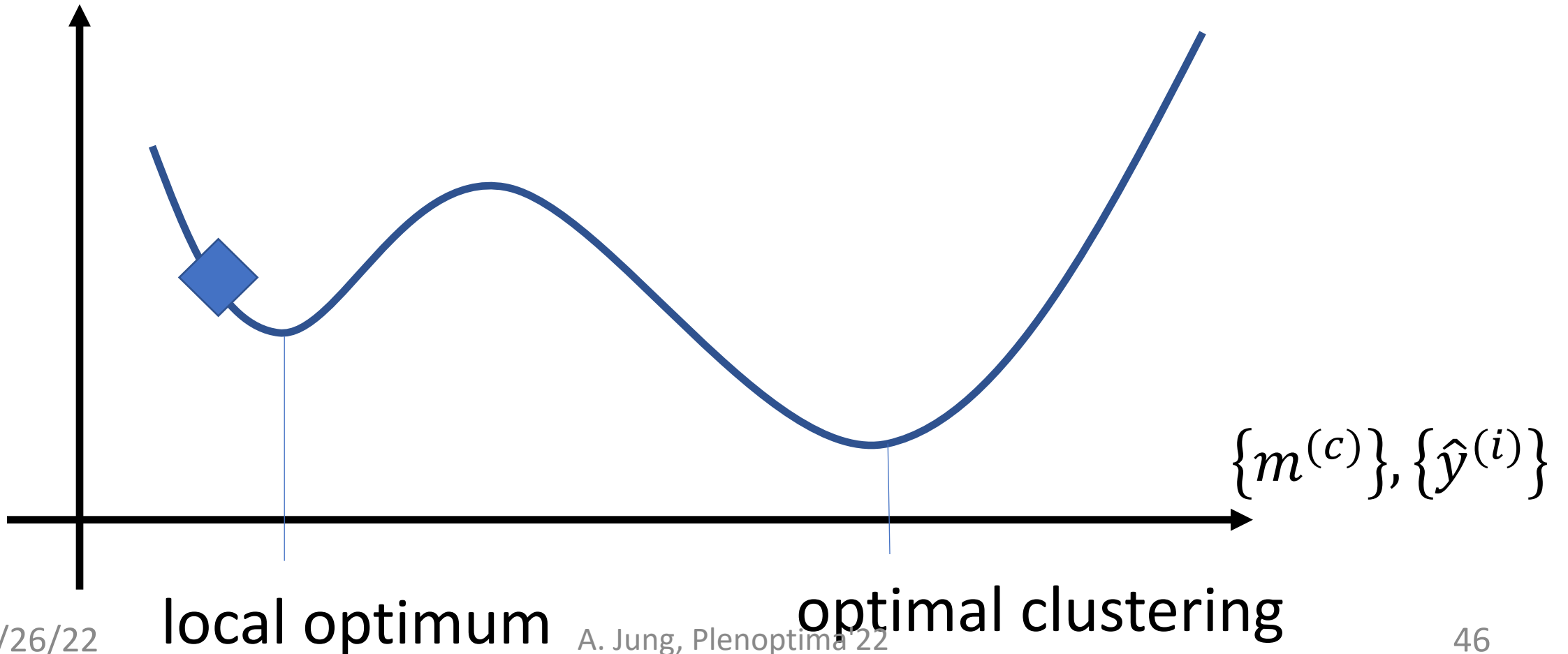
Bad Initialization



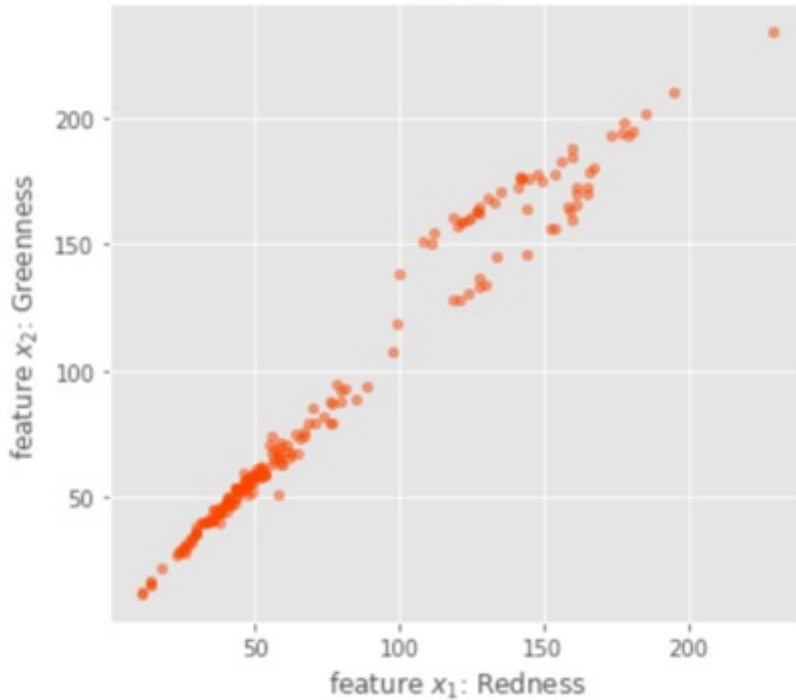
Bad Initialization

◆ initial cluster means

$$\mathcal{E}(\{m^{(c)}\}, \{\hat{y}^{(i)}\})$$



How to choose number k of clusters?



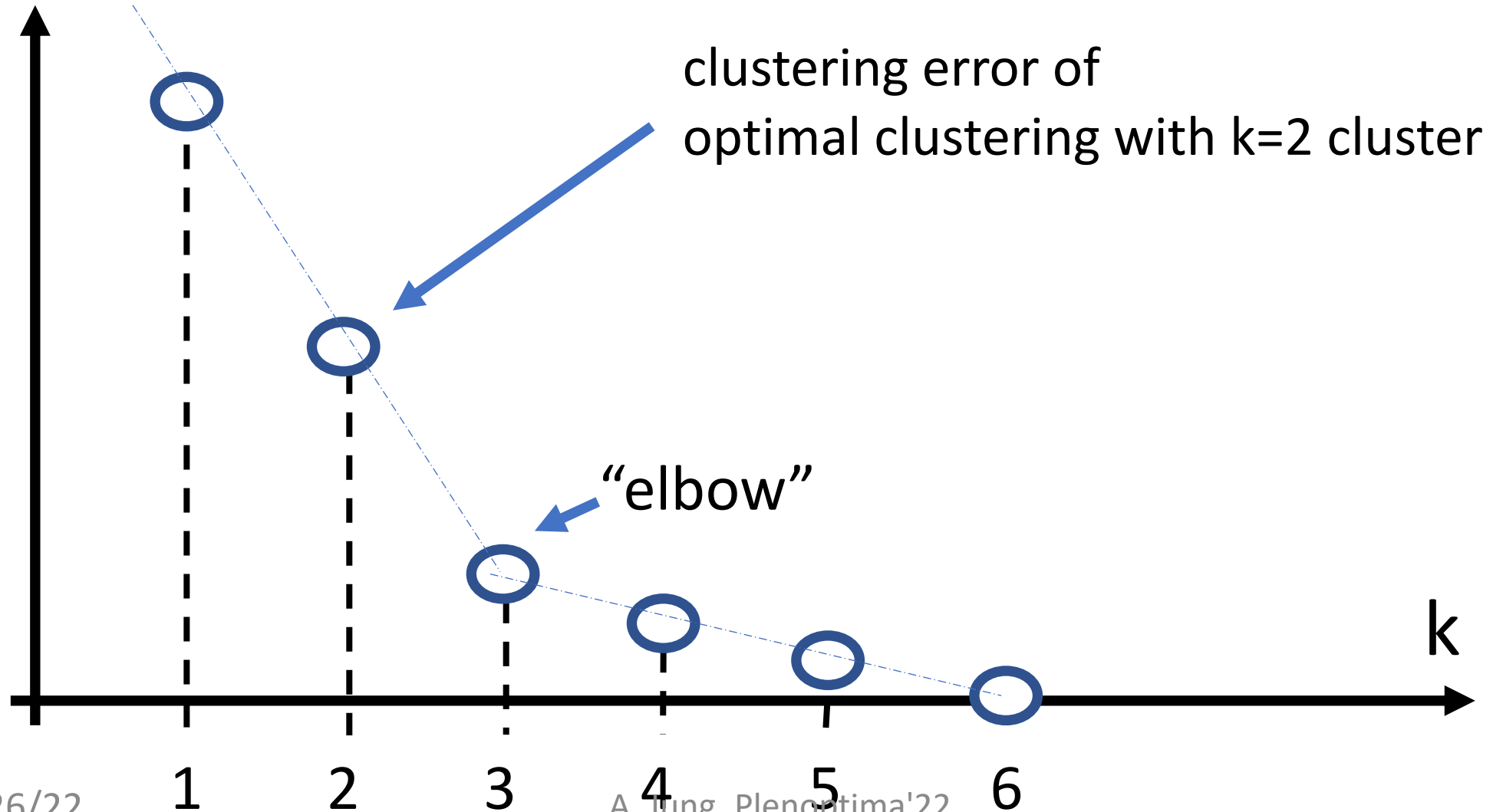
- defined by application (img. seg.)
- desired compression rate
- “elbow-method”

For/Background Segmentation $k=2$

Cluster 1 = Background, Cluster 2=Foreground



Elbow Method



Choose k by Validation Error

- clustering can be used as pre-processing for follow-up regression method
- try different values of k and pick the one resulting in smallest validation error

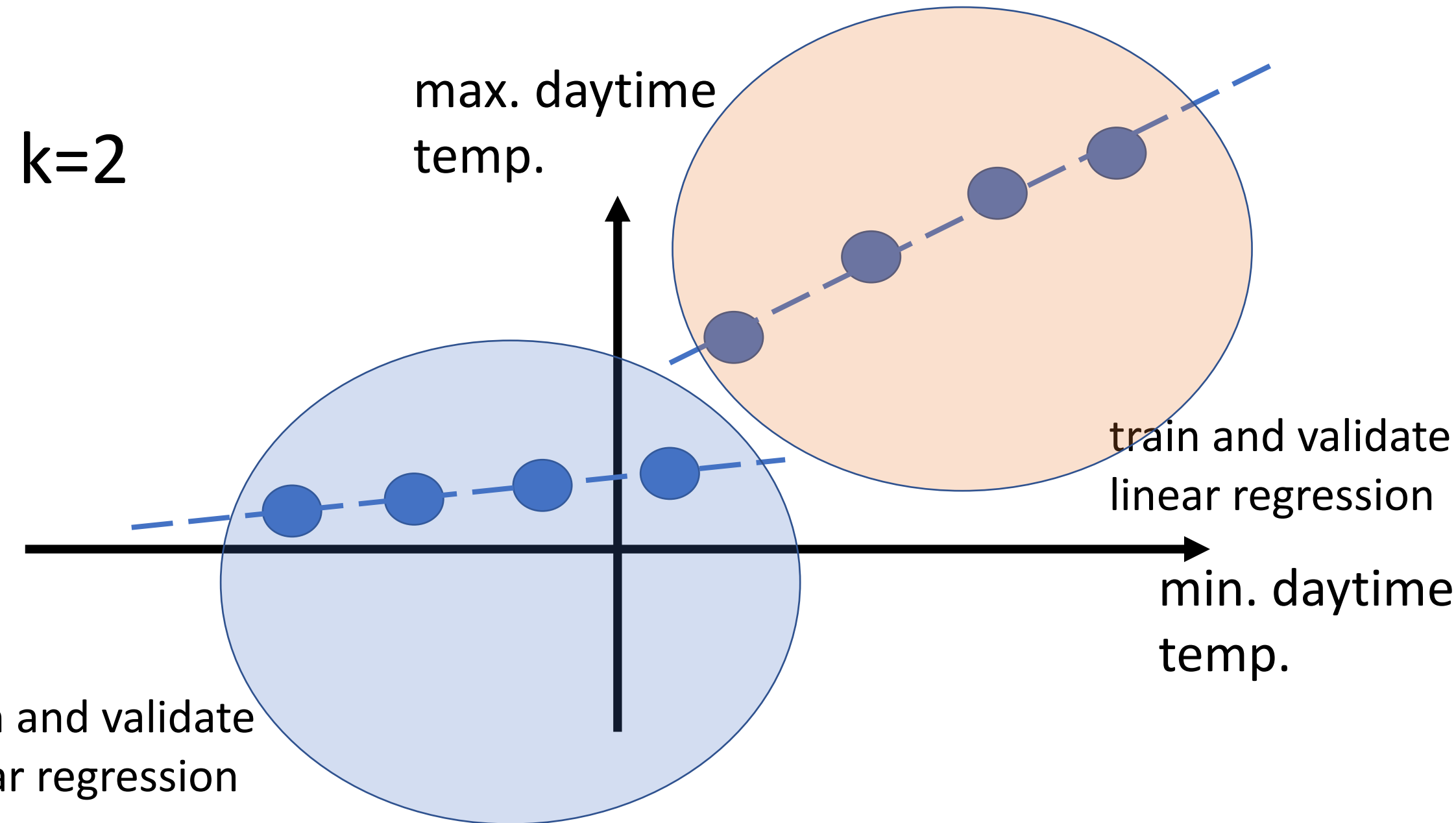
$k=2$

max. daytime
temp.

train and validate
linear regression

train and validate
linear regression

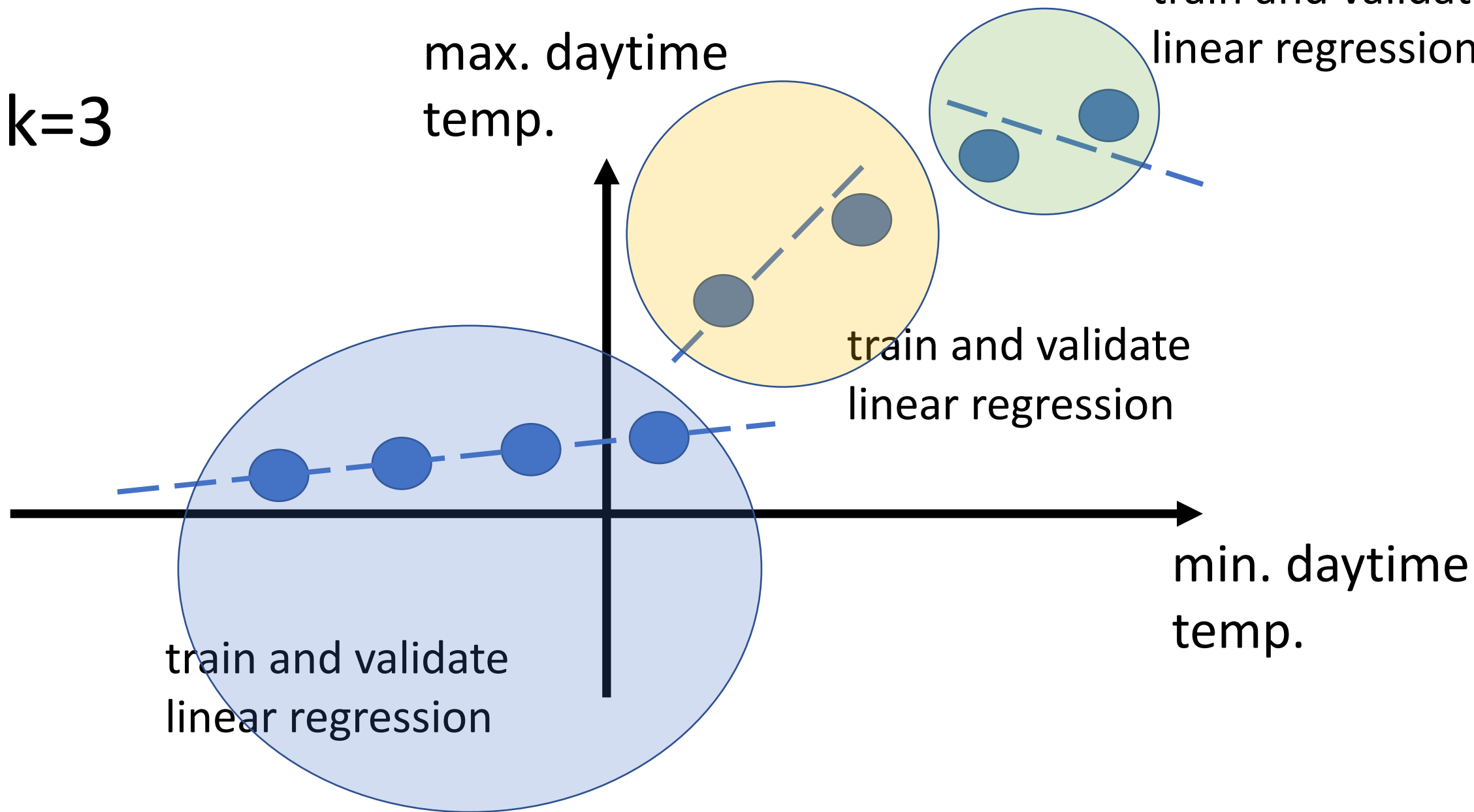
min. daytime
temp.



$k=3$

max. daytime
temp.

train and validate
linear regression



To Sum Up

- k-means partitions dataset into k clusters
- k-means iteratively minimizes clustering error
- k-means might deliver sub-optimal clustering
- repeat k-means with different initial cluster means
- number k of clusters needs to be given