

# Foundations of Data Science: Introduction to unsupervised learning

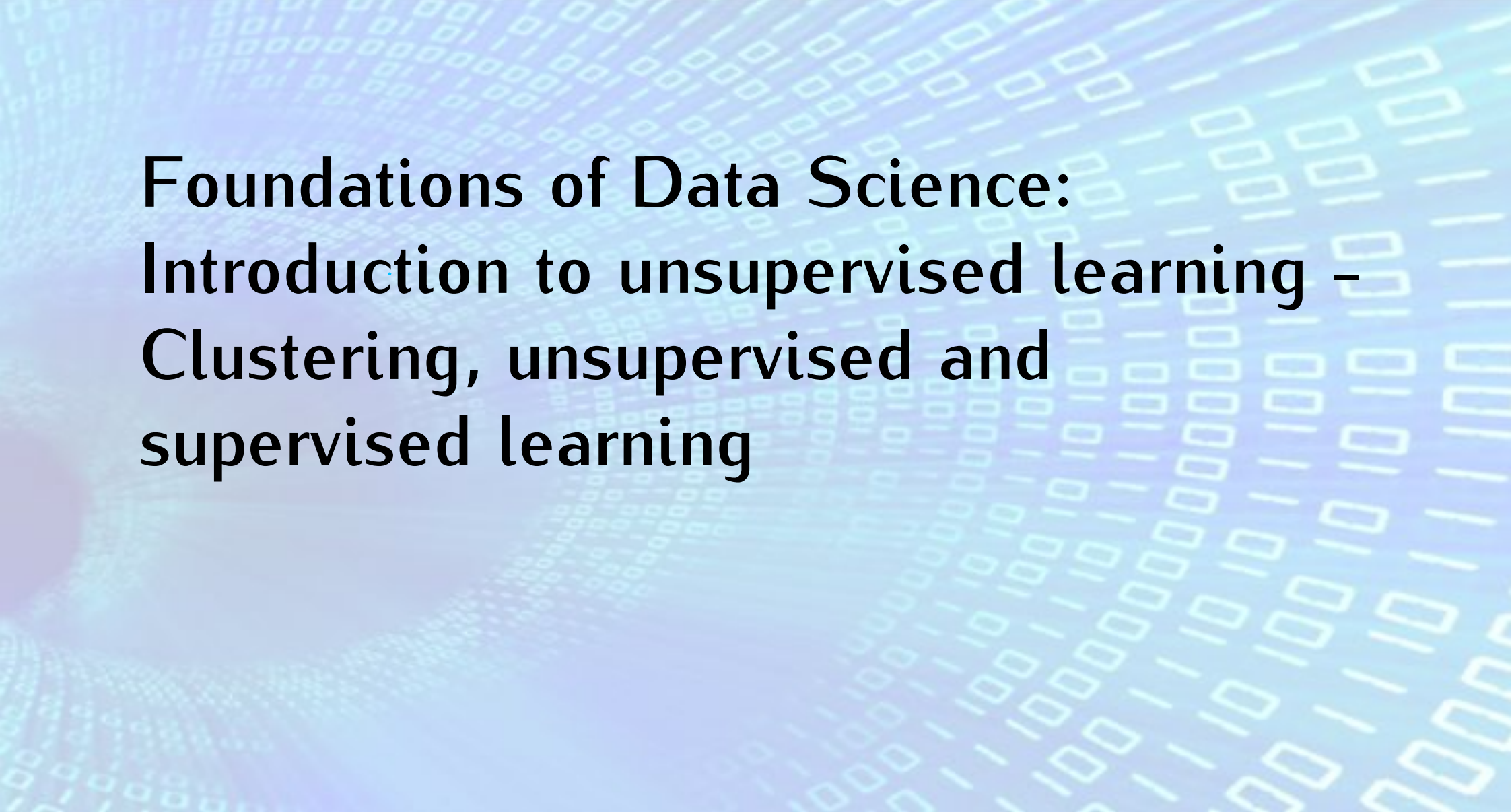


THE UNIVERSITY *of* EDINBURGH  
**informatics**

**FOUNDATIONS**  
**OF**  
**DATA**  
**SCIENCE**

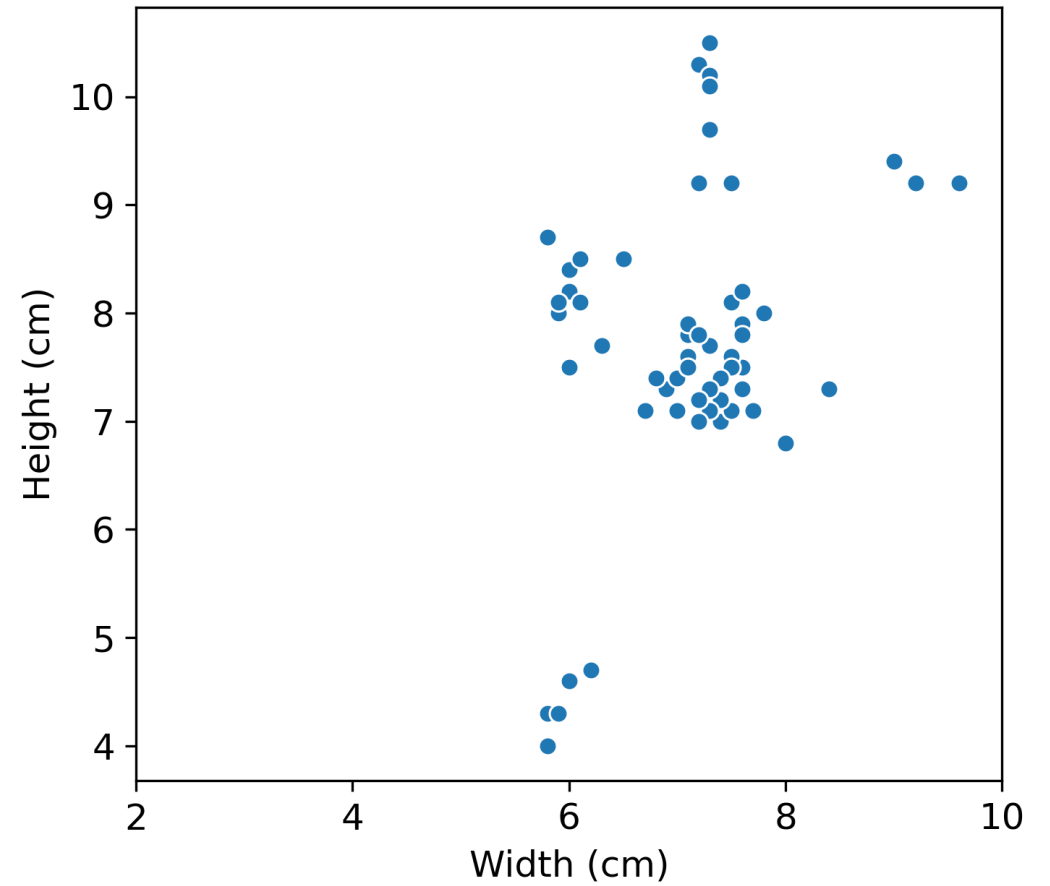
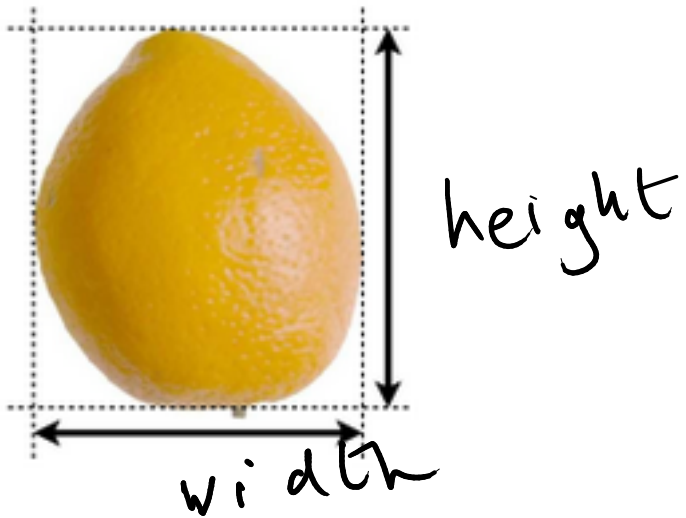
# Overview

1. Unsupervised learning, supervised learning, clustering
2. Partitional versus hierarchical clustering
3. K-means
4. Evaluation of K-means



**Foundations of Data Science:  
Introduction to unsupervised learning -  
Clustering, unsupervised and  
supervised learning**

# Clustering



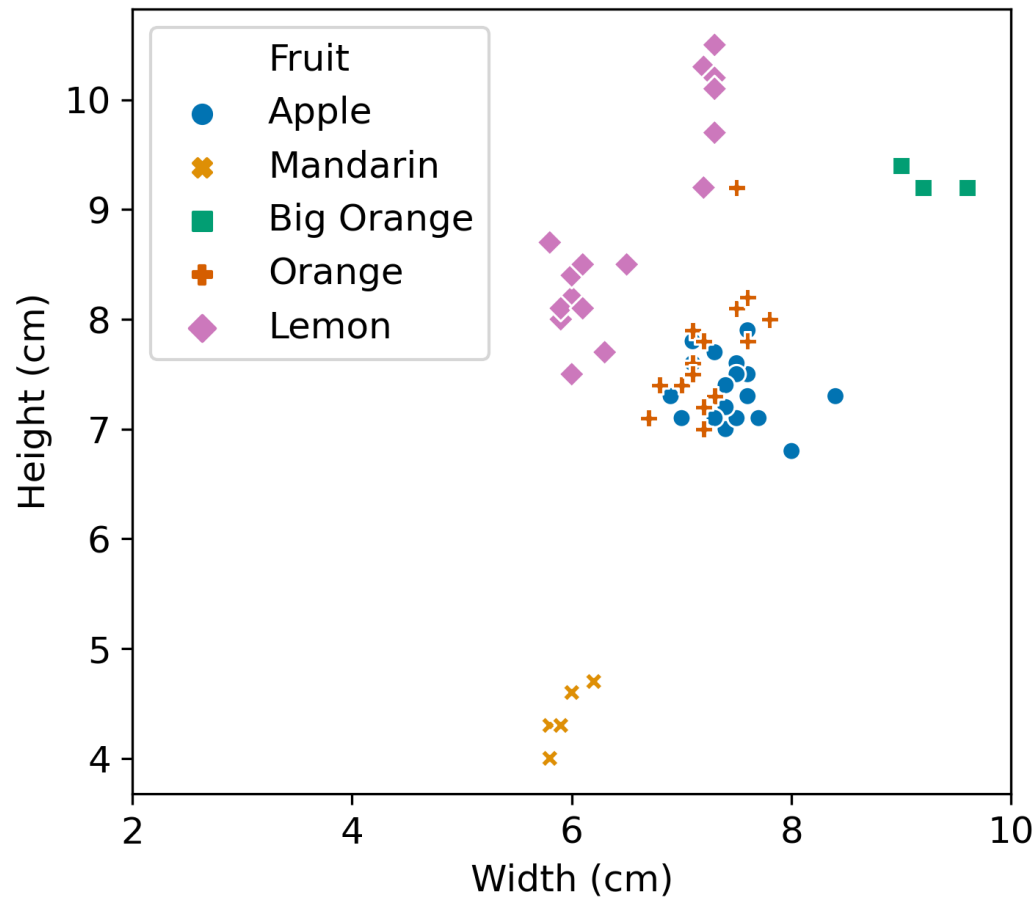
# Supervised learning process

For example, classification:

Training set:

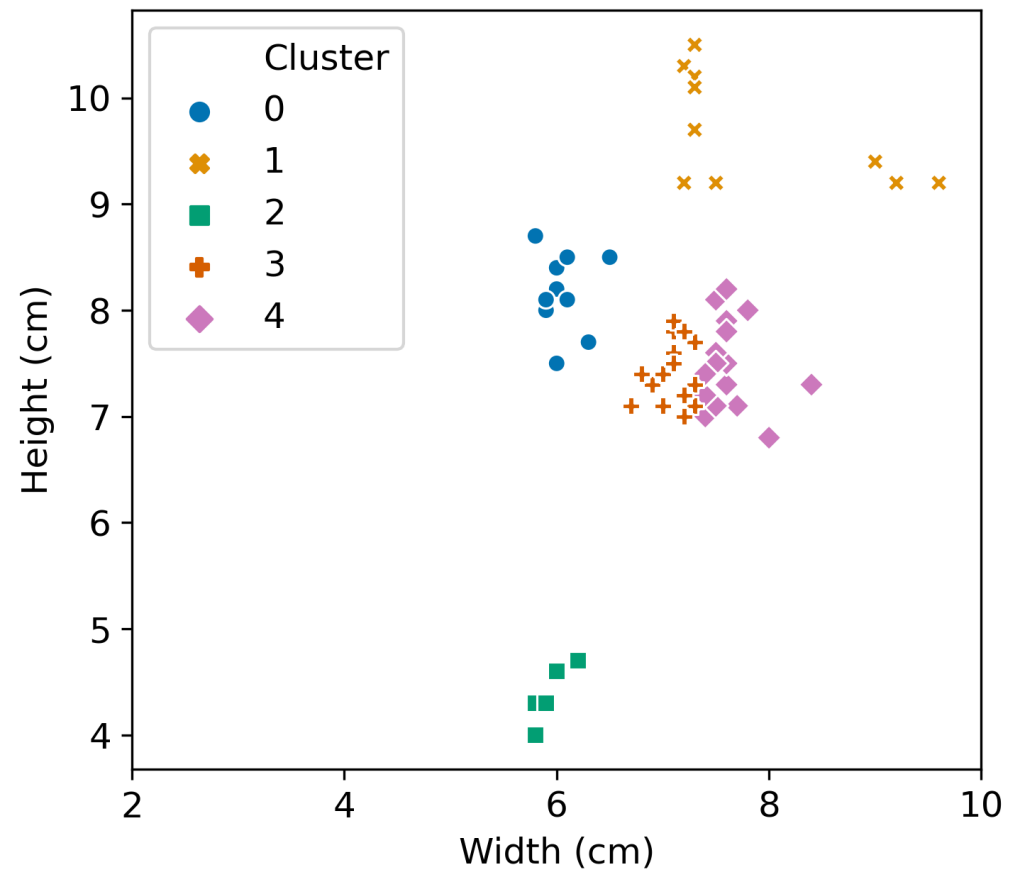
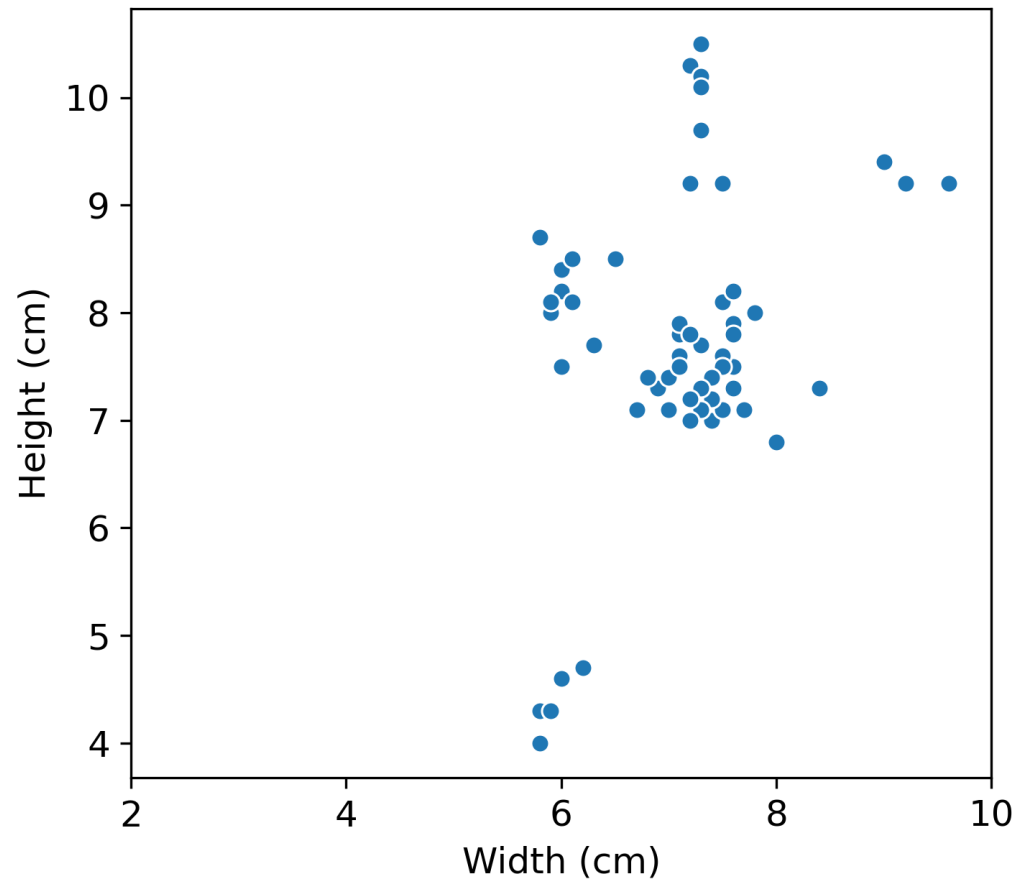
Features: width, height

Label: fruit

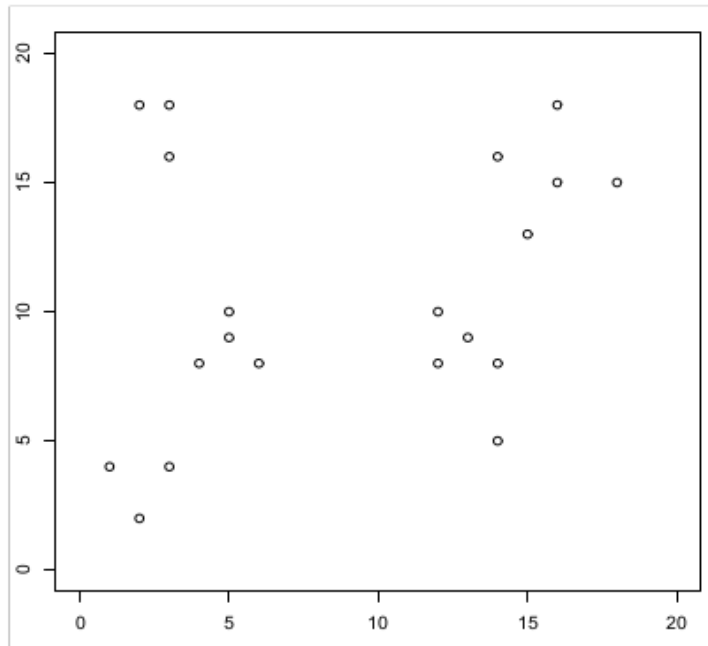


# Unsupervised learning process

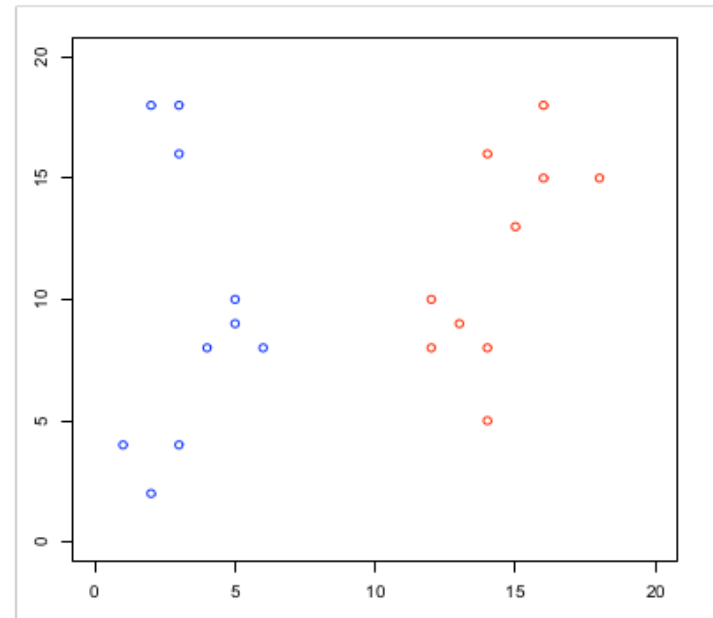
For example, clustering



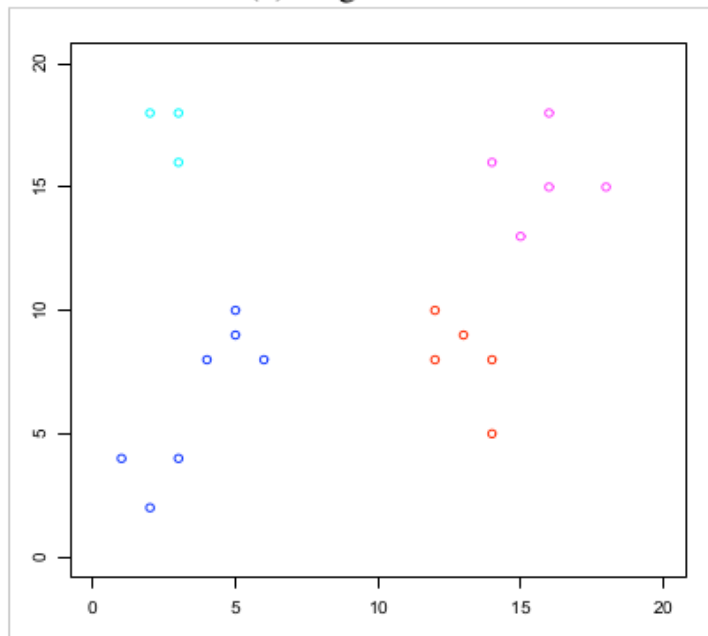
# How many clusters are there?



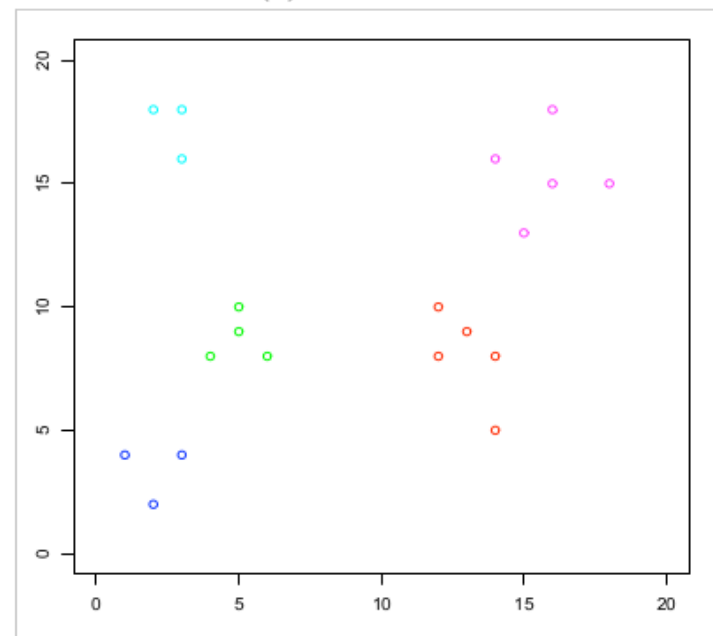
(a) original data



(b) two clusters



(c) four clusters

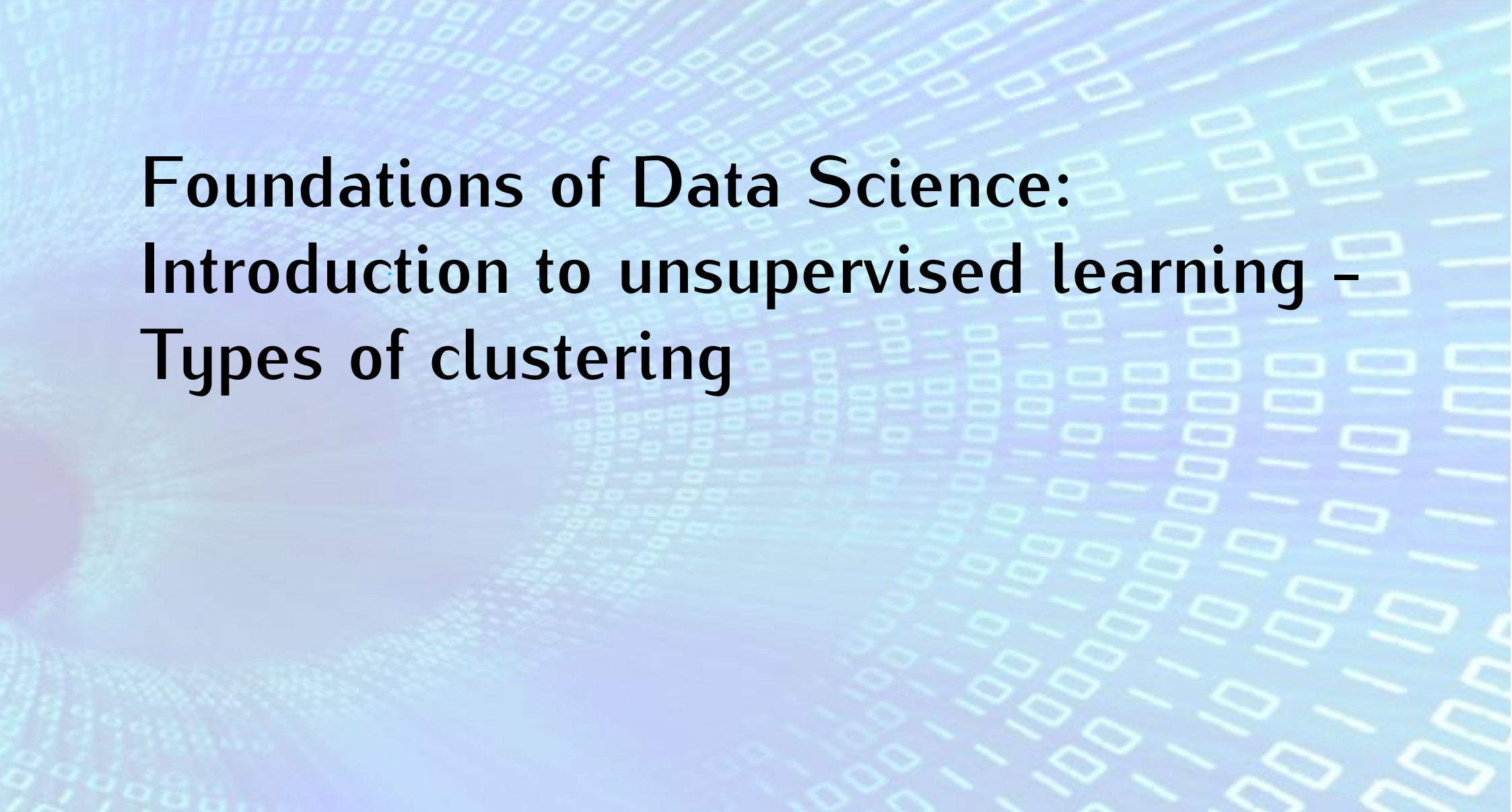


(d) five clusters

# Why cluster?

1. Interpretation
2. Data compression

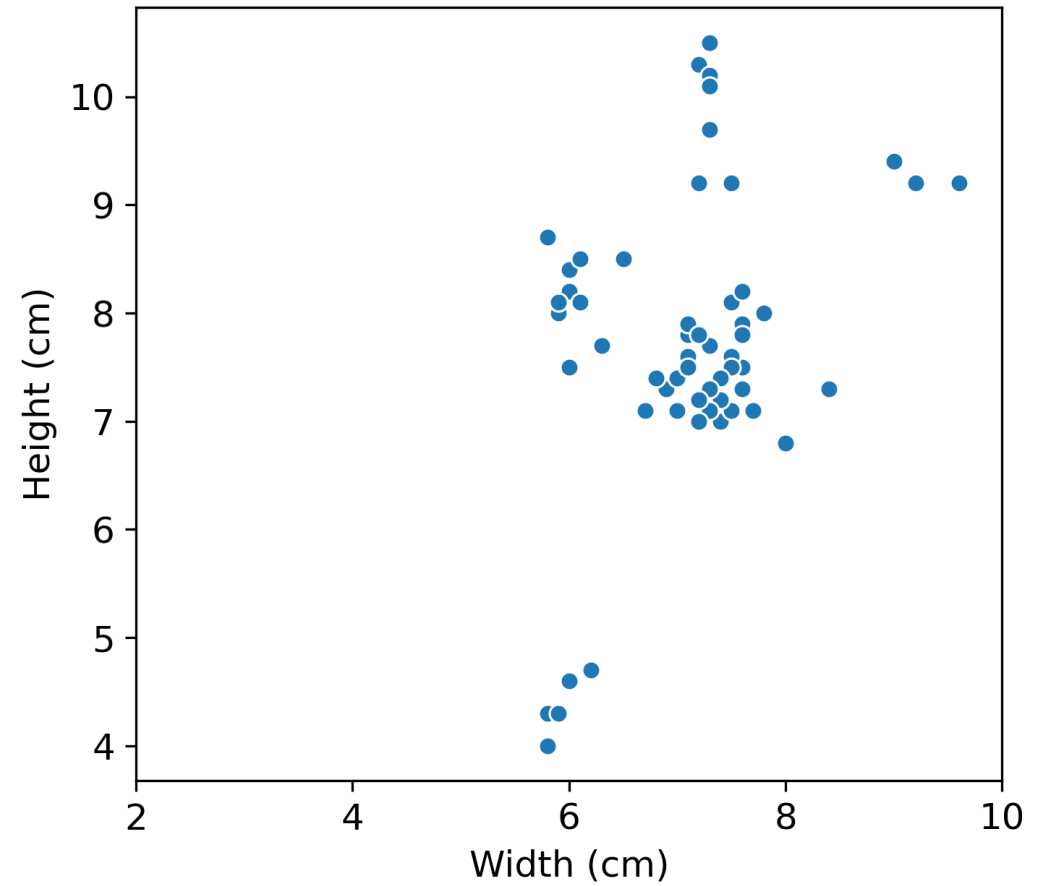
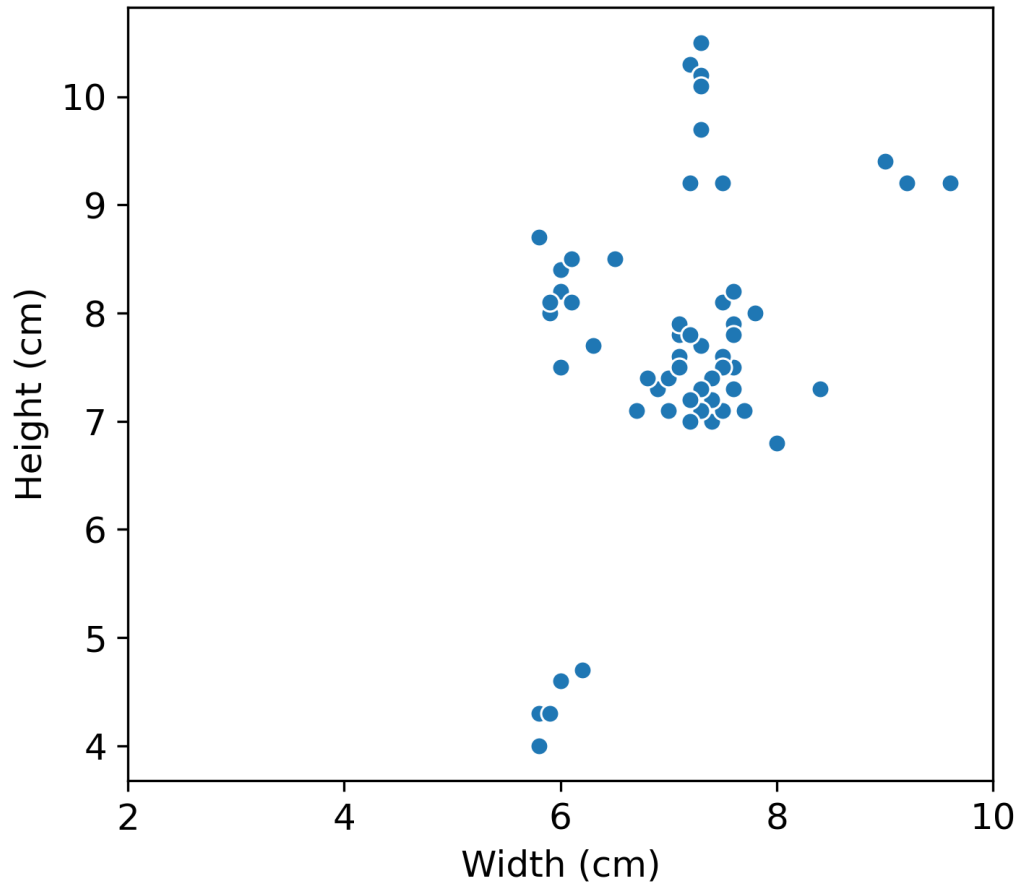




**Foundations of Data Science:  
Introduction to unsupervised learning -  
Types of clustering**

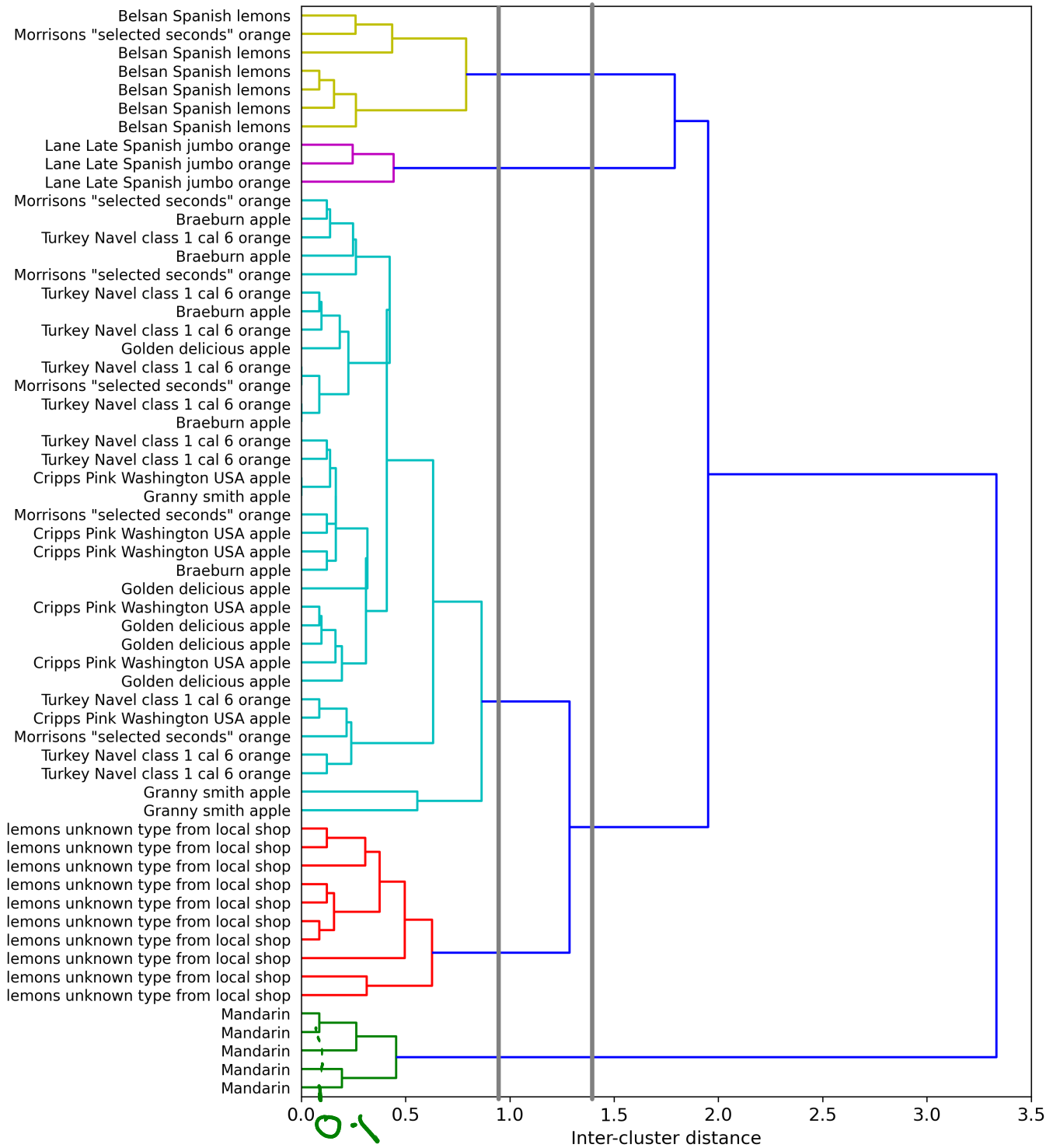
# Types of clustering

## Partitional

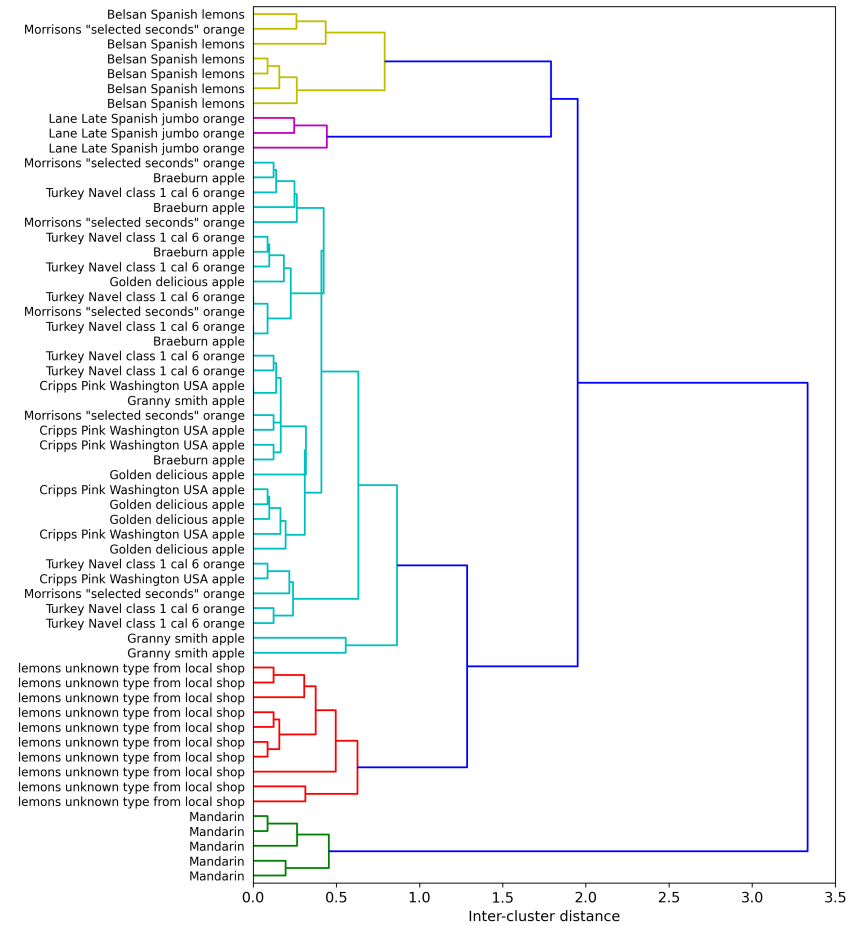
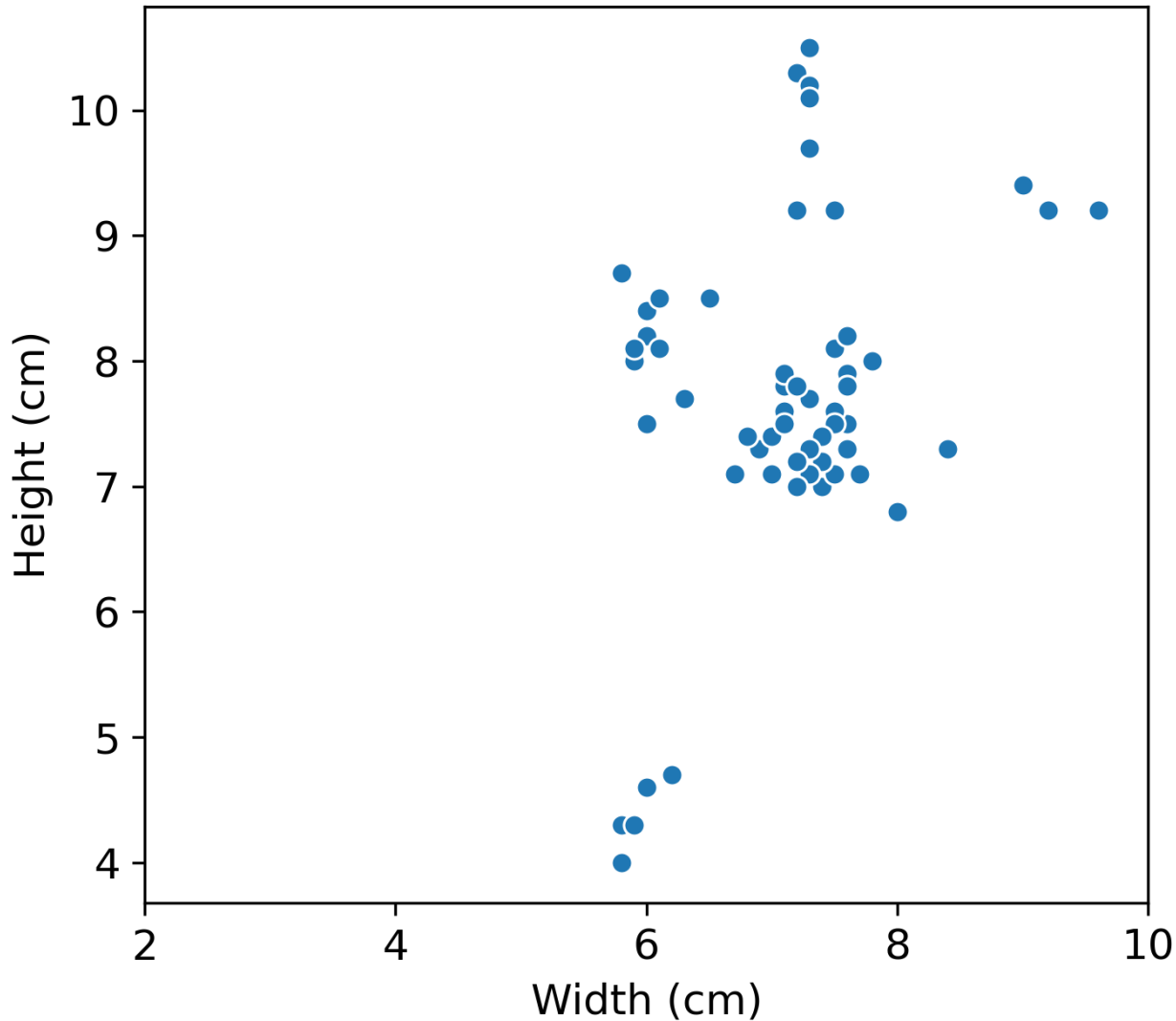


# Hierarchical clustering

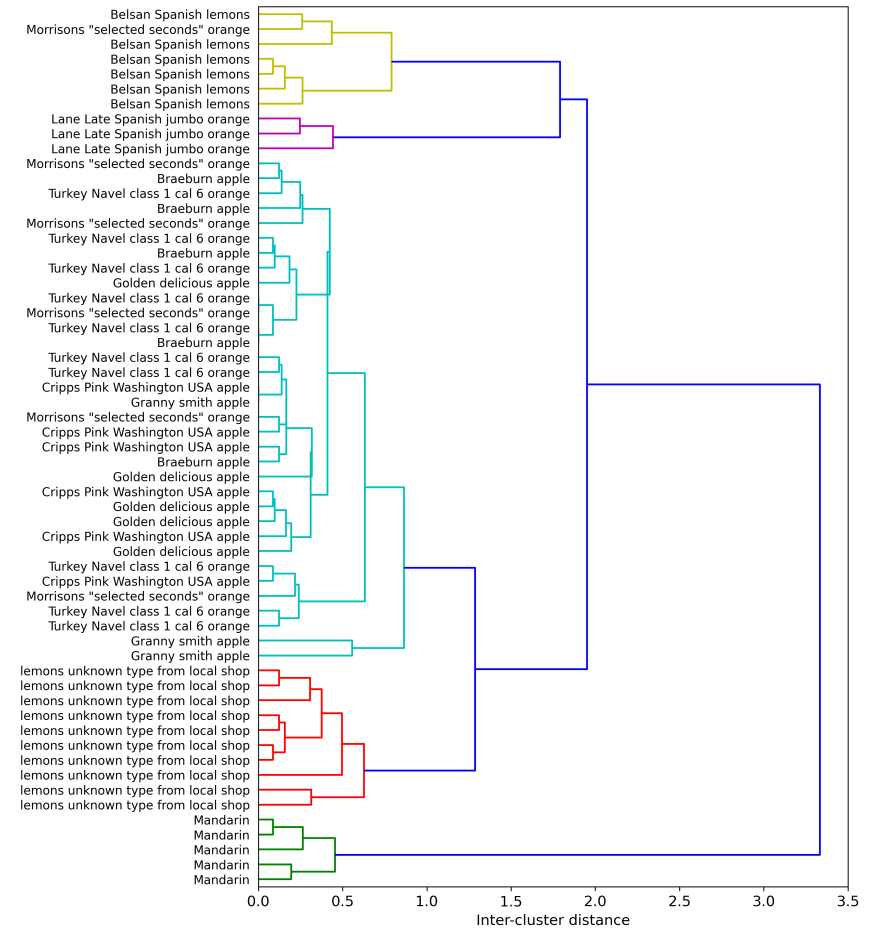
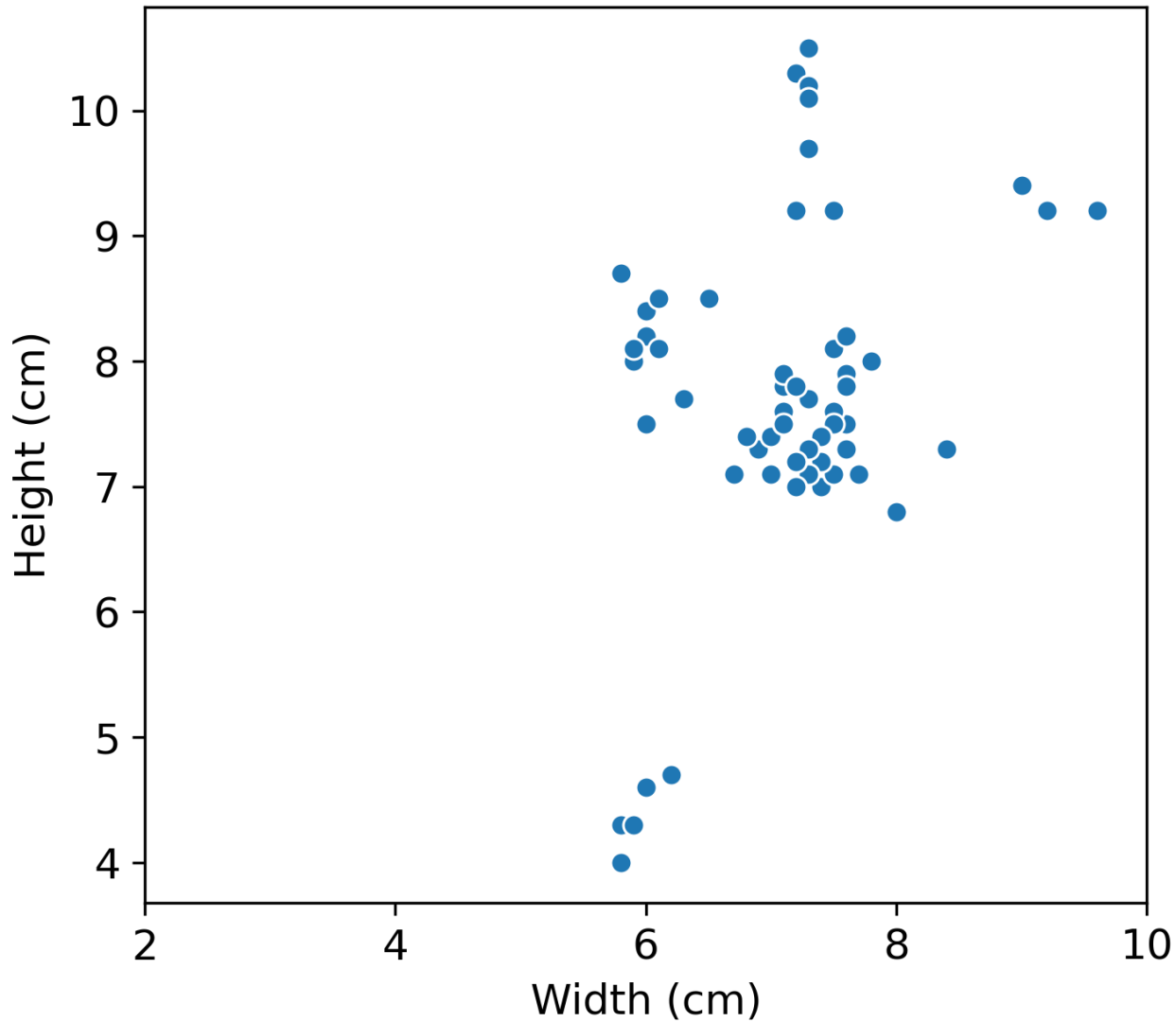
## Dendrogram

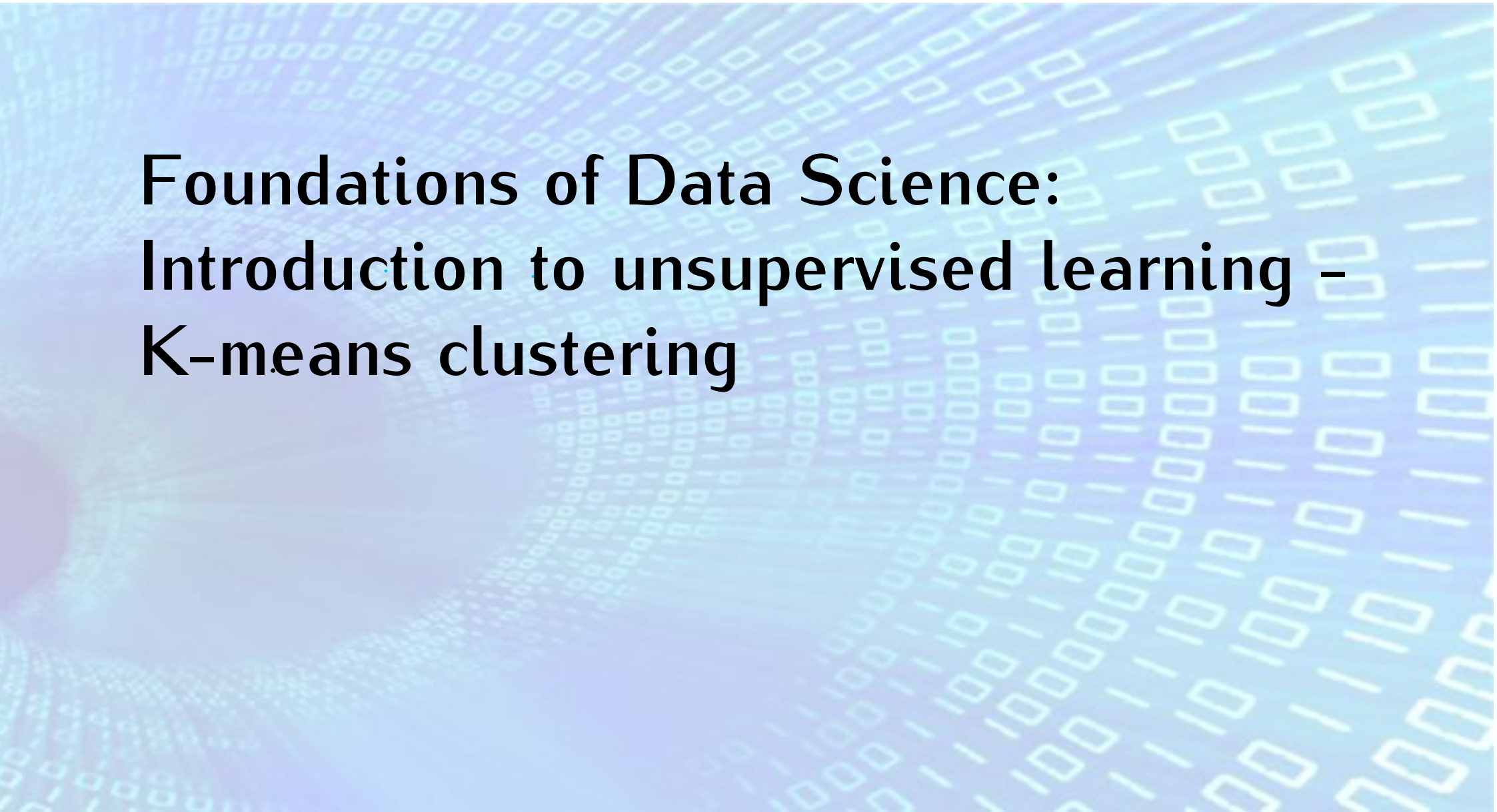


# Top-down hierarchical clustering



# Agglomerative (bottom-up) hierarchical clustering





**Foundations of Data Science:  
Introduction to unsupervised learning -  
K-means clustering**

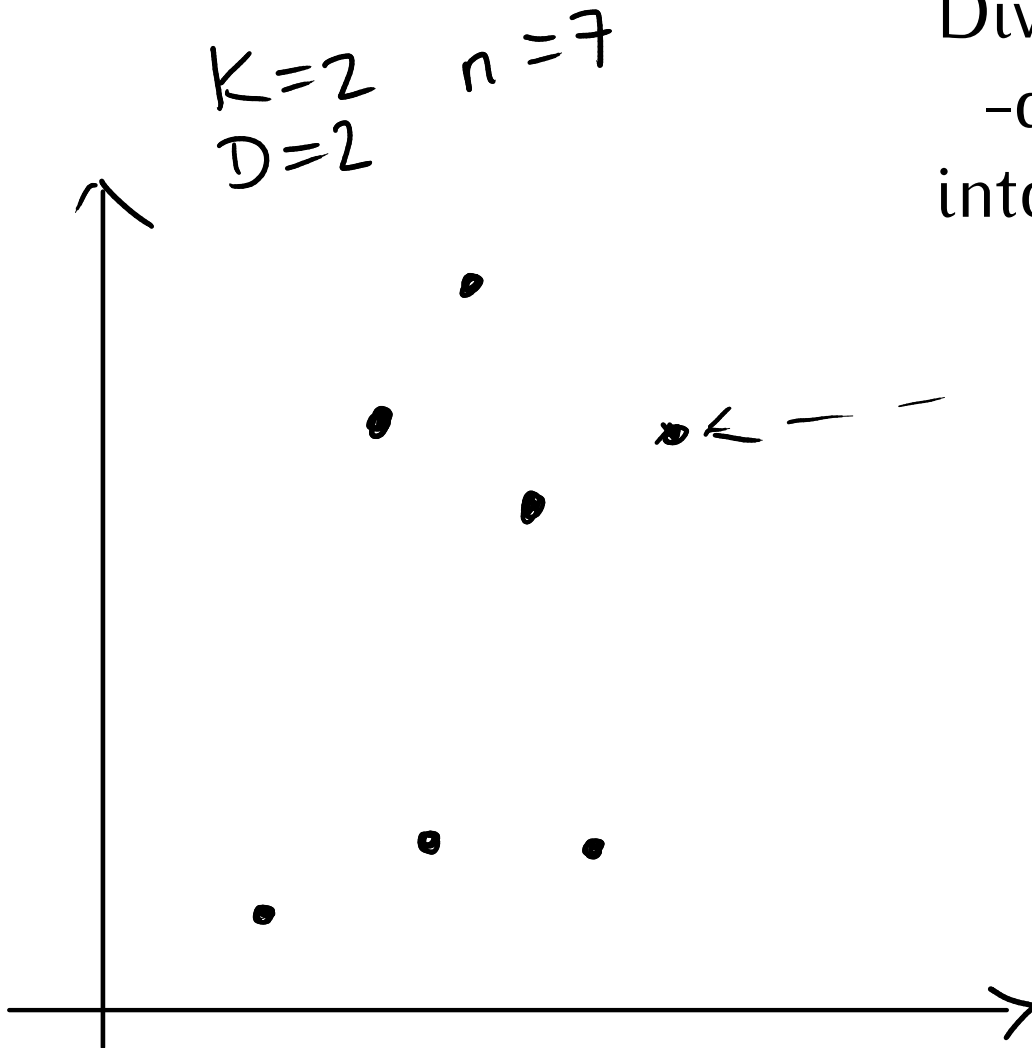
# K-means algorithm

Problem:

Divide a set of  
-dimensional data points  
into clusters

Algorithm:

1. Initialise cluster centres
2. While not converged
  - Assign
  - Recompute



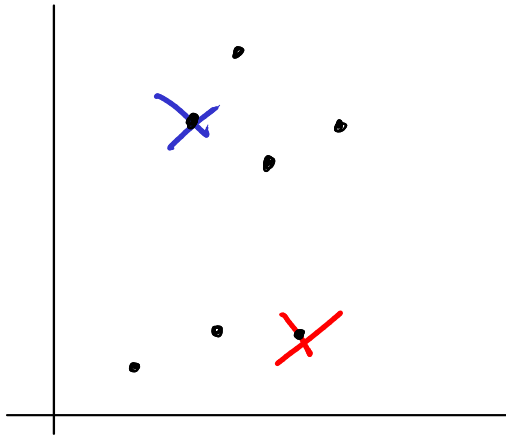
# Distance measure



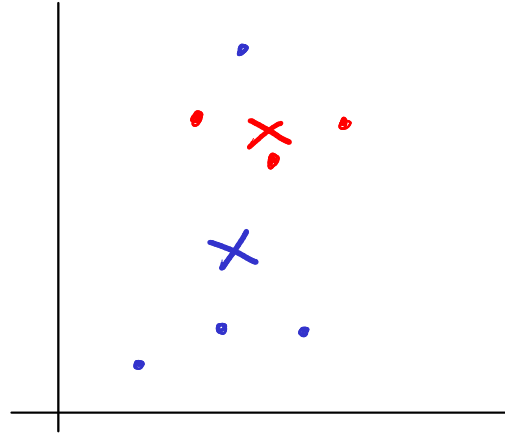
$$d(\underline{x}, \underline{y}) = \|\underline{x} - \underline{y}\| = \sqrt{\sum_{j=1}^D (x_j - y_j)^2}$$



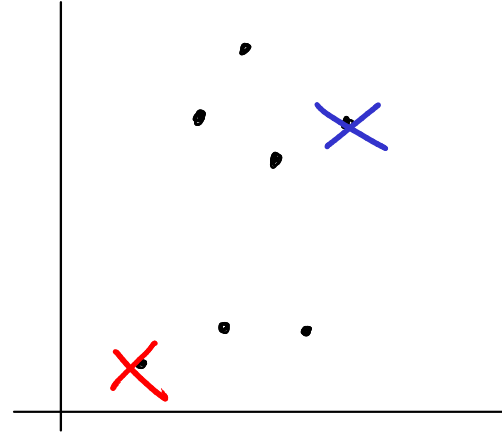
# Initialisation methods



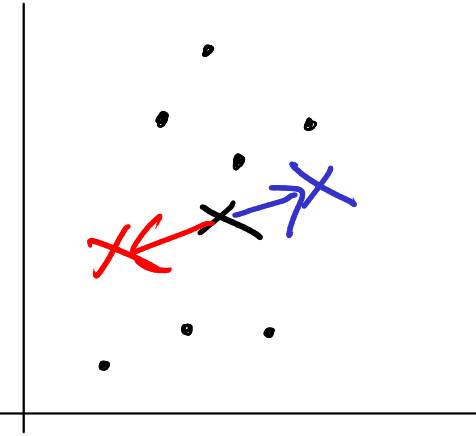
Random data points as cluster centres



Random assignment to clusters



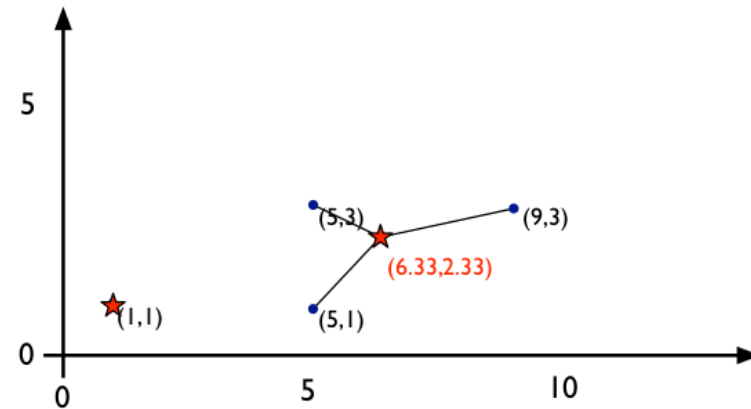
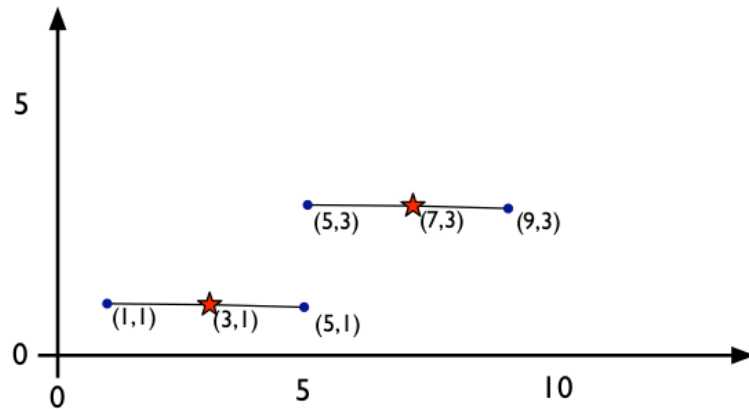
Data points with extreme values



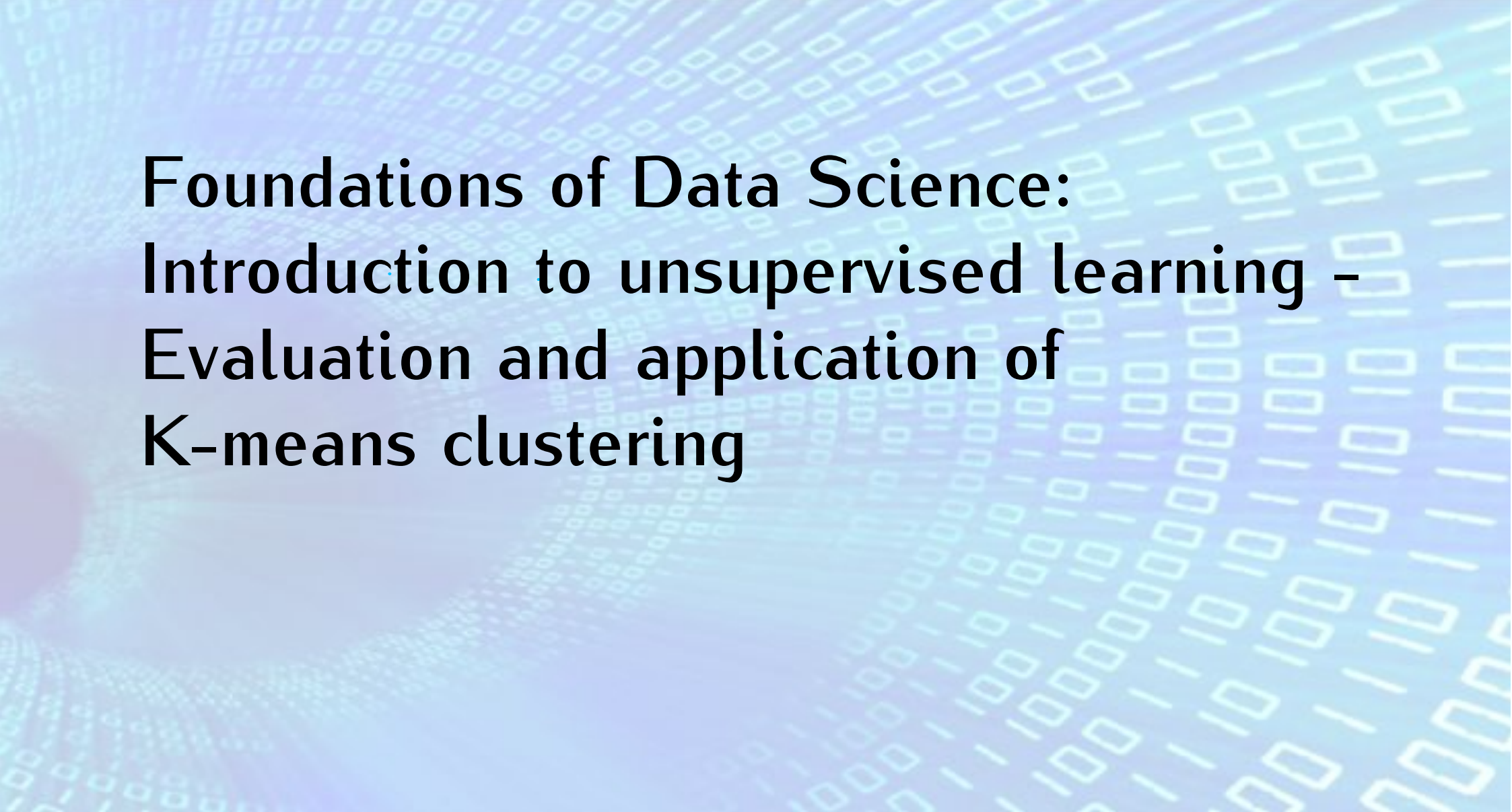
Mean for whole dataset and perturbation

# Convergence

Convergence is guaranteed...

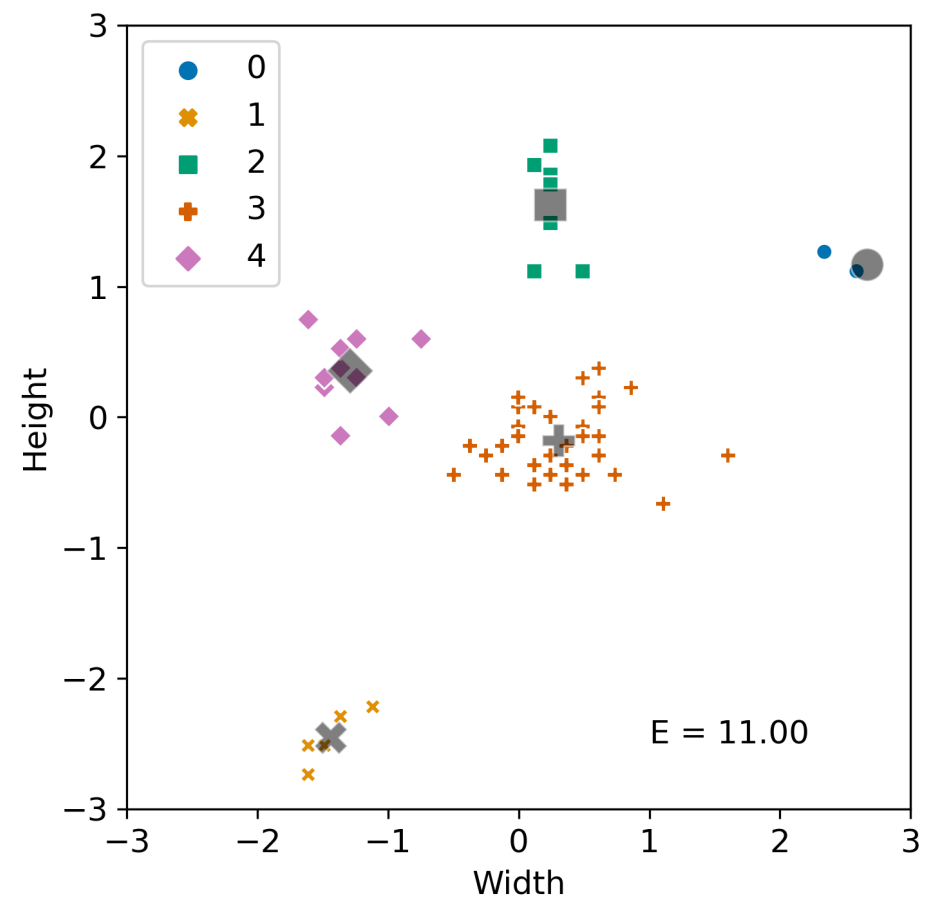
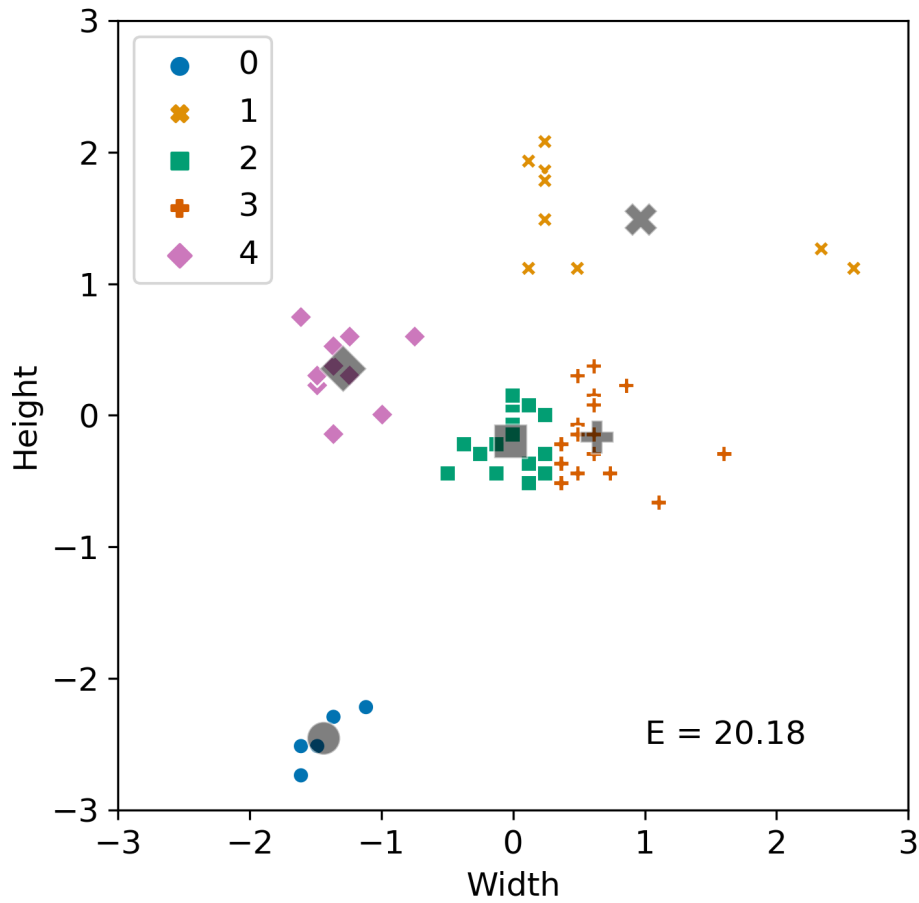


... but unique solutions are not.

The background of the slide features a blue-toned digital aesthetic. It includes a faint, glowing globe on the left side and a pattern of binary code (0s and 1s) that appears to be receding into the distance, creating a sense of depth and data flow.

**Foundations of Data Science:  
Introduction to unsupervised learning -  
Evaluation and application of  
K-means clustering**

# Multiple solutions



# Mean squared error function

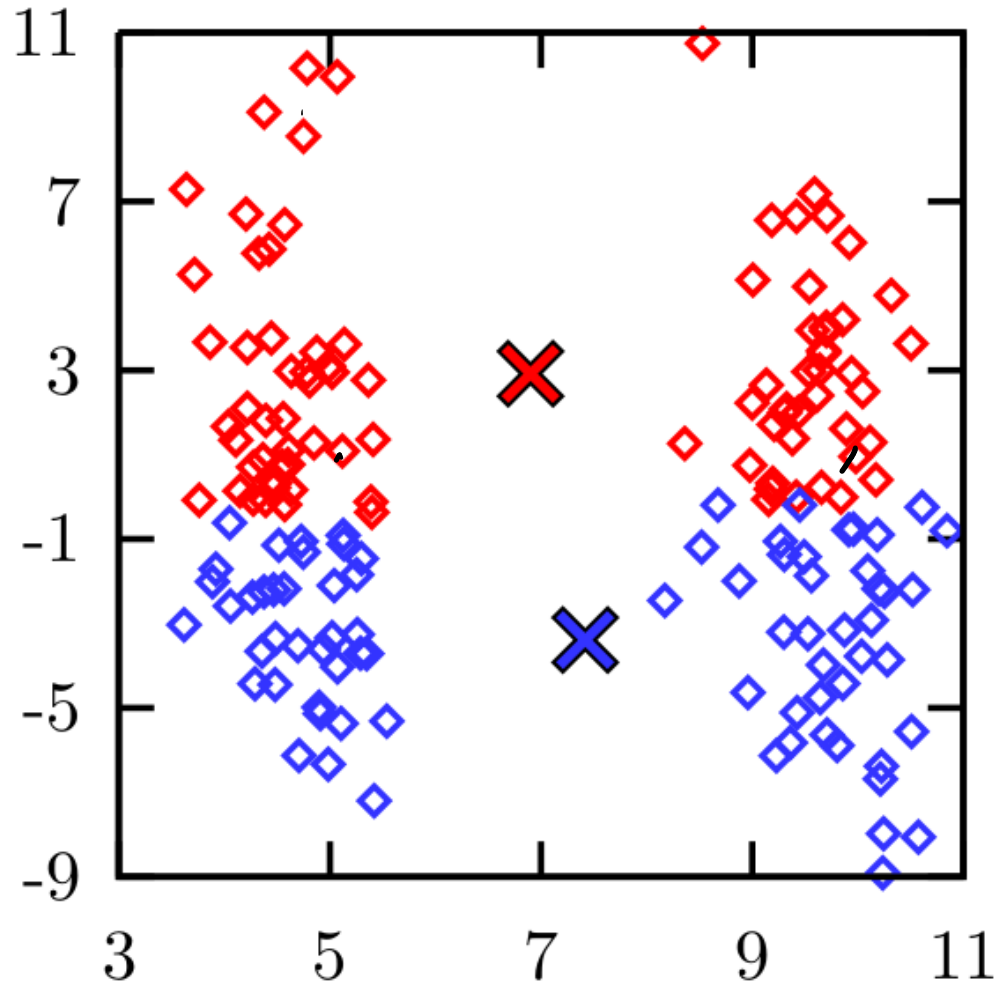
AKA inertia

If point  $i$  belongs to cluster  $k$   
 $i \in C_k$

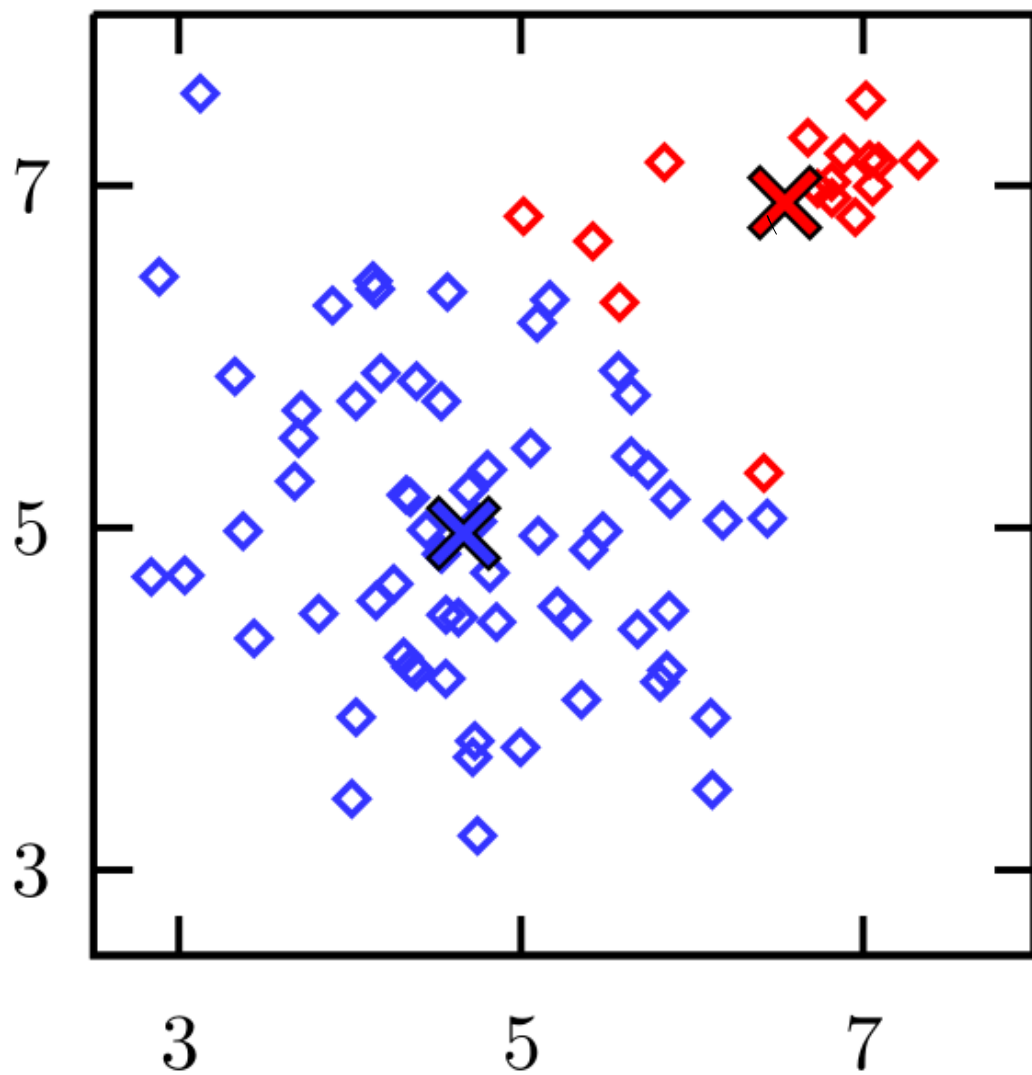
$$E = \frac{1}{n} \sum_{k=1}^K \sum_{i \in C_k} \| \underline{x}_i - \underline{m}_k \|^2$$

Minimum variance clustering

# Failures of K-means (1)

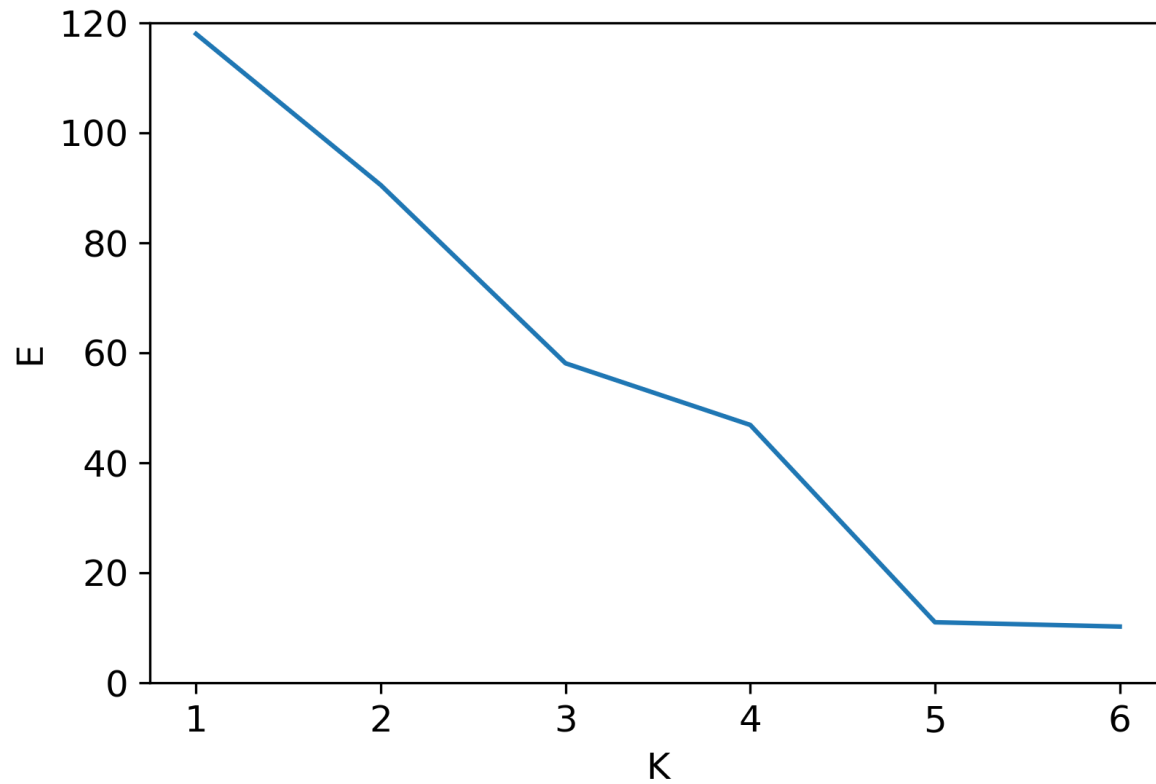


# Failure of K-means (2)



# How to choose K?

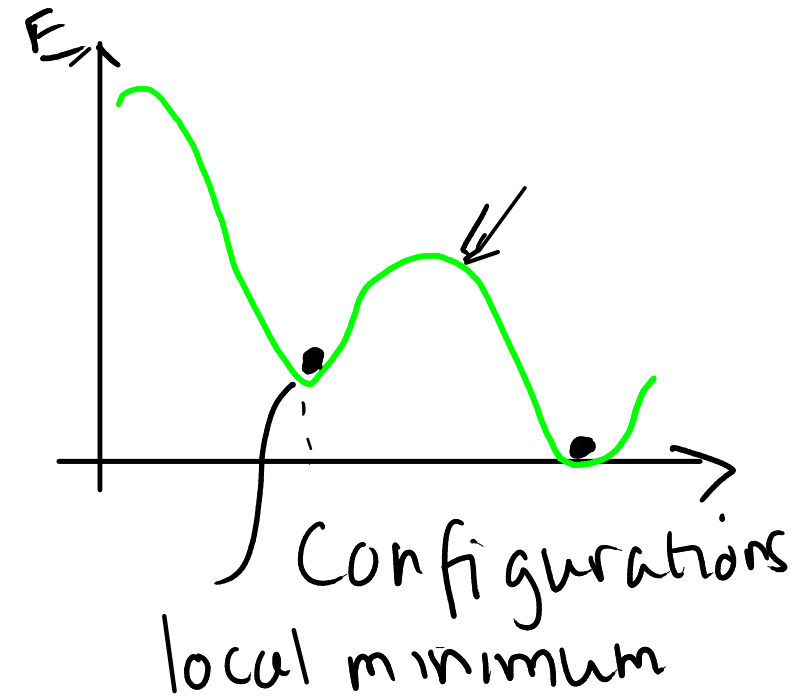
Scree plot





# Batch versus online

Online versus batch



# The curse of Dimensionality

⇒ Use dimensionality  
reduction to  
overcome the curse.

