# Clustering

# Unsupervised Learning



- *Supervised learning*: There is a labeled training/validation set that can be used to tune the algorithm parameters
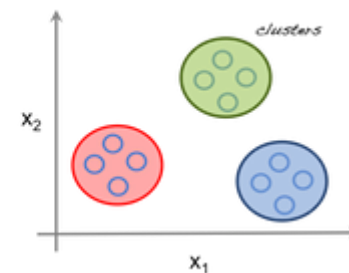- *Unsupervised Learning*: Training data is not labeled

*Unsupervised Learning*:
- We are interested in finding some interesting structure in the data, or, equivalently, to organize it in some meaningful way
- *Target:* find a function that describes the structure of "*unlabeled*" data (i.e., data that has not been classified or categorized)
- Several approaches are based on the idea that the data is the realization of a *hidden probability density function*, i.e., unsupervised learning is linked to the density estimation of a hidden PDF producing the data

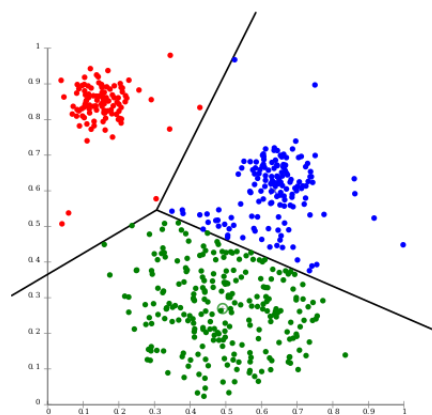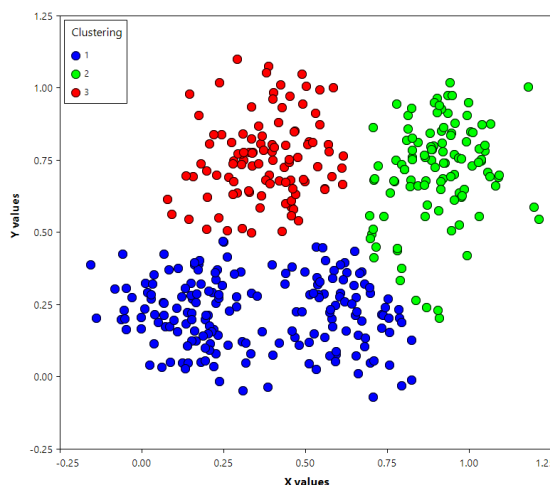We are going to see only a couple of unsupervised learning techniques and a few very simple and commonly used methods
1. *Clustering*
   o *K-means*
   o *Linkage-based clustering*
2. Dimensionality reduction
   o Principal Component Analysis (PCA) →later in the course


There are many other techniques (*not part of this course*)
- Mean shift clustering, spectral clustering….
- Compressive sensing

**Clustering**

Idea: Divide a set of objects represented by *N*-dimensional vectors into groups (*clusters*) of similar objects

➢ Key target: *identifying meaningful groups among data points*

➢ The definition is not rigorous and may be ambiguous, different definitions have been proposed leading to different algorithms

**Formal Definition:**

Clustering is the task of grouping a set of objects such that *similar objects end up in the same group* and *dissimilar objects are separated into different groups*

**DIPARTIMENTO DI INGEGNERIA DELL'INFORMAZIONE**

*Similarity is not transitive:*

"*similar objects in same group*" and
"*dissimilar objects into different groups*"
may contradict each other...

2 clusters

"similar objects in same group":

"dissimilar objects into different groups"

## *There is no ground truth:*

How To Evaluate Performances ?

Divide in
2 clusters

# Clustering: Applications

**Many Applications:**

- Business and marketing
  - Market research
  - Grouping of shopping items
- World Wide Web
  - Search engines
  - Social network analysis
- Image segmentation
- Medicine, Medical imaging
- Biology and bioinformatics
- Recommender systems
- Anomaly detection
- Natural language processing

# Example (1): Customer Segmentation



- Data: features (e.g., products bought, demographic info, etc..) for a large number of customers
- Goal: customers segmentation → identify groups of homogeneous customers
- Useful for advertising, product development, …

# Example (2): Image Segmentation



Medical Imaging

Movies/Special Effects (chroma keying)

Features extraction/detection

Object Recognition

3D Reconstruction

- *Data*: one data sample for each pixel
  - ➤ *Samples*: vectors containing features (e.g., color or spatial position of pixels)
- *Goal*: divide the image into regions (*clusters*) with uniform properties
- Useful for medical imaging, image analysis, background segmentation in movies, object recognition, ….

# Clustering Model

**Input:**

- Set of elements $x \in \mathcal{X}$
- Distance function $d: \mathcal{X} x \mathcal{X} \rightarrow \mathbb{R}_+$
    1. symmetric, i.e.,: $d(x, y) = d(y, x) \; \forall \; x, y$
    2. $d(x, y) \geq 0 \; \forall \; x, y$ and $d(x, x) = 0$
    3. Triangle inequality $d(x, z) \leq d(x, y) + d(y, z)$

**Output:**

A partition $C = (C_1, C_2, .., C_k)$ of set $\mathcal{X}$ into $k$ clusters

- $\bigcup_{i=1}^{k} C_i = \mathcal{X}$
- $\forall i \neq j: C_i \cap C_j = \emptyset$

- $k$ (# of clusters): sometimes given in input, sometimes computed by the algorithm

x                y

z < x + y

x                y

z

x                y

z ≈ x + y

**Dendrogram**: tree, with input points $x \in \mathcal{X}$ as leaves, that shows the arrangement/relation between clusters.

*Sometimes, the output is a dendrogram (from Greek Dendron = tree, gramma = drawing), a tree diagram showing the arrangement of the clusters*

Very Common approach in clustering:

- Define a cost function over possible partitions of the objects
- Find the partition (→clustering) of minimal cost

Assumptions:

- Data points come from a larger space $\mathcal{X}'$ (typically $\mathbb{R}^n$)
- Distance function $d(x, x')\ for\ x, x' \in \mathcal{X}$
- For simplicity: assume $\mathcal{X}' = \mathbb{R}^n$ and $d(x, x') = \|x - x'\|_2$
  - I.e., use Euclidean distance

# K-Means

- ❑ The simplest distance-based clustering algorithm
- ❑ Proposed in 1957, also known as Lloyd algorithm
- ❑ Choose a fixed number of clusters
- ❑ Find cluster centers and point-cluster allocations in order to minimize the error made by approximating the points with the cluster centers
- ❑ Can't do this by exhaustive search, because there are too many possible allocations
- ❑ Iterative algorithm
  - o fix cluster centers; assign each point to the closest cluster
  - o fix allocation; compute best cluster centers
- ❑ Vectors $x$ can be any set of features for which we can compute a distance (*careful about scaling for non-homogenous data*)

*Some material on K-Means D.A. Forsyth*

| | |
|---|---|
| $\mathcal{X} \subset \mathbb{R}^n$ | Set of vectors to be clustered |
| $\boldsymbol{x} \in \mathcal{X}$ | Vector to be clustered |
| $k$ | Number of clusters (*parameter of the algorithm*) |
| $C_i \quad i = 1, \dots, k$ | Clusters (each vector $\boldsymbol{x}$ is associated to a cluster) |
| $\boldsymbol{\mu}_i \quad i = 1, \dots, k$ | Centroids of the clusters |

Find cluster centers and allocations in order to minimize the error made by approximating the points with the cluster centers:

centroid of $C_i$

use squared distance

if using euclidean distance

$$\boldsymbol{\mu_i} = \operatorname*{argmin}_{\boldsymbol{\mu}} \sum_{\boldsymbol{x} \in C_i} d(\boldsymbol{x}, \boldsymbol{\mu})^2 = \operatorname*{argmin}_{\boldsymbol{\mu}} \sum_{\boldsymbol{x} \in C_i} \|\boldsymbol{x} - \boldsymbol{\mu}\|^2$$

$$\mathrm{G_{km}}\big((\mathcal{X}, d), (C_1, \dots, C_k)\big) = \sum_{i=1}^{k} \sum_{\boldsymbol{x} \in C_i} d(\boldsymbol{x}, \boldsymbol{\mu}_i)^2$$

---

**Theorem:**

Given a cluster $C_i$, the center $\boldsymbol{\mu_i}$ that minimizes $\sum_{\boldsymbol{x} \in C_i} d(\boldsymbol{x}, \boldsymbol{\mu_i})^2$ is

$$\boldsymbol{\mu_i} = \frac{1}{|C_i|} \sum_{\boldsymbol{x} \in C_i} \boldsymbol{x}$$

---

- ❑ Demonstration: compute gradient and set to 0

$$\frac{\partial}{\partial \boldsymbol{\mu_i}} \left( \sum_{\boldsymbol{x} \in C_i} \|\boldsymbol{x} - \boldsymbol{\mu_i}\|^2 \right) = \sum_{\boldsymbol{x} \in C_i} 2(\boldsymbol{x} - \boldsymbol{\mu_i}) = \boldsymbol{0} \ \rightarrow \ \sum_{\boldsymbol{x} \in C_i} \boldsymbol{x} = |C_i| \boldsymbol{\mu_i} \ \rightarrow \ \boldsymbol{\mu_i} = \frac{1}{|C_i|} \sum_{\boldsymbol{x} \in C_i} \boldsymbol{x}$$

- ❑ Naive (brute-force) algorithm to solve K-Means Clustering?
- ❑ Try all possible partitions of the $m$ points into $k$ clusters, evaluate each partition, and find the best one
- ❑ Is it efficient?
  - o Number of possible partitions is exponential in $m$
  - o NP-Hard problem

# K-Means: Algorithm

Procedure:

1. Select $k$ random centroids (*or use some more advanced initialization strategy*)

2. Each point is associated to the closest centroid (according to the distance measure)

$$\forall i: C_i = \{x \in \mathcal{X} : i = \operatorname*{argmin}_j \|x - \mu_j\|\}$$

3. Compute the new centroids (each centroid is the barycentre of the associated points)

$$\forall i : \mu_i = \frac{\sum_{x \in C_i} x}{|C_i|}$$

4. Repeat step 2 and 3 until the algorithm converges

*Theorem:*

*At each iteration the value of the objective function $G_{km}$ does not increase*

*Theorem: at each iteration the value of the objective function $G_{km}$ does not increase*

1. Consider K-means objective func. (simplified notation) $G(C_1, \ldots, C_k) = \min\limits_{\mu_1,\ldots,\mu_k \in \mathbb{R}^n} \sum_{i=1}^{k} \sum_{x \in C_i} \|x - \mu_i\|^2$

2. Centroids minimize distance w.r.t points in the associated cluster

$$\mu(C_i) \stackrel{\text{def}}{=} \frac{1}{|C_i|} \sum_{x \in C_i} x = \operatorname*{argmin}_{\mu \in \mathbb{R}^n} \sum_{x \in C_i} \|x - \mu_i\|^2$$

3. We can rewrite the objective function as $G(C_1, \ldots, C_k) = \sum_{i=1}^{k} \sum_{x \in C_i} \|x - \mu(C_i)\|^2$

4. Define with $C_i^{(t)}$ the *i-th* cluster at time *t* and with $\mu_i^{(t)}$ the value of $\mu(C_i)$ at step t

5. *Centroid computation:* The new centroids minimize the distance w.r.t the points in the cluster

$$G\left(C_1^{(t)}, \ldots, C_k^{(t)}\right) = \sum_{i=1}^{k} \sum_{x \in C_i^{(t)}} \left\|x - \mu_i^{(t)}\right\|^2 \leq \sum_{i=1}^{k} \sum_{x \in C_i^{(t)}} \left\|x - \mu_i^{(t-1)}\right\|^2$$

6. *Points allocation:* Each point is assigned to the closest centroid

$$\sum_{i=1}^{k} \sum_{x \in C_i^{(t)}} \left\|x - \mu_i^{(t-1)}\right\|^2 \leq \sum_{i=1}^{k} \sum_{x \in C_i^{(t-1)}} \left\|x - \mu_i^{(t-1)}\right\|^2$$

7. By placing all together:

From 5

From 6

$$G\left(C_1^{(t)}, \ldots, C_k^{(t)}\right) = \sum_{i=1}^{k} \sum_{x \in C_i^{(t)}} \left\|x - \mu_i^{(t)}\right\|^2 \leq \sum_{i=1}^{k} \sum_{x \in C_i^{(t)}} \left\|x - \mu_i^{(t-1)}\right\|^2 \leq \sum_{i=1}^{k} \sum_{x \in C_i^{(t-1)}} \left\|x - \mu_i^{(t-1)}\right\|^2 = G\left(C_1^{(t-1)}, \ldots, C_k^{(t-1)}\right)$$

*Note: monotonic not decreasing, but no guarantees on # iterations to converge and could fall in local min.*

1. The centroids positions and allocations do not change any more
2. Error improvement below threshold in 2 consecutive iterations ($\Delta G_{km} < T_1$)
3. Maximum number of iterations
4. Reached a target value for $G$ ($G < T_2$)

Complexity:

- Assignment of $m$ points in $\mathbb{R}^n$ to $k$ clusters : time $O(kmn)$
- Computation of centers: time $O(mn)$
- If convergence after $t$ iterations: $O(tkmn)$
- In practice convergence after a few iterations but can be very long on some critical cases

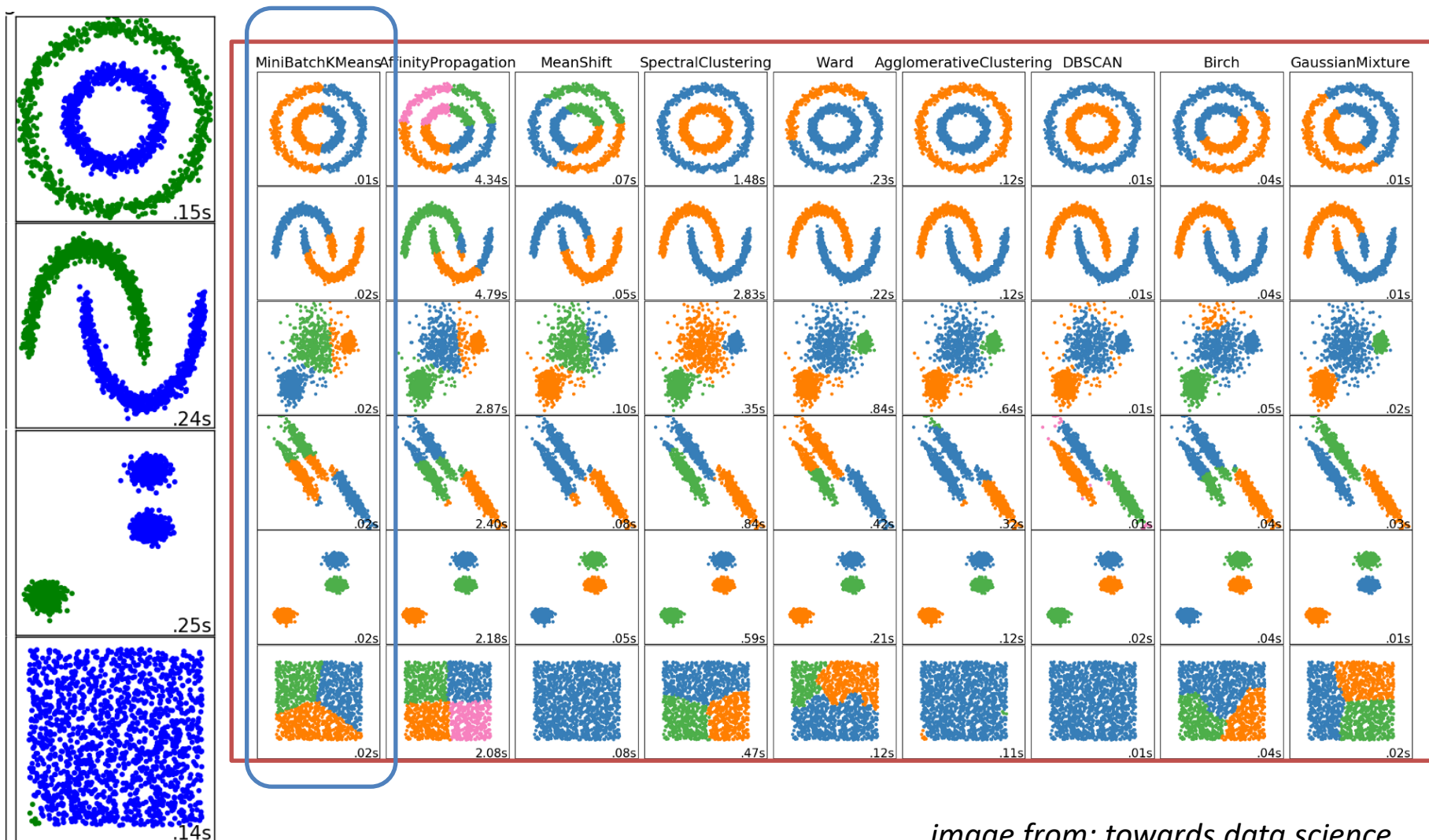*Notation: m: #samples, k: # clusters, n: dimensionality of data*

*Pros*

- ☑ Fast and simple
  - ○ True in practice, in theory it is a NP-hard problem
- ☑ Always converges and typically also very fast
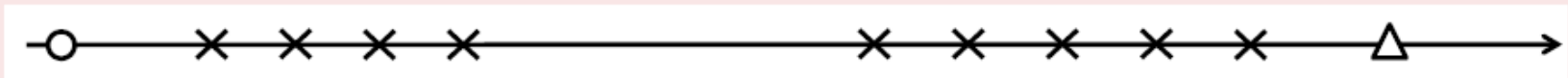
*Cons*

- ☒ It does not guarantee an optimal solution
- ☒ The solution depends on the initial centroids
- ☒ K must be known a priori
- ☒ Forces spherical symmetry of clusters (in the $n$-dimensional space)

# Examples



*image from: towards data science*

# Example/Exercise

Draw (approximately) the solution (clusters and centers) found by Lloyd algorithm for the 2 clusters ($K = 2$) problem, when the data ($x_i \in \mathbb{R}$) are the crosses in the figure below and the algorithm is initialised with center values indicated with the circle ($\circ$, cluster 1) and triangle ($\triangle$, cluster 2) shown in the figure.



The lab of Friday 1/12 will be on clustering and the K-Means algorithm

# Linkage-based Clustering

General class of algorithms that follow the general scheme below

1. Start from the trivial clustering: each data sample/point is a (single-point) cluster

2. Until "*termination condition*": repeatedly merge the "*closest*" clusters of the previous clustering

Two "*parameters*":

1. How to define distance *D(A,B)* between two clusters *A* and *B*

   o *Need cluster-to-cluster distance (not point-to-point)*
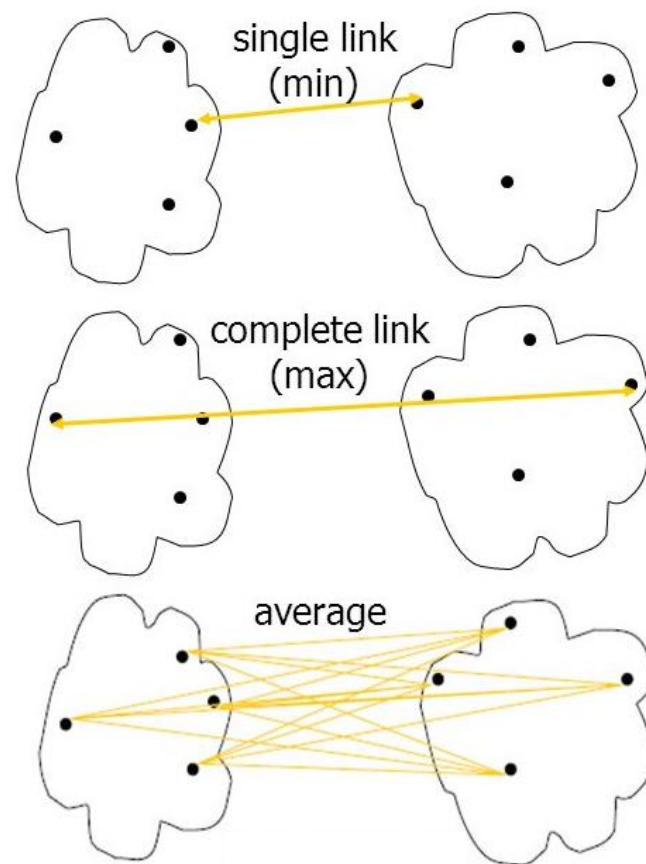
2. Termination condition

Different distances $D(A, B)$ between two clusters $A$ and $B$ can be used, resulting into different linkage methods:

- **single linkage**: $D(A, B) = \min\{d(\mathbf{x}, \mathbf{x}') : \mathbf{x} \in A, \mathbf{x}' \in B\}$
- **average linkage**: $D(A, B) = \frac{1}{|A||B|} \sum_{\mathbf{x} \in A, \mathbf{x}' \in B} d(\mathbf{x}, \mathbf{x}')$
- **max linkage**: $D(A, B) = \max\{d(\mathbf{x}, \mathbf{x}') : \mathbf{x} \in A, \mathbf{x}' \in B\}$
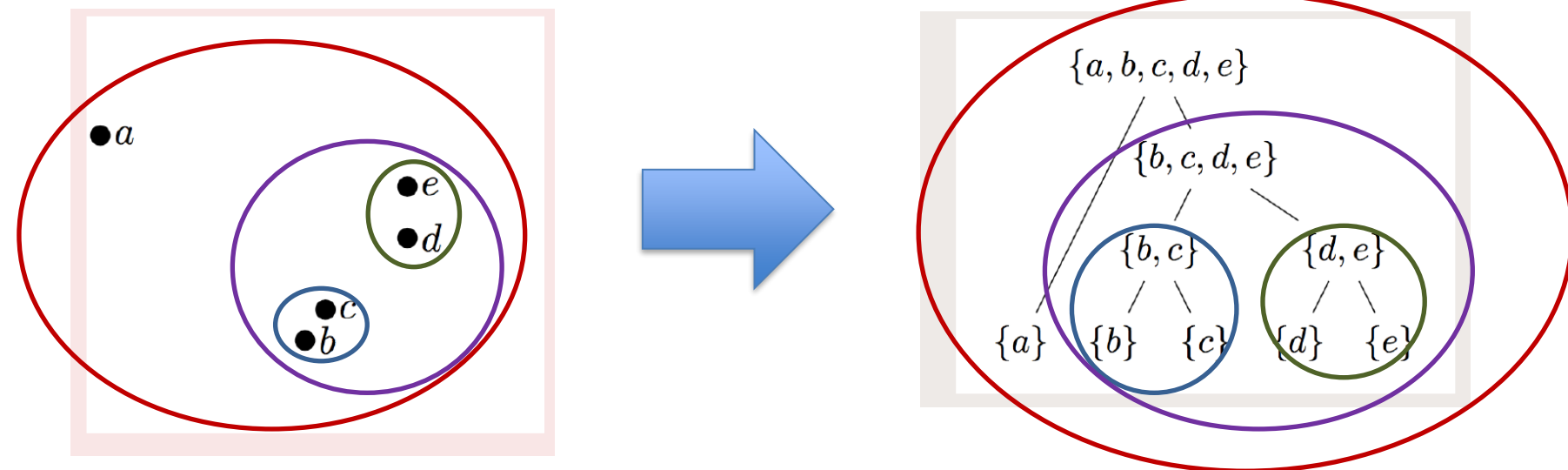
Common termination condition:

- data points are partitioned into $k$ clusters
- minimum distance between pairs of clusters is $> r$, where $r$ is a parameter provided in input
- all points are in a cluster $\Rightarrow$ output is a dendrogram

See next slide



single link (min)

complete link (max)

average

image from Univ. Manchester

- ❑ Single linkage (use minimum distance between points in the cluster)
- ❑ End when all points are in a single cluster → output is a dendogram
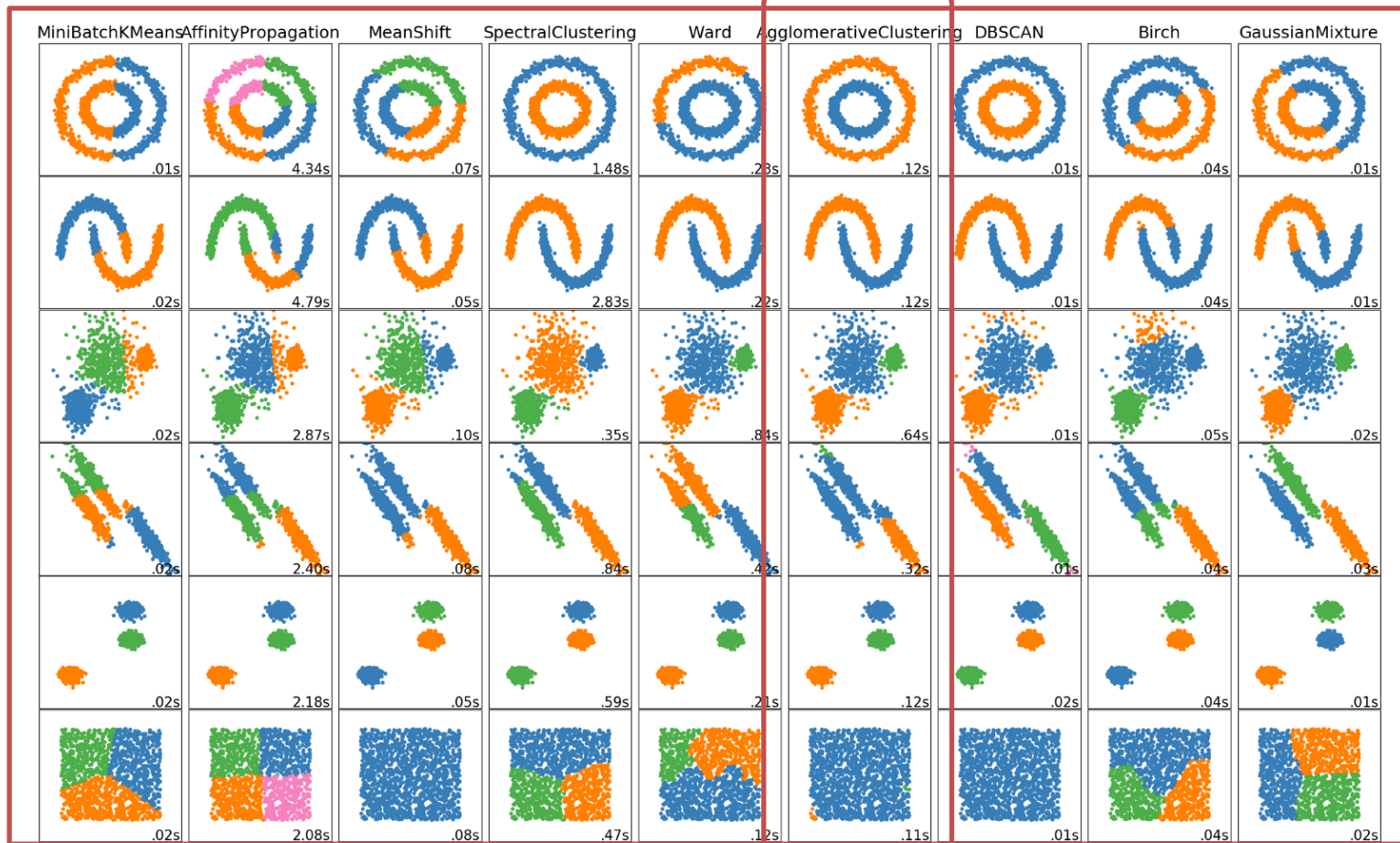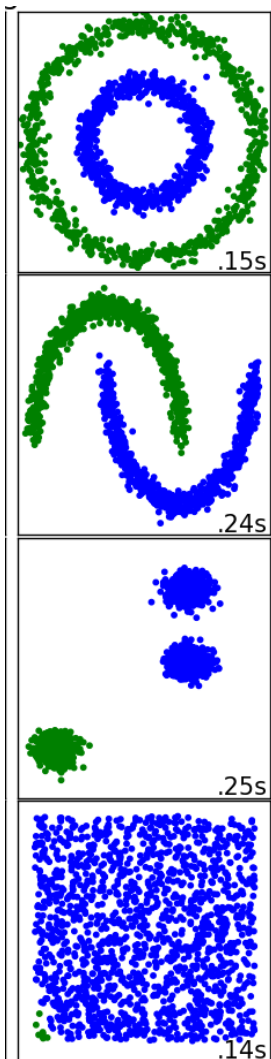  - ➢ from the dendogram various clusterings can be extracted

# Examples



*image from: towards data science*