# Lecture 20 - Mar 18

Clustering
Expectation-Maximization (EM) Algorithm

<u>References</u>

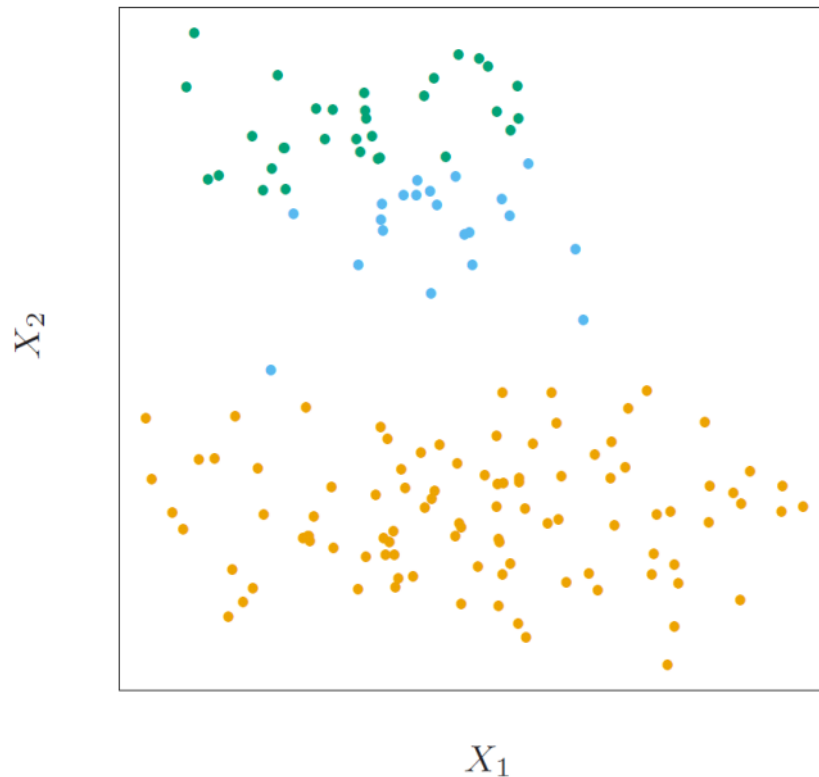*Data Mining and Machine Learning*
Ch 13 - Representation-based Clustering

*Elements of Statistical Learning*
14.3.6 K-means
14.3.7 Gaussian Mixtures as Soft K-Means Clustering

AKA data segmentation has a variety of goals.

All group or segment objects into "clusters" where points inside are similar but different from points in other clusters
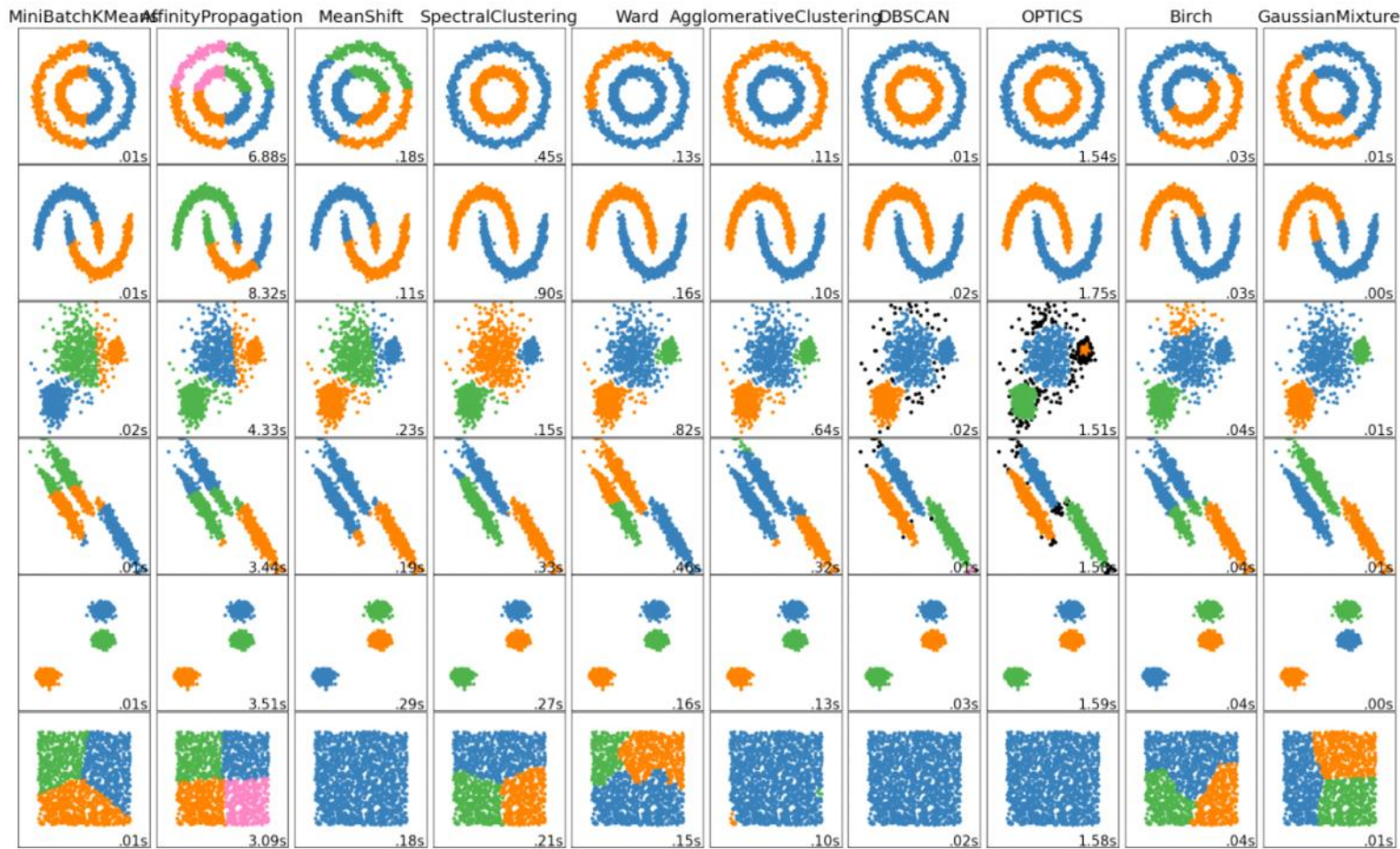


- Sometimes the goal is to find a natural hierarchy of clusters

- Also used for descriptive statistics regarding whether or not there exist distinct subgroups with substantially different properties

# Clustering Algorithms



A comparison of the clustering algorithms in scikit-learn

https://scikit-learn.org/stable/modules/clustering.html

Types: Combinatorial, Mixture modeling, mode seeking

work directly on
the observed data

assumes data is an i.i.d.
sample from a population
described by a pdf
characterized by a
parameterized model
taken as a mixture
of component density
functions: each describing
one cluster

"bump hunters" attempt to
directly estimate distinct
modes of the pdf as
centers of clusters

Let $X \in \mathbb{R}^{n \times d}$ be a dataset of points $x_i \in \mathbb{R}^d$

Representative-based clustering partitions $X$ into $k$ clusters...

$$C = \{C_1, \ldots, C_k\}$$

For each cluster $C_i$, there exists a representative point that summarizes the cluster

$\hookrightarrow$ e.g. the Centroid (mean) $M_i$ of points in $C_i$

$$M_i = \frac{1}{n_i} \sum_{x_j \in C_i} x_j \quad \text{where} \quad n_i = |C_i|$$

Two common methods: K-means + Expectation-Maximization (EM) algorithms

For a given clustering $C$, we need a scoring function to evaluate it

The SSE scoring function is

$$SSE(C) = \sum_{i=1}^{k} \sum_{x_j \in C_i} \|x_j - m_i\|^2$$

Sum over all clusters

SS of difference between cluster pts + cluster centroid

$$C^* = \arg\min_{C} \left\{ SSE(C) \right\}$$

K-means minimizs SSE via a greedy iterative approach (risks reaching local min)

K-means initializes centroids randomly + iterates 2 steps
① Assign points to nearest centroid's cluster
② Compute centroid as mean of points in each cluster

⌐ ② Compute Centroid as mean of points in each ⸳⸳⸳

Stop when centroids stop moving.

Often, people run K-means several times due to random initialization + keep the best result (acct. SSE)

K-means creates Convex-shaped clusters.

## Algorithm 13.1

Input: $X, k, \varepsilon$

often, sample uniformly in each feature

1. [Randomly initialize] $k$ centroids $\mu_1^0, \ldots, \mu_k^0 \in \mathbb{R}^d$

while $\sum_{i=1}^{k} \| \mu_i^t - \mu_i^{t-1} \| > \varepsilon$ :

while centroids move substantially in re-calculations

    for $x_j \in X$:

        $i^* = \arg\min_i \{ \| x_j - \mu_i^{t-1} \|^2 \}$

        $C_{i^*} = C_{i^*} \cup \{ x_j \}$

assign points to cluster with nearest centroid

    for $C_i \in C$

        $\mu_i^t = \frac{1}{n_i} \sum_{x_j \in C_i} x_j$

re-calculate centroids with new clusters

# K-Means in 1D



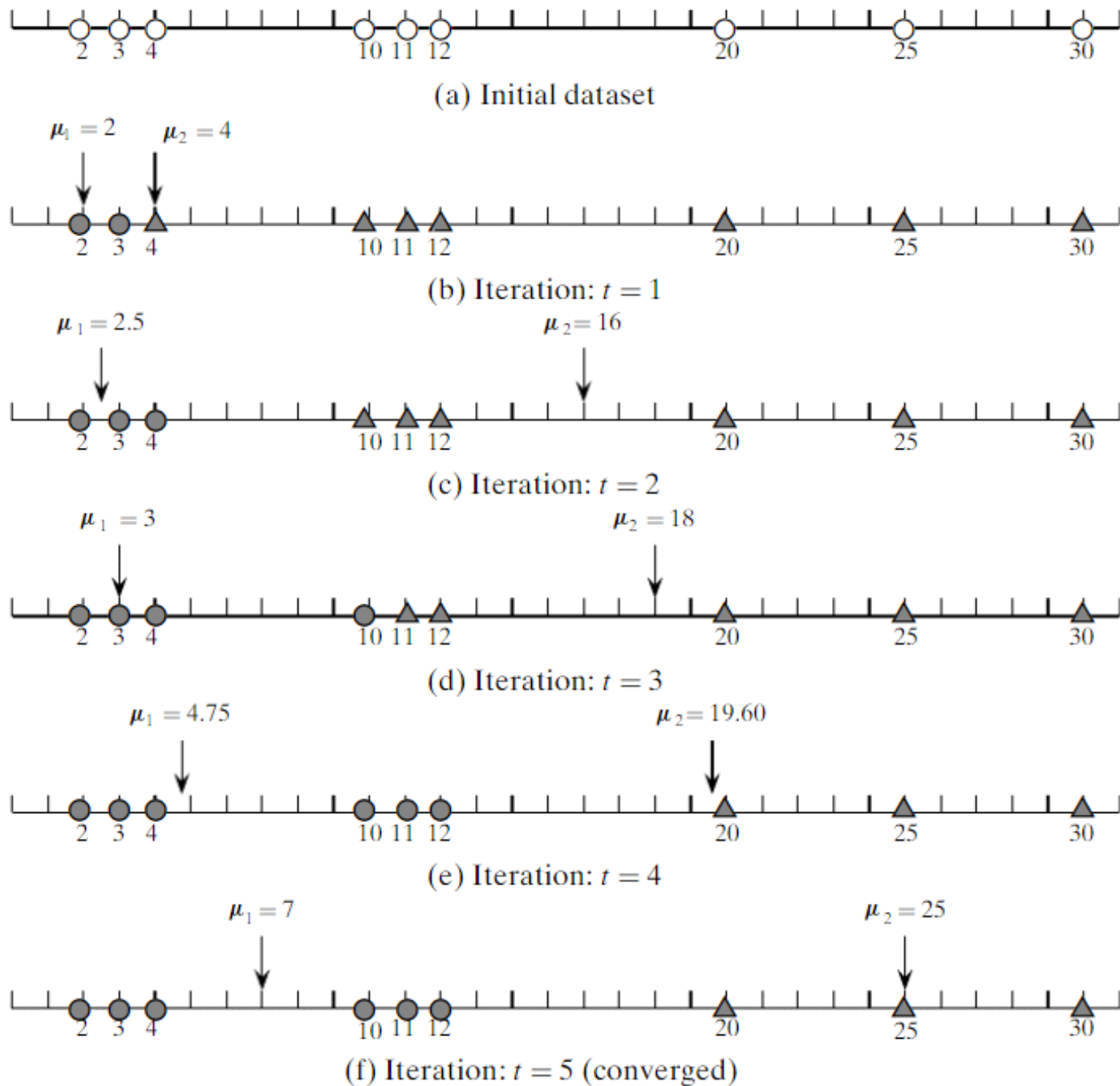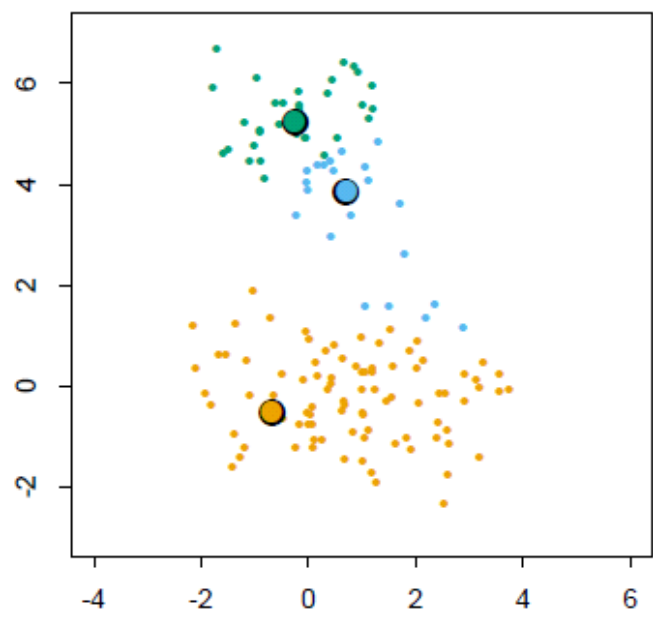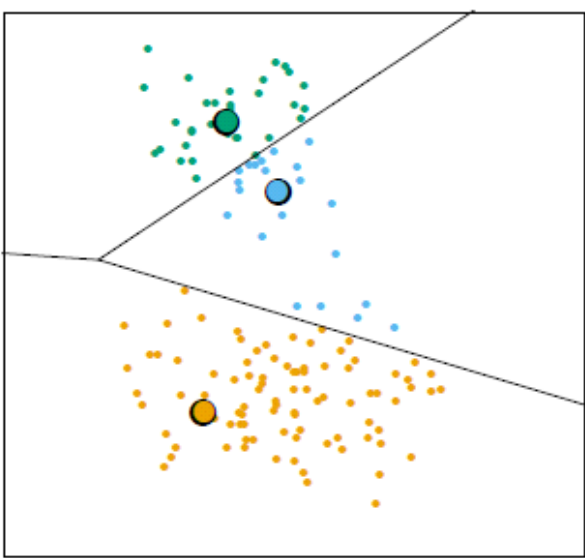(a) Initial dataset

(b) Iteration: $t = 1$

$\mu_1 = 2$ $\mu_2 = 4$

(c) Iteration: $t = 2$

$\mu_1 = 2.5$ $\mu_2 = 16$

(d) Iteration: $t = 3$

$\mu_1 = 3$ $\mu_2 = 18$

(e) Iteration: $t = 4$

$\mu_1 = 4.75$ $\mu_2 = 19.60$

(f) Iteration: $t = 5$ (converged)

$\mu_1 = 7$ $\mu_2 = 25$

Figure 13.1.  K-means in one dimension.

# K-Means in 2D