

Representative based clustering

Chapter 13

Clustering: what and why?

- Clustering is a method of partitioning a set of data instances into k different groups, such that similar data instances belong to the same group
- For clustering task, the label of data instances are not available, so clustering is also called unsupervised data analysis
- Clustering is important because it provides an insight about the data by finding subgroups of similar data points.
 - After clustering, the properties of each group can be studied independently
- In scenarios, where we are mostly interested in classification, we sometimes perform clustering because the labeling of data instances is costly.
- Clustering methods:
 - Representative Based clustering
 - Density based Clustering
 - Graph or spectral clustering

Representative Based Clustering

- Given
 - a dataset of n points in a d -dimensional space, $D = \{\mathbf{x}_i\}_{i=1}^n$
 - the number of desired clusters k ,
- Goal of representative based clustering is
 - partition the data into k groups $C = \{C_1, C_2, \dots, C_k\}$ so that similar objects are in the same clusters
 - For each cluster, C_i , find a representative data objects (say, $\boldsymbol{\mu}_i$)
 - We typically try to minimize the following objective function, called sum of square errors
$$SSE(C) = \sum_{i=1}^k \sum_{\mathbf{x}_j \in C_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2, \mu_i = \frac{1}{|C_i|} \sum_{\mathbf{x}_j \in C_i} \mathbf{x}_j$$
- This is a combinatorial optimization problem, which is difficult to solve
- Brute-force solution is also not possible, because of exponentially many partitioning.
 - The number of ways n objects can be partitioned into k groups is $S(n, k) = \frac{1}{k!} \sum_{t=0}^k (-1)^t \binom{k}{t} (k-t)^n$
 - Brute-force method is not possible

k -means Algorithm

- A greedy iterative algorithm for representative based clustering
- Very efficient, if the dataset are points in \mathbb{R}^d
- Does not obtain global optimal solution, but provides a local optimal guaranty
 - Reassigning the cluster membership of exactly one data instance (keeping the membership of the remaining data instances unchanged), does not yield a better solution
- Algorithm
 1. Randomly choose k points in the data space as initial cluster means (also called cluster seeds)
 2. Assign each data points to the closest seed, points assigned to the same seed makes a cluster
 3. Re-compute a new mean for each cluster, as the centroid of the data points in that cluster
 4. Repeat between step 2 and 3 until the convergence is reached

k -means Pseudocode

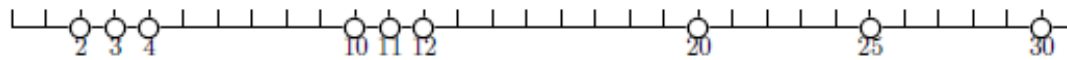
Algorithm 13.1: K-means Algorithm

K-MEANS (D, k, ϵ):

```
1  $t = 0$ 
2 Randomly initialize  $k$  centroids:  $\mu_1^t, \mu_2^t, \dots, \mu_k^t$ 
3 repeat
4    $t = t + 1$ 
5   // Cluster Assignment Step
6   foreach  $x_j \in D$  do
7      $j^* = \arg \min_i \{\|x_j - \mu_i^t\|^2\}$  // Assign  $x_j$  to closest centroid
8      $C_{j^*} = C_{j^*} \cup \{x_j\}$ 
9   // Centroid Update Step
10  foreach  $i = 1$  to  $k$  do
11     $\mu_i^t = \frac{1}{|C_i|} \sum_{x_j \in C_i} x_j$ 
12 until  $\sum_{i=1}^k \|\mu_i^t - \mu_i^{t-1}\| \leq \epsilon$ 
```

Complexity: is $O(nkdt)$, n : data points, d : dimension size, k : number of clusters, t : number of iterations

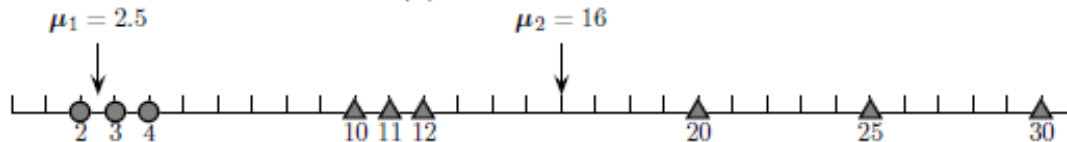
Example (1-D)



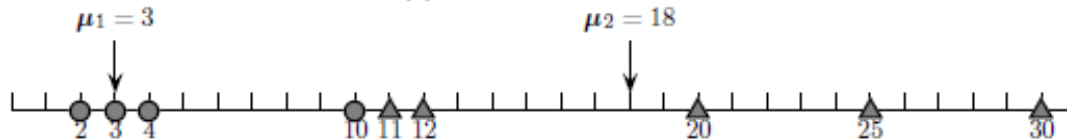
(a) Initial dataset



(b) Iteration: $t = 1$



(c) Iteration: $t = 2$



(d) Iteration: $t = 3$

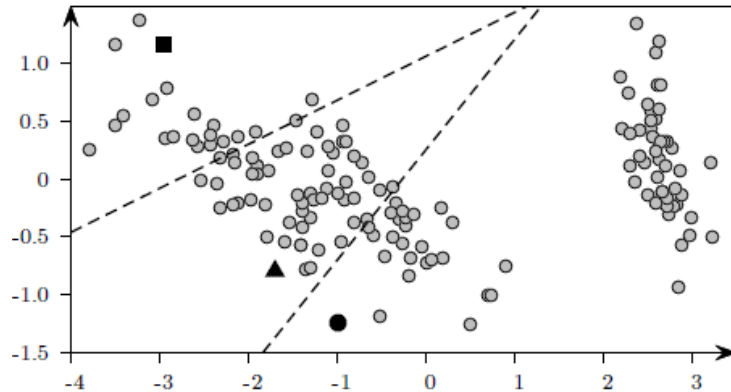


(e) Iteration: $t = 4$

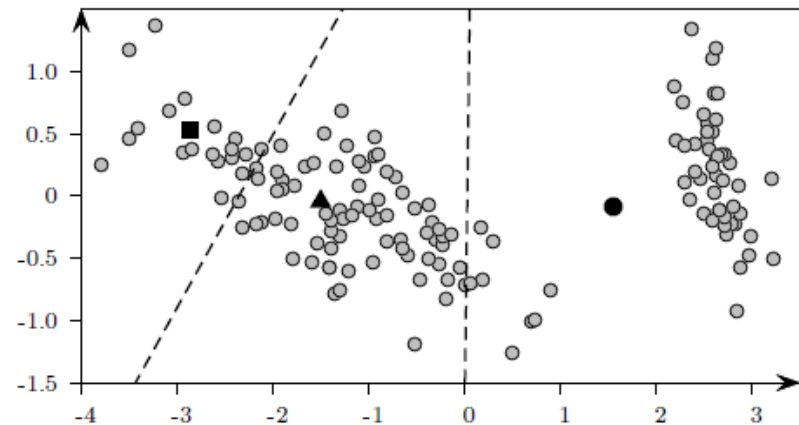


(f) Iteration: $t = 5$ (converged)

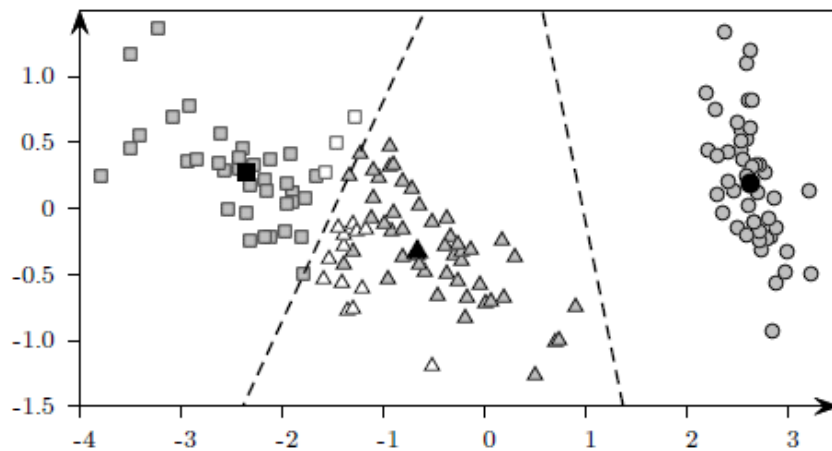
Example (2 –Dimension)



(a) Random Initialization: $t = 0$



(b) Iteration: $t = 1$



(c) Iteration: $t = 8$ (converged)

Expectation Maximization (EM)

Clustering

- The k -means approach is an example of hard assignment, where each point can belong to only one cluster
- EM clustering allow soft cluster assignment, where each point has a distinct probability of belonging to each of the clusters
- EM assumes that a cluster is modeled as a multivariate normal distribution with parameters μ_i and Σ_i .
- It also assumes that each of the data points is generated by a mixture of k clusters, denoted by $P(C_i)$ such that $\sum_{i=1}^k P(C_i) = 1$
- Thus solving the clustering is to find the model parameters $\theta = \{\mu_1, \Sigma_1, P(C_1), \dots, \mu_k, \Sigma_k, P(C_k)\}$
- Assuming each column as a random variable X_j , and $X = (X_1, X_2, \dots, X_d)$ is a vector-random variable. A data point, x can be assumed to be as an instance of the variable X
- Thus the probability density function of the dataset X is given as a mixture model of the k cluster normals, i.e., $f(x) = \sum_{i=1}^k f_i(x)P(C_i) = \sum_{i=1}^k f(x|\mu_i, \Sigma_i)P(C_i)$

Maximum Likelihood Estimation (MLE)

- For a parametric model, MLE process chooses the parameter θ that maximizes the likelihood ($P(D|\theta)$) of the observed data. $\theta^* = \arg \max_{\theta} P(D|\theta)$
- Since each of the n points are assumed to be chosen from an iid distribution, $P(D|\theta) = \prod_{j=1}^n f(x_j)$. Instead of maximizing the likelihood, we can also maximize the log-likelihood, $\log P(D|\theta) = \log \prod_{j=1}^n f(x_j) = \sum_{j=1}^n \log(\sum_{i=1}^k f(x_j|\mu_i, \Sigma_i)P(C_i))$
- Directly maximizing log-likelihood is a difficult task, so we can use the Expectation-Maximization (EM) approach for finding the MLE estimate of θ , which is θ^* .
- EM starts from an initial guess of θ and improve θ iteratively.
 - Expectation step computes the posterior probability $P(C_i|x_j) = \frac{P(x_j|C_i) \cdot P(C_i)}{\sum_{a=1}^k P(x_j|C_a) \cdot P(C_a)} = \frac{f(x_j|\pi_i, \Sigma_i) \cdot P(C_i)}{\sum_{a=1}^k f(x_j|\pi_a, \Sigma_a) \cdot P(C_a)}$
 - Maximization step uses the weights $P(C_i|x_j)$ to re-estimate θ
 - The process continue until a convergence criteria is met

EM in one Dimension

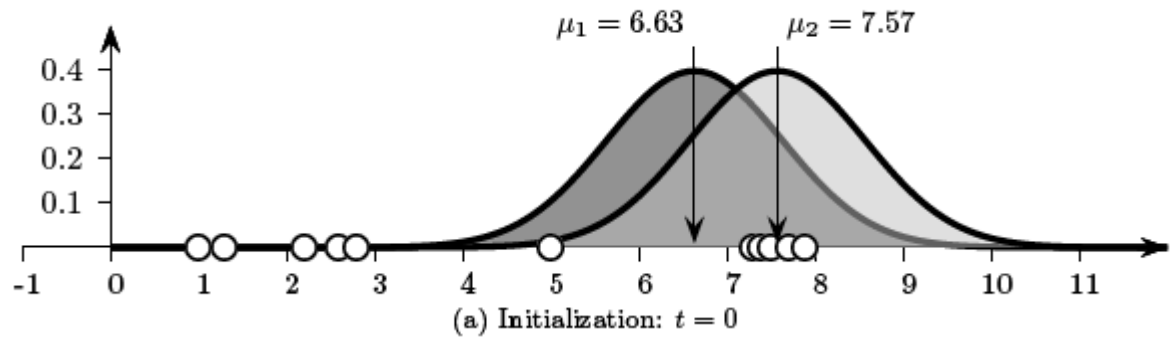
- Consider D consists of points in one dimension, $f_i(x) = f(x|\mu_i, \sigma_i^2) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left\{-\frac{(x - \mu_i)^2}{2\sigma_i^2}\right\}$
- Initialization: For each cluster C_i we initialize μ_i, σ_i^2 , and $P(C_i)$ randomly. Mean is selected uniformly at random from the range of possible values, σ_i^2 is typically assumed to be 1, and $P(C_i) = 1/k$
- Expectation step: This step computes the posterior probabilities: $P(C_i|x_j) = \frac{f(x_j|\mu_i, \sigma_i^2) \cdot P(C_i)}{\sum_{a=1}^k f(x_j|\mu_a, \sigma_a^2) \cdot P(C_a)}$
- For convenience, we assume that $P(C_i|x_j) = w_{ij}$. Also assume that $w_i = (w_{i1}, w_{i2}, \dots, w_{in})^T$
- The new estimate of cluster mean is simply the weighted mean of all the points: $\mu_i = \frac{\sum_{j=1}^n w_{ij} \cdot x_j}{\sum_{j=1}^n w_{ij}}$
- The new estimate of the cluster variance is the weighted variance across all the points

$$\sigma_i^2 = \frac{\sum_{j=1}^n w_{ij} (x_j - \mu_i)^2}{\sum_{j=1}^n w_{ij}}$$

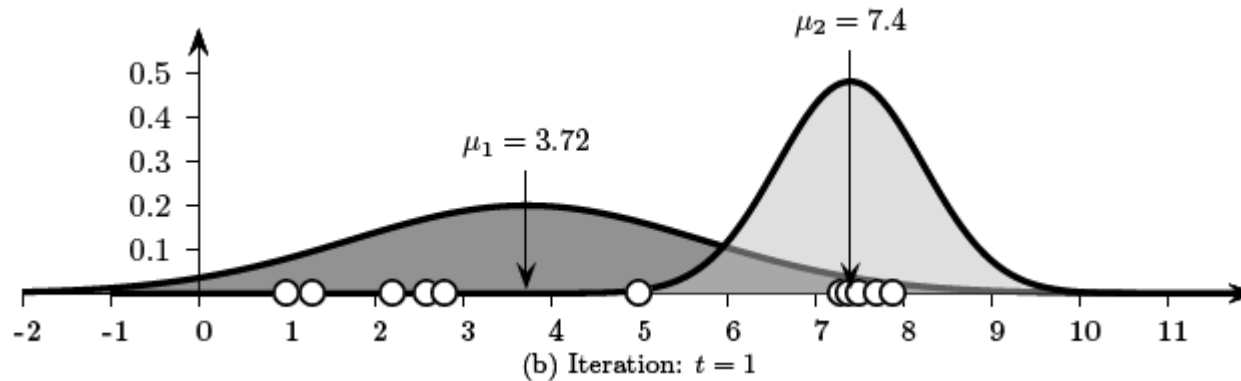
- Finally the probability $P(C_i)$ is simply the fraction of total weight belonging to the cluster C_i

$$P(C_i) = \frac{\sum_{j=1}^n w_{ij}}{\sum_{a=1}^k \sum_{j=1}^n w_{aj}} = \frac{\sum_{j=1}^n w_{ij}}{\sum_{j=1}^n 1} = \frac{\sum_{j=1}^n w_{ij}}{n}$$

Example

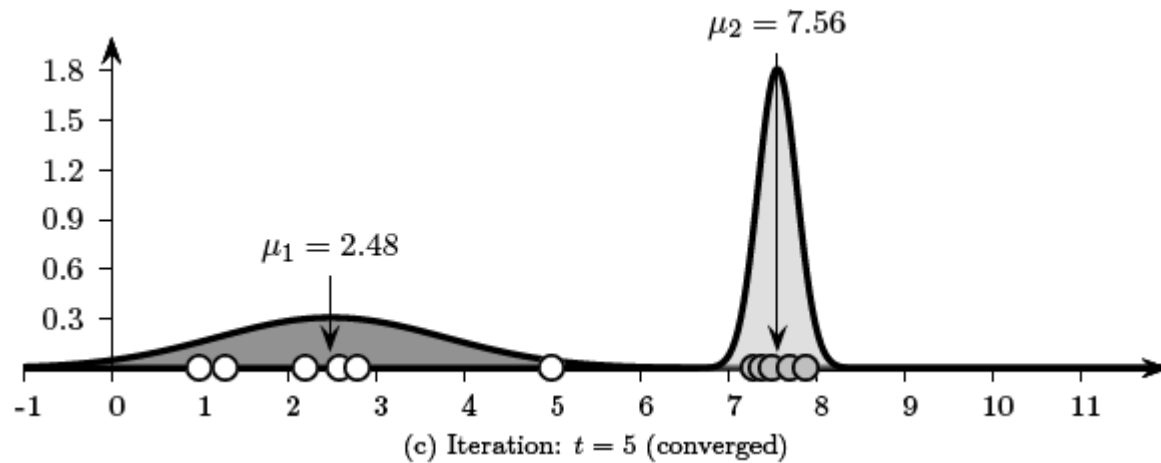


$\mu_1 = 6.63$	$\sigma_1^2 = 1$	$P(C_2) = 0.5$
$\mu_2 = 7.57$	$\sigma_2^2 = 1$	$P(C_2) = 0.5$



$\mu_1 = 3.72$	$\sigma_1^2 = 6.13$	$P(C_1) = 0.71$
$\mu_2 = 7.4$	$\sigma_2^2 = 0.69$	$P(C_2) = 0.29$

Example



$$\mu_1 = 2.48$$

$$\sigma_1^2 = 1.69$$

$$P(C_1) = 0.55$$

$$\mu_2 = 7.56$$

$$\sigma_2^2 = 0.05$$

$$P(C_2) = 0.45$$

EM algorithm

Algorithm 13.3: Expectation-Maximization (EM) Algorithm

EXPECTATION-MAXIMIZATION (\mathbf{D}, k, ϵ):

```
1  $t = 0$ 
   // Random Initialization
2 Randomly initialize  $\mu_1^t, \dots, \mu_k^t$ 
3  $\Sigma_i^t = \mathbf{I}, \forall i = 1, \dots, k$ 
4  $P^t(C_i) = \frac{1}{k}, \forall i = 1, \dots, k$ 
5 repeat
6    $t = t + 1$ 
   // Expectation Step
7   for  $i = 1, \dots, k$  and  $j = 1, \dots, n$  do
8      $w_{ij}^t = P^t(C_i | \mathbf{x}_j) = \frac{f(\mathbf{x}_j | \mu_i, \Sigma_i) \cdot P(C_i)}{\sum_{a=1}^k f(\mathbf{x}_j | \mu_a, \Sigma_a) \cdot P(C_a)}$ 
   // Maximization Step
9   for  $i = 1, \dots, k$  do
10     $\mu_i^t = \frac{\sum_{j=1}^n w_{ij} \cdot \mathbf{x}_j}{\sum_{j=1}^n w_{ij}}$ 
11     $\Sigma_i^t = \frac{\sum_{j=1}^n w_{ij} (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^T}{\sum_{j=1}^n w_{ij}}$ 
12     $P^t(C_i) = \frac{\sum_{j=1}^n w_{ij}}{n}$ 
13 until  $\sum_{i=1}^k \|\mu_i^t - \mu_i^{t-1}\| \leq \epsilon$ 
```

EM in d-dimension

- Like 1-D, we need to estimate the parameters, $\theta = \{\mu_i, \Sigma_i, P(C_i)\}_{i=1..k}$

$$\Sigma_i = \begin{pmatrix} (\sigma_1^i)^2 & \sigma_{12}^i & \dots & \sigma_{1d}^i \\ \sigma_{21}^i & (\sigma_2^i)^2 & \dots & \sigma_{2d}^i \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1}^i & \sigma_{d2}^i & \dots & (\sigma_d^i)^2 \end{pmatrix} \Rightarrow \Sigma_i = \begin{pmatrix} (\sigma_1^i)^2 & 0 & \dots & 0 \\ 0 & (\sigma_2^i)^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & (\sigma_d^i)^2 \end{pmatrix}$$

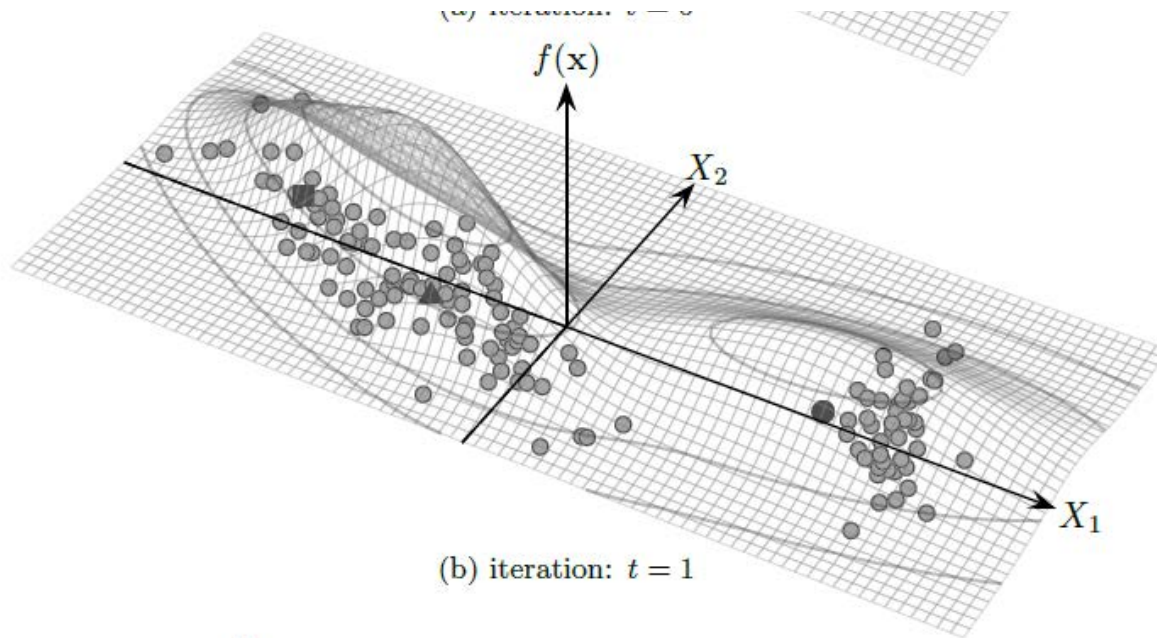
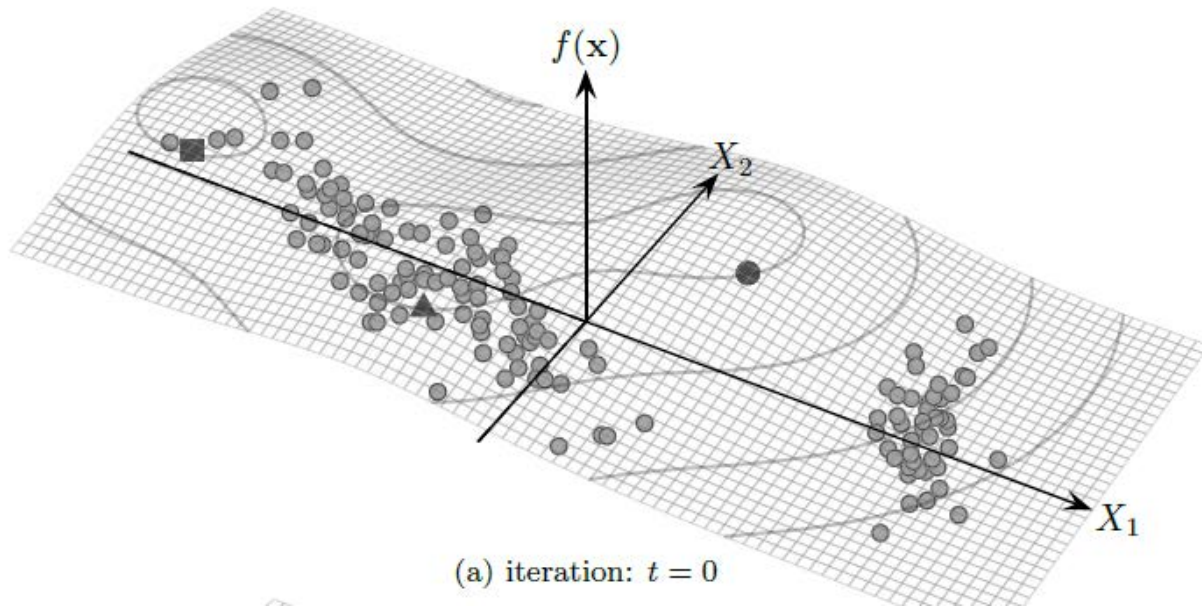
Update rules

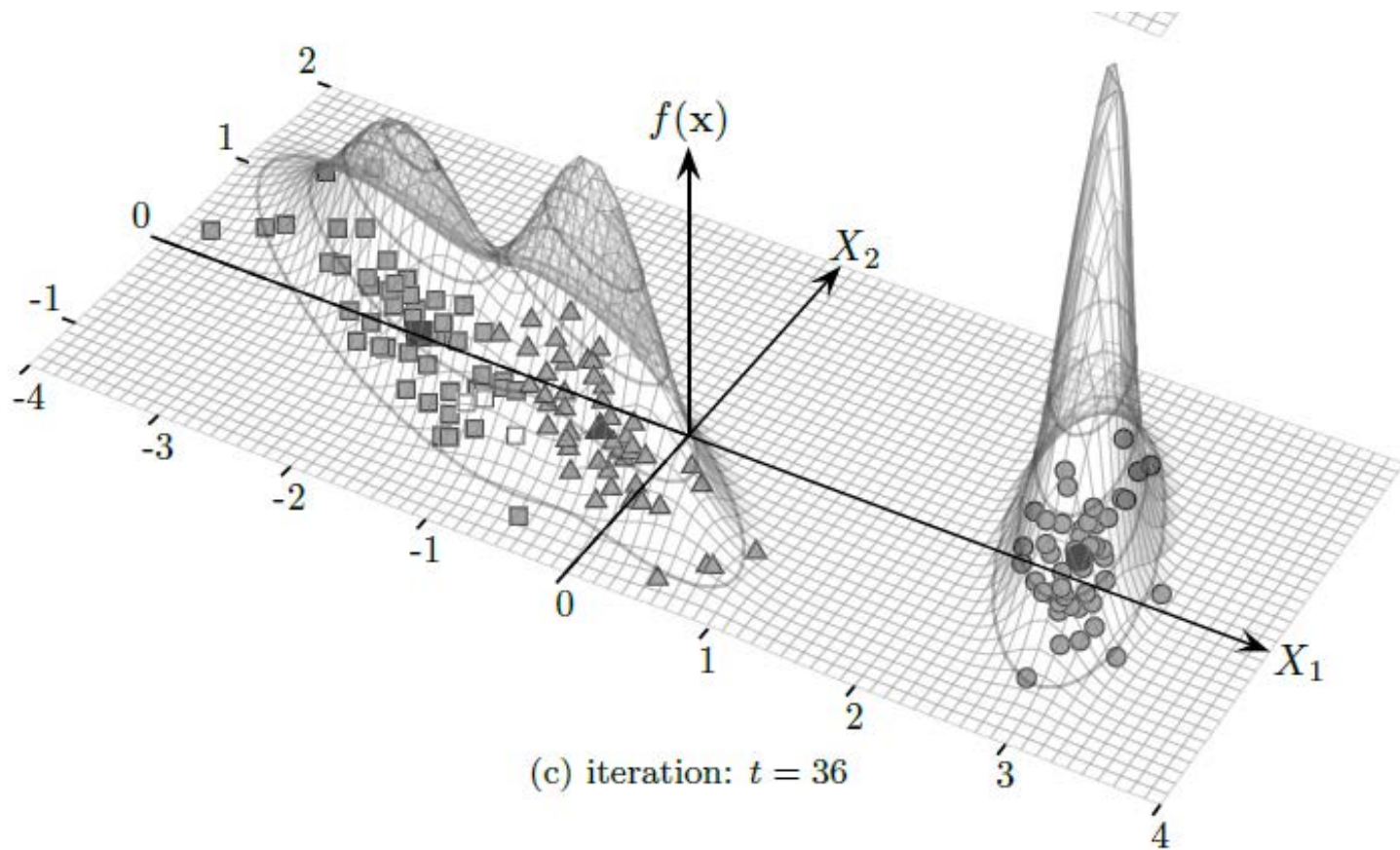
$$\mu_i = \frac{\sum_{j=1}^n w_{ij} \cdot \mathbf{x}_j}{\sum_{j=1}^n w_{ij}} \quad (13.11)$$

$$\Sigma_i = \frac{\sum_{j=1}^n w_{ij} \mathbf{z}_{ji} \mathbf{z}_{ji}^T}{\mathbf{w}_i^T \mathbf{1}} \quad (13.12)$$

$$P(C_i) = \frac{\sum_{j=1}^n w_{ij}}{n} = \frac{\mathbf{w}_i^T \mathbf{1}}{n}$$

Example





K-Means is a special case of EM

- K-Means is obtained from the EM algorithm as follows:

$$P(\mathbf{x}_j|C_i) = \begin{cases} 1 & \text{if } C_i = \arg \min_a \{ \|\mathbf{x}_j - \boldsymbol{\mu}_a\|^2 \} \\ 0 & \text{otherwise} \end{cases}$$

$$P(C_i|\mathbf{x}_j) = \begin{cases} 1 & \text{if } \mathbf{x}_j \in C_i, \text{ i.e., if } C_i = \arg \min_{C_a} \{ \|\mathbf{x}_j - \boldsymbol{\mu}_a\|^2 \} \\ 0 & \text{otherwise} \end{cases}$$