# Machine Learning

## B. Unsupervised Learning
## B.1 Cluster Analysis

Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL)
Institute for Computer Science
University of Hildesheim, Germany

# Outline

1. k-means & k-medoids

2. Gaussian Mixture Models

3. Hierarchical Cluster Analysis

# Syllabus

# Outline

## 1. k-means & k-medoids

## 2. Gaussian Mixture Models

## 3. Hierarchical Cluster Analysis

# Partitions

Let $X$ be a set. A set $P \subseteq \mathcal{P}(X)$ of subsets of $X$ is called
a **partition of** $X$ if the subsets

   1. are pairwise disjoint:                   $A \cap B = \emptyset, \quad A, B \in P, A \neq B$

   2. cover $X$:                             $\displaystyle\bigcup_{A \in P} A = X$, and

   3. do not contain the empty set:    $\emptyset \notin P$.

# Partitions

Let $X := \{x_1, \ldots, x_N\}$ be a finite set. A set $P := \{X_1, \ldots, X_K\}$ of subsets $X_k \subseteq X$ is called a **partition of** $X$ if the subsets

1. are **pairwise disjoint**:          $X_k \cap X_j = \emptyset, \quad k, j \in \{1, \ldots, K\}, k \neq j$

2. **cover** $X$:                                  $\displaystyle\bigcup_{k=1}^{K} X_k = X,$ and

3. do **not contain the empty set**:   $X_k \neq \emptyset, \quad k \in \{1, \ldots, K\}.$

The sets $X_k$ are also called **clusters**, a partition $P$ a **clustering**.
$K \in \mathbb{N}$ is called **number of clusters**.

$\text{Part}(X)$ denotes the set of all partitions of $X$.

# Partitions

Let $X$ be a finite set. A **surjective** function

$$p : \{1, \ldots, |X|\} \to \{1, \ldots, K\}$$

is called a **partition function of $X$**.

The sets $X_k := p^{-1}(k)$ form a partition $P := \{X_1, \ldots, X_K\}$.

# Partitions

Let $X := \{x_1, \ldots, x_N\}$ be a finite set. A binary $N \times K$ matrix

$$P \in \{0, 1\}^{N \times K}$$

is called a **partition matrix of $X$** if it

1. is **row-stochastic**:
$$\sum_{k=1}^{K} P_{i,k} = 1, \quad i \in \{1, \ldots, N\}, k \in \{$$

2. does **not contain a zero column**: $X_{i,k} \neq (0, \ldots, 0)^T, \quad k \in \{1, \ldots, K\}.$

The sets $X_k := \{i \in \{1, \ldots, N\} \mid P_{i,k} = 1\}$ form a partition
$P := \{X_1, \ldots, X_K\}$.

$P_{.,k}$ is called **membership vector of class $k$**.

# The Cluster Analysis Problem

Given

- a set $\mathcal{X}$ called **data space**, e.g., $\mathcal{X} := \mathbb{R}^m$,
- a set $X \subseteq \mathcal{X}$ called **data**, and
- a function

$$D : \bigcup_{X \subseteq \mathcal{X}} \text{Part}(X) \to \mathbb{R}_0^+$$

  called **distortion measure** where $D(P)$ measures how bad a partition $P \in \text{Part}(X)$ for a data set $X \subseteq \mathcal{X}$ is,

find a partition $P = \{X_1, X_2, \ldots X_K\} \in \text{Part}(X)$ with minimal distortion $D(P)$.

# The Cluster Analysis Problem (given K)

Given

- a set $\mathcal{X}$ called **data space**, e.g., $\mathcal{X} := \mathbb{R}^m$,
- a set $X \subseteq \mathcal{X}$ called **data**,
- a function

$$D : \bigcup_{X \subseteq \mathcal{X}} \text{Part}(X) \to \mathbb{R}_0^+$$

  called **distortion measure** where $D(P)$ measures how bad a partition $P \in \text{Part}(X)$ for a data set $X \subseteq \mathcal{X}$ is, and

- a number $K \in \mathbb{N}$ of clusters,

find a partition $P = \{X_1, X_2, \ldots X_K\} \in \text{Part}_K(X)$ with $K$ clusters with minimal distortion $D(P)$.

# k-means: Distortion Sum of Distances to Cluster Centers

Sum of squared distances to cluster centers:

$$D(P) := \sum_{k=1}^{K} \sum_{\substack{i=1: \\ P_{i,k}=1}}^{n} ||x_i - \mu_k||^2$$

with

$$\mu_k := \text{mean } \{x_i \mid P_{i,k} = 1, i = 1, \ldots, n\}$$

# k-means: Distortion Sum of Distances to Cluster Centers

Sum of squared distances to cluster centers:

$$D(P) := \sum_{i=1}^{n} \sum_{k=1}^{K} P_{i,k} ||x_i - \mu_k||^2 = \sum_{k=1}^{K} \sum_{\substack{i=1: \\ P_{i,k}=1}}^{n} ||x_i - \mu_k||^2$$

with

$$\mu_k := \frac{\sum_{i=1}^{n} P_{i,k} x_i}{\sum_{i=1}^{n} P_{i,k}} = \text{mean } \{x_i \mid P_{i,k} = 1, i = 1, \dots, n\}$$

# k-means: Distortion Sum of Distances to Cluster Centers

Sum of squared distances to cluster centers:

$$D(P) := \sum_{i=1}^{n} \sum_{k=1}^{K} P_{i,k} ||x_i - \mu_k||^2 = \sum_{k=1}^{K} \sum_{\substack{i=1: \\ P_{i.k}=1}}^{n} ||x_i - \mu_k||^2$$

with

$$\mu_k := \text{mean } \{x_i \mid P_{i,k} = 1, i = 1, \ldots, n\}$$

Minimizing $D$ over partitions with varying number of clusters leads to singleton clustering with distortion 0; only the cluster analysis problem with given $K$ makes sense.

Minimizing $D$ is not easy as reassigning a point to a different cluster also shifts the cluster centers.

# k-means: Minimizing Distances to Cluster Centers

Add cluster centers $\mu$ as auxiliary optimization variables:

$$D(P, \mu) := \sum_{i=1}^{n} \sum_{k=1}^{K} P_{i,k} ||x_i - \mu_k||^2$$

# k-means: Minimizing Distances to Cluster Centers

Add cluster centers $\mu$ as auxiliary optimization variables:

$$D(P, \mu) := \sum_{i=1}^{n} \sum_{k=1}^{K} P_{i,k} ||x_i - \mu_k||^2$$

Block coordinate descent:

1. fix $\mu$, optimize $P \rightsquigarrow$ reassign data points to clusters:

$$P_{i,k} := \delta(k = \ell_i), \quad \ell_i := \underset{k \in \{1,...,K\}}{\arg \min} \; ||x_i - \mu_k||^2$$

# k-means: Minimizing Distances to Cluster Centers

Add cluster centers $\mu$ as auxiliary optimization variables:

$$D(P, \mu) := \sum_{i=1}^{n} \sum_{k=1}^{K} P_{i,k} ||x_i - \mu_k||^2$$

Block coordinate descent:

1. fix $\mu$, optimize $P$ ⤳ reassign data points to clusters:

$$P_{i,k} := \delta(k = \ell_i), \quad \ell_i := \underset{k \in \{1,...,K\}}{\arg \min} \; ||x_i - \mu_k||^2$$

2. fix $P$, optimize $\mu$ ⤳ recompute cluster centers:

$$\mu_k := \frac{\sum_{i=1}^{n} P_{i,k} x_i}{\sum_{i=1}^{n} P_{i,k}}$$

Iterate until partition is stable.

# k-means: Initialization

k-means is usually initialized by picking $K$ data points as cluster centers at random:

1. pick the first cluster center $\mu_1$ out of the data points at random and then

2. sequentially select the data point with the largest sum of distances to already choosen cluster centers as next cluster center

$$\mu_k := x_i, \quad i := \arg\max_{i \in \{1, \ldots, n\}} \sum_{\ell=1}^{k-1} ||x_i - \mu_\ell||^2, \quad k = 2, \ldots, K$$

# k-means: Initialization

k-means is usually initialized by picking $K$ data points as cluster centers at random:

1. pick the first cluster center $\mu_1$ out of the data points at random and then

2. sequentially select the data point with the largest sum of distances to already choosen cluster centers as next cluster center

$$\mu_k := x_i, \quad i := \underset{i \in \{1,...,n\}}{\arg \max} \sum_{\ell=1}^{k-1} ||x_i - \mu_\ell||^2, \quad k = 2, \ldots, K$$

Different initializations may lead to different local minima.

▶ run k-means with different random initializations and

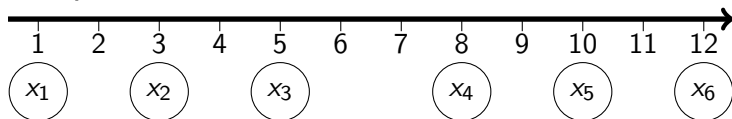▶ keep only the one with the smallest distortion (**random restarts**).

# k-means Algorithm

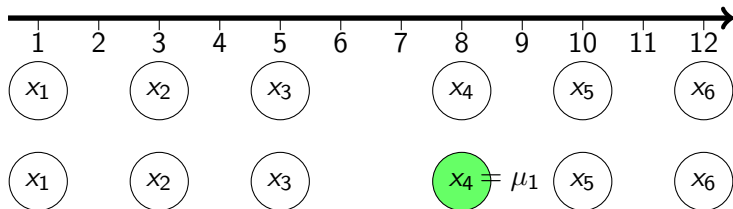1: **procedure** CLUSTER-KMEANS($\mathcal{D} := \{x_1, \ldots, x_N\} \subseteq \mathbb{R}^M, K \in \mathbb{N}, \epsilon \in \mathbb{R}^+$)
2:     $i_1 \sim \text{unif}(\{1, \ldots, N\}), \mu_1 := x_{i_1}$
3:     **for** $k := 2, \ldots, K$ **do**
4:        $i_k := \arg\max_{n \in \{1, \ldots, N\}} \sum_{\ell=1}^{k-1} ||x_n - \mu_\ell||, \mu_i := x_{i_k},$.
5:     **repeat**
6:        $\mu^{\text{old}} := \mu$
7:        **for** $n := 1, \ldots, N$ **do**
8:           $P_n := \arg\min_{k \in \{1, \ldots, K\}} ||x_n - \mu_k||$
9:        **for** $k := 1, \ldots, K$ **do**
10:        $\mu_k := \text{mean} \{x_n \mid P_n = k\}$
11:     **until** $\frac{1}{K} \sum_{k=1}^{K} ||\mu_k - \mu_k^{\text{old}}|| < \epsilon$
12:     **return** $P$

Note: In implementations, the two loops over the data (lines 6 and 9) can be combined in one loop.
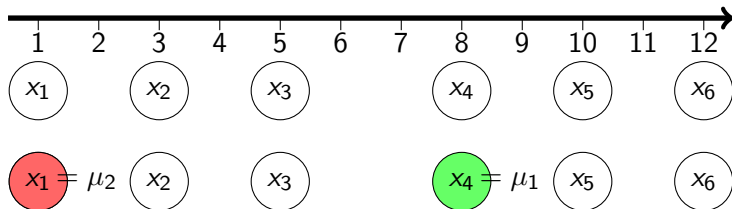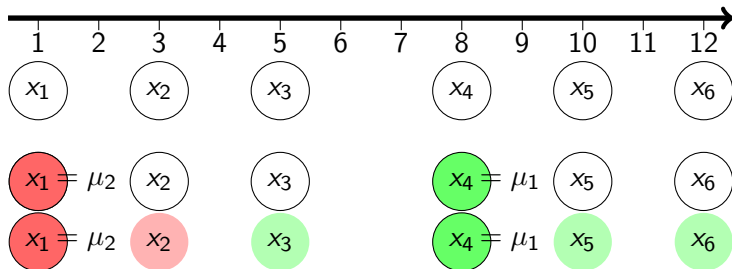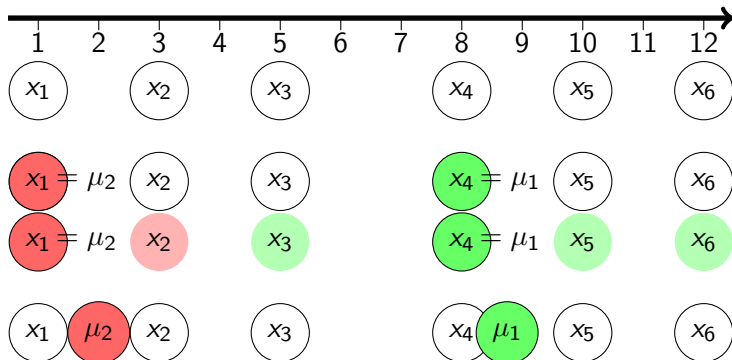
# Example

# Example

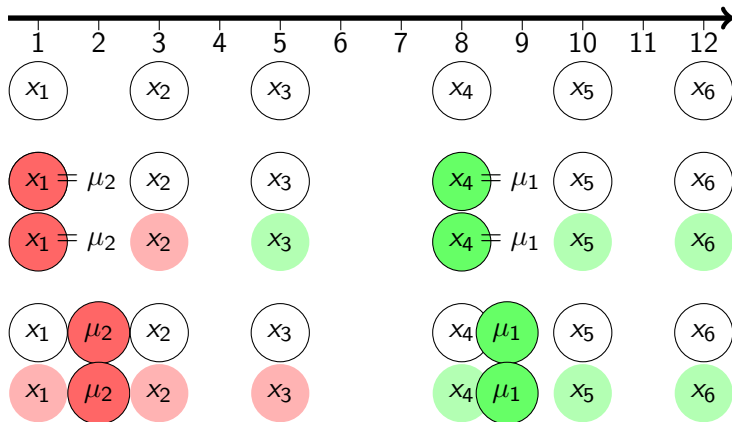# Example

# Example

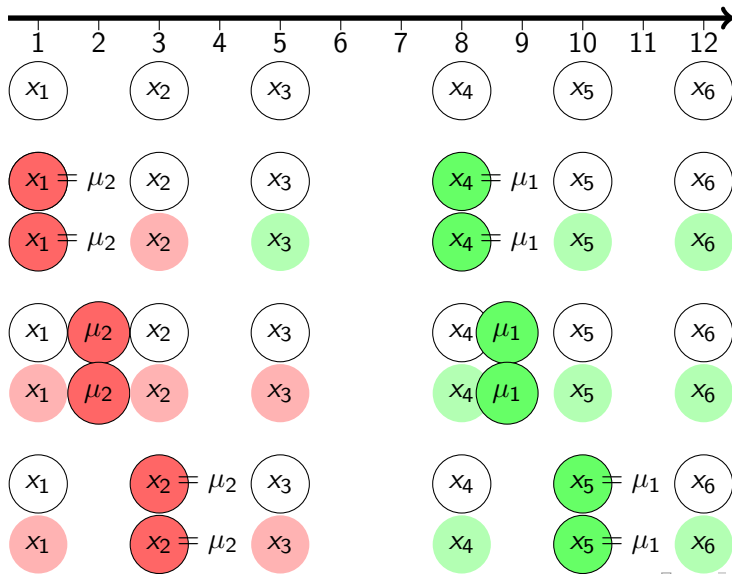# Example

# Example

# Example

# Example

# Example

# How Many Clusters $K$?

# How Many Clusters $K$?

# k-medoids: k-means for General Distances

One can generalize k-means to general distances $d$:

$$D(P, \mu) := \sum_{i=1}^{n} \sum_{k=1}^{K} P_{i,k} d(x_i, \mu_k)$$

# k-medoids: k-means for General Distances

One can generalize k-means to general distances $d$:

$$D(P, \mu) := \sum_{i=1}^{n} \sum_{k=1}^{K} P_{i,k} d(x_i, \mu_k)$$

▶ step 1 assigning data points to clusters remains the same

$$P_{i,k} := \underset{k \in \{1, \dots, K\}}{\arg \min} \; d(x_i, \mu_k)$$

▶ but step 2 finding the best **cluster representatives** $\mu_k$ is not solved by the mean and may be difficult in general.

# k-medoids: k-means for General Distances

One can generalize k-means to general distances $d$:

$$D(P, \mu) := \sum_{i=1}^{n} \sum_{k=1}^{K} P_{i,k} d(x_i, \mu_k)$$

- step 1 assigning data points to clusters remains the same

$$P_{i,k} := \arg\min_{k \in \{1,\dots,K\}} d(x_i, \mu_k)$$

- but step 2 finding the best **cluster representatives** $\mu_k$ is not solved by the mean and may be difficult in general.

idea **k-medoids**: choose cluster representatives out of cluster data points:

$$\mu_k := x_j, \quad j := \arg\min_{j \in \{1,\dots,n\}: P_{j,k}=1} \sum_{i=1}^{n} P_{i,k} d(x_i, x_j)$$

# k-medoids: k-means for General Distances

k-medoids is a "kernel method": it requires no access to the variables, just to the distance measure.

For the **Manhattan distance/$L_1$ distance**, step 2 finding the best cluster representatives $\mu_k$ can be solved without restriction to cluster data points:

$$(\mu_k)_j := \text{median}\{(x_i)_j \mid P_{i,k} = 1, i = 1, \ldots, n\}, \quad j = 1, \ldots, m$$

# Outline

# Soft Partitions: Row Stochastic Matrices

Let $X := \{x_1, \ldots, x_N\}$ be a finite set. A $N \times K$ matrix

$$P \in [0, 1]^{N \times K}$$

is called a **soft partition matrix of $X$** if it

1. is row-stochastic:
$$\sum_{k=1}^{K} P_{i,k} = 1, \quad i \in \{1, \ldots, N\}, k \in \{1, \ldots, K\}$$

2. does not contain a zero column: $X_{i,k} \neq (0, \ldots, 0)^T, \quad k \in \{1, \ldots, K\}$.

$P_{i,k}$ is called the **membership degree of instance $i$ in class $k$** or the **cluster weight of instance $i$ in cluster $k$**.

$P_{.,k}$ is called **membership vector of class $k$**.

SoftPart($X$) denotes the set of all soft partitions of $X$.

Note: Soft partitions are also called **soft clusterings** and **fuzzy clusterings**.

# The Soft Clustering Problem

Given

- a set $\mathcal{X}$ called **data space**, e.g., $\mathcal{X} := \mathbb{R}^m$,
- a set $X \subseteq \mathcal{X}$ called **data**, and
- a function

$$D : \bigcup_{X \subseteq \mathcal{X}} \mathrm{SoftPart}(X) \to \mathbb{R}_0^+$$

called **distortion measure** where $D(P)$ measures how bad a soft partition $P \in \mathrm{SoftPart}(X)$ for a data set $X \subseteq \mathcal{X}$ is,

find a soft partition $P \in \mathrm{SoftPart}(X)$ with minimal distortion $D(P)$.

# The Soft Clustering Problem (with given $K$)

Given

- a set $\mathcal{X}$ called **data space**, e.g., $\mathcal{X} := \mathbb{R}^m$,
- a set $X \subseteq \mathcal{X}$ called **data**,
- a function

$$D : \bigcup_{X \subseteq \mathcal{X}} \mathsf{SoftPart}(X) \to \mathbb{R}_0^+$$

  called **distortion measure** where $D(P)$ measures how bad a soft partition $P \in \mathsf{SoftPart}(X)$ for a data set $X \subseteq \mathcal{X}$ is, and

- a number $K \in \mathbb{N}$ of clusters,

find a soft partition $P \in \mathsf{SoftPart}_K(X) \subseteq [0,1]^{|X| \times K}$ with $K$ clusters with minimal distortion $D(P)$.

# Mixture Models

Mixture models assume that there exists an **unobserved nominal variable** $Z$ with $K$ levels:

$$p(X, Z) = p(Z)p(X \mid Z) = \prod_{k=1}^{K} (\pi_k p(X \mid Z = k)^{\delta(Z=k)}$$

The **complete data loglikelihood** of the **completed data** $(X, Z)$ then is

$$\ell(\Theta; X, Z) := \sum_{i=1}^{n} \sum_{k=1}^{K} \delta(Z_i = k)(\ln \pi_k + \ln p(X = x_i \mid Z = k; \theta_k)$$

$$\text{with } \Theta := (\pi_1, \ldots, \pi_K, \theta_1, \ldots, \theta_K)$$

$\ell$ cannot be computed because $z_i$'s are unobserved.

# Mixture Models: Expected Loglikelihood

Given an estimate $\Theta^{(t-1)}$ of the parameters, mixtures aim to optimize the **expected complete data loglikelihood**:

$$
\begin{aligned}
Q(\Theta;\Theta^{(t-1)}) &:= \mathbb{E}[\ell(\Theta; X, Z) \mid \Theta^{(t-1)}] \\
&= \sum_{i=1}^{n} \sum_{k=1}^{K} \mathbb{E}[\delta(Z_i = k) \mid x_i, \Theta^{(t-1)}](\ln \pi_k + \ln p(X = x_i \mid Z = k; \theta_k))
\end{aligned}
$$

which is relaxed to

$$
\begin{aligned}
Q(\Theta, r; \Theta^{(t-1)}) = \sum_{i=1}^{n} \sum_{k=1}^{K} & r_{i,k}(\ln \pi_k + \ln p(X = x_i \mid Z = k; \theta_k)) \\
&+ (r_{i,k} - \mathbb{E}[\delta(Z_i = k) \mid x_i, \Theta^{(t-1)}])^2
\end{aligned}
$$

# Mixture Models: Expected Loglikelihood

Block coordinate descent (**EM algorithm**): alternate until convergence

1. **expectation step**:

$$
\begin{aligned}
r_{i,k}^{(t-1)} &:= \mathbb{E}[\delta(Z_i = k) \mid x_i, \Theta^{(t-1)}] = p(Z = k \mid X = x_i; \Theta^{(t-1)}) \\
&= \frac{p(X = x_i \mid Z = k; \Theta^{(t-1)})p(Z = k; \Theta^{(t-1)})}{\sum_{k'=1}^{K} p(X = x_i \mid Z = k'; \Theta^{(t-1)})p(Z = k'; \Theta^{(t-1)})} \\
&= \frac{p(X = x_i \mid Z = k; \theta_k^{(t-1)})\pi_k^{(t-1)}}{\sum_{k'=1}^{K} p(X = x_i \mid Z = k'; \theta_k^{(t-1)})\pi_k^{(t-1)}}
\end{aligned}
\tag{0}
$$

2. **maximization step**:

$$
\begin{aligned}
\Theta^{(t)} &:= \underset{\Theta}{\arg\max}\, Q(\Theta, r^{(t-1)}; \Theta^{(t-1)}) \\
&= \underset{\pi_1,\ldots,\pi_K,\theta_1,\ldots,\theta_K}{\arg\max} \sum_{i=1}^{n} \sum_{k=1}^{K} r_{i,k}(\ln \pi_k + \ln p(X = x_i \mid Z = k; \theta_k))
\end{aligned}
$$

# Mixture Models: Expected Loglikelihood

2. **maximization step**:

$$\Theta^{(t)} = \underset{\pi_1,\ldots,\pi_K,\theta_1,\ldots,\theta_K}{\arg\max} \sum_{i=1}^{n}\sum_{k=1}^{K} r_{i,k}(\ln \pi_k + \ln p(X = x_i \mid Z = k; \theta_k))$$

$$\leadsto \quad \pi_k^{(t)} = \frac{\sum_{i=1}^{n} r_{i,k}}{n} \tag{1}$$

$$\sum_{i=1}^{n} \frac{r_{i,k}}{p(X = x_i \mid Z = k; \theta_k)} \frac{\partial p(X = x_i \mid Z = k; \theta_k)}{\partial \theta_k} = 0, \quad \forall k \tag{*}$$

(*) needs to be solved for specific cluster specific distributions $p(X|Z)$.

## Gaussian Mixtures

Gaussian mixtures:

► use Gaussians for $p(X|Z)$:

$$p(X = x \mid Z = k) = \frac{1}{\sqrt{(2\pi)^m |\Sigma_k|}} e^{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1}(x-\mu_k)}, \quad \theta_k := (\mu_k, \Sigma_k)$$

$$\rightsquigarrow \quad \mu_k^{(t)} = \frac{\sum_{i=1}^n r_{i,k}^{(t-1)} x_i}{\sum_{i=1}^k r_{i,k}^{(t-1)}} \tag{2}$$

$$\Sigma_k^{(t)} = \frac{\sum_{i=1}^n r_{i,k}^{(t-1)} (x_i - \mu_k^{(t)})^T (x_i - \mu_k^{(t)})}{\sum_{i=1}^n r_{i,k}^{(t-1)}}$$

$$= \frac{\sum_{i=1}^n r_{i,k}^{(t-1)} x_i^T x_i - \mu_k^{(t)T} \mu_k^{(t)}}{\sum_{i=1}^n r_{i,k}^{(t-1)}} \tag{3}$$

# Gaussian Mixtures: EM Algorithm, Summary

1. **expectation step**: $\forall i, k$

$$\tilde{r}_{i,k}^{(t-1)} = \frac{1}{\sqrt{(2\pi)^m |\Sigma_k^{(t-1)}|}} e^{-\frac{1}{2}(x_i - \mu_k^{(t-1)})^T \Sigma_k^{(t-1)-1}(x_i - \mu_k^{(t-1)})} \tag{0a}$$

$$r_{i,k}^{(t-1)} = \frac{\tilde{r}_{i,k}^{(t-1)}}{\sum_{k'=1}^{K} \tilde{r}_{i,k'}^{(t-1)}} \tag{0b}$$

2. **maximization step**: $\forall k$

$$\pi_k^{(t)} = \frac{\sum_{i=1}^{n} r_{i,k}^{(t-1)}}{n} \tag{1}$$

$$\mu_k^{(t)} = \frac{\sum_{i=1}^{n} r_{i,k}^{(t-1)} x_i}{\sum_{i=1}^{n} r_{i,k}^{(t-1)}} \tag{2}$$

$$\Sigma_k^{(t)} = \frac{\sum_{i=1}^{n} r_{i,k}^{(t-1)} x_i^T x_i - \mu_k^{(t)T} \mu_k^{(t)}}{\sum_{i=1}^{n} r_{i,k}^{(t-1)}} \tag{3}$$

# Gaussian Mixtures for Soft Clustering

▶ The **responsibilities** $r \in [0,1]^{N \times K}$ are a soft partition.

$$P := r$$

▶ The negative expected loglikelihood can be used as cluster distortion:

$$D(P) := -\max_{\Theta} Q(\Theta, r)$$

▶ To optimize $D$, we simply can run EM.

# Gaussian Mixtures for Soft Clustering

▶ The **responsibilities** $r \in [0,1]^{N \times K}$ are a soft partition.

$$P := r$$

▶ The negative expected loglikelihood can be used as cluster distortion:
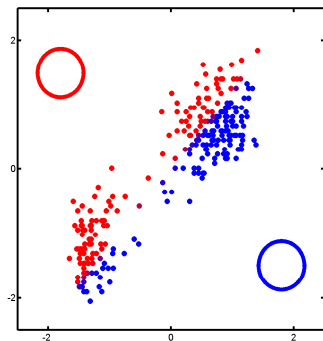
$$D(P) := -\max_{\Theta} Q(\Theta, r)$$

▶ To optimize $D$, we simply can run EM.

For hard clustering:

▶ assign points to the cluster with highest responsibility (**hard EM**):

$$r_{i,k}^{(t-1)} = \delta(k = \arg\max_{k'=1,...,K} \tilde{r}_{i,k'}^{(t-1)}) \qquad (0\text{b}')$$
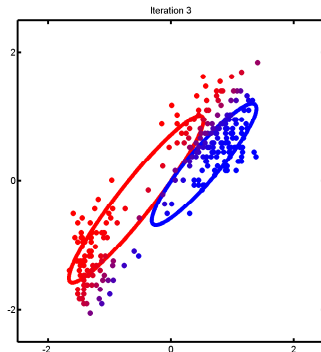
# Gaussian Mixtures for Soft Clustering / Example



[Mur12, fig. 11.11]

# Gaussian Mixtures for Soft Clustering / Example



[Mur12, fig. 11.11]

# Gaussian Mixtures for Soft Clustering / Example



[Mur12, fig. 11.11]

# Model-based Cluster Analysis

Different parametrizations of the covariance matrices $\Sigma_k$ restrict possible **cluster shapes**:

- full $\Sigma$:
  all sorts of ellipsoid clusters.
- diagonal $\Sigma$:
  ellipsoid clusters with axis-parallel axes
- unit $\Sigma$:
  spherical clusters.

One also distinguishes

- cluster-specific $\Sigma_k$:
  each cluster can have its own shape.
- shared $\Sigma_k = \Sigma$:
  all clusters have the same shape.

# k-means: Hard EM with spherical clusters

1. **expectation step**: $\forall i, k$

$$\tilde{r}_{i,k}^{(t-1)} = \frac{1}{\sqrt{(2\pi)^m |\Sigma_k^{(t-1)}|}} e^{-\frac{1}{2}(x_i - \mu_k^{(t-1)})^T \Sigma_k^{(t-1)-1}(x_i - \mu_k^{(t-1)})} \quad (0a)$$

$$= \frac{1}{\sqrt{(2\pi)^m}} e^{-\frac{1}{2}(x_i - \mu_k^{(t-1)})^T (x_i - \mu_k^{(t-1)})}$$

$$r_{i,k}^{(t-1)} = \delta(k = \underset{k'=1,\ldots,K}{\arg\max}\, \tilde{r}_{i,k'}^{(t-1)}) \quad (0b')$$

$$\underset{k'=1,\ldots,K}{\arg\max}\, \tilde{r}_{i,k'}^{(t-1)} = \underset{k'=1,\ldots,K}{\arg\max}\, \frac{1}{\sqrt{(2\pi)^m}} e^{-\frac{1}{2}(x_i - \mu_k^{(t-1)})^T (x_i - \mu_k^{(t-1)})}$$

$$= \underset{k'=1,\ldots,K}{\arg\max}\, -(x_i - \mu_k^{(t-1)})^T (x_i - \mu_k^{(t-1)})$$

$$= \underset{k'=1,\ldots,K}{\arg\min}\, ||x_i - \mu_k^{(t-1)}||^2$$

# Outline

# Hierarchies

Let $X$ be a set.

A tree $(H, E)$, $E \subseteq H \times H$ edges pointing towards root

- with leaf nodes $h$ corresponding bijectively to elements $x_h \in X$
- plus a surjective map $L : H \to \{0, \ldots, d\}, d \in \mathbb{N}$ with
  - $L(\text{root}) = 0$ and
  - $L(h) = d$ for all leaves $h \in H$ and
  - $L(h) \leq L(g)$ for all $(g, h) \in E$

  called **level map**

is called an **hierarchy over** $X$.

# Hierarchies

Let $X$ be a set.

A tree $(H, E)$, $E \subseteq H \times H$ edges pointing towards root

- with leaf nodes $h$ corresponding bijectively to elements $x_h \in X$
- plus a surjective map $L : H \to \{0, \ldots, d\}, d \in \mathbb{N}$ with
  - $L(\text{root}) = 0$ and
  - $L(h) = d$ for all leaves $h \in H$ and
  - $L(h) \leq L(g)$ for all $(g, h) \in E$

  called **level map**

is called an **hierarchy over** $X$.

$d$ is called the **depth** of the hierarchy.

$\text{Hier}(X)$ denotes the set of all hierarchies over $X$.

# Hierarchies / Example

$X$ :          $x_1$          $x_2$          $x_3$          $x_4$          $x_5$          $x_6$

# Hierarchies / Example



$X$ :    $x_1$    $x_2$    $x_3$    $x_4$    $x_5$    $x_6$

# Hierarchies / Example

# Hierarchies / Example

# Hierarchies: Nodes Correspond to Subsets

Let $(H, E)$ be such an hierarchy:
- ► nodes of an hierarchy correspond to subsets of $X$:
    - ► leaf nodes $h$ correspond to a singleton subset by definition.

$$\text{subset}(h) := \{x_h\}, \quad x_h \in X \text{ corresponding to leaf } h$$

    - ► interior nodes $h$ correspond to the union of the subsets of their children:

$$\text{subset}(h) := \bigcup_{\substack{g \in H \\ (g,h) \in E}} \text{subset}(g)$$

- ► thus the root node $h$ corresponds to the full set $X$:

$$\text{subset}(h) = X$$

# Hierarchies: Nodes Correspond to Subsets

# Hierarchies: Levels Correspond to Partitions

Let $(H, E)$ be such an hierarchy:

- levels $\ell \in \{0, \ldots, d\}$ correspond to partitions

$$P_\ell(H, L) := \{h \in H \mid L(h) \geq \ell, \nexists g \in H : L(g) \geq \ell, h \subsetneq g\}$$

# Hierarchies: Levels Correspond to Partitions

$$\{x_1, x_3, x_4, x_2, x_5, x_6\} \qquad\qquad \{\{x_1, x_3, x_4, x_2, x_5, x_6\}\}$$

$$\{x_1, x_3, x_4\} \qquad\qquad\qquad \{\{x_1, x_3, x_4\}, \{x_2, x_5, x_6\}\}$$

$$\{x_2, x_5, x_6\} \qquad \{\{x_1, x_3\}, \{x_4\}, \{x_2, x_5, x_6\}\}$$

$$\{x_2, x_5\} \qquad \{\{x_1, x_3\}, \{x_4\}, \{x_2, x_5\}, \{x_6\}\}$$

$$\{x_1, x_3\} \qquad\qquad\qquad\qquad \{\{x_1, x_3\}, \{x_4\}, \{x_2\}, \{x_5\}, \{x_6\}\}$$

$$\{x_1\} \quad \{x_3\} \quad \{x_4\} \quad \{x_2\} \quad \{x_5\} \quad \{x_6\} \; \{\{x_1\}, \{x_3\}, \{x_4\}, \{x_2\}, \{x_5\}, \{x_6\}\}$$

# The Hierarchical Cluster Analysis Problem

Given

- a set $\mathcal{X}$ called **data space**, e.g., $\mathcal{X} := \mathbb{R}^m$,
- a set $X \subseteq \mathcal{X}$ called **data** and
- a function

$$D : \bigcup_{X \subseteq \mathcal{X}} \text{Hier}(X) \to \mathbb{R}_0^+$$

called **distortion measure** where $D(P)$ measures how bad a hierarchy $H \in \text{Hier}(X)$ for a data set $X \subseteq \mathcal{X}$ is,

find a hierarchy $H \in \text{Hier}(X)$ with minimal distortion $D(H)$.

# Distortions for Hierarchies

Examples for distortions for hierarchies:

$$D(H) := \sum_{K=1}^{n} \tilde{D}(P_K(H))$$

where

- $P_K(H)$ denotes the partition at level $K-1$ (with $K$ classes) and
- $\tilde{D}$ denotes a distortion for partitions.

# Agglomerative and Divisive Hierarchical Clustering

Hierarchies are usually learned by greedy search level by level:

- **agglomerative clustering:**

  1. start with the singleton partition $P_n$:

  $$P_n := \{X_k \mid k = 1, \ldots, n\}, \quad X_k := \{x_k\}, \quad k = 1, \ldots, n$$

  2. in each step $K = n, \ldots, 2$ build $P_{K-1}$ by joining the two clusters $k, \ell \in \{1, \ldots, K\}$ that lead to the minimal distortion

  $$D(\{X_1, \ldots, \widehat{X_k}, \ldots, \widehat{X_\ell}, \ldots, X_K, X_k \cup X_\ell)$$

Note: $\widehat{X_k}$ denotes that the class $X_k$ is omitted from the partition.

# Agglomerative and Divisive Hierarchical Clustering

Hierarchies are usually learned by greedy search level by level:

- **divisive clustering:**
    1. start with the all partition $P_1$:

    $$P_1 := \{X\}$$

    2. in each step $K = 1, n - 1$ build $P_{K+1}$ by splitting one cluster $X_k$ in two clusters $X'_k, X'_\ell$ that lead to the minimal distortion

    $$D(\{X_1, \ldots, \widehat{X_k}, \ldots, X_K, X'_k, X'_\ell\}), \quad X_k = X'_k \cup X'_\ell$$

Note: $\widehat{X_k}$ denotes that the class $X_k$ is omitted from the partition.

# Class-wise Defined Partition Distortions

If the partition distortion can be written as a sum of distortions of its classes,

$$D(\{X_1, \ldots, X_K\}) = \sum_{k=1}^{K} \tilde{D}(X_k)$$

then the optimal pair does only depend on $X_k, X_\ell$:

$$D(\{X_1, \ldots, \widehat{X_k}, \ldots, \widehat{X_\ell}, \ldots, X_K, X_k \cup X_\ell\}) = \tilde{D}(X_k \cup X_\ell) - (\tilde{D}(X_k) + \tilde{D}(X_\ell))$$

# Closest Cluster Pair Partition Distortions

For a **cluster distance**

$$\tilde{d} : \mathcal{P}(X) \times \mathcal{P}(X) \to \mathbb{R}_0^+$$
$$\text{with} \quad \tilde{d}(A \cup B, C) \geq \min\{\tilde{d}(A, C), \tilde{d}(B, C)\}, \quad A, B, C \subseteq X$$

a partition can be judged by the closest cluster pair it contains:

$$D(\{X_1, \ldots, X_K\}) = \min_{\substack{k, \ell = 1, K \\ k \neq \ell}} \tilde{d}(X_k, X_\ell)$$

Such a distortion has to be maximized.

To increase it, the closest cluster pair has to be joined.

# Single Link Clustering

$$d_{\mathsf{sl}}(A, B) := \min_{x \in A, y \in B} d(x, y), \quad A, B \subseteq X$$

# Complete Link Clustering

$$d_{\mathrm{cl}}(A, B) := \max_{x \in A, y \in B} d(x, y), \quad A, B \subseteq X$$

# Average Link Clustering

$$d_{\text{al}}(A, B) := \frac{1}{|A||B|} \sum_{x \in A, y \in B} d(x, y), \quad A, B \subseteq X$$

# Recursion Formulas for Cluster Distances

$$d_{\text{sl}}(X_i \cup X_j, X_k) := \min_{x \in X_i \cup X_j, y \in X_k} d(x, y)$$

$$= \min\{\min_{x \in X_i, y \in X_k} d(x, y), \min_{x \in X_j, y \in X_k} d(x, y)\}$$

$$= \min\{d_{\text{sl}}(X_i, X_k), d_{\text{sl}}(X_j, X_k)\}$$

# Recursion Formulas for Cluster Distances

$$d_{\mathsf{sl}}(X_i \cup X_j, X_k) = \min\{d_{\mathsf{sl}}(X_i, X_k), d_{\mathsf{sl}}(X_j, X_k)\}$$

$$d_{\mathsf{cl}}(X_i \cup X_j, X_k) := \max_{x \in X_i \cup X_j, y \in X_k} d(x, y)$$

$$= \max\{\max_{x \in X_i, y \in X_k} d(x, y), \max_{x \in X_j, y \in X_k} d(x, y)\}$$

$$= \max\{d_{\mathsf{cl}}(X_i, X_k), d_{\mathsf{cl}}(X_j, X_k)\}$$

# Recursion Formulas for Cluster Distances

$$d_{\text{sl}}(X_i \cup X_j, X_k) = \min\{d_{\text{sl}}(X_i, X_k), d_{\text{sl}}(X_j, X_k)\}$$

$$d_{\text{cl}}(X_i \cup X_j, X_k) = \max\{d_{\text{cl}}(X_i, X_k), d_{\text{cl}}(X_j, X_k)\}$$

$$d_{\text{al}}(X_i \cup X_j, X_k) := \frac{1}{|X_i \cup X_j||X_k|} \sum_{x \in X_i \cup X_j, y \in X_k} d(x, y)$$

$$= \frac{|X_i|}{|X_i \cup X_j|} \frac{1}{|X_i||X_k|} \sum_{x \in X_i, y \in X_k} d(x, y)$$

$$+ \frac{|X_j|}{|X_i \cup X_j|} \frac{1}{|X_j||X_k|} \sum_{x \in X_j, y \in X_k} d(x, y)$$

$$= \frac{|X_i|}{|X_i| + |X_j|} d_{\text{al}}(X_i, X_k) + \frac{|X_j|}{|X_i| + |X_j|} d_{\text{al}}(X_j, X_k)$$

# Recursion Formulas for Cluster Distances

$$d_{\mathsf{sl}}(X_i \cup X_j, X_k) = \min\{d_{\mathsf{sl}}(X_i, X_k), d_{\mathsf{sl}}(X_j, X_k)\}$$

$$d_{\mathsf{cl}}(X_i \cup X_j, X_k) = \max\{d_{\mathsf{cl}}(X_i, X_k), d_{\mathsf{cl}}(X_j, X_k)\}$$

$$d_{\mathsf{al}}(X_i \cup X_j, X_k) = \frac{|X_i|}{|X_i| + |X_j|} d_{\mathsf{al}}(X_i, X_k) + \frac{|X_j|}{|X_i| + |X_j|} d_{\mathsf{al}}(X_j, X_k)$$

$\rightsquigarrow$ agglomerative hierarchical clustering requires to compute the
**distance matrix $D \in \mathbb{R}^{n \times n}$** only once:

$$D_{i,j} := d(x_i, x_j), \quad i, j = 1, \ldots, K$$

# Conclusion (1/2)

- Cluster analysis aims at **detecting latent groups** in data, without labeled examples ($\leftrightarrow$ **record linkage**).

- Latent groups can be described in three different granularities:
    - **partitions** segment data into $K$ subsets (**hard clustering**).
    - **hierarchies** structure data into an hierarchy, in a sequence of consistent partitions (**hierarchical clustering**).
    - **soft clusterings / row-stochastic matrices** build overlapping groups to which data points can belong with some **membership degree** (**soft clustering**).

- **k-means** finds a $K$-partition by finding $K$ **cluster centers** with smallest **Euclidean distance** to all their cluster points.

- **k-medoids** generalizes k-means to **general distances**; it finds a $K$-partition by selecting $K$ data points as **cluster representatives** with smallest distance to all their cluster points.

# Conclusion (2/2)

- **hierarchical single link, complete link and average link methods**
    - find a hierarchy by greedy search over consistent partitions,
    - starting from the singleton parition (**agglomerative**)
    - being efficient due to **recursion formulas**,
    - requiring only a distance matrix.

- **Gaussian Mixture Models** find soft clusterings by modeling data by a class-specific multivariate Gaussian distribution $p(X \mid Z)$ and estimating expected class memberships (**expected likelihood**).

- The **Expectation Maximiation Algorithm (EM)** can be used to learn Gaussian Mixture Models via block coordinate descent.

- k-means is a special case of a Gaussian Mixture Model
    - with hard/binary cluster memberships (**hard EM**) and
    - **spherical cluster shapes**.

# Readings

- ▶ k-means:
  - ▶ [HTFF05], ch. 14.3.6, 13.2.3, 8.5 [Bis06], ch. 9.1, [Mur12], ch. 11.4.2
- ▶ hierarchical cluster analysis:
  - ▶ [HTFF05], ch. 14.3.12, [Mur12], ch. 25.5. [PTVF07], ch. 16.4.
- ▶ Gaussian mixtures:
  - ▶ [HTFF05], ch. 14.3.7, [Bis06], ch. 9.2, [Mur12], ch. 11.2.3, [PTVF07], ch. 16.1.

# References

Christopher M. Bishop.
*Pattern Recognition and Machine Learning*.
Springer, 2006.

Trevor Hastie, Robert Tibshirani, Jerome Friedman, and James Franklin.
*The elements of statistical learning: data mining, inference and prediction*, volume 27.
2005.

Kevin P. Murphy.
*Machine learning: a probabilistic perspective*.
The MIT Press, 2012.

William H. Press, Saul A. Teukolsky, William T. Vetterling, and Brian P. Flannery.
*Numerical Recipes*.
Cambridge University Press, 3rd edition, 2007.