

Syllabus

Fri. 21.10. (1) 0. Introduction

A. Supervised Learning: Linear Models & Fundamentals

Fri. 27.10. (2) A.1 Linear Regression

Fri. 3.11. (3) A.2 Linear Classification

Fri. 10.11. (4) A.3 Regularization

Fri. 17.11. (5) A.4 High-dimensional Data

B. Supervised Learning: Nonlinear Models

Fri. 24.11. (6) B.1 Nearest-Neighbor Models

Fri. 1.12. (7) B.2 Neural Networks

Fri. 8.12. (8) B.3 Decision Trees

Fri. 15.12. (9) B.4 Support Vector Machines

Fri. 12.1. (10) B.5 A First Look at Bayesian and Markov Networks

C. Unsupervised Learning

Fri. 19.1. (11) C.1 Clustering

Fri. 26.1. (12) C.2 Dimensionality Reduction

Fri. 2.2. (13) C.3 Frequent Pattern Mining

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

1 / 37

Machine Learning

Outline

1. k-means & k-medoids

2. Gaussian Mixture Models

3. Hierarchical Cluster Analysis

Outline

1. k-means & k-medoids
2. Gaussian Mixture Models
3. Hierarchical Cluster Analysis

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

1 / 37

Machine Learning 1. k-means & k-medoids

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 : $\bigcup_{A \in P} A = X, \text{ and}$
3. do not contain the empty set: $\emptyset \notin P.$

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

1. are **pairwise disjoint**: $X_k \cap X_j = \emptyset, \quad k, j \in \{1, \dots, K\}, k \neq j$
2. **cover X** : $\bigcup_{k=1}^K X_k = X, \text{ and}$
3. do **not contain the empty set**: $X_k \neq \emptyset, \quad k \in \{1, \dots, K\}.$

A set X_k is also called a **cluster**, a partition P a **clustering**.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

1 / 37

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**, and
- ▶ a function

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

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

k-means: Distortion Sum of Distances to Cluster Centers

Sum of squared distances to cluster centers:

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

with

$$\mu_k := \frac{\sum_{n=1}^N P_{n,k} x_n}{\sum_{n=1}^N P_{n,k}} = \text{mean} \{x_n \mid P_{n,k} = 1, n \in \{1, \dots, N\}\}$$

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 μ as auxiliary optimization variables:

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

Block coordinate descent:

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

$$P_{n,k} := \delta(k = \ell_n), \quad \ell_n := \arg \min_{k \in \{1, \dots, K\}} \|x_n - \mu_k\|^2$$

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

$$\mu_k := \frac{\sum_{n=1}^N P_{n,k} x_n}{\sum_{n=1}^N P_{n,k}}$$

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 μ_1 out of the data points at random and then
2. sequentially select the data point with the largest sum of distances to already chosen cluster centers as next cluster center

$$\mu_k := x_n, \quad n := \arg \max_{n \in \{1, \dots, N\}} \sum_{\ell=1}^{k-1} \|x_n - \mu_\ell\|^2, \quad k = 2, \dots, K$$

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( $\mathcal{D} := \{x_1, \dots, x_N\} \subseteq \mathbb{R}^M, K \in \mathbb{N}, \epsilon \in \mathbb{R}^+$ )
2:    $n_1 \sim \text{unif}(\{1, \dots, N\}), \mu_1 := x_{n_1}$ 
3:   for  $k := 2, \dots, K$  do
4:      $n_k := \arg \max_{n \in \{1, \dots, N\}} \sum_{j=1}^{k-1} \|x_n - \mu_j\|, \mu_k := x_{n_k}$ 
5:   repeat
6:      $\mu^{\text{old}} := \mu$ 
7:     for  $n := 1, \dots, N$  do
8:        $P_n := \arg \min_{k \in \{1, \dots, K\}} \|x_n - \mu_k\|$ 
9:     for  $k := 1, \dots, K$  do
10:       $\mu_k := \text{mean} \{x_n \mid P_n = k, n \in \{1, \dots, N\}\}$ 
11:   until  $\frac{1}{K} \sum_{k=1}^K \|\mu_k - \mu_k^{\text{old}}\| < \epsilon$ 
12:   return  $P$ 

```

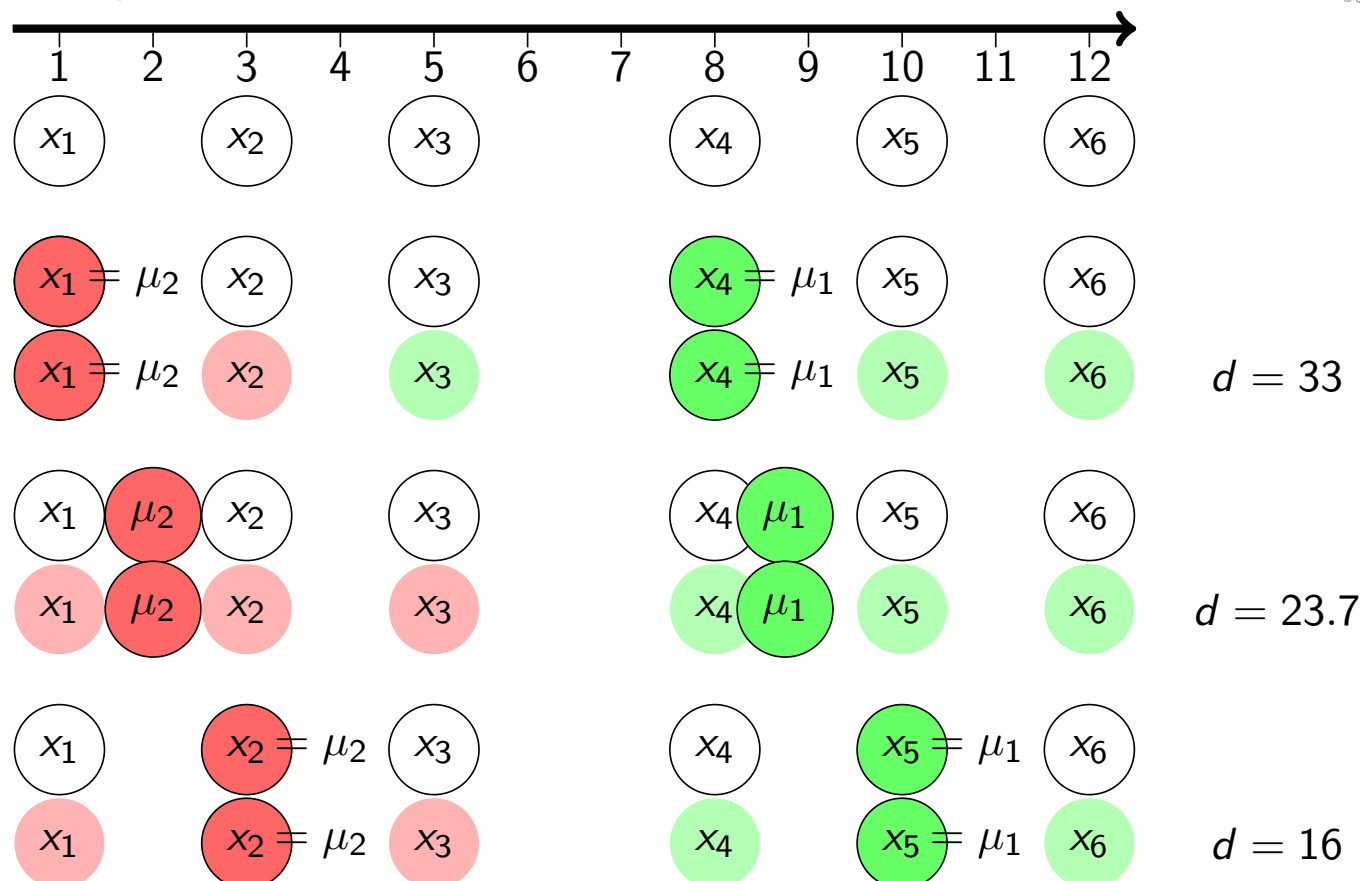
Note: In implementations, the two loops over the data (lines 7 and 10) can be combined in one loop.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

6 / 37

Machine Learning 1. k-means & k-medoids

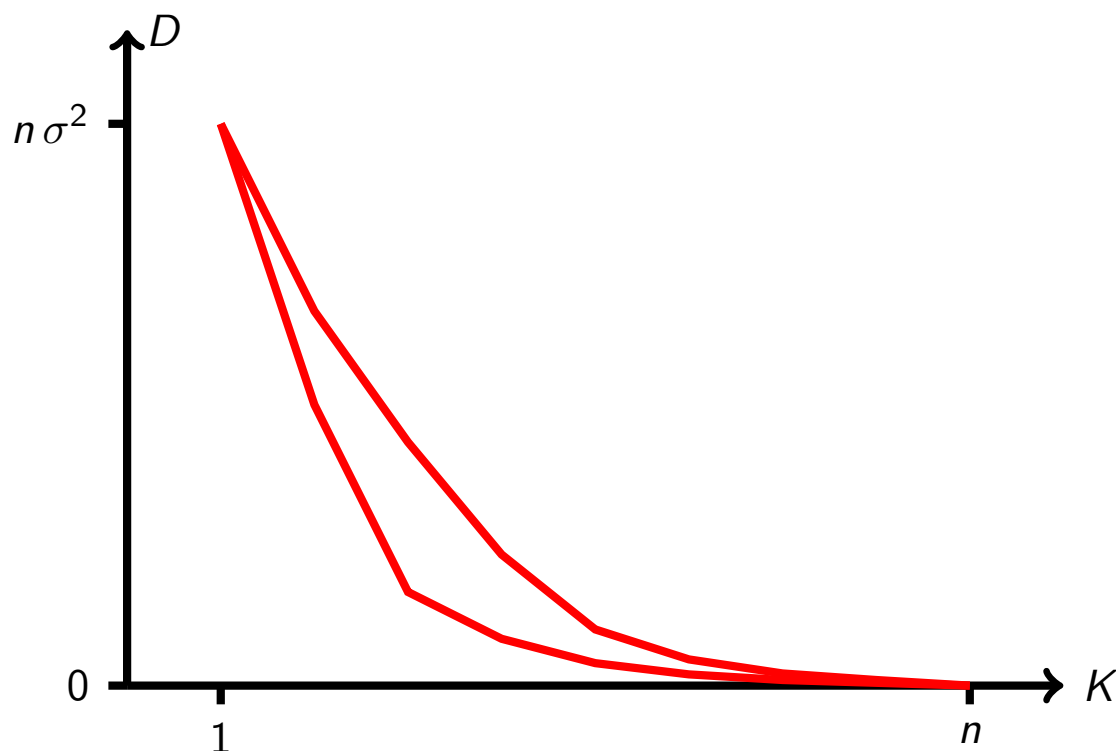
Example



Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

7 / 37

How Many Clusters K ?



Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

8 / 37

Machine Learning 1. k-means & k-medoids

k-medoids: k-means for General Distances

One can generalize k-means to general distances d :

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

- step 1 assigning data points to clusters remains the same

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

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

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

$$\mu_k := x_n, \quad n := \arg \min_{n \in \{1, \dots, N\}: P_{n,k}=1} \sum_{\ell=1}^N P_{\ell,k} d(x_\ell, x_n)$$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

9 / 37

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/L₁ distance**, step 2 finding the best cluster representatives μ_k can be solved without restriction to cluster data points:

$$(\mu_k)_m := \text{median}\{(x_n)_m \mid P_{n,k} = 1, n \in \{1, \dots, N\}\}, \quad m = 1, \dots, M$$

Outline

1. k-means & k-medoids
2. Gaussian Mixture Models
3. Hierarchical Cluster Analysis

Soft Partitions: Row Stochastic Matrices

Let $X := \{x_1, \dots, x_N\}$ 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:
$$\sum_{k=1}^K P_{n,k} = 1, \quad n \in \{1, \dots, N\}$$
2. does not contain a zero column:
$$X_{:,k} \neq (0, \dots, 0)^T, \quad k \in \{1, \dots, K\}$$

$P_{n,k}$ is called the

- ▶ **membership degree of instance n in class k** or the
- ▶ **cluster weight of instance n in cluster k .**

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

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

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

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

11 / 37

Machine Learning 2. Gaussian Mixture Models

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**, and
- ▶ a function

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

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

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

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

Mixture Models

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

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

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

$$\ell(\Theta; X, Z) := \sum_{n=1}^N \sum_{k=1}^K \delta(Z_n = k) (\ln \pi_k + \ln p(X = x_n | Z = k; \theta_k))$$

with $\Theta := (\pi_1, \dots, \pi_K, \theta_1, \dots, \theta_K)$

ℓ cannot be computed because z_n '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**:

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

which is relaxed to

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

Mixture Models: Expected Loglikelihood

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

1. expectation step:

$$\begin{aligned}
 r_{n,k}^{(t-1)} &:= \mathbb{E}[\delta(Z_n = k) \mid x_n, \Theta^{(t-1)}] = p(Z = k \mid X = x_n; \Theta^{(t-1)}) \\
 &= \frac{p(X = x_n \mid Z = k; \Theta^{(t-1)})p(Z = k; \Theta^{(t-1)})}{\sum_{k'=1}^K p(X = x_n \mid Z = k'; \Theta^{(t-1)})p(Z = k'; \Theta^{(t-1)})} \\
 &= \frac{p(X = x_n \mid Z = k; \theta_k^{(t-1)})\pi_k^{(t-1)}}{\sum_{k'=1}^K p(X = x_n \mid Z = k'; \theta_{k'}^{(t-1)})\pi_{k'}^{(t-1)}} \quad (0)
 \end{aligned}$$

2. maximization step:

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

Mixture Models: Expected Loglikelihood

2. maximization step:

$$\begin{aligned}
 \Theta^{(t)} &= \arg \max_{\pi_1, \dots, \pi_K, \theta_1, \dots, \theta_K} \sum_{n=1}^N \sum_{k=1}^K r_{n,k} (\ln \pi_k + \ln p(X = x_n \mid Z = k; \theta_k)) \\
 \rightsquigarrow \pi_k^{(t)} &= \frac{\sum_{n=1}^N r_{n,k}}{N} \quad (1) \\
 \sum_{n=1}^N \frac{r_{n,k}}{p(X = x_n \mid Z = k; \theta_k)} \frac{\partial p(X = x_n \mid Z = k; \theta_k)}{\partial \theta_k} &= 0, \quad \forall k \quad (*)
 \end{aligned}$$

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

Gaussian Mixtures

Gaussian mixtures:

- use Gaussians for $p(X|Z)$:

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

$$\leadsto \mu_k^{(t)} = \frac{\sum_{n=1}^N r_{n,k}^{(t-1)} x_n}{\sum_{n=1}^N r_{n,k}^{(t-1)}} \quad (2)$$

$$\begin{aligned} \Sigma_k^{(t)} &= \frac{\sum_{n=1}^N r_{n,k}^{(t-1)} (x_n - \mu_k^{(t)})^T (x_n - \mu_k^{(t)})}{\sum_{n=1}^N r_{n,k}^{(t-1)}} \\ &= \frac{\sum_{n=1}^N r_{n,k}^{(t-1)} x_n^T x_n - \mu_k^{(t)T} \mu_k^{(t)}}{\sum_{n=1}^N r_{n,k}^{(t-1)}} \end{aligned} \quad (3)$$

Gaussian Mixtures: EM Algorithm, Summary

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

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

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

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

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

$$\mu_k^{(t)} = \frac{\sum_{n=1}^N r_{n,k}^{(t-1)} x_n}{\sum_{n=1}^N r_{n,k}^{(t-1)}} \quad (2)$$

$$\Sigma_k^{(t)} = \frac{\sum_{n=1}^N r_{n,k}^{(t-1)} x_n^T x_n - \mu_k^{(t)T} \mu_k^{(t)}}{\sum_{n=1}^N r_{n,k}^{(t-1)}} \quad (3)$$

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.

For hard clustering:

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

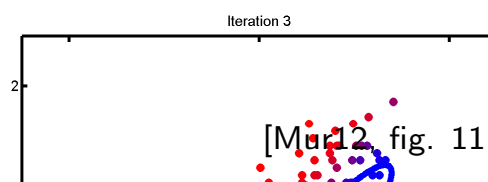
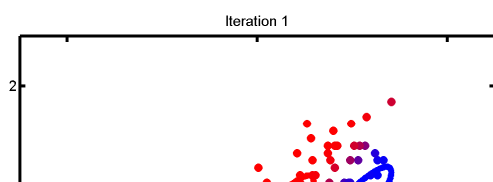
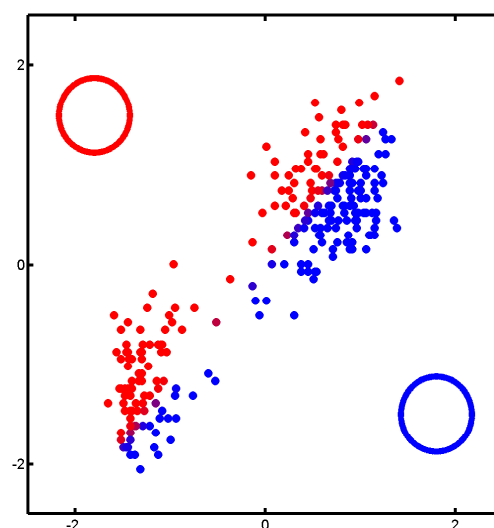
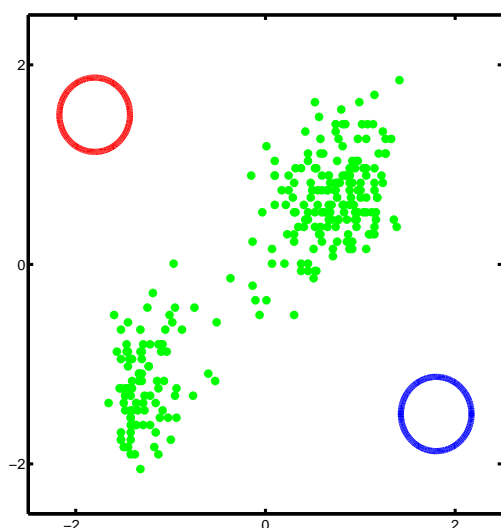
$$r_{n,k}^{(t-1)} = \delta(k = \arg \max_{k'=1,\dots,K} \tilde{r}_{n,k'}^{(t-1)}) \quad (0b')$$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

19 / 37

Machine Learning 2. Gaussian Mixture Models

Gaussian Mixtures for Soft Clustering / Example



[Mur12, fig. 11.11]

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

20 / 37

Model-based Cluster Analysis

Different parametrizations of the covariance matrices Σ_k restrict possible **cluster shapes**:

- ▶ full Σ :
all sorts of ellipsoid clusters.
- ▶ diagonal Σ :
ellipsoid clusters with axis-parallel axes
- ▶ unit Σ :
spherical clusters.

One also distinguishes

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

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

21 / 37

Machine Learning 2. Gaussian Mixture Models

k-means: Hard EM with spherical clusters

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

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

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

$$r_{n,k}^{(t-1)} = \delta(k = \arg \max_{k'=1,\dots,K} \tilde{r}_{n,k'}^{(t-1)}) \quad (0b')$$

$$\arg \max_{k'=1,\dots,K} \tilde{r}_{n,k'}^{(t-1)} = \arg \max_{k'=1,\dots,K} \frac{1}{\sqrt{(2\pi)^M}} e^{-\frac{1}{2}(x_n - \mu_k^{(t-1)})^T (x_n - \mu_k^{(t-1)})}$$

$$= \arg \max_{k'=1,\dots,K} -(x_n - \mu_k^{(t-1)})^T (x_n - \mu_k^{(t-1)})$$

$$= \arg \min_{k'=1,\dots,K} \|x_n - \mu_k^{(t-1)}\|^2$$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

22 / 37

Outline

1. k-means & k-medoids
2. Gaussian Mixture Models
3. Hierarchical Cluster Analysis

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

23 / 37

Machine Learning 3. Hierarchical Cluster Analysis

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 $L : H \rightarrow \{0, \dots, d\}$, $d \in \mathbb{N}$ with
 - ▶ $L(\text{root}) = 0$ and
 - ▶ $L(h) = d$ for all leaves $h \in H$ and
 - ▶ $L(h) \leq L(g)$ for all $(g, h) \in E$

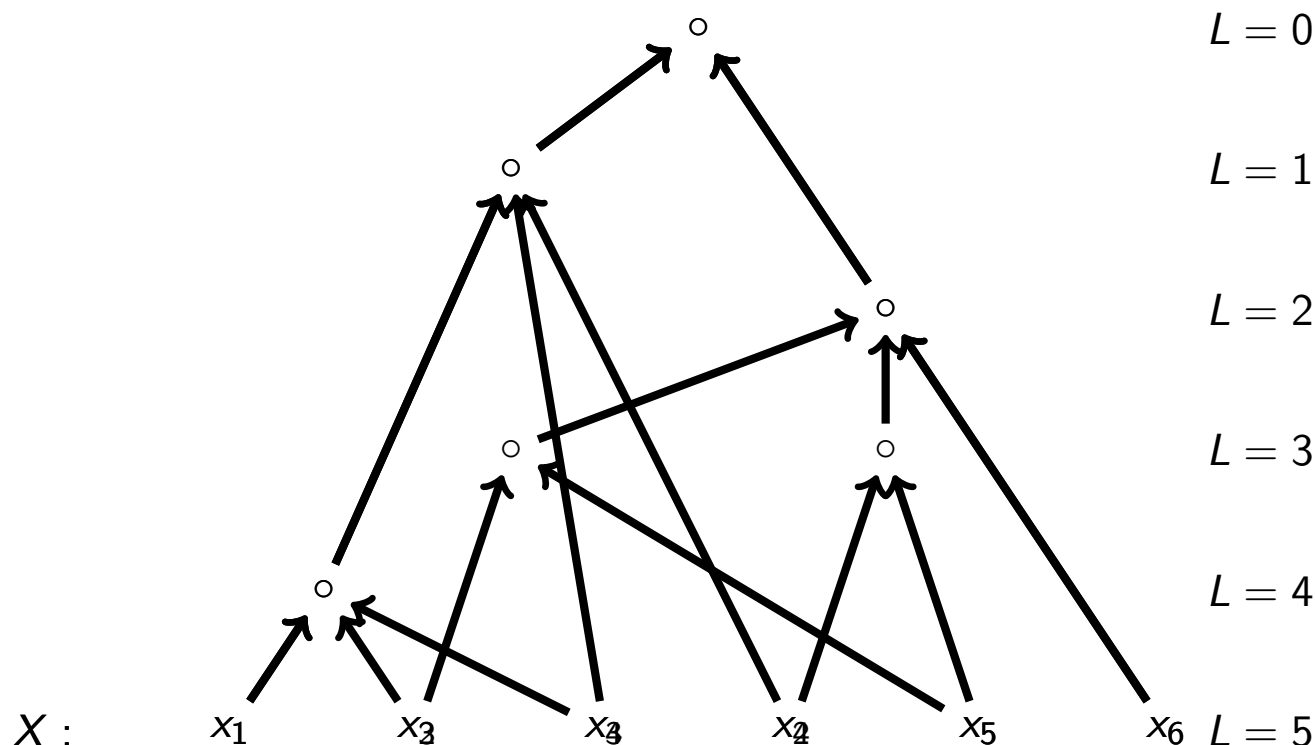
called **level map**

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



Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

24 / 37

Machine Learning 3. Hierarchical Cluster Analysis

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.

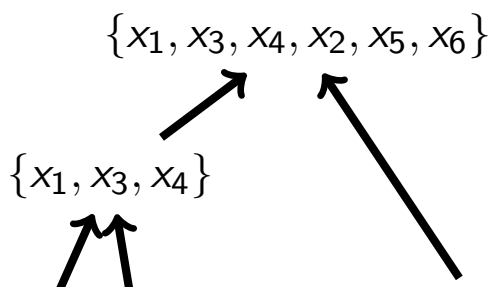
$$\text{subset}(h) := \{x_h\}, \quad x_h \in X \text{ corresponding to leaf } h$$

- 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 :

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



Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

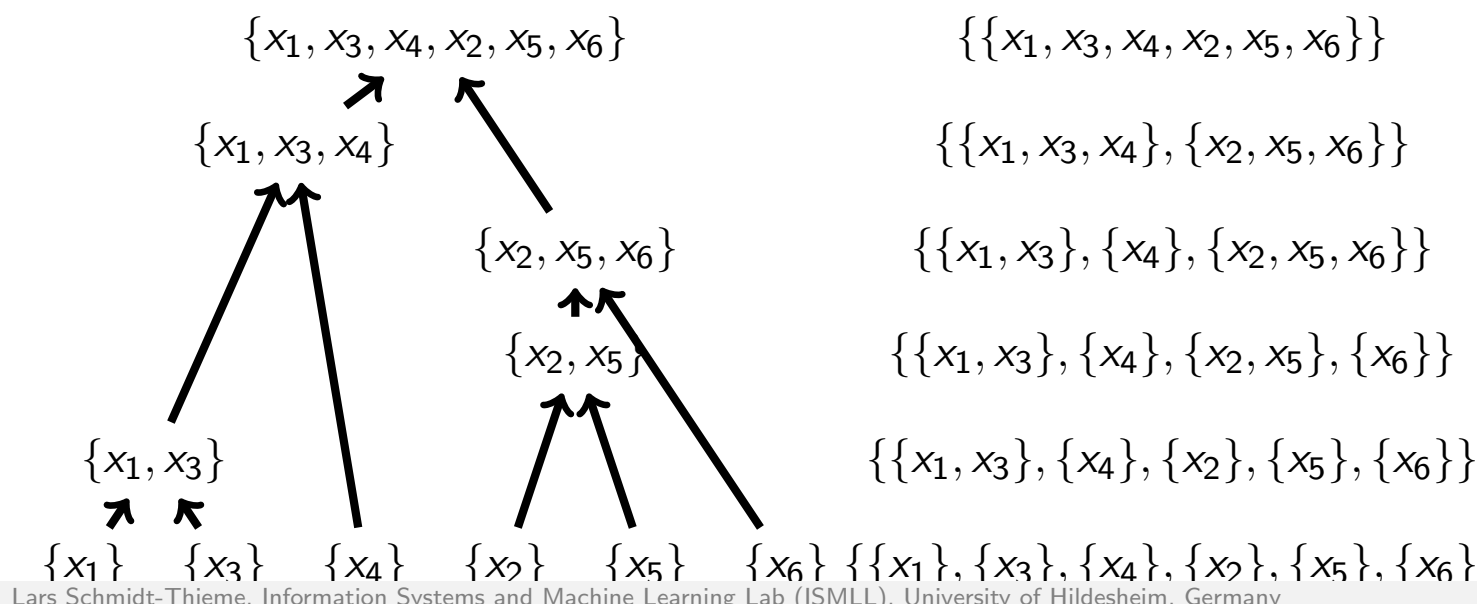
25 / 37

Hierarchies: Levels Correspond to Partitions

Let (H, E) be such an hierarchy:

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

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



Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

26 / 37

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

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

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 $D(H)$.

27 / 37

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 :

$$P_N := \{X_k \mid k = 1, \dots, N\}, \quad X_k := \{x_k\}, \quad k = 1, \dots, N$$

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

$$D(\{X_1, \dots, \widehat{X}_k, \dots, \widehat{X}_\ell, \dots, X_K, X_k \cup X_\ell\})$$

▶ **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'_ℓ that lead to the minimal distortion

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

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,

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

then the optimal pair does only depend on X_k, X_ℓ :

$$D(\{X_1, \dots, \widehat{X}_k, \dots, \widehat{X}_\ell, \dots, X_K, X_k \cup X_\ell\}) = \tilde{D}(X_k \cup X_\ell) - (\tilde{D}(X_k) + \tilde{D}(X_\ell))$$

Closest Cluster Pair Partition Distortions

For a **cluster distance**

$$\tilde{d} : \mathcal{P}(X) \times \mathcal{P}(X) \rightarrow \mathbb{R}_0^+$$

$$\text{with } \tilde{d}(A \cup B, C) \geq \min\{\tilde{d}(A, C), \tilde{d}(B, C)\}, \quad A, B, C \subseteq X$$

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

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

Such a distortion has to be maximized.

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

Single Link Clustering

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

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

32 / 37

Machine Learning 3. Hierarchical Cluster Analysis

Complete Link Clustering

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

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

33 / 37

Average Link Clustering

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

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

34 / 37

Machine Learning 3. Hierarchical Cluster Analysis

Recursion Formulas for Cluster Distances

$$\begin{aligned}
 d_{sl}(X_i \cup X_j, X_k) &:= \min_{x \in X_i \cup X_j, y \in X_k} d(x, y) \\
 &= \min \left\{ \min_{x \in X_i, y \in X_k} d(x, y), \min_{x \in X_j, y \in X_k} d(x, y) \right\} \\
 &= \min \{ d_{sl}(X_i, X_k), d_{sl}(X_j, X_k) \} \\
 d_{cl}(X_i \cup X_j, X_k) &:= \max_{x \in X_i \cup X_j, y \in X_k} d(x, y) \\
 &= \max \left\{ \max_{x \in X_i, y \in X_k} d(x, y), \max_{x \in X_j, y \in X_k} d(x, y) \right\} \\
 &= \max \{ d_{cl}(X_i, X_k), d_{cl}(X_j, X_k) \} \\
 d_{al}(X_i \cup X_j, X_k) &:= \frac{1}{|X_i \cup X_j||X_k|} \sum_{x \in X_i \cup X_j, y \in X_k} d(x, y) \\
 &= \frac{|X_i|}{|X_i \cup X_j|} \frac{1}{|X_i||X_k|} \sum_{x \in X_i, y \in X_k} d(x, y)
 \end{aligned}$$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

35 / 37

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.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

36 / 37

Machine Learning 3. Hierarchical Cluster Analysis

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 partition (**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 $p(X | Z)$ and estimating expected class memberships (**expected likelihood**).
- ▶ The **Expectation Maximisation 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**.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

37 / 37

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.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

38 / 37

Machine Learning

References



Christopher M. Bishop.

Pattern recognition and machine learning, volume 1.
springer New York, 2006.



Trevor Hastie, Robert Tibshirani, Jerome Friedman, and James Franklin.

The elements of statistical learning: data mining, inference and prediction, volume 27.
Springer, 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.