

# Computational Intelligence Laboratory

## Lecture 6

### Data Clustering and Mixture Models

**Thomas Hofmann**

ETH Zurich – `cil.inf.ethz.ch`

April 7, 2017

# Section 1

## Motivation

# Motivation: Data Clustering

- ▶ Given: set of data points  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$
- ▶ Goal: find a **meaningful partition** of the data
  - ▶ i.e. an **assignment** of each data point to a cluster

$$\pi : \{1, \dots, N\} \rightarrow \{1, \dots, K\} \quad \text{or}$$

$$\pi : \mathbb{R}^d \rightarrow \{1, \dots, K\}$$

- ▶ note: numbering of clusters is arbitrary
- ▶  $j$ -th cluster recovered by

$$\pi^{-1}(j) \subseteq \{1, \dots, n\} \quad \text{or} \quad \subseteq \mathbb{R}^d$$

# Motivation: Data Clustering

- ▶ Clustering via **similarity**:
  - ▶ group together similar data points, avoid grouping together dissimilar ones
  - ▶ uncover hidden group structure of data
  - ▶ learn a data density model
  - ▶ may give rise to data compression schemes

# Clustering Example

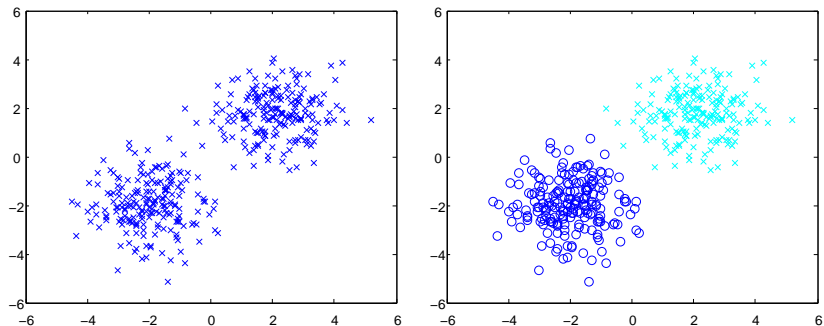


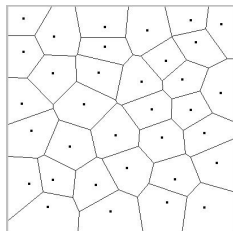
Figure: A simple clustering example. Left: 1 cluster, right: 2 clusters.

# Vector Quantization

- ▶ Partitioning of the space  $\mathbb{R}^d$
- ▶ Clusters represented by **centroids**  $\mathbf{u}_j \in \mathbb{R}^d$ .
- ▶ Mapping induced via nearest centroid rule

$$\pi(\mathbf{x}) = \arg \min_{j=1,\dots,k} \|\mathbf{u}_j - \mathbf{x}\|$$

- ▶ Voronoi (or Dirichlet) tessellation of  $\mathbb{R}^d$



# Color Reduction by Vector Quantization



**Figure:** Top: original images, Bottom: image represented with 10 colors, selected by clustering color vectors in RGB space.

## Section 2

### $K$ -Means



# Encoding via Indicators

- ▶ Formalize clustering problem as optimization problem
  - ▶ find centroids  $\mathbf{u}_j \in \mathbb{R}^d$  and assignment  $\pi$  minimizing ...
  - ▶ **loss function** or **distortion**, e.g. squared Euclidean norm
- ▶ Encode  $\pi$  via indicator matrix  $\mathbf{Z} \in \{0, 1\}^{n \times k}$

$$z_{ij} := \begin{cases} 1 & \text{if } \pi(\mathbf{x}_i) = j \\ 0 & \text{otherwise} \end{cases}$$

- ▶ note that

$$\sum_{j=1}^k z_{ij} = 1 \quad (\forall i)$$

# Objective Function

- $K$ -means objective function

$$\begin{aligned} J(\mathbf{U}, \mathbf{Z}) &= \sum_{i=1}^n \sum_{j=1}^k z_{ij} \|\mathbf{x}_i - \mathbf{u}_j\|^2 \\ &= \|\mathbf{X} - \mathbf{U}\mathbf{Z}^\top\|_F^2 \end{aligned}$$

- where

$\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ , data matrix

$\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_k] \in \mathbb{R}^{d \times k}$ , centroid matrix.

# $K$ -means Algorithm: Idea

- ▶ How do we minimize the  $K$ -means objective?
- ▶ Simple observation:
  - ▶ determining optimal centroids given assignments is easy (continuous variables)
  - ▶ determining optimal assignments given centroids is easy (integer variables)
- ▶ Computational strategy: alternating minimization

# *K*-means Algorithm: Optimal Assignment

- ▶ Compute optimal assignment  $\mathbf{Z}$ , given centroids  $\mathbf{U}$ 
  - ▶ each data point contributes to exactly one term in outer sum
  - ▶ minimize assignment of each data point separately

$$z_{ij}^* = \begin{cases} 1 & \text{if } j = \arg \min_s \|\mathbf{x}_i - \mathbf{u}_s\|^2 \\ 0 & \text{otherwise} \end{cases}$$

- ▶ map each data point to the closest centroid

# *K*-means Algorithm: Optimal Centroids

- ▶ Compute optimal choice of  $\mathbf{U}$ , given assignments  $\mathbf{Z}$ 
  - ▶ continuous variables: compute gradient and set to zero (1st order optimality condition)
  - ▶ look at (partial) gradient for every centroid  $\mathbf{u}_j$

$$\nabla_{\mathbf{u}_j} J(\mathbf{U}, \mathbf{Z}) = \sum_{i=1}^n z_{ij} \underbrace{\frac{1}{2} \nabla_{\mathbf{u}_j} \|\mathbf{x}_i - \mathbf{u}_j\|^2}_{=\mathbf{u}_j - \mathbf{x}_i}$$

- ▶ setting gradient to zero

$$\nabla_{\mathbf{U}} J(\mathbf{U}, \mathbf{Z}) \stackrel{!}{=} 0 \implies \mathbf{u}_j^* = \frac{\sum_{i=1}^n z_{ij} \mathbf{x}_i}{\sum_{i=1}^n z_{ij}}, \quad \text{if } \sum_{i=1}^n z_{ij} \geq 1$$

- ▶ **centroid condition** (center of mass of assigned data points)

# *K*-means Algorithm: Summary

initialize  $\mathbf{U}$  on  $k$  distinct random data points

initialize  $\mathbf{Z} \leftarrow \mathbf{Z}^*(\mathbf{U})$

**repeat**

$\mathbf{U} \leftarrow \mathbf{U}^*(\mathbf{Z})$  (see above)

$\mathbf{Z}^{\text{new}} \leftarrow \mathbf{Z}^*(\mathbf{U})$  (see above)

same = ( $\mathbf{Z}^{\text{new}} == \mathbf{Z}$ )

$\mathbf{Z} \leftarrow \mathbf{Z}^{\text{new}}$

**until** (same)

- ▶ different initialization strategies, here: random points
- ▶ better handling of empty clusters: random re-initialization

# *K*-means Algorithm:

- ▶ Computational cost of each iteration is  $O(knd)$
- ▶ *K*-means convergence is guaranteed
  - ▶ non-increasing objective, bounded from below by 0
- ▶ *K*-means optimizes a non-convex objective
  - ▶ we are not guaranteed to find the global optimum

# Illustration of the $K$ -means Algorithm

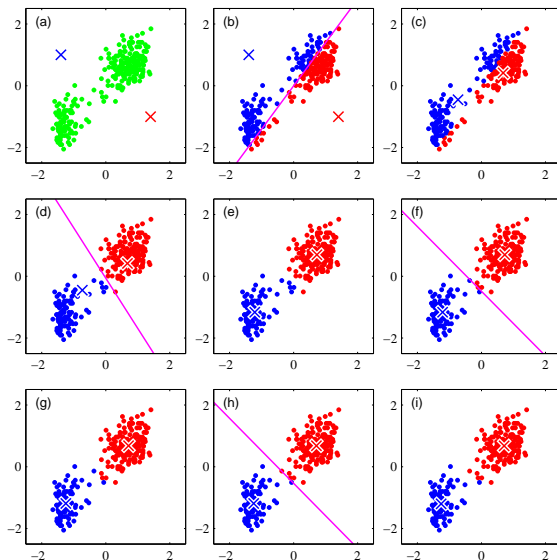


Figure: Bishop, *Pattern Recognition & Machine Learning*, Springer (2006).



## Section 3

### Mixture Models

# Probabilistic Clustering

From hard to probabilistic assignments

- ▶  $K$ -means: each data point assigned to exactly one cluster
- ▶ probabilistic or "soft" assignments: assign  $\mathbf{x}_i$  to each cluster  $j$  with some probability  $z_{ij}$
- ▶ generalize (relax) constraints on  $\mathbf{Z}$

$$z_{ij} \in [0; 1] \ (\forall i, j), \quad \sum_{j=1}^k z_{ij} = 1 \ (\forall i)$$

# Cluster Conditional Probability Distributions

Model each cluster by a **probability distribution**

- ▶ Simplest choice: multivariate normal distribution
- ▶ Univariate Gaussian with mean  $\mu$  and variance  $\sigma^2$ , density:

$$p(x; \mu; \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[ -\frac{(x - \mu)^2}{2\sigma^2} \right]$$

- ▶ Isotropic multivariate normal distribution with mean  $\boldsymbol{\mu}$ , density:

$$p(\mathbf{x}; \boldsymbol{\mu}; \sigma) = \prod_{i=1}^d \frac{1}{\sigma\sqrt{2\pi}} \exp \left[ -\frac{(x_i - \mu_i)^2}{2\sigma^2} \right]$$

# Cluster Conditional Probability Distributions

- ▶ Multivariate normal distribution with covariance matrix  $\Sigma$ , density:

$$p(\mathbf{x}; \boldsymbol{\mu}; \Sigma) = \frac{1}{|\Sigma|^{\frac{1}{2}} (2\pi)^{\frac{d}{2}}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

- ▶  $\Sigma$ : symmetric, positive definite
- ▶ generally difficult to estimate for large  $d$ ,  $(d+1)d/2$  parameters

# Probabilistic Clustering Model

- ▶ Finite **Mixture Model**

$$p(\mathbf{x}; \theta) = \sum_{j=1}^k \pi_j p(\mathbf{x}; \theta_j), \quad \theta = (\pi, \theta_1, \dots, \theta_k) \in \mathbb{R}^{k+k \cdot m}$$

- ▶ mixing proportions  $\pi \geq 0, \sum_{j=1}^k \pi_j = 1$
- ▶ component density functions  $p(\mathbf{x}; \theta_j)$  with  $\theta_j \in \mathbb{R}^m$
- ▶ Mixture models for clustering
  - ▶ relative cluster sizes =  $\pi_j$  ( $j = 1, \dots, k$ )
  - ▶ location & "shape" of clusters = specific form of  $p(\mathbf{x}; \theta_j)$
  - ▶ special case: Gaussian densities with,  $\theta_j = (\underbrace{\mu_j}_{\text{location}}, \underbrace{\Sigma_j}_{\text{shape}})$

# Gaussian Mixture Model

- ▶ Gaussian Mixture Model (GMM):

$$p(\mathbf{x}; \theta) = \sum_{j=1}^k \pi_j p(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \quad (\text{normal densities})$$

- ▶ Two-stage generative model: generate a data point as follows
  - ▶ sample cluster index from categorical distribution  
 $j \sim \text{Categorical}(\boldsymbol{\pi})$
  - ▶ given  $j$ , sample a data point  $\mathbf{x}$  from the  $j$ -th component  
 $\mathbf{x} \sim \text{Normal}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$
- ▶ Cluster index  $j$ : **latent variable**; final outcome  $\mathbf{x}$ : observed
- ▶ Probabilistic clustering: compute posteriors of latent cluster memberships ...

# Complete Data Distribution

- ▶ Explicitly introduce latent variables into generative model
- ▶ Assignment variable (for a generic data point)

$$\mathbf{z} \in \{0, 1\}^k, \quad \sum_{j=1}^k z_j = 1.$$

- ▶ Categorical distribution

$$\Pr(z_j = 1) = \pi_j \quad \text{or} \quad p_{\pi}(\mathbf{z}) = \prod_{j=1}^k \pi_j^{z_j}$$

- ▶ Joint distribution over  $(\mathbf{x}, \mathbf{z})$  (**complete data** distribution)

$$p(\mathbf{x}, \mathbf{z}; \theta) = \prod_{j=1}^k [\pi_j p(\mathbf{x}; \theta_j)]^{z_j}$$

# Posterior Assignments

- ▶ **Generation**: given  $\mathbf{z}$ , generate  $\mathbf{x}$ ; **Inference**: given  $\mathbf{x}$ , infer  $\mathbf{z}$
- ▶ Bayes rule
  - ▶ reminder, posterior  $p(A|B) = \frac{p(B|A)p(A)}{p(B)}$
  - ▶ here:  $p(A)$  prior,  $p(B|A)$  likelihood and  $p(B)$  evidence
- ▶ Posterior probabilities for assignments

$$\Pr(z_j = 1 \mid \mathbf{x}) = \frac{\Pr(z_j = 1)p(\mathbf{x} \mid z_j = 1)}{\sum_{l=1}^k \Pr(z_l = 1)p(\mathbf{x} \mid z_l = 1)} = \frac{\pi_j p(\mathbf{x}; \theta_j)}{\sum_{l=1}^k \pi_l p(\mathbf{x}; \theta_l)}$$

- ▶ assumes access to parameters  $\pi, \{\theta_j = (\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)\}$



# Maximum Likelihood: Mixture Model

- ▶ MLE requires to optimize

$$\hat{\theta} = \arg \max_{\theta} \sum_{i=1}^n \log \left[ \sum_{j=1}^k \pi_j p(\mathbf{x}_i; \theta_j) \right]$$

- ▶ Challenge: summation over  $j$  inside the logarithm

⇒ MLE has **no closed-form solution**

# Lower Bounding the Log-Likelihood

- ▶ Expectation Maximization
  - ▶ maximize a lower bound on the log-likelihood
  - ▶ based on complete data distribution

- ▶ Specifically:

$$\begin{aligned}\log p(\mathbf{x}; \theta) &= \log \left[ \sum_{j=1}^k \pi_j p(\mathbf{x}; \theta_j) \right] = \log \left[ \sum_{j=1}^k q_j \frac{\pi_j p(\mathbf{x}; \theta_j)}{q_j} \right] \\ &\geq \sum_{j=1}^k q_j [\log p(\mathbf{x}; \theta_j) + \log \pi_j - \log q_j]\end{aligned}$$

- ▶ follows from Jensen's inequality (concavity of logarithm)
- ▶ can be done for the contribution of each data point (additive)

# Mixture Model: Expectation Step

- ▶ Optimize bound with regard to the distribution  $q$ 
  - ▶ formulate Lagrangian (decoupled for each data point)

$$\max_q \left\{ \sum_{j=1}^k q_j [\log p(\mathbf{x}; \theta_j) + \log \pi_j - \log q_j] + \lambda \left( \sum_{j=1}^k q_j - 1 \right) \right\}$$

- ▶ first order optimality condition (setting gradient to zero):

$$\log p(\mathbf{x}; \theta_j) + \log \pi_j - \log q_j - 1 + \lambda \stackrel{!}{=} 0 \iff$$
$$q_j^* = \frac{\pi_j p(\mathbf{x}; \theta_j)}{\sum_{l=1}^k \pi_l p(\mathbf{x}; \theta_l)} = \Pr(z_j = 1 \mid \mathbf{x})$$

- ▶ optimal  $q$ -distribution equals posterior (given the parameters)
- ▶ E-step selects the best lower bound on the log-likelihood

# Mixture Model: Maximization Step

- ▶ Optimize expected complete data log-likelihood with regard to the model parameters
  - ▶ problem decouples for each cluster and with regard to  $\pi$
  - ▶ solution for mixing proportions  $\pi$

$$\pi_j^* = \frac{1}{n} \sum_{i=1}^n q_{ij}$$

- ▶ solution for  $\theta_j = (\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$

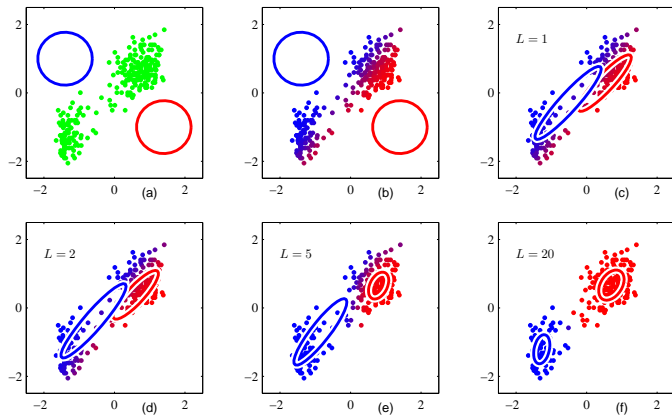
$$\boldsymbol{\mu}_j^* = \frac{\sum_{i=1}^n q_{ij} \mathbf{x}_i}{\sum_{i=1}^n q_{ij}}, \quad \boldsymbol{\Sigma}_j^* = \frac{\sum_{i=1}^n q_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j)(\mathbf{x}_i - \boldsymbol{\mu}_j)^\top}{\sum_{i=1}^n q_{ij}}$$

# Expectation Maximization Algorithm

- ▶ Alternate E-step and M-step
  - ▶ both E- and M-step maximize the same (bounded) objective
  - ▶ guaranteed convergence towards a point  $\theta^*$
  - ▶ like in  $K$ -means:  $\theta^*$  may not be the global maximizer
  - ▶ convergence criterion (e.g. change in objective)
- ▶ E-step: compute probabilistic assignments of points to clusters (keeping their location and shape fixed)
- ▶ M-step: recompute optimal cluster locations and shapes, given probabilistic assignments

# Example of EM for Gaussian Mixtures

Illustration of the EM algorithm using the Old Faithful data set.



**Figure:** Gaussian mixture model fitting via EM for two clusters. Remark: here the covariance is also estimated (illustrated by the two ellipsoids).

# Comparison with $K$ -means

- ▶ Assignments
  - ▶  $K$ -means algorithm: hard assignment points to clusters
  - ▶ EM algorithm: soft assignment based on posteriors
- ▶ Shapes
  - ▶  $K$ -means: spherical cluster shapes, uniform spread
  - ▶ EM algorithm: can learn covariance matrix
- ▶  $K$ -means as a special case
  - ▶ Gaussian mixture model with (fixed) covariances  $\Sigma_j = \sigma^2 \mathbf{I}$
  - ▶ in the limit of  $\sigma \rightarrow 0$ , recover  $K$ -means (hard assignments)
  - ▶ can be more formally derived (EM objective  $\rightarrow K$ -means objective)

# Practical Points about $K$ -means and EM

- ▶ EM algorithm
  - ▶ takes many more iterations to reach convergence
  - ▶ each cycle requires significantly more computation.
- ▶  $K$ -means algorithm can be used to find a good initialization
  - ▶ covariance matrices can be initialized to the sample covariances of the clusters found by the  $K$ -means algorithm.
  - ▶ mixing coefficients can be set to the fractions of data points assigned to the respective clusters



## Section 4

### Model Selection

# Occam's Razor

William Occam:

*Entities must not be multiplied  
beyond necessity.*



# Model order selection: General principle

Trade-off between two conflicting goals:

**Data fit:** We want to predict the data well, e.g., maximize the likelihood. The likelihood usually increases with increasing number of clusters.

**Complexity:** Choose a model that is not very complex which is often measured by the number of free parameters.

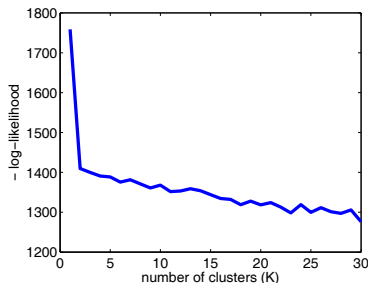
Find a compromise between these two goals!

# Better fit with increasing $K$

Negative Log-Likelihood of data for  $K$  mixture Gaussians:

$$-\log p(\mathbf{X}; \theta) = -\sum_{i=1}^n \log \left[ \sum_{j=1}^k \pi_j p(\mathbf{x}_i; \theta_j) \right].$$

- ▶ smaller negative log-likelihood = better fit
- ▶ decreasing with  $k$  (some noise due to local minima)



# AIC and BIC

- ▶ Model complexity: can be measured by the number of free parameters  $\kappa(\cdot)$ .
- ▶ Different Heuristics for choosing  $k$ 
  - ▶ Akaike Information Criterion (AIC)

$$AIC(\theta|\mathbf{X}) = -\log p(\mathbf{X}; \theta) + \kappa(\theta)$$

- ▶ Bayesian Information Criterion (BIC)

$$BIC(\theta|\mathbf{X}) = -\log p(\mathbf{X}; \theta) + \frac{1}{2}\kappa(\theta) \log n$$

- ▶ Generally speaking, the BIC criterion penalizes complexity more than the AIC criterion.

# AIC and BIC: Remarks and Example

## Analysis

A single AIC (BIC) result is meaningless. One has to repeat the analysis for different  $k$ s and compare the differences: the most suitable number of clusters corresponds to the smallest AIC (BIC) value.

## Example (Mixture of Gaussians)

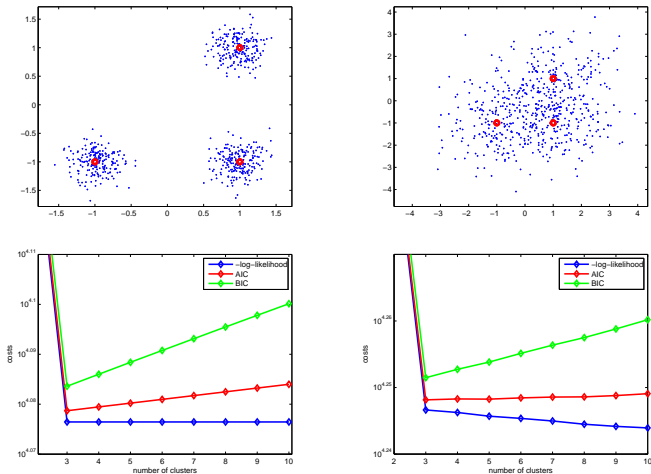
- Number of free parameters (with fixed covariance matrices)

$$\kappa(\theta) = k \cdot d + (k - 1).$$

- Number of free parameters (with full covariance matrices)

$$\kappa(\theta) = k \cdot \left( d + \frac{d(d+1)}{2} \right) + (k - 1).$$

# AIC and BIC example: 3 clusters



**Figure:** Information criteria for a synthetic dataset with 3 clusters. Synthetic data has smaller variance on the left than on the right.

# AIC and BIC example: 5 clusters

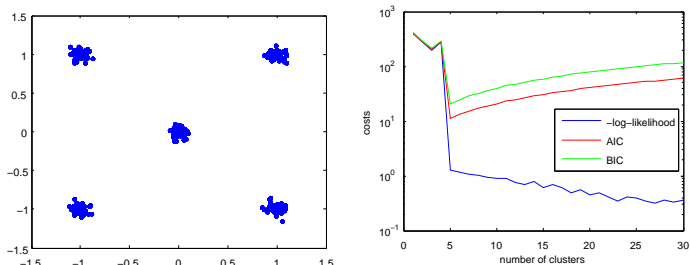


Figure: Information criteria for a synthetic dataset with 5 clusters.