

Machine Learning

Lecture 11: Clustering

Prof. Dr. Stephan Günnemann

Data Mining and Analytics
Technical University of Munich

21.01.2019

Reading material

Reading material

- Bishop: chapters 9.1, 9.2, 9.3.0, 9.3.1

Unsupervised learning

Main idea

- Given some **unlabeled** data $\{x_i\}$ we want to discover **latent** structure in it.

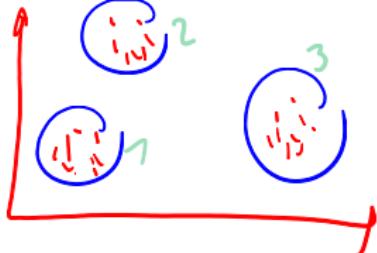
Last week: Dimensionality reduction

- Given high-dimensional data $x_i \in \mathbb{R}^D$, find a **continuous** representation $z_i \in \mathbb{R}^K$ with $K \ll D$, such that the structure is preserved.

This week: Clustering

- Group objects x_i into K groups (**clusters**) based on their similarity.

Problem definition



Given

- Collection of data points $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subset \mathcal{X}$.

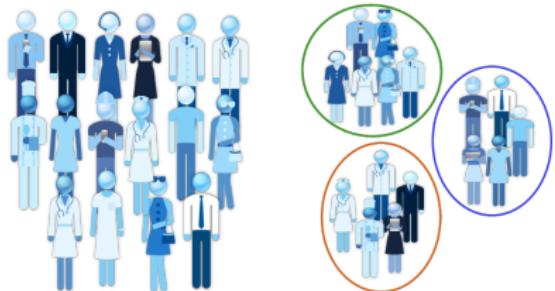
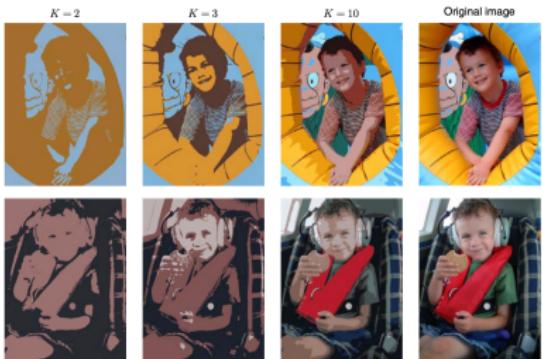
Goal

- For each object \mathbf{x}_i find the corresponding assignment $\mathbf{z}_i \in \mathbb{R}^K$ to one of the K groups (**clusters**), such that:
 - objects within the same group are **similar** to each other;
 - objects between different groups are **dissimilar** from each other;

\mathbf{z}_i is a one-hot indicator vector, with $z_{ik} = 1 \iff$ point i belongs to cluster k .

Clustering: Examples

- Image segmentation
- User profiling
- Data compression
- Visualization
- ...



Sources: <https://blogs.sas.com/content/subconsciousmusings/2016/05/26/data-mining-clustering/>
and C. Bishop - "Pattern Recognition and Machine Learning"

Section 1

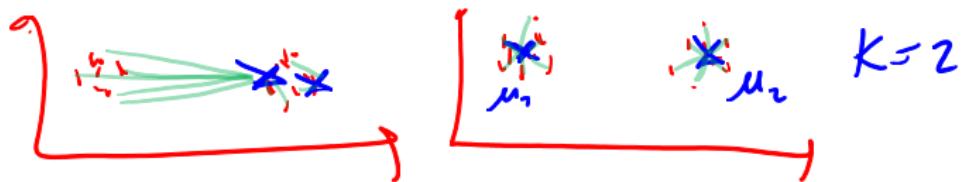
K-means Algorithm

Distance-based clustering

- If we want to group **similar** objects together, we need some notion of **similarity** to begin with.
- Typically defined as **similarity function** $s(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$.
- Alternatively: minimize dissimilarity, usually provided by a **distance function** $d(\cdot, \cdot)$.
- Typical distance functions include
 - Manhattan distance: $d(\mathbf{x}_i, \mathbf{x}_j) = \sum_d |x_{id} - x_{jd}|$
 - Euclidean distance: $d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_d (x_{id} - x_{jd})^2}$
 - Mahalanobis distance: $d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \mathbf{x}_j)}$

K-means

Model



- Consider the standard \mathbb{R}^D space with Euclidean distance as metric.
- Each of K clusters is defined by its **centroid** $\mu_k \in \mathbb{R}^D$.
- **Cluster indicator** $z_i \in \{0, 1\}^K$ with $\|z_i\| = 1$ (i.e. one-hot vector) shows which cluster point i belongs to.
- The **K-means objective** (also called **distortion measure**) is defined as

$$J(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}) = \sum_{i=1}^N \sum_{k=1}^K z_{ik} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2$$

// CONSTRAINT
REGULARIZATION
// IDEA:
// $z_i \in [0, 1]^K$
// $\sum z_{ik} = 1$

Goal

- We need to find the "best" centroids $\boldsymbol{\mu}$ and cluster assignments \mathbf{Z}

$$\mathbf{Z}^*, \boldsymbol{\mu}^* = \arg \min_{\mathbf{Z}, \boldsymbol{\mu}} J(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu})$$

Minimizing the K -means objective

- The joint optimization problem is difficult to solve (discrete optim.; even if Z is assumed to be continuous it is non-convex)

$$\min_{Z, \mu} J(X, Z, \mu) = \min_{Z, \mu} \sum_{i=1}^N \sum_{k=1}^K z_{ik} \|x_i - \mu_k\|_2^2$$

- However, the two subproblems

$$\min_Z J(X, Z, \mu) \quad \text{and} \quad \min_\mu J(X, Z, \mu)$$

have simple closed form solutions!

- Idea:** alternating optimization - update Z and μ in turns!
- This method for solving K -means is called Lloyd's algorithm.

Lloyd's algorithm for K -means

1. Initialize the centroids $\mu = \{\mu_1, \dots, \mu_K\}$.
2. Update cluster indicators (solve $\min_{\mathbf{Z}} J(\mathbf{X}, \mathbf{Z}, \mu)$)

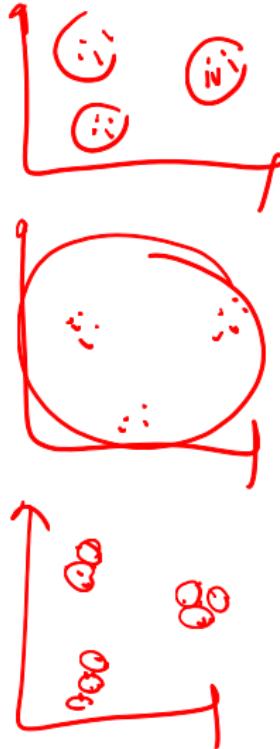
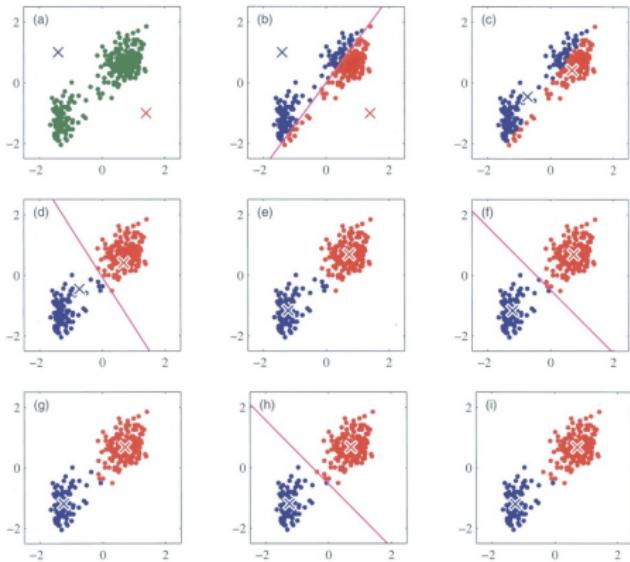
$$z_{ik} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_i - \mu_j\|^2, \\ 0 & \text{else.} \end{cases}$$

3. Update centroids (solve $\min_{\mu} J(\mathbf{X}, \mathbf{Z}, \mu)$)

$$\mu_k = \frac{1}{N_k} \sum_{i=1}^N z_{ik} \mathbf{x}_i, \quad \text{where } N_k = \sum_{i=1}^N z_{ik}$$

4. If objective $J(\mathbf{X}, \mathbf{Z}, \mu)$ has not converged, return to step 2.

K-means in action



An interactive demo is available at

- <http://stanford.edu/class/ee103/visualizations/kmeans/kmeans.html>

Source: Bishop, Figure. 9.1

Initialization of the centroids



We saw in the interactive demo that the initial locations of the centroids μ_k greatly affect the performance.

So, how do we initialize the centroids in a good way?

K-means++ algorithm

1	20
1	30
1.5	25
2	
2.1	27

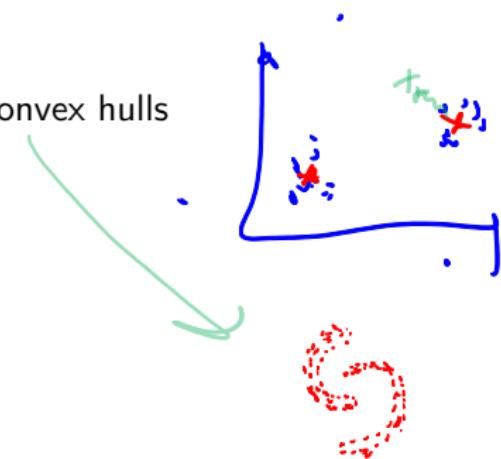
1. Choose the first centroid μ_1 uniformly at random among the data points.
2. For each point x_i compute the distance $D_i^2 = \|x_i - \mu_1\|_2^2$.
3. Sample the next centroid μ_k from ~~remaining~~ points with probability proportional to D_i^2 .
4. Recompute the distances $D_i^2 = \min\{\|x_i - \mu_1\|_2^2, \dots, \|x_i - \mu_k\|_2^2\}$.
5. Continue steps 3 and 4 until K initial centroids have been chosen.

Limitations of K -means

It is important to distinguish between the **model** (distortion $J(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu})$ in our case) and the **algorithm** used to fit it (here Lloyd's algorithm).

Modeling issues

- Underlying assumptions are not explicit
- Can't detect clusters with overlapping convex hulls
- Sensitivity to outliers
- No uncertainty measure



Algorithmic issues

- Extreme sensitivity to initialization

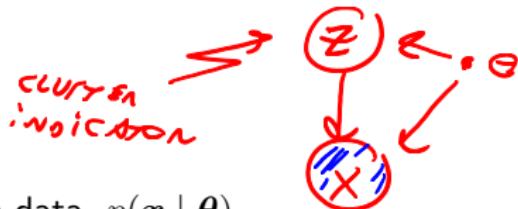
How can we address (some of) these issues?

Switch to the more principled probabilistic framework!

Section 2

Gaussian Mixture Model

Probabilistic unsupervised learning



- Goal: Design a probabilistic model for the data, $p(x | \theta)$, parametrized by θ , which we have to estimate.
- It's difficult to come up with a useful model $p(x | \theta)$ without making any assumptions about the data.
- When doing clustering, we implicitly assume that each point x belongs to some cluster (indicated by z).
⇒ Let's model $p(x, z | \theta)$ instead! $p(x|z, \theta) \cdot p(z|\theta)$
- Let's recall what tools we developed during our discussion of supervised learning, and see if any of them can help us here.
- For this, assume for now, that we know the true cluster indicators z , and are only interested in estimating θ .

Parameter estimation

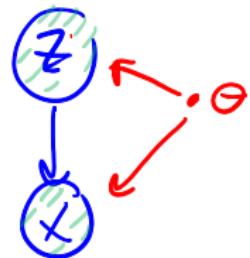
$$\Rightarrow p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\theta}) \cdot p(\mathbf{Z} | \boldsymbol{\theta})$$

How can we estimate the parameters $\boldsymbol{\theta}$ of a given model $p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$?

Supervised learning

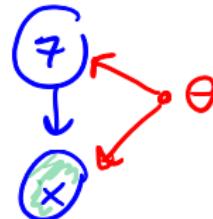
- Both $\mathbf{X} \in \mathbb{R}^{N \times D}$ and $\mathbf{Z} \in \{0, 1\}^{N \times K}$ are known.
- Simple solution - maximum likelihood estimation

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$$



Unsupervised learning

- Only the data $\mathbf{X} \in \mathbb{R}^{N \times D}$ is given.
- Is there any way to estimate $\boldsymbol{\theta}$ in this case?



Latent variables

How can we estimate θ when Z is unknown?

We can compute the joint probability $p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$, so we can obtain $p(\mathbf{X} | \boldsymbol{\theta})$ by marginalization

$$p(\mathbf{X} | \boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\theta}) \cdot p(\mathbf{Z} | \boldsymbol{\theta})$$

and simply optimize w.r.t $\boldsymbol{\theta}$

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} p(\mathbf{X} | \boldsymbol{\theta})$$



Since true Z are never observed, and are only our abstraction used to simplify the model, they are called **latent variables**.¹

¹Sometimes also called hidden or explanatory variables.

Gaussian mixture model (GMM)

We decompose the joint probability $p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})$ using the product rule ²

$$p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) = p(\mathbf{z} | \boldsymbol{\theta}) \cdot p(\mathbf{x} | \mathbf{z}, \boldsymbol{\theta})$$

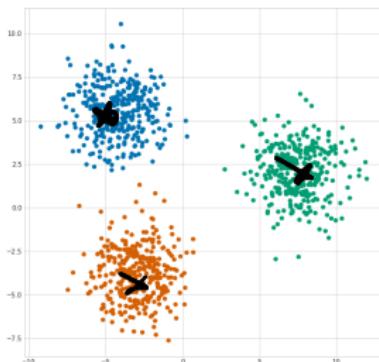
In Gaussian mixture model we assume:

- Cluster prior is **categorical**

$$p(\mathbf{z} | \boldsymbol{\theta}) = \text{Cat}(\boldsymbol{\pi})$$

- Each cluster has its own **multivariate normal** distribution

$$p(\mathbf{x} | z_k = 1, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$



²We assume that samples are drawn i.i.d., so consider only a single (\mathbf{x}, \mathbf{z}) .

Generative process of GMM

We can draw data from a GMM as following.

Setup

We have K Gaussian components (clusters), each with its own $\{\mu_k, \Sigma_k\}$.

Prior probability of each cluster k is π_k .

Generative process

1. Draw a one-hot cluster indicator $z \sim \text{Cat}(\boldsymbol{\pi})$.
 $z_k = 1 \iff$ current point belongs to cluster k .
2. Draw the sample $x \sim \mathcal{N}(\mu_k, \Sigma_k)$ if $z_k = 1$.

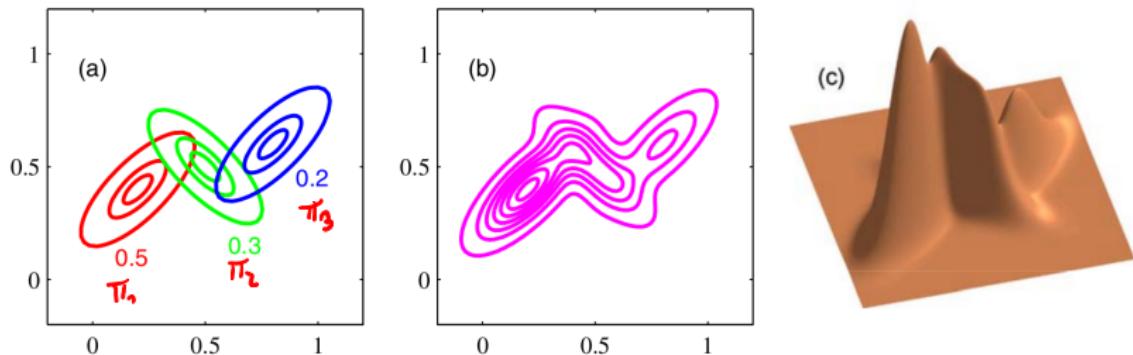
We only observe the sample x , and don't know which k it came from.

Using marginalization over z , the probability of a single sample \mathbf{x} under a Gaussian mixture model can be computed as

$$\begin{aligned} p(\mathbf{x} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \sum_{k=1}^K p(z_k = 1 \mid \boldsymbol{\pi}) \cdot p(\mathbf{x} \mid z_k = 1, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ &= \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{aligned}$$

Hence, for the entire dataset $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ the log-likelihood is

$$\log p(\mathbf{X} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^N \log \left(\sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)$$



- (a) Isocontours of each component $p(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ with respective mixing coefficients π_k .
- (b) Isocontours of the GMM distribution

$$p(\mathbf{x} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- (c) Surface plot of $p(\mathbf{x} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$.

Source: Bishop, Figure. 2.23

Learning and inference in GMM

Remember, that we are interested in two things

Learning

- We want to determine the "optimal" values of the parameters (e.g., find values that maximize the log-likelihood)

$$\boldsymbol{\pi}^*, \boldsymbol{\mu}^*, \boldsymbol{\Sigma}^* = \arg \max_{\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}} \log p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

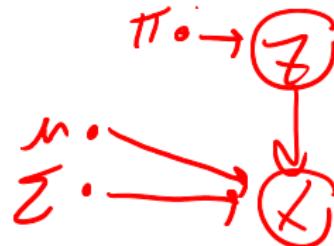
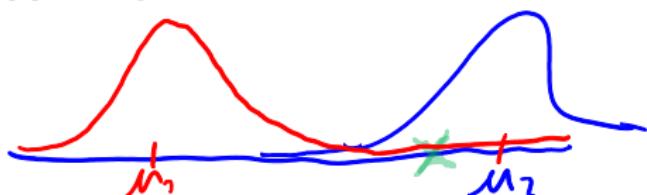
Inference

- Given the parameters $\{\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\}$, we want to assign the points to clusters, i.e. compute the posterior distribution of cluster indicators

$$p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

(This is the actual goal of clustering.)

Inference in GMM



By straightforward application of the Bayes formula, we obtain the posterior over the entries of Z

$$p(z_{ik} = 1 \mid \mathbf{x}_i, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_i \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \rightrightarrows (k)$$

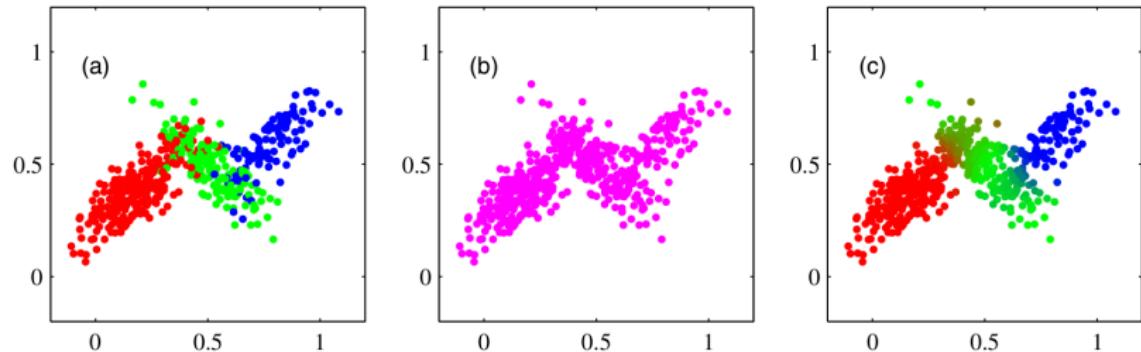
We will call this quantity the **responsibility** of component k for the observation i .

We also introduce shorthand notation for it

$$\gamma(z_{ik}) = p(z_{ik} = 1 \mid \mathbf{x}_i, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$(*) \propto \pi(z_{ik}=1|\boldsymbol{\pi}) \cdot p(\mathbf{x}_i | z_{ik}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Inference: Example



- (a) Data colored according to the true underlying Z .
- (b) Observed data X .
- (c) Data colored according to the inferred $\gamma(Z)$ (with π, μ, Σ known).

$$\hat{p}(z|x, \pi, \mu, \Sigma)$$

Source: Bishop, Figure. 9.5

Learning in GMM

Maximization of the log-likelihood w.r.t. the parameters appears to be more challenging. Recall, that the log-likelihood is

$$\log p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^N \log \left(\sum_{k=1}^K \pi_k \cdot \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)$$

Due to the summation over k , the logarithm doesn't act directly on the Gaussian density.

This happens to be a serious problem, as this means that the derivatives w.r.t. the model parameters $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$ have no analytical roots.³

Is there anything that we can do in this case?

³I.e. we cannot obtain closed form expressions for the optimal values of $\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$.

Learning in GMM

LinEnn

Recall from Linear Classification lecture (Gaussian discriminant analysis), that parameter estimation is rather straightforward if Z is known.

That is, we know how to solve the "easy" problem

$$\pi^*, \mu^*, \Sigma^* = \arg \max_{\pi, \mu, \Sigma} \log p(\mathbf{X}, \mathbf{Z} \mid \pi, \mu, \Sigma) \quad (*)$$

Learning in GMM

Recall from Linear Classification lecture (Gaussian discriminant analysis), that parameter estimation is rather straightforward if Z is known.

That is, we know how to solve the "easy" problem

$$\pi^*, \mu^*, \Sigma^* = \arg \max_{\pi, \mu, \Sigma} \log p(\mathbf{X}, \mathbf{Z} \mid \pi, \mu, \Sigma) \quad (*)$$

Too bad we don't know anything about Z ... Or do we?

Learning in GMM

..

Recall from Linear Classification lecture (Gaussian discriminant analysis), that parameter estimation is rather straightforward if Z is known.

That is, we know how to solve the "easy" problem

$$\pi^*, \mu^*, \Sigma^* = \arg \max_{\pi, \mu, \Sigma} \log p(X, Z | \pi, \mu, \Sigma) \quad (*)$$

Too bad we don't know anything about Z ... Or do we?

Idea: use our current beliefs $p(Z | X, \pi^*, \mu^*, \Sigma^*)$ and solve the easier problem (*).

(
n)

Important: As we are dealing with a distribution over Z , we have to take the expectation to account for the uncertainty

$$\pi^*, \mu^*, \Sigma^* = \arg \max_{\pi, \mu, \Sigma} \mathbb{E}_{Z \sim p(Z)} [\log p(X, Z | \pi, \mu, \Sigma)]$$

$$= \underset{\pi, \mu, \Sigma}{\text{AUC MAX}} \sum_z p(z | X, \pi^*, \mu^*, \Sigma^*) \cdot \log p(x, z | \pi^*, \mu^*, \Sigma^*)$$

Expectation-Maximization algorithm

This two-step procedure that we just derived is called **Expectation-Maximization algorithm (EM)**.

In its most general form, the two steps are

- **E step** - evaluate the posterior $\gamma(\mathbf{Z}) = p(\mathbf{Z} \mid \mathbf{X}, \boldsymbol{\theta})$.
- **M step** - maximize the current estimate of $p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\theta})$ w.r.t. $\boldsymbol{\theta}$

$$\underset{\boldsymbol{\theta}}{\text{arg max}} \mathbb{E}_{\mathbf{Z} \sim \gamma(\mathbf{Z})} [\log p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\theta})]$$

()*

at t+1 $\boldsymbol{\theta}^{new} = \boldsymbol{\theta}^{(t)}$

Note: we do **not** get a closed form solution neither for $\gamma(\mathbf{Z})$, nor for the parameters, as $\gamma(\mathbf{Z})$ and $\boldsymbol{\theta}$ recursively depend on each other.

Therefore, we have to iterate between these two steps until the objective $\mathbb{E}_{\mathbf{Z} \sim \gamma(\mathbf{Z})} [\log p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\theta})]$ converges.

EM algorithm for GMM

$$N_k = \sum_{i=1}^n \gamma(z_{ik})$$

1. Initialize model parameters $\theta = \{\pi, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K\}$.
2. **E step.** Evaluate the responsibilities for current θ

$$\gamma(z_{ik}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}.$$

3. **M step.** Re-estimate the parameters using current $\gamma(z_{ik})$'s

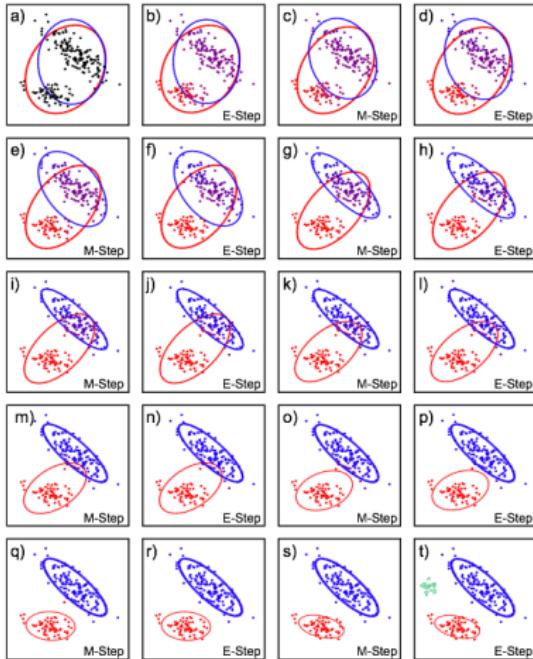
$$\boldsymbol{\mu}_k^{new} = \frac{1}{N_k} \sum_{i=1}^N \gamma(z_{ik}) \mathbf{x}_i \quad \text{(WEIGHTED MEAN)}$$

$$\boldsymbol{\Sigma}_k^{new} = \frac{1}{N_k} \sum_{i=1}^N \gamma(z_{ik}) (\mathbf{x}_i - \boldsymbol{\mu}_k^{new})(\mathbf{x}_i - \boldsymbol{\mu}_k^{new})^T$$

$$\pi_k^{new} = \frac{N_k}{N}.$$

4. If $\mathbb{E}_{\mathbf{Z} \sim \gamma(\mathbf{Z})} [\log p(\mathbf{X}, \mathbf{Z} | \theta)]$ has not converged, return to step 2.

EM algorithm for GMM: Demo



Gif: https://upload.wikimedia.org/wikipedia/commons/6/69/EM_Clustering_of_Old_Faithful_data.gif

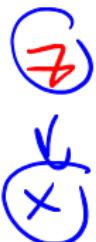
Source: Computer Vision: Models, Learning, and Inference - S.J.Prince

EM algorithm: Summary

When should we use EM?

- Intractable to optimize the log-likelihood

$$\log p(\mathbf{X} \mid \boldsymbol{\theta}) = \log \left(\sum_{\mathbf{Z}} p(\mathbf{Z} \mid \boldsymbol{\theta}) p(\mathbf{X} \mid \mathbf{Z}, \boldsymbol{\theta}) \right)$$



- Easy to optimize the joint probability

$$\log p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\theta}) = \log p(\mathbf{Z} \mid \boldsymbol{\theta}) + \log p(\mathbf{X} \mid \mathbf{Z}, \boldsymbol{\theta})$$

The main idea of EM

- Use our current beliefs $p(\mathbf{Z} \mid \mathbf{X}, \boldsymbol{\theta})$ to "pretend" that we know \mathbf{Z} .
- E step - update the responsibilities $\gamma(\mathbf{Z}) = p(\mathbf{Z} \mid \mathbf{X}, \boldsymbol{\theta})$.
- M step - update parameters

$$\boldsymbol{\theta}^{new} = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{Z} \sim \gamma(\mathbf{Z})} [\log p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\theta})]$$