**Machine Learning Course - CS-433**

# Gaussian Mixture Models
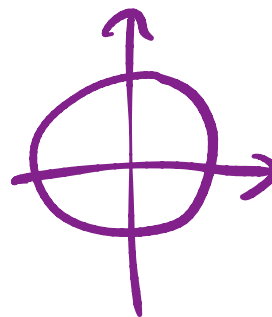
Nov 29, 2022

Martin Jaggi
Last updated on: November 28, 2022

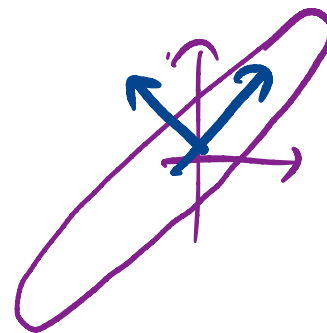credits to Mohammad Emtiyaz Khan & Rüdiger Urbanke

**EPFL**

# Motivation

K-means forces the clusters to be *spherical,* but sometimes it is desirable to have *elliptical* clusters. Another issue is that, in K-means, each example can only belong to one cluster, but this may not always be a good choice, e.g. for data points that are near the "border". Both of these problems are solved by using Gaussian Mixture Models.

$\Sigma = \mathbb{1}$

$\Sigma$ general

# Clustering with Gaussians

The first issue is resolved by using full covariance matrices $\Sigma_k$ instead of *isotropic* covariances.

① $$p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{z}) = \prod_{n=1}^{N} \prod_{k=1}^{K} [\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_{nk}}$$

parameters:

$\boldsymbol{\mu} \in \mathbb{R}^{D \cdot K}$

$\boldsymbol{\Sigma} \in \mathbb{R}^{D \cdot D \cdot K}$

$\boldsymbol{\pi} \in \mathbb{R}^{K}$

# Soft-clustering

The second issue is resolved by defining $z_n$ to be a random variable. Specifically, define $z_n \in \{1, 2, \ldots, K\}$ that follows a multinomial distribution.

② random variable (vector 0... 1... 0)

$$p(z_n = k) = \pi_k \text{ where } \pi_k > 0, \forall k \text{ and } \sum_{k=1}^{K} \pi_k = 1$$
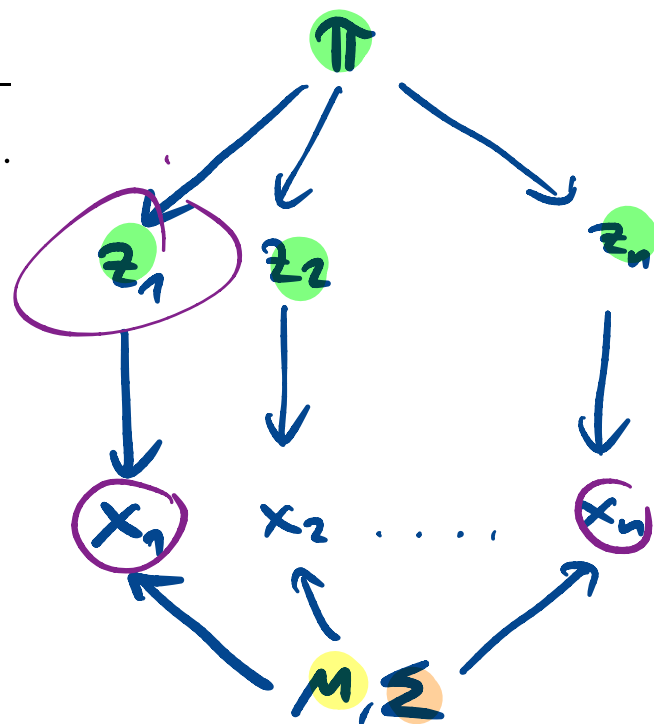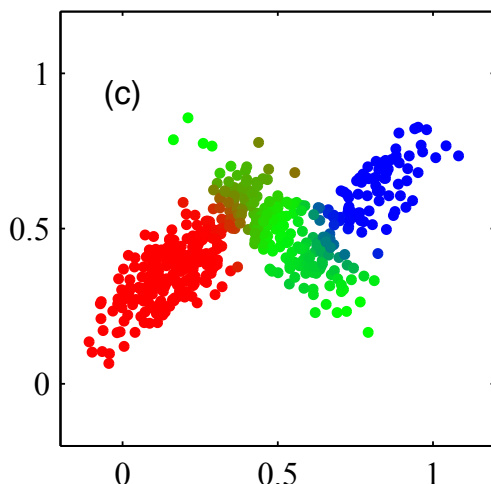
$0, \ldots 1 \ldots 0$
$\uparrow$ k-th entry

importance of cluster k

$z_n = (0 \dots 1 \dots 0)$

$z_{nk} = \begin{cases} 1 \\ 0 \end{cases}$   $\tau_k$

This leads to soft-clustering as opposed to having "hard" assignments.



(c)

## Gaussian mixture model

Together, the likelihood and the prior define the joint distribution of Gaussian mixture model (GMM):

joint

$$p(\mathbf{X}, \mathbf{z} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})$$

$$= \prod_{n=1}^{N} \underbrace{p(\mathbf{x}_n \mid z_n, \boldsymbol{\mu}, \boldsymbol{\Sigma})}_{\text{likelihood}} \underbrace{p(z_n \mid \boldsymbol{\pi})}_{\text{prior}}$$

$$= \prod_{n=1}^{N} \prod_{k=1}^{K} [\mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_{nk}} \prod_{k=1}^{K} [\pi_k]^{z_{nk}}$$
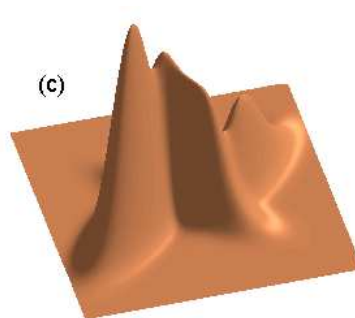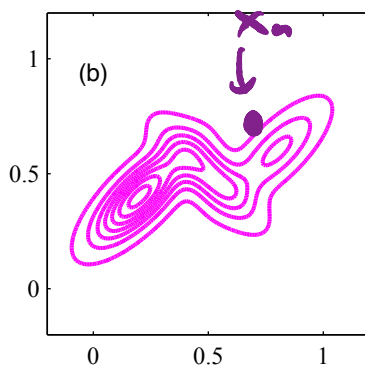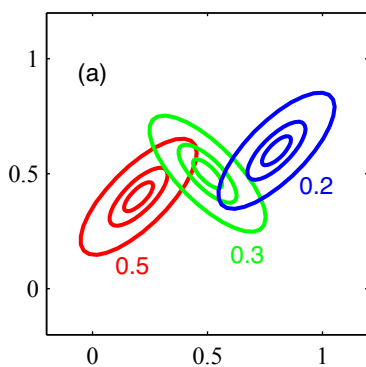
Here, $\mathbf{x}_n$ are observed data vectors, $z_n$ are *latent* unobserved variables, and the unknown *parameters* are given by $\boldsymbol{\theta} := \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K, \boldsymbol{\pi}\}$.

Bayes Rule:

$$P(a,b \mid) = P(a \mid b)\, P(b)$$

# Marginal likelihood

GMM is a latent variable model with $z_n$ being the unobserved (latent) variables. An advantage of treating $z_n$ as latent variables instead of *parameters* is that we can *marginalize* them out to get a cost function that does not depend on $z_n$, i.e. as if $z_n$ never existed.

Specifically, we get the following marginal likelihood by marginalizing $z_n$ out from the likelihood:

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

joint:
$$p(x_n, z_n)$$

marginal:
$$p(x_n) = \sum_{k=1}^{K} p(x_n, z_n = k)$$
$$= \sum_k p(x_n | \textcircled{z}) \, p(\textcircled{z})$$

N of cluster k     $\pi_k$

$\mu, \Sigma, \pi$



(a)   0.5   0.3   0.2

(b)   $x_n$

(c)

Deriving cost functions this way is good for *statistical efficiency*. Without a latent variable model, the number of parameters grows at rate $\mathcal{O}(N)$. After marginalization, the growth is reduced to $\mathcal{O}(D^2 K)$ (assuming $D, K \ll N$).

$z: \quad N \cdot K$

$\theta: \mu \quad K \cdot D$
$\quad \Sigma \quad K \cdot D^2$
$\quad \pi \quad K$

# Maximum likelihood

To get a maximum (marginal) likelihood estimate of $\boldsymbol{\theta}$, we maximize the following: $\log(P(X|\theta))$

$$\max_{\boldsymbol{\theta}} \sum_{n=1}^{N} \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

marginal    likelihood

$$\log \left/ p\left(\begin{smallmatrix} \text{all} \\ x_n \end{smallmatrix} \middle| \theta\right) \right.$$

$$= \prod_{n=1}^{N} p(x_n | \theta)$$

$$\sum_k \pi_k \mathcal{N}(x_n | M_k, \Sigma_k)$$

Is this cost convex? Identifiable? Bounded?

① non-convex (see k-means)

② non-unique optima

permutation of $1 \ldots k$

$k \to k' \quad \pi_k \to \pi_{k'}$
$M_k \to M_{k'}$
$\Sigma_k \to \Sigma_{k'}$

③ unbounded

$$\Sigma_k = \sigma_k \mathbf{I}$$

↑ scalar width



$p(x)$

$M_1$     $M_2$

width $\sigma_k \to 0$