Annotated Version

**Machine Learning Course - CS-433**

# Gaussian Mixture Models

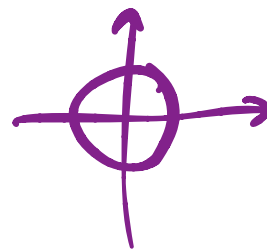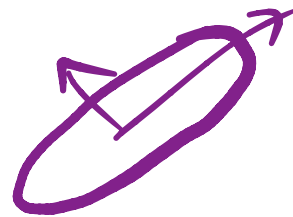Nov 30, 2021

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.

*spherical* $\Sigma = 1$

*general* $\Sigma$ (Ellipse)

# Clustering with Gaussians

The first issue is resolved by using full covariance matrices $\boldsymbol{\Sigma}_k$ instead of *isotropic* covariances.

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

# 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.
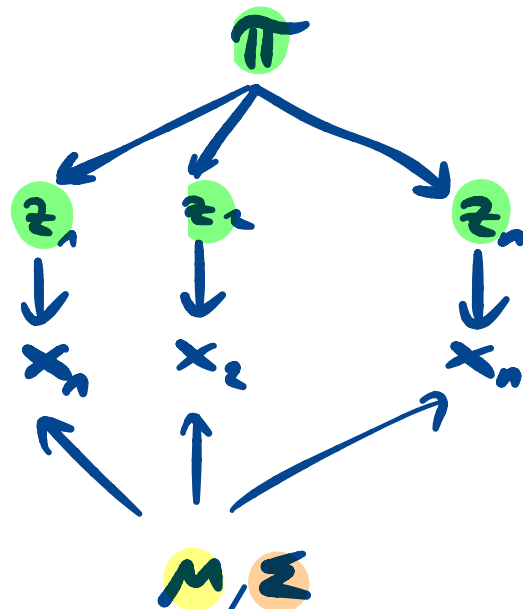
*parameters*
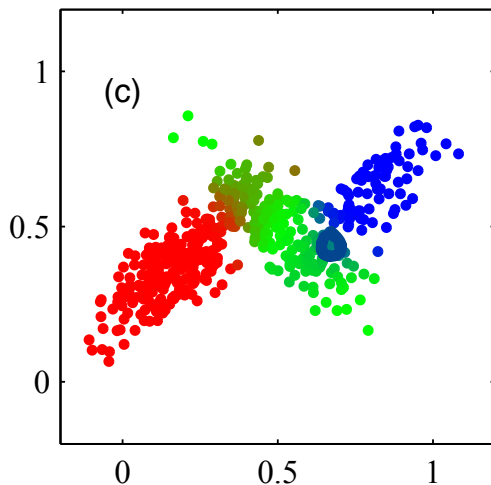- $\boldsymbol{\mu} \in \mathbb{R}^{D \cdot K}$
- $\boldsymbol{\Sigma} \in \mathbb{R}^{D^2 K}$
- $\boldsymbol{\pi} \in \mathbb{R}^{K}$

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

importance of cluster $k$

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

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

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

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

Here, $\mathbf{x}_n$ are observed data vectors, $z_n$ are *latent* unobserved variables, and the unknown *parameters* are given by $\boldsymbol{\theta}$ := $\{\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \ldots, \boldsymbol{\Sigma}_K, \boldsymbol{\pi}\}$.

Bayes Rule

$$p(a, b) = p(a|b) \cdot p(b)$$

$$z_n = (0, \ldots, 1, \ldots 0)$$

$$z_{nk} = \begin{cases} 0 & \ldots \\ 1 & \ldots \end{cases} \underset{k}{\sim}$$
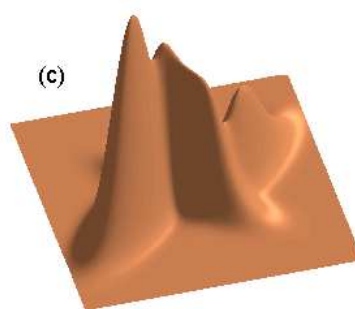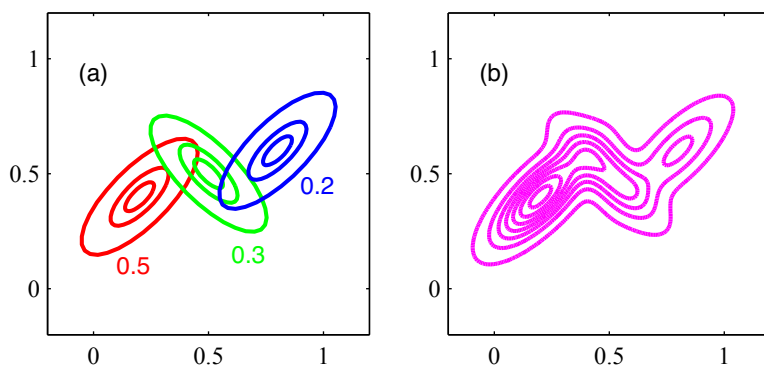
# Marginal likelihood

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.

*orig. likelihood* $\quad p(x_n, z_n \mid \theta)$

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*

$p(x_n) = \sum_{k=1}^{K} p(x, z_n = k)$

$= \sum p(x|z) \cdot p(z=k)$

$\quad\quad\quad \pi_k$

*Bayes*



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 $D, K \ll N$).

$z : \quad \mathcal{O}(N)$

$\theta : \quad \mu : KD$

$\quad\quad\quad \Sigma : KD^2$

$\quad\quad\quad \pi : K$

# Maximum likelihood

To get a maximum (marginal) like-lihood estimate of $\boldsymbol{\theta}$, we maximize the following:

$$\text{marginal likelihood}$$

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

$\theta = (\mu, \Sigma, \pi)$

$$\log\left( p(\overset{all}{x_n} | \theta) = \prod_{n=1}^{N} p(x_n | \theta) \right)$$

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

$$= \sum_{n=1}^{N} \log \sum_{k}^{K} \pi \mathcal{N}(\ldots)$$

Is this cost convex?  Identifiable?
Bounded?
no                      no

no

$p(x)$



$x$

$\mathcal{N}_1$

$\mathcal{N}_2$

$\sigma$ width
$\Sigma$

① non-convex
(see k-means)

② non-unique optima

permutation of 1..K

$k \to k'$     $\pi_k \to \pi_{k'}$
$\Sigma_k \to \Sigma_{k'}$
$\mu_k \to \mu_{k'}$

③ unbounded
$-\mathcal{L} \to \infty$

if $\Sigma = \sigma \cdot \mathbf{1}$
                ↑
              width
and $\sigma \to 0$