

Annotated  
Version

Machine Learning Course - CS-433

# Gaussian Mixture Models

Nov 9, 2017

©Mohammad Emtiyaz Khan 2015

minor changes by Martin Jaggi 2016

minor changes by Martin Jaggi 2017

Last updated: November 9, 2017



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

# Motivation

- ① 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.
- ②

## Clustering with Gaussians

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

for k-means:  $\Sigma_k = I$

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

## Soft-clustering vs hard assignment

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

	parameters
$\mathbb{R}^{D \times K}$	$\mu = (\mu_1, \dots, \mu_K)$
$\mathbb{R}^{D \times D \times K}$	$\Sigma = (\Sigma_1, \dots, \Sigma_K)$
$\mathbb{R}^K$	$\pi = (\pi_1, \dots, \pi_K)$

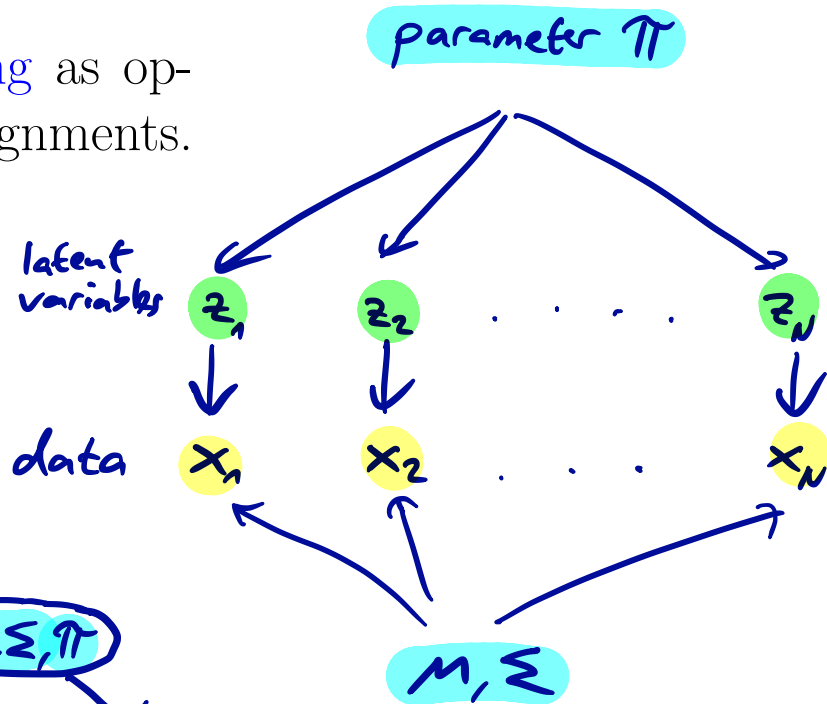
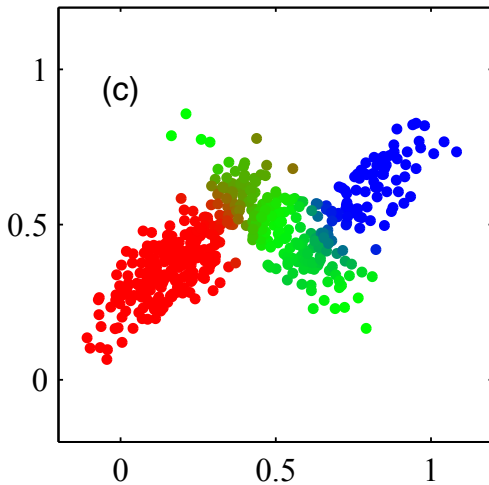
importance weight of group k

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

$$\mathbf{z}_n = (0, \dots, \overset{k}{1}, \dots, 0)$$

with  $1P = \pi_k$

This leads to **soft-clustering** as opposed to having “hard” assignments.



## Gaussian mixture model

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

$$\begin{aligned}
 \rightarrow p(\underline{\mathbf{X}}, \underline{\mathbf{z}} | \underline{\mu}, \underline{\Sigma}, \underline{\pi}) &= \prod_n p(\mathbf{x}_n, \mathbf{z}_n | \underline{\mu}, \underline{\Sigma}, \underline{\pi}) \quad \text{Bayes Rule} \\
 &= \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \underline{\mu}, \underline{\Sigma}) p(\mathbf{z}_n | \underline{\pi}) \\
 &= \prod_{n=1}^N \prod_{k=1}^K [\mathcal{N}(\mathbf{x}_n | \underline{\mu}_k, \underline{\Sigma}_k)]^{z_{nk}} \prod_{k=1}^K [\pi_k]^{z_{nk}}
 \end{aligned}$$

for each datapoint  $n$

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

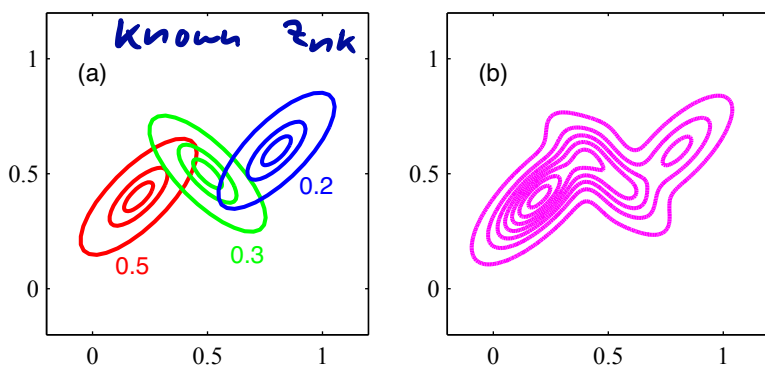
$$\begin{aligned}
 \underline{\pi} &\rightarrow \mathbf{z}_n \\
 (0.3, 0.6, 0.1) &\rightarrow (0, 1, 0) \\
 &\quad \uparrow k
 \end{aligned}$$

# 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:

$$p(\mathbf{x}_n | \theta) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)$$



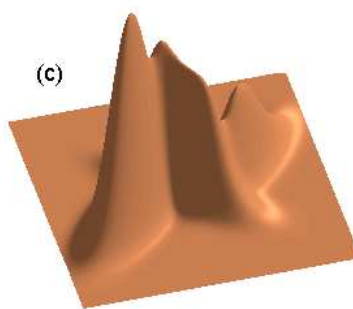
Deriving cost functions this way, is good for *statistical efficiency*. Without a latent variable model, the number of parameters grow at rate  $O(N)$ . After marginalization, the growth is reduced to  $O(D^2K)$  (assuming  $D, K \ll N$ ).

$$z_n = \begin{cases} (1 & 0 & 0 & \dots) & \text{if } k=1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ (0 & 0 & 0 & \dots & 1) & \text{if } k=K \end{cases}$$

joint  
 $= p(\mathbf{x}_n, z_n)$

marginal  
 $= p(\mathbf{x}_n)$   
 $= \sum_{k=1}^K p(\mathbf{x}_n, z_n=k)$   
 $= \sum_{k=1}^K p(\mathbf{x}_n | z_n=k) \cdot p(z_n=k)$

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



Marginalization  
 ~~$z \cdot N$~~

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

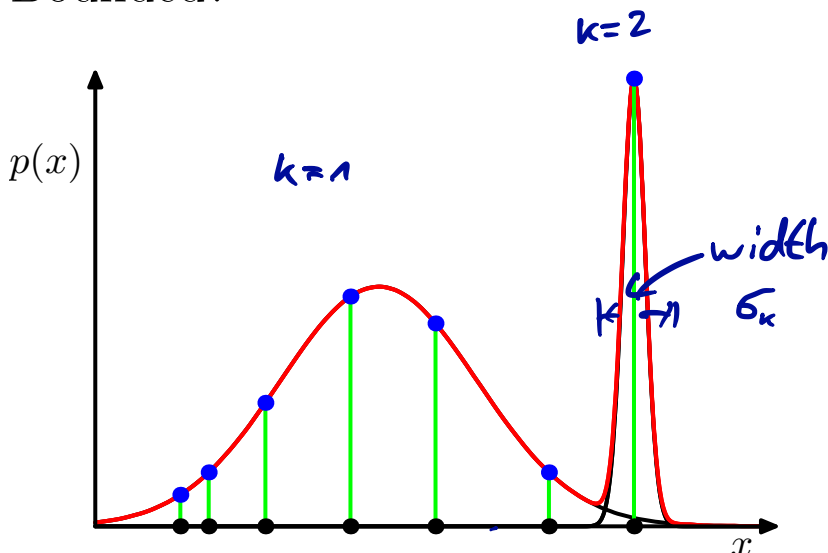
# Maximum likelihood

To get a maximum (marginal) likelihood estimate of  $\theta$ , we maximize the following:

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

$\mathcal{L}(\theta)$

Is this cost convex? Identifiable? Bounded?



$$\begin{aligned} \log(p(\mathbf{x} | \theta)) \\ &= \log \left( \prod_{n=1}^N \underbrace{p(\mathbf{x}_n | \theta)}_{\sum_{k=1}^K \dots} \right) \\ &= \sum_{n=1}^N \log(\dots) \end{aligned}$$

( $\mathcal{L}$  not concave)

① non-convex in  $\theta$

② non-unique optimum  $\theta$   
permutation of  $k$

$$k \rightarrow k'$$

$$\begin{aligned} \pi_k &\rightarrow \pi_{k'} \\ \mu_k &\rightarrow \mu_{k'} \\ \Sigma_k &\rightarrow \Sigma_{k'} \end{aligned}$$

③ non-bounded

$$\mathcal{L}(\theta) \rightarrow \infty$$

if  $\Sigma_k = \sigma_k I$   
in the limit

$$\sigma_k \rightarrow 0$$

## Exercises

1. Understand K-means extension to GMM. Why do we treat  $z_n$  as a random variable? Identify the joint, likelihood, prior, and marginal distributions, respectively.