

*Annotated  
version*

Machine Learning Course - CS-433

# Gaussian Mixture Models

Nov 10, 2016

©Mohammad Emtiyaz Khan 2015

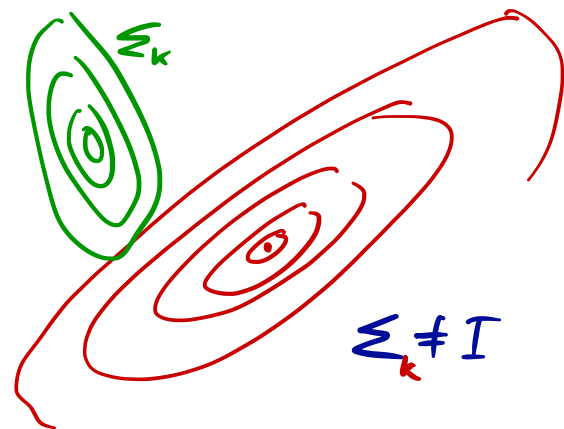
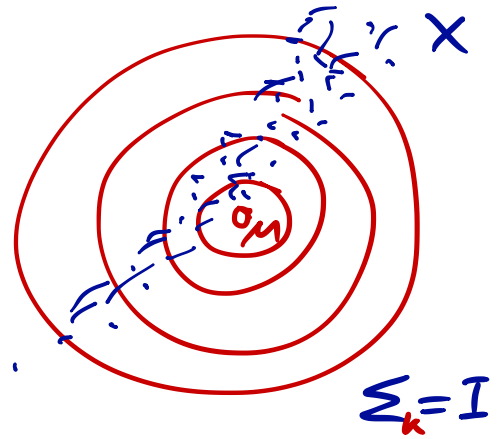
minor changes by Martin Jaggi 2016



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

# 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 Model.



## Clustering with Gaussians

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

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

$I^{D \times D}$  for K-means

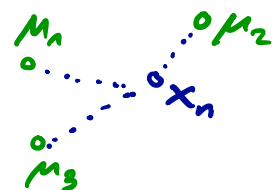
any  $\Sigma_k$  (p.s.d.)

$$\Sigma = (\Sigma_1, \dots, \Sigma_K)$$

$$\mu = (\mu_1, \dots, \mu_K)$$

## Soft-clustering

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.



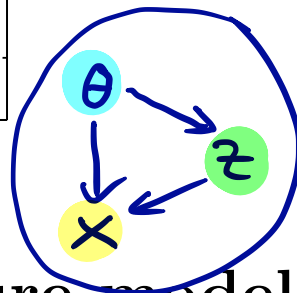
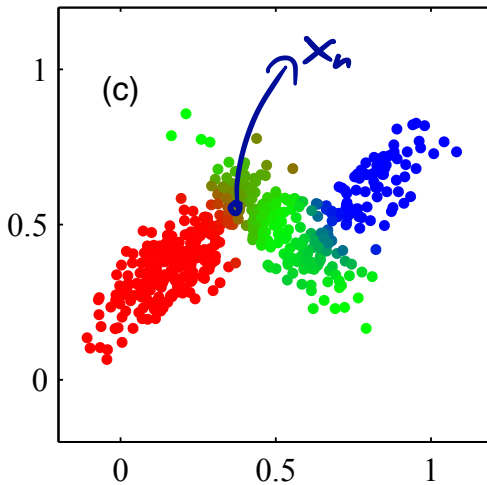
$$\sum_{k=1}^K \pi_k = 1$$

$n=1$	$\pi=0.1$
$n=2$	$\pi=0.7$
$n=9$	$\pi=0.2$

$z_n$

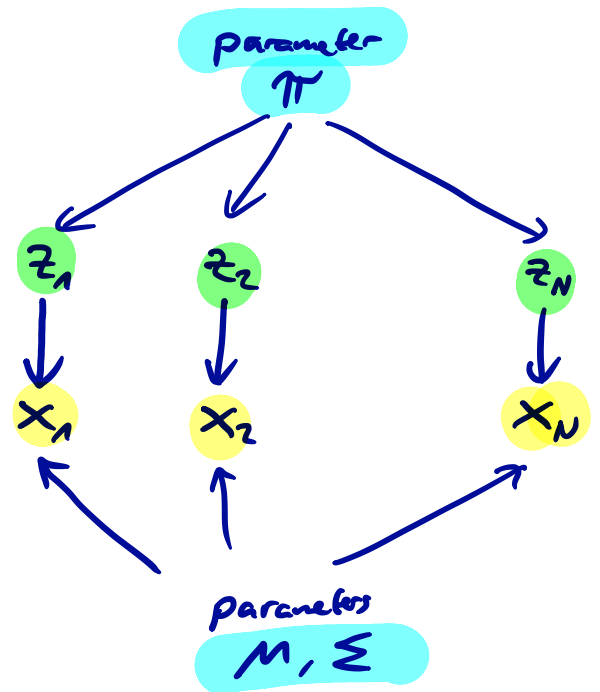
$$p(z_n = k) \doteq \pi_k \text{ where } \pi_k > 0, \forall k \text{ and } \sum_{k=1}^K \pi_k = 1$$

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



latent variables

data



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

## Gaussian mixture model

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

Bayes Rule

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

$$p(\mathbf{X}, \mathbf{z} | \mu, \Sigma, \pi) = p(\mathbf{X} | \mathbf{z}, \mu, \Sigma, \pi) \cdot p(\mathbf{z} | \mu, \Sigma, \pi)$$

$$= \prod_{n=1}^N p(\mathbf{x}_n | z_n, \mu, \Sigma) p(z_n | \pi)$$

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

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

$$\pi = (0.1, 0.7, 0.2)$$

$\pi_1 \quad \pi_2 \quad \pi_3$

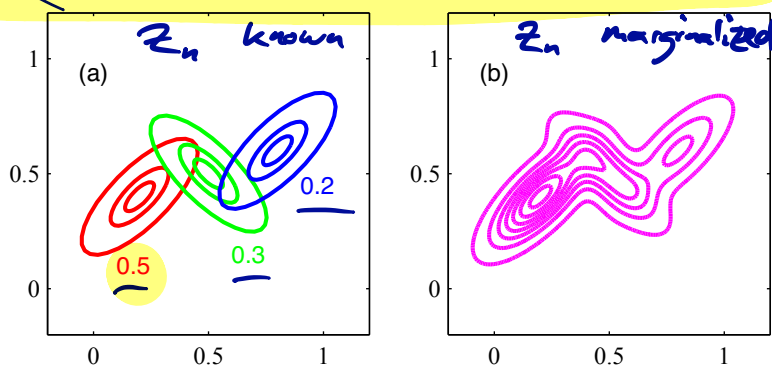
Here,  $\mathbf{x}_n$  are observed **data** vectors,  $z_n$  are **latent** unobserved variables, and the unknown **parameters** are given by  $\theta := \{\mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K, \pi\}$ .

# 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  $D, K \ll N$ ).

joint:

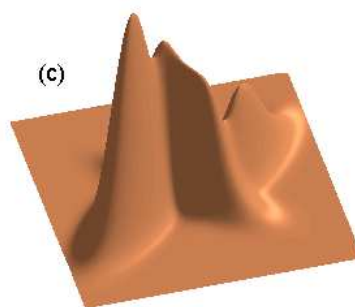
$$p(\mathbf{x}, \mathbf{z})$$

marginal:

$$p(\mathbf{x}) = \sum_{k=1}^K p(\mathbf{x}, \mathbf{z}=k)$$

$$= \sum_k p(\mathbf{x} | \mathbf{z}=k) \cdot p(\mathbf{z}=k)$$

if  $z_n = k$   
 $(\Leftrightarrow z_{nk} = 1 \quad z_{nj} = 0 \quad j \neq k)$   
 then  $p(\mathbf{x}_n) = \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)$



$$\theta = (\pi_{1:K}, \mu_{1:K}, \Sigma_{1:K})$$

$z$	$N$
$\pi$	$K$
$\mu$	$K \cdot D$
$\Sigma$	$K \cdot D^2$

# Maximum likelihood

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

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

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

$\mathcal{L}(\theta)$

Is this cost convex? Identifiable? Bounded?

• non-convex in  $\theta$

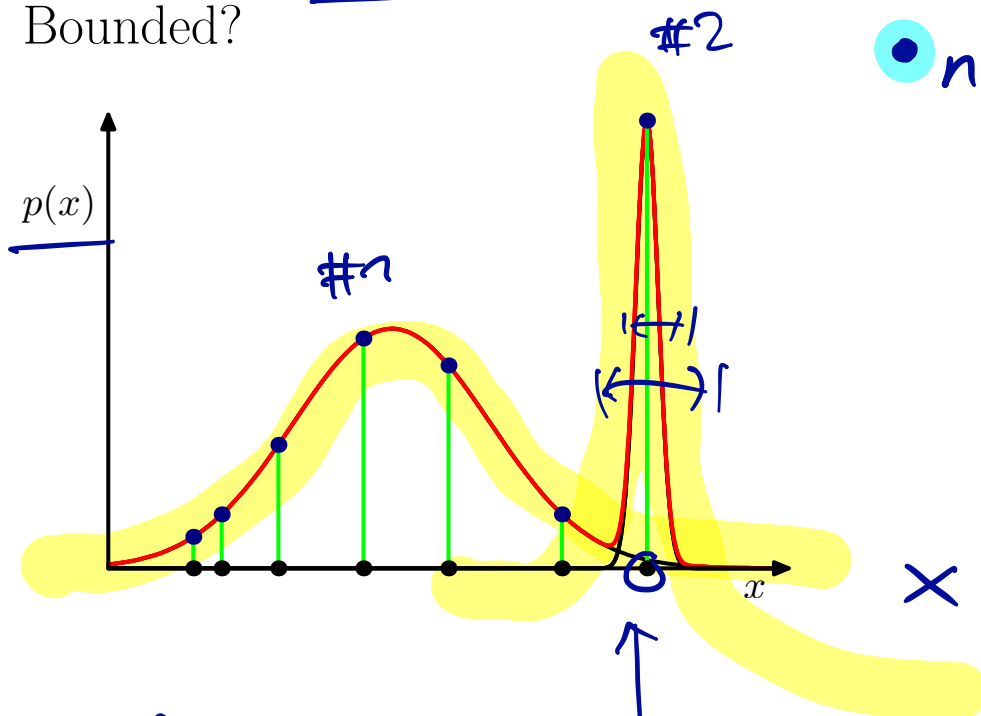
• not identifiable

permutation,

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

$$\theta_k \leftrightarrow \theta_{k'}$$

$\hookrightarrow K!$



• degenerate

$$\Sigma_2 = \sigma_2 I \rightarrow 0$$

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

$$\begin{aligned} &\mathcal{L}(\hat{\pi}_1=7, \hat{\pi}_2=3, \hat{\pi}_3=1) \\ &\quad \swarrow \quad \downarrow \\ &= \mathcal{L}(\pi_1=3, \pi_2=7, \pi_3=1) \end{aligned}$$

ToDo

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