

# Gaussian Mixture Models

Mohammad Emtiyaz Khan  
EPFL

Nov 5, 2015



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

©Mohammad Emtiyaz Khan 2015

# 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

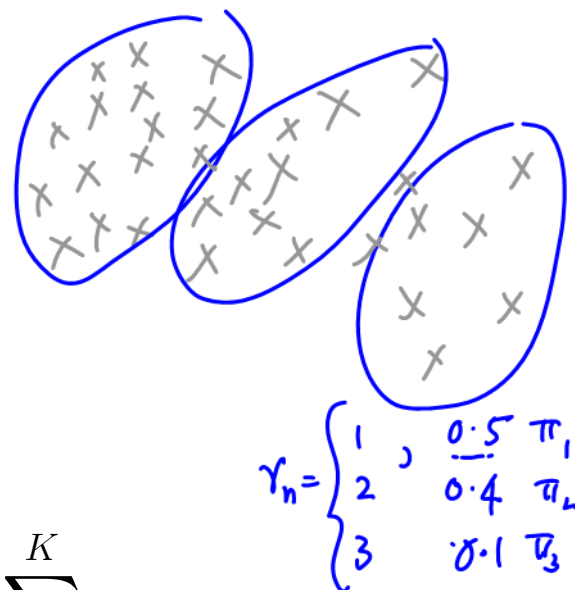
The first issue is resolved by using full covariance matrices  $\Sigma_k$  instead of *isotropic* covariances.

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

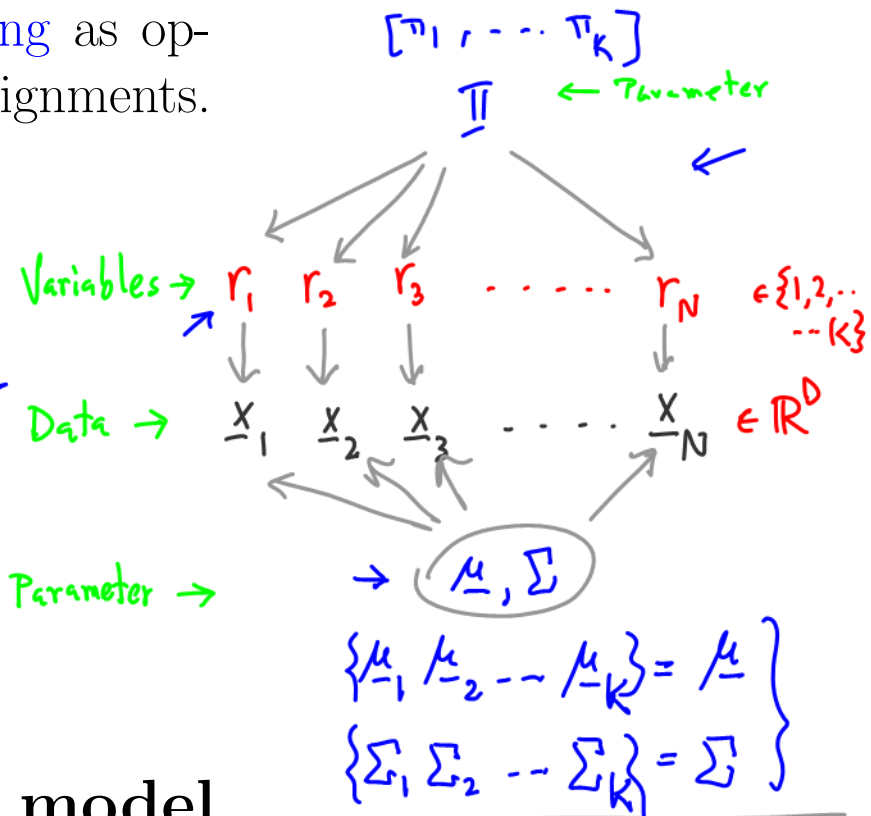
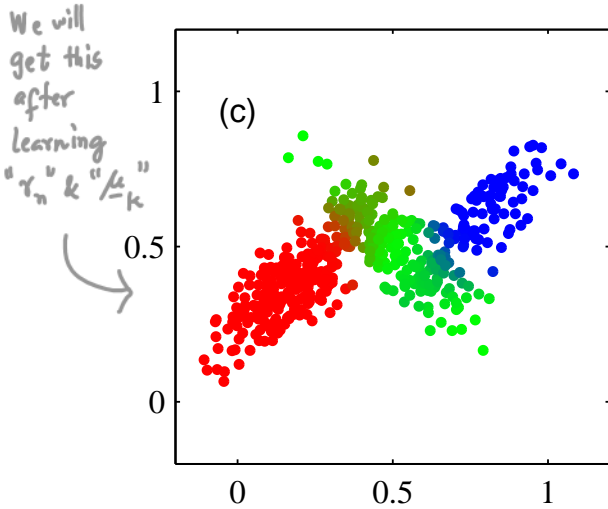
## Soft-clustering

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

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



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



## Gaussian mixture model

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

$$\theta = \{\mu, \Sigma, \pi\}$$

$$p(x_1, x_2, \dots, x_N, r_1, r_2, \dots, r_N / \theta)$$

$$= p(x, r / \theta)$$

**Joint**  $p(\mathbf{X}, \mathbf{r} | \mu, \Sigma, \pi)$  . **prior**

$$= \prod_{n=1}^N \left[ p(\mathbf{x}_n | r_n, \mu, \Sigma) p(r_n | \pi) \right]$$

**likelihood**

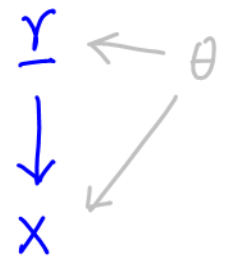
$$= \prod_{n=1}^N \left[ \prod_{k=1}^K [\mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)]^{r_{nk}} \right] \prod_{k=1}^K [\pi_k]^{r_{nk}}$$

Joint = Likelihood x Prior

$$P_{\theta}(D, r) = P_{\theta}(D | r) P_{\theta}(r)$$

Latent variable models

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



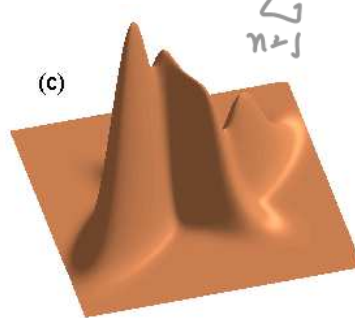
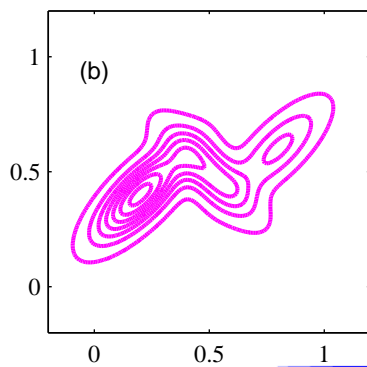
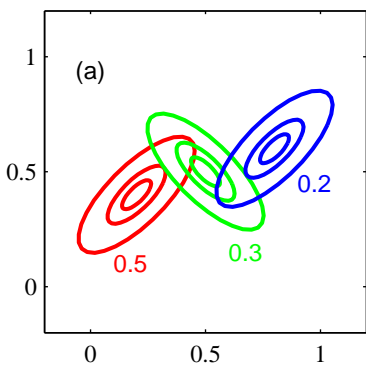
This is a generative model

# Marginal likelihood

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

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

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

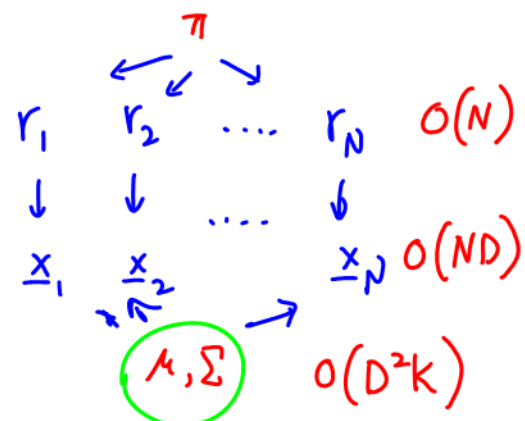


$$\begin{aligned} p(\mathbf{x}_n | \boldsymbol{\theta}) &= \sum_{k=1}^K p(\mathbf{x}_n, r_n = k | \boldsymbol{\theta}) \\ &= \sum_{k=1}^K \underbrace{p(\mathbf{x}_n | r_n = k, \boldsymbol{\theta})}_{\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \underbrace{p(r_n = k | \boldsymbol{\theta})}_{\pi_k} \end{aligned}$$

$$\begin{aligned} \log p(\mathbf{X} | \boldsymbol{\theta}) &= \sum_{n=1}^N \log p(\mathbf{x}_n | \boldsymbol{\theta}) \\ &= \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{aligned}$$

Maximum Marginal Likelihood

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  $D, K \ll N$ ).



# Maximum likelihood

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

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

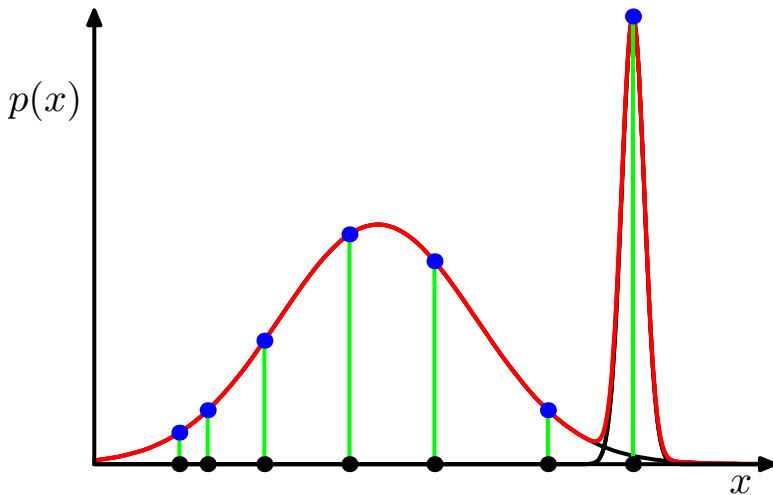
w.r.t.  $\theta$  \*

Is this cost convex? Identifiable?  
Bounded?

$$\theta^* = \max_{\theta} \log \mathcal{P}_{\theta}(x)$$

$$\theta \rightarrow \mathcal{P}_{\theta}(x)$$

is one-to-one



## To do

1. Understand K-means extension to GMM. Why do we need to treat  $r_n$  as a random variable? Identify the joint, likelihood, prior, and marginal distributions, respectively.
2. Understand identifiability and the difficulty with the maximum-likelihood estimation.