Machine Learning Course - CS-433

# Gaussian Mixture Models
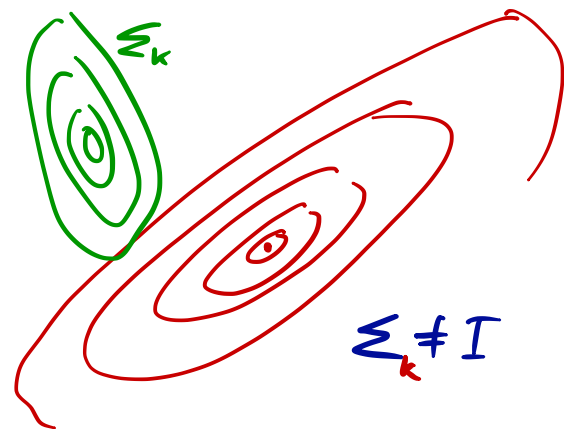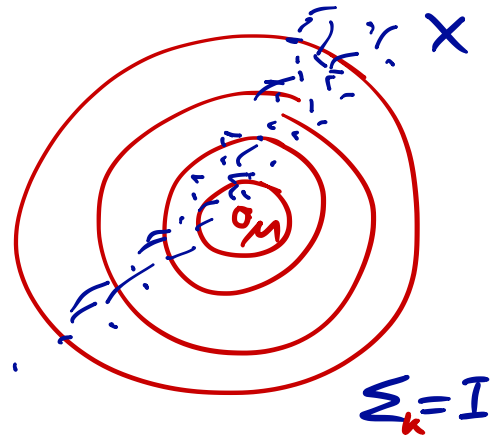
Nov 10, 2016

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.

$\Sigma_k = I$

# Clustering with Gaussians

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

$\Sigma_k$

$\Sigma_k \neq I$

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

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

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

# 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.
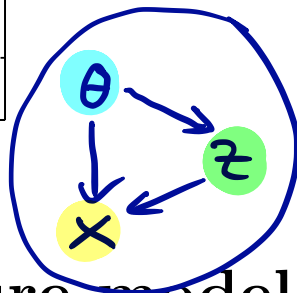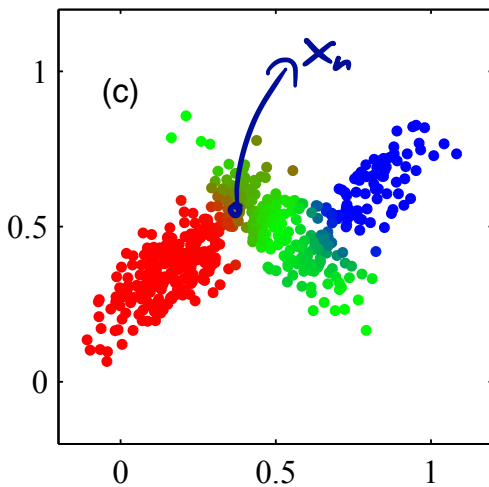
$\Sigma = (\Sigma_1, \ldots, \Sigma_k)$
$M = (M_1, \ldots, M_K)$

$M_1 \quad \circ \mu_2$
$\circ x_n$
$M_3$

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

| | |
|---|---|
| $n \to 1$ | $\mathbb{P} = 0.1$ |
| $n \to 2$ | $\mathbb{P} = 0.7$ |
| $n \to 3$ | $\mathbb{P} = 0.2$ |

$z_n$

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



(c)



Parameter $\pi$

latent variables

$z_1$ $z_2$ $z_N$

data

$x_1$ $x_2$ $x_N$

parameters $M, \Sigma$

## Gaussian mixture model

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

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

Bayes Rule
$$p(a,b) = p(a|b)\,p(b)$$

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

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

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

$z_{n:} = (0, 1, 0)$

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

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

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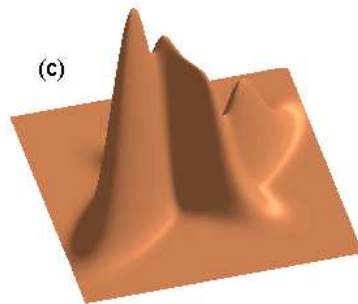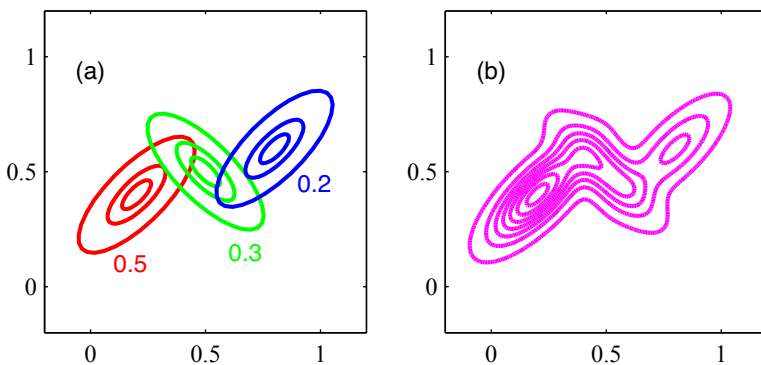
# Marginal likelihood

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:

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

marginal

$$p(a,b)$$

joint

marginal

$$p(a) = \sum_{k} p(a, b=k)$$

$$= \sum_{k} p(a|b=k)$$

$$\cdot p(b=k)$$

(a) 1 0.5 0 with values 0.5, 0.3, 0.2
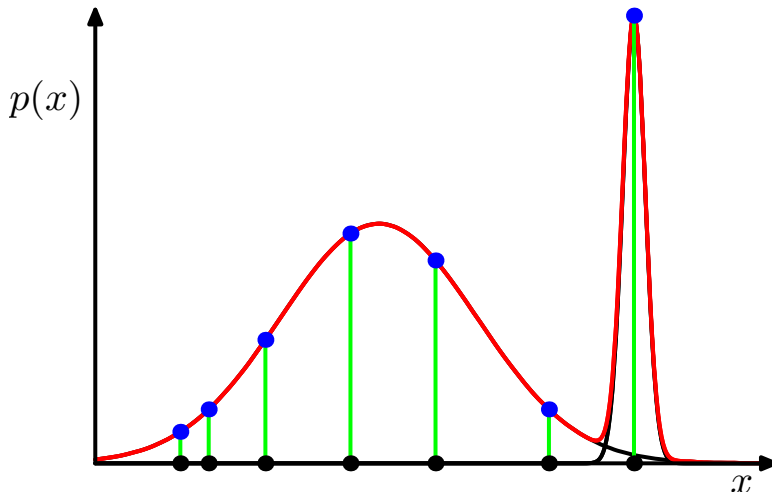
(b) 1 0.5 0

(c)

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 $O(D^2 K)$ (assuming $D, K \ll N$).

3

# Maximum likelihood

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

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

Is this cost convex? Identifiable? Bounded?



## 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.