

annotated
version

Machine Learning Course - CS-433

Gaussian Mixture Models

Nov 28, 2023

Martin Jaggi

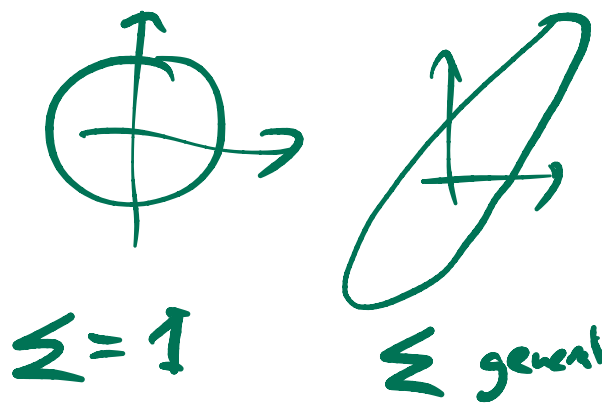
Last updated on: November 28, 2023

credits to Mohammad Emtiyaz Khan & Rüdiger Urbanke

EPFL

Motivation

- ① *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 Σ_k instead of *isotropic* covariances.

parameters:
 $\mu \in \mathbb{R}^{D \cdot K}$
 $\Sigma \in \mathbb{R}^{D^2 \cdot K}$
 $\pi \in \mathbb{R}^K$

- ①

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

Soft-clustering

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

$z_n = \text{vector } 0, \dots, 1, \dots, 0$
 $z_{nk} = \begin{cases} 1 & \text{if assigned to } k \\ 0 & \end{cases}$

- ②
- random variable: vector

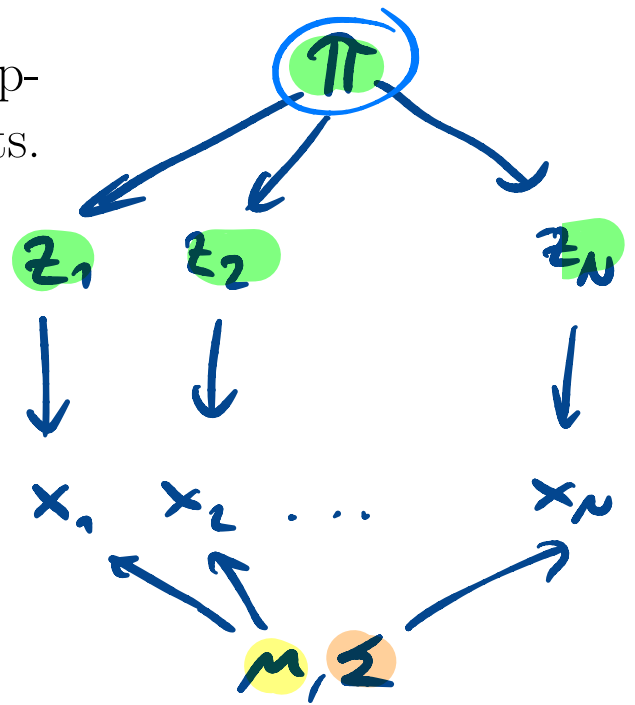
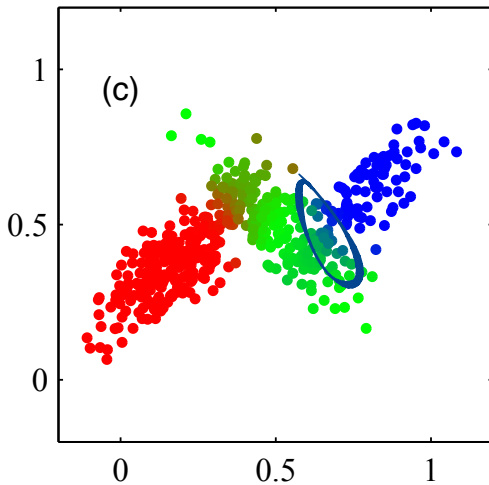
$$p(z_n = k) = \pi_k \quad \text{where } \pi_k > 0, \forall k \text{ and } \sum_{k=1}^K \pi_k = 1$$

def

$z_{nk} = 1$

cluster k
 $0, \dots, 1, \dots, 0$

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):

Bayes Rule

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

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

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

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

$p(z | \pi)$
 $\approx \pi_{k'}$
 $k' \text{ is where } z_{k'n} = 1$

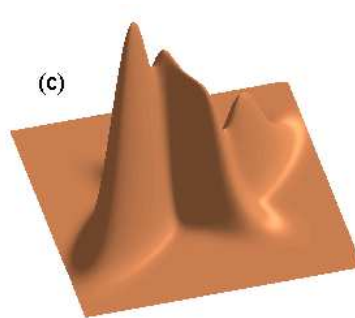
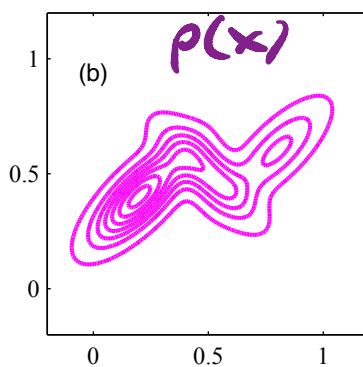
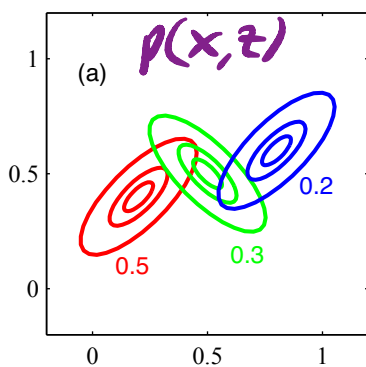
Here, \mathbf{x}_n are observed data vectors, z_n are latent unobserved variables, and the unknown parameters are given by $\theta := \{\mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K, \pi\}$.

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 grows at rate $\mathcal{O}(N)$. After marginalization, the growth is reduced to $\mathcal{O}(D^2 K)$ (assuming $D, K \ll N$).

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

marginal

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

$$= \sum_k \underbrace{p(\mathbf{x}_n | \mathbf{z}_n)}_{\text{likelihood}} \underbrace{p(\mathbf{z}_n)}_{\text{prior}}$$

$$\mathbf{z} \in \{0, 1\}^{N \cdot K}$$

Maximum likelihood

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

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

$$\log p(\mathbf{x}^n | \theta)$$

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

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

Is this cost convex? Identifiable?
Bounded?

① non-convex
(see k-means)

② non-unique optimals

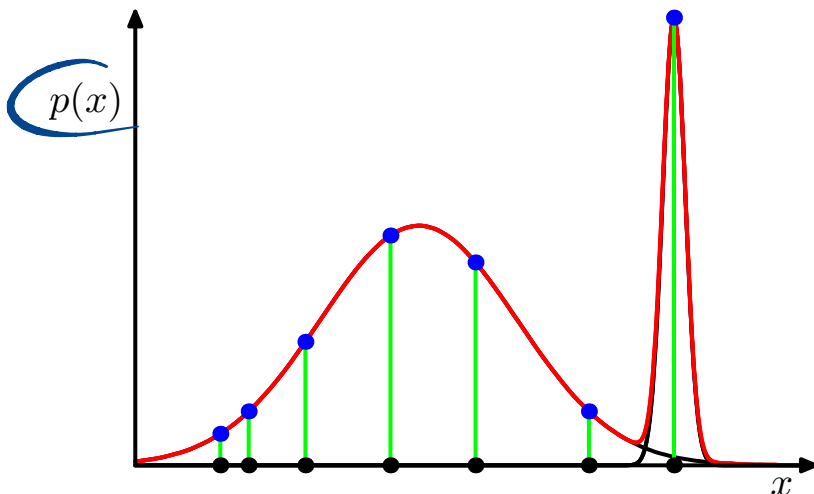
$$\begin{aligned} \pi_k &\leftrightarrow \pi_{k'} \\ \theta: \mu_k &\leftrightarrow \mu_{k'} \\ \Sigma_k &\leftrightarrow \Sigma_{k'} \end{aligned}$$

permutation of K

③ unbounded

$$\Sigma_k = \sigma_k \mathbf{I}$$

↑ scalar



width $\sigma_k \rightarrow 0$