

Gaussian mixtures

Nathanaël Carraz Rakotonirina

Mathématiques Informatique et Statistique Appliquées (MISA)
Université d'Antananarivo

Mixture models

It is an unsupervised learning approach. A mixture model is a combination of multiple simple distributions.

$$p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^K \pi_i p_i(\mathbf{x})$$

where p_i is the i th component of the mixture and π_i the mixture weights ($0 \leq \pi_i \leq 1$ and $\sum_{i=1}^K \pi_i = 1$)

We introduce the latent (hidden) variable z which can take values $1, \dots, K$. It specifies which distribution to use to generate the output. We note $p(z = i) = \pi_i$ and $p(\mathbf{x}|z = i) = p_i(\mathbf{x}) = p(\mathbf{x}; \boldsymbol{\theta}_i)$.

$$p(z; \boldsymbol{\theta}) = \text{Cat}(z; \boldsymbol{\pi})$$
$$p(\mathbf{x}|z = i; \boldsymbol{\theta}) = p(\mathbf{x}; \boldsymbol{\theta}_i)$$

According to the model, the data is generated as follows : we first sample a specific component z , and then depending on z , we generate the data \mathbf{x} (using the corresponding parameter).

We obtain the form of the model by marginalizing out z :

$$p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^K p(z = i; \boldsymbol{\theta}) p(\mathbf{x} | z = i; \boldsymbol{\theta}) = \sum_{i=1}^K \pi_i p(\mathbf{x}; \boldsymbol{\theta}_i)$$

This is difficult because we do not know the value of z .

Gaussian mixture models

A Gaussian mixture model (GMM) is a mixture model with a gaussian base distribution.

$$p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^K \pi_i \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

For large enough K a GMM can approximate any smooth distribution over R^D .

Applications of GMM

- ▶ Clustering
- ▶ Prior to regularize an inverse problem (image restoration tasks like image denoising, deblurring, inpainting, super-resolution ...)

GMM for clustering

Suppose we know the MLE of the parameters $\hat{\theta} = (\pi, \{\mu_i, \Sigma_i\})$, we use Bayes rule to compute the **responsibility** r_{ji} of cluster k for data point \mathbf{x}_j :

$$r_{ji} = p(z_j = i | \mathbf{x}_j; \theta) = \frac{p(z_j = i; \theta)p(\mathbf{x}_j | z_j = i; \theta)}{\sum_{i'=1}^K p(z_j = i'; \theta)p(\mathbf{x}_j | z_j = i'; \theta)} = \frac{\pi_i p(\mathbf{x}_j; \theta_i)}{\sum_{i'=1}^K \pi_{i'} p(\mathbf{x}_j; \theta_{i'})}$$

Then, we can compute the most probable cluster for \mathbf{x}_j as

$$\hat{z}_j = \arg \max_i r_{ji} = \arg \max_i (\log p(z_j = i; \pi) + \log p(\mathbf{x}_j | z_j = i; \theta))$$

K-means clustering

K-means clustering is a special case where $\Sigma_i = \mathbf{I}$ and $\pi_i = 1/K$ so we just have to estimate the means μ_i .

EM algorithm

The **expectation maximization (EM)** algorithm is designed to compute the MLE or MAP parameter estimate for probability models that have missing data and/or hidden variables.

We alternate between two steps:

- ▶ **E step**: expectation step, we estimate the hidden variables
- ▶ **M step**: maximization step, we compute the MLE

We need to iterate this process since the steps depend on each other. It may converge to a local maximum.

MLE using the EM algorithm

- **E step:** It just computes the responsibility of cluster i for generating \mathbf{x}_j using the current parameter estimates:

$$r_{ji}^{(t)} = p(z_j = i | \mathbf{x}_j; \boldsymbol{\theta}^{(t)}) = \frac{\pi_i p(\mathbf{x}_j; \boldsymbol{\theta}_i^{(t)})}{\sum_{i'=1}^K \pi_{i'} p(\mathbf{x}_j; \boldsymbol{\theta}_{i'}^{(t)})}$$

- **M step:** The M step maximizes the expected complete data log likelihood :

$$\begin{aligned}\mathcal{L}^{(t)}(\boldsymbol{\theta}) &= \mathbb{E} \left[\sum_j \log p(z_j; \boldsymbol{\pi}) + \log p(\mathbf{x}_j | z_j; \boldsymbol{\theta}) \right] \\ &= \mathbb{E} \left[\sum_j \log \left(\prod_i \pi_i^{z_{ij}} \right) + \log \left(\prod_i \mathcal{N}(\mathbf{x}_j; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)^{z_{ij}} \right) \right] \\ &= \sum_j \sum_i \mathbb{E}[z_{ij}] \log \pi_i + \sum_j \sum_i \mathbb{E}[z_{ij}] \mathcal{N}(\mathbf{x}_j; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \\ &= \sum_j \sum_i r_{ji}^{(t)} \log \pi_i + \sum_j \sum_i r_{ji}^{(t)} \mathcal{N}(\mathbf{x}_j; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)\end{aligned}$$

The parameter estimates are the following:

$$\mu_i^{(t+1)} = \frac{\sum_j r_{ji}^{(t)} \mathbf{x}_j}{r_k^{(t)}}$$
$$\Sigma_i^{(t+1)} = \frac{\sum_j r_{ji}^{(t)} (\mathbf{x}_j - \mu_i^{(t+1)})(\mathbf{x}_j - \mu_i^{(t+1)})^\top}{r_k^{(t)}}$$

where $r_i^{(t)} = \sum_j r_{ji}^{(t)}$ weighted number of points assigned to cluster k.
The mixture weights during the M step:

$$\pi_i^{(t+1)} = \frac{1}{N} \sum_j r_{ji}^{(t)} = \frac{r_i^{(t)}}{N}$$

What's next ?

Explore

- ▶ Unidentifiability and label switching
- ▶ Bayesian model selection to select K
- ▶ MAP estimate of a GMM
- ▶ Mixtures of Bernoullis
- ▶ Gaussian scale mixtures