

As recompiled from https://pillowlab.princeton.edu/teaching/statneuro2024/notes/notes17_MoG.pdf

EM for the Mixture of Gaussians model

1 General form of EM

$$L = \overbrace{\log p(x | \theta)}^{\text{log likelihood}} = \log \overbrace{\int p(x, z | \theta) dz}^{\text{total-data likelihood}} \geq \int \underbrace{q(z | \phi)}_{\text{variational distribution}} \log \left[\frac{p(x, z | \theta)}{q(z | \phi)} \right] dz \triangleq \overbrace{F(\phi, \theta)}^{\text{ELBO}} \quad (1)$$

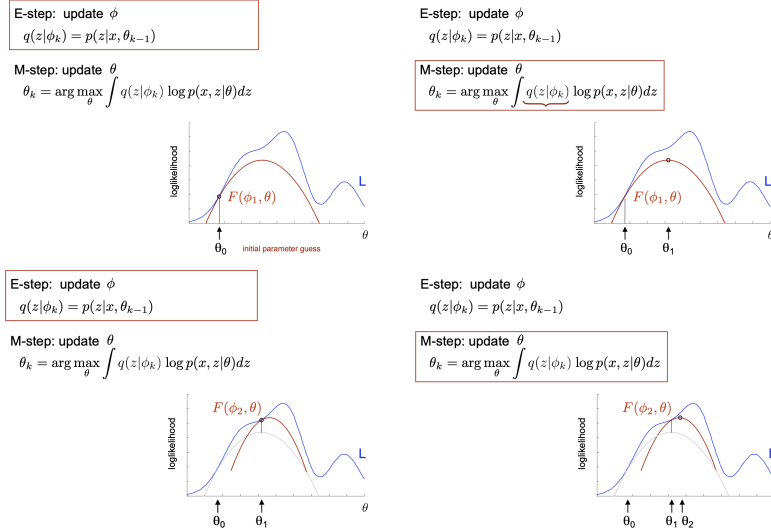
The ELBO, or negative free energy, has two convenient forms, which we exploit in the two alternating phases of EM:

$$F(\phi, \theta) = \log p(x | \theta) - KL(q(z | \phi) || p(z | x, \theta)) \quad (\text{used in E-step}) \quad (2)$$

$$F(\phi, \theta) = \int q(z | \phi) \log p(x, z | \theta) dz + \underbrace{H[q(z | \phi)]}_{\text{indep of } \theta} \quad (\text{used in M-step}) \quad (3)$$

Specifically, EM involves alternating between:

- E-step: Update ϕ by setting $q(z | \phi) = p(z | x, \theta)$, with θ held fixed.
- M-step: Update θ by maximizing $\int q(z | \phi) \log p(x, z | \theta) dz$, with ϕ held fixed.



Note that for discrete latent variable models, where the latent z takes on finite or countably infinitely many discrete values, the integral over z is replaced by a sum:

$$F(\phi, \theta) = \sum_{j=1}^m q(z = z_j | \phi) \log p(x, z = z_j | \theta) + \underbrace{H[q(z | \phi)]}_{\text{indep of } \theta} \quad (4)$$

where $\{z_1, \dots, z_m\}$ are the possible values of z .

1.1 Motivating the variational posterior, $q(z | \phi)$

Read this if the above is confusing

We sample some data, and think it emerges from a distribution with a latent variable, z – e.g. GMMs. We want to see what this distribution looks like, but we can only draw x . So what we need to find is $p(z|x)$. Our immediate instinct should be to use:

$$p(z | x) = \frac{p(x | z)p(z)}{p(x)}$$

Unfortunately, $p(x)$ is an intractable integral as it marginalises x over all possible z . Bayes by itself is therefore no good. We instead resort to directly approximating the posterior $p(z | x)$ using EM. To do this, we assume there exists a distribution $q(z | \phi)$ that approximates the posterior well. Notice that we use ϕ to signal that we are looking at the particular distribution with parameters that does this best. We know that this distribution should closely resemble the posterior, so as to minimise the KL-divergence below.

$$q^*(z | \phi) = \arg \min \text{KL}(q(z | \phi) || p(z | x))$$

The intractable $p(x)$ again appears, hidden in the posterior. We can use Bayes on $p(x)$ to re-express the KL-divergence as:

$$\text{KL}(q(z | \phi) || p(z | x)) = \underbrace{\mathbb{E}_q[\log p(x, z)] - \mathbb{E}_q[\log q(z | \phi)]}_{\text{ELBO}} + \log p(x)$$

Notice that we have aggregated the messy expectations into a term we call ELBO. We can show using Jensen's inequality that $\log p(x)$ is always greater than ELBO. So if we wish to minimise KL, we can only decrease ELBO, as $\log p(x)$ is invariant of ϕ . The conclusion is that minimising ELBO is equivalent to minimising KL.

Looking at EM in the context of Mixture of Gaussians may clarify things further.

2 EM for Mixture of Gaussians

The model:

$$z \sim \text{Ber}(p) \quad (5)$$

$$x \mid z \sim \begin{cases} \mathcal{N}(\mu_0, C_0), & \text{if } z = 0 \\ \mathcal{N}(\mu_1, C_1), & \text{if } z = 1 \end{cases} \quad (6)$$

Suppose we have a dataset consisting of N samples $\{x_i\}_{i=1}^N$. Our model seeks to model these in terms of a set of pairs of iid random variables $\{(z_i, x_i)\}_{i=1}^N$, each consisting of a latent and an observation. These samples are independent under the model, meaning that the log-likelihood is given by a sum of independent terms, each of which is bounded by its ELBO:

$$\log p(\{x_i\}_{i=1}^N \mid \theta) = \sum_{i=1}^N \log p(x_i \mid \theta) \quad (7)$$

$$\geq \sum_{i=1}^N \text{ELBO}(x_i \mid \theta) \quad (8)$$

Each of these ELBO terms is given by a sum over the discrete values of z :

$$\text{ELBO}(x_i \mid \theta) = \sum_{z_i=0}^1 q(z_i \mid \phi_i) \frac{\log p(x_i, z_i \mid \theta)}{q(z_i \mid \phi)} \quad (9)$$

$$= q(z_i = 0 \mid \phi_i) \log \frac{p(x_i, z_i = 0 \mid \theta)}{q(z_i = 0 \mid \phi_i)} + q(z_i = 1 \mid \phi_i) \log \frac{p(x_i, z_i = 1 \mid \theta)}{q(z_i = 1 \mid \phi_i)} \quad (10)$$

Thus the ELBO for the entire model can be written as the sum over all datapoints:

$$\text{ELBO} = \sum_{i=1}^N \left[q(z_i = 0 \mid \phi_i) \log \frac{p(x_i, z_i = 0 \mid \theta)}{q(z_i = 0 \mid \phi_i)} + q(z_i = 1 \mid \phi_i) \log \frac{p(x_i, z_i = 1 \mid \theta)}{q(z_i = 1 \mid \phi_i)} \right] \quad (11)$$

Here ϕ_i is the variational parameter associated with the i 'th latent variable z_i .

2.1 E-step

The E step involves setting $q(z_i \mid \phi_i)$ equal to the conditional distribution of z_i given the data and current parameters θ . We will denote these binary probabilities by ϕ_{i0} and ϕ_{i1} , given by the recognition distribution of the z_i under the model:

$$\phi_{i0} = p(z_i = 0 \mid x_i, \theta) = \frac{(1-p)\mathcal{N}_0(x_i)}{(1-p)\mathcal{N}_0(x_i) + p\mathcal{N}_1(x_i)} \quad (12)$$

$$\phi_{i1} = p(z_i = 1 \mid x_i, \theta) = \frac{p\mathcal{N}_1(x_i)}{(1-p)\mathcal{N}_0(x_i) + p\mathcal{N}_1(x_i)} \quad (13)$$

where $\mathcal{N}_0(x_i) = \mathcal{N}(x_i \mid \mu_0, C_0)$ and $\mathcal{N}_1(x_i) = \mathcal{N}(x_i \mid \mu_1, C_1)$, and note that $\phi_{i0} + \phi_{i1} = 1$.

At the end of the E-step we have a pair of these probabilities for each sample, which can be represented as an $N \times 2$ matrix

$$\phi = \begin{bmatrix} \phi_{10} & \phi_{11} \\ \phi_{20} & \phi_{21} \\ \vdots & \vdots \\ \phi_{N0} & \phi_{N1} \end{bmatrix} \quad (14)$$

where each row sums to 1 . The values ϕ_{ij} are often described as the "soft assignments" of datapoints to classes, since they describe how likely it is that datapoint i came from class j , given the current model parameters θ .

2.2 M-step

The M-step involves updating the parameters $\theta = \{p, \mu_0, \mu_1, C_0, C_1\}$ using the current variational distribution $q(z | \phi)$. To do this, we plug in the assignment probabilities $\{\phi_{i0}, \phi_{i1}\}$ from the E-step into the ELBO (eq. 11) to obtain:

$$F = \sum_{i=1}^N [\phi_{i0} \log p(x_i, z_i = 0 | \theta) + \phi_{i1} \log p(x_i, z_i = 1 | \theta)] \quad (15)$$

$$= \sum_{i=1}^N [\phi_{i0} (\log(1-p) + \log \mathcal{N}(x_i | \mu_0, C_0)) + \phi_{i1} (\log p + \log \mathcal{N}(x_i | \mu_1, C_1))] \quad (16)$$

Maximizing this expression for the model parameters (see next section for derivations) gives updates:

$$\hat{\mu}_0 = \left(\frac{1}{\sum \phi_{i0}} \right) \sum \phi_{i0} x_i \quad (17)$$

$$\hat{\mu}_1 = \left(\frac{1}{\sum \phi_{i1}} \right) \sum \phi_{i1} x_i \quad (18)$$

$$\hat{C}_0 = \left(\frac{1}{\sum \phi_{i0}} \right) \sum \phi_{i0} (x_i - \hat{\mu}_0) (x_i - \hat{\mu}_0)^\top \quad (19)$$

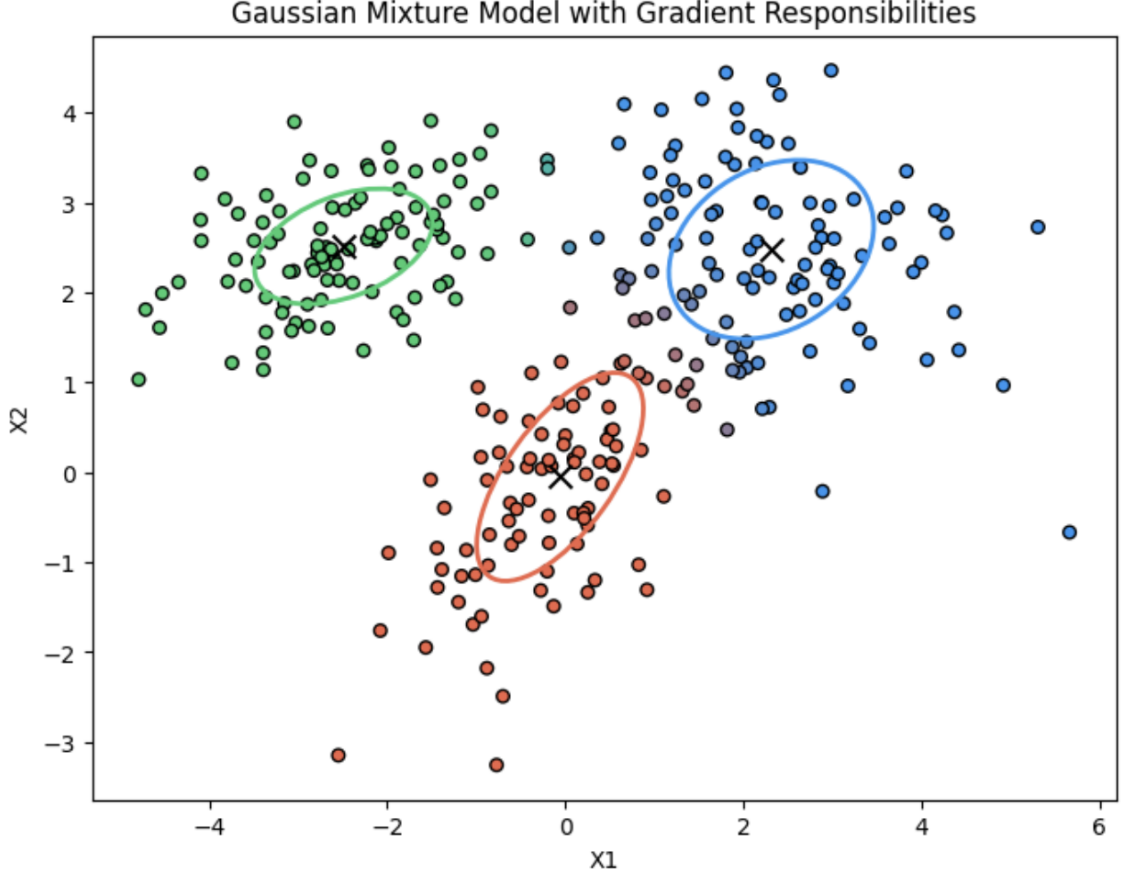
$$\hat{C}_1 = \left(\frac{1}{\sum \phi_{i1}} \right) \sum \phi_{i1} (x_i - \hat{\mu}_1) (x_i - \hat{\mu}_1)^\top \quad (20)$$

$$p = \frac{1}{N} \sum \phi_{i1} \quad (21)$$

Note that the mean and covariance updates are formed by taking the weighted average and weighted covariance of the samples, with weights given by the assignment probabilities ϕ_{i0} and ϕ_{i1} .

3 GMM EM Algorithm

```
1 import numpy as np
2
3 # Initialise parameters
4 def initialise_parameters(X, K):
5     N, D = X.shape
6     pi = np.ones(K) / K # Initialise mixing coefficients
7     mu = np.random.randn(K, D) # Initialise means
8     sigma = np.array([np.eye(D) for _ in range(K)]) # Initialise covariances
9     return pi, mu, sigma
10
11 # E-step: Calculate responsibilities (posterior probabilities)
12 def e_step(X, pi, mu, sigma):
13     N, K = X.shape[0], pi.shape[0]
14     gamma = np.zeros((N, K))
15     for k in range(K):
16         gamma[:, k] = pi[k] * multivariate_gaussian(X, mu[k], sigma[k])
17     return gamma / gamma.sum(axis=1, keepdims=True) # Normalise
18
19 # M-step: Update parameters based on responsibilities
20 def m_step(X, gamma):
21     N, D = X.shape
22     K = gamma.shape[1]
23     N_k = gamma.sum(axis=0) # Effective number of samples for each component
24     pi = N_k / N # Update mixing coefficients
25     mu = np.dot(gamma.T, X) / N_k[:, np.newaxis] # Update means
26     sigma = np.array([(gamma[:, k][:, np.newaxis] * (X - mu[k])).T @ (X - mu[k]) / N_k[k]
27                        for k in range(K)]) # Update covariances
28     return pi, mu, sigma
29
30 # Gaussian likelihood function
31 def multivariate_gaussian(X, mu, sigma):
32     D = X.shape[1]
33     norm_const = 1 / (np.sqrt((2 * np.pi) ** D * np.linalg.det(sigma)))
34     X_centered = X - mu
35     return norm_const * np.exp(-0.5 * np.sum(X_centered @ np.linalg.inv(sigma) *
36                                               X_centered, axis=1))
37
38 # Log likelihood calculation
39 def log_likelihood(X, pi, mu, sigma):
40     return np.sum(np.log(np.sum([pi[k] * multivariate_gaussian(X, mu[k], sigma[k]) for k
41                               in range(len(pi))], axis=0)))
42
43 # EM Algorithm
44 def em_algorithm(X, K, max_iters=100, tol=1e-4):
45     pi, mu, sigma = initialise_parameters(X, K)
46     log_likelihood_values = []
47
48     for _ in range(max_iters):
49         # E-step
50         gamma = e_step(X, pi, mu, sigma)
51
52         # M-step
53         pi, mu, sigma = m_step(X, gamma)
54
55         # Log likelihood
56         ll = log_likelihood(X, pi, mu, sigma)
57         log_likelihood_values.append(ll)
58
59         # Check for convergence
60         if len(log_likelihood_values) > 1 and np.abs(log_likelihood_values[-1] -
61                                                       log_likelihood_values[-2]) < tol:
62             break
63
64     return pi, mu, sigma, log_likelihood_values
```



4 Derivation of M-step updates

4.1 Updates for μ_0, μ_1

To derive the updates for μ_0 , we collect the terms from the ELBO (eq. 16) that involve μ_0 , giving:

$$F(\mu_0) = \sum_{i=1}^N \phi_{i0} \log \mathcal{N}(x_i | \mu_0, C_0) + \text{const} \quad (22)$$

$$= -\frac{1}{2} \sum_{i=1}^N \phi_{i0} (x_i - \mu_0)^\top C_0^{-1} (x_i - \mu_0) + \text{const} \quad (23)$$

$$- \frac{1}{2} \sum_{i=1}^N \phi_{i0} (-2\mu_0^\top C_0^{-1} x_i + \mu_0^\top C_0^{-1} \mu_0) + \text{const} \quad (24)$$

$$= \mu_0^\top C_0^{-1} \left(\sum_{i=1}^N \phi_{i0} x_i \right) - \frac{1}{2} \left(\sum_{i=1}^N \phi_{i0} \right) \mu_0^\top C_0^{-1} \mu_0 + \text{const.} \quad (25)$$

Differentiating with respect to μ_0 and setting to zero gives:

$$\frac{\partial}{\partial \mu_0} F = C_0^{-1} \left(\sum_{i=1}^N \phi_{i0} x_i \right) - \left(\sum_{i=1}^N \phi_{i0} \right) C_0^{-1} \mu_0 = 0 \quad (26)$$

$$\implies \left(\sum_{i=1}^N \phi_{i0} x_i \right) = \left(\sum_{i=1}^N \phi_{i0} \right) \mu_0 \quad (27)$$

$$\implies \hat{\mu}_0 = \frac{\sum_{i=1}^N \phi_{i0} x_i}{\sum_{i=1}^N \phi_{i0}} \quad (28)$$

A similar approach leads to the update for μ_1 , with weights ϕ_{i1} instead of ϕ_{i0} .

4.2 Updates for C_0, C_1

Matrix derivative identities:

Assume C is a symmetric, positive definite matrix. We have the following identities ([1]):

- log-determinant:

$$\frac{\partial}{\partial C} \log |C| = C^{-1} \quad (29)$$

- quadratic form:

$$\frac{\partial}{\partial C} x^\top C x = x x^\top \quad (30)$$

Derivation:

The simplest approach for deriving updates for C_0 is to differentiate the ELBO F with respect to C_0^{-1} and then solve for C_0 . We assume we already have the updated mean $\hat{\mu}_0$ (which did not depend on C_0 or any other parameters).

The ELBO as a function of C_0 can be written:

$$F(C_0) = \sum_{i=1}^N \phi_{i0} \log \mathcal{N}(x_i | \hat{\mu}_0, C_0) + \text{const} \quad (31)$$

$$= \sum_{i=1}^N \phi_{i0} \left(+\frac{1}{2} \log |C_0^{-1}| - \frac{1}{2} (x_i - \hat{\mu}_0)^\top C_0^{-1} (x_i - \hat{\mu}_0) \right) + \text{const} \quad (32)$$

Differentiating with respect to C_0^{-1} gives us:

$$\frac{\partial}{\partial C_0^{-1}} F = \frac{1}{2} \left(\sum_{i=1}^N \phi_{i0} \right) C_0 + \frac{1}{2} \left(\sum_{i=1}^N \phi_{i0} (x_i - \hat{\mu}_0) (x_i - \hat{\mu}_0)^\top \right) = 0 \quad (33)$$

$$\implies \hat{C}_0 = \frac{1}{\left(\sum_{i=1}^N \phi_{i0} \right)} \left(\sum_{i=1}^N \phi_{i0} (x_i - \hat{\mu}_0) (x_i - \hat{\mu}_0)^\top \right) \quad (34)$$

which as noted above is simply the covariance matrix of all stimuli weighted by their recognition weights. The same derivation can be used for C_1 .

4.3 Mixing probability p update

Finally, updates for p are obtained by collecting terms involving p :

$$F(p) = \sum_{i=1}^N (\phi_{i0} \log(1-p) + \phi_{i1} \log p) + \text{const} \quad (35)$$

$$= \log(1-p) \left(\sum_{i=1}^N \phi_{i0} \right) + (\log p) \left(\sum_{i=1}^N \phi_{i1} \right) + \text{const} \quad (36)$$

Differentiating and setting to zero gives

$$\frac{\partial}{\partial p} F = \frac{1}{p-1} \left(\sum_{i=1}^N \phi_{i0} \right) + \frac{1}{p} \left(\sum_{i=1}^N \phi_{i1} \right) = 0 \quad (37)$$

$$\implies p \left(\sum_{i=1}^N \phi_{i0} \right) + (p-1) \left(\sum_{i=1}^N \phi_{i1} \right) = 0 \quad (38)$$

$$\implies p \left(\sum_{i=1}^N \phi_{i0} + \phi_{i1} \right) = \left(\sum_{i=1}^N \phi_{i1} \right) \quad (39)$$

$$\implies p = \frac{1}{N} \sum_{i=1}^N \phi_{i1} \quad (40)$$

where note that we have used $p_{i0} + p_{i1} = 1$ for all i . Thus the m-step estimate for p is simply the average probability assigned to cluster 1 .

5 K-means clustering

The K-means clustering algorithm can be seen as applying the EM algorithm to a mixture-of Gaussians latent variable with covariances $C_0 = C_1 = \epsilon I$ in the limit where $\epsilon \rightarrow 0$. Note that in this limit the recognition probabilities go to 0 or 1 :

$$p(z=1 | x) = \frac{pN(x | \mu_1, \epsilon I)}{pN(x | \mu_0, \epsilon I) + (1-p)N(x | \mu_0, \epsilon I)} \quad (41)$$

$$= \frac{1}{1 + \frac{1-p}{p} \exp \left(\frac{1}{2\epsilon} \left(\|x - \mu_1\|^2 - \|x - \mu_0\|^2 \right) \right)} \quad (42)$$

$$= \begin{cases} 0, & \text{if } \|x - \mu_1\|^2 > \|x - \mu_0\|^2 \\ 1, & \text{if } \|x - \mu_1\|^2 < \|x - \mu_0\|^2. \end{cases} \quad (43)$$

The E-step for this model results in "hard assignments", since each datapoint is assigned definitively to one cluster or the other, and the M-step involves updating the means μ_0 and μ_1 to be the sample means of the points assigned to each cluster.

Note that the recognition distribution is independent of p , and we can therefore drop that parameter from the model. Thus, the only parameters of the K-means model are the means μ_0 and μ_1 .

References

[1] Kaare Brandt Petersen and Michael Syskind Pedersen. The matrix cookbook. Technical University of Denmark, 7(15):510, 2008.