

Gaussian Mixture Models

Andrea Casalino

December 25, 2019

1 What is a Gaussian Mixture model?

A Gaussian Mixture Model (GMM) is a class of probability distribution functions, often adopted to approximate unknown distributions. GMM is basically a mixture (Section 3) of Gaussians (Section 2).

The generic probability density function (aka pdf) f of a continuous random variable X , having $\mathcal{X} \in \mathbb{R}^n$ as domain, is a function defined as follows:

$$f : \mathcal{X} \rightarrow [0, 1] \quad (1)$$

$$\int_{\mathcal{X}} f(X) dx_1 \cdots dx_n = 1 \quad (2)$$

Prior to discuss the GMM properties, the expectation \mathbb{E} operator must be introduced. The expectation \mathbb{E} of a distribution f w.r.t a function $g(X)$ defined over the same domain is equal to:

$$\mathbb{E}_f[g(X)] = \int_{\mathcal{X}} f(X)g(X)dx_1, \dots, n \quad (3)$$

2 Gaussian distributions

Among all the possible distributions, Gaussians are ones of the most important. The mono-variate Gaussian distribution is a pdf defined as follows:

$$\phi_{(\mu, \Sigma)} : \mathcal{X} \rightarrow [0, 1] \quad (4)$$

$$\phi_{(\mu, \Sigma)}(x) = \frac{1}{\sqrt{2\pi\Sigma}} \exp\left(-\frac{1}{2}(x - \mu) \cdot \frac{1}{\Sigma} \cdot (x - \mu)\right) \quad (5)$$

where $X \subset \mathbb{R}$. Σ and μ are the covariance and mean of the distribution and are obtained through the following expectations:

$$\mu = \mathbb{E}_{\phi}[x] = \int_{\mathcal{X}} \phi_{(\mu, \Sigma)}(x) x dx \quad (6)$$

$$\Sigma = \mathbb{E}_{\phi}[(x - \mu)^2] = \int_{\mathcal{X}} \phi_{(\mu, \Sigma)}(x - \mu)^2 dx \quad (7)$$

Figure 1 reports examples of Gaussian distributions.

The multivariate version of a Gaussian is built generalizing the uni-variate one. Indeed, every multivariate Gaussian distribution can be seen, in a proper space, as a union of n independent Gaussian distributions. Indeed, suppose to have n independent Gaussians having a null means and the variances equal to $\Sigma_1, \dots, \Sigma_n$. They form a multivariate Gaussian distribution $\Phi_{(0, \Sigma)}$, with a covariance matrix Σ equal to the following diagonal matrix:

$$\Sigma = \begin{bmatrix} \Sigma_1 & & \\ & \ddots & \\ & & \Sigma_n \end{bmatrix} \quad (8)$$

Since all the Gaussians are independent, the expression of the density function is obtained as the following product:

$$\Phi_{(0, \Sigma)}(Y) = \phi_{(0, \Sigma_1)}(y_1) \cdots \phi_{(0, \Sigma_n)}(y_n) \quad (9)$$

$$\Phi_{(0, \Sigma)}(Y) = \frac{1}{\sqrt{2\pi\Sigma_1}} \exp\left(-\frac{1}{2} \frac{y_1^2}{\Sigma_1}\right) \cdots \frac{1}{\sqrt{2\pi\Sigma_n}} \exp\left(-\frac{1}{2} \frac{y_n^2}{\Sigma_n}\right) \quad (10)$$

$$\Phi_{(0, \Sigma)}(Y) = \frac{1}{\sqrt{(2\pi)^n \cdot \Sigma_1 \cdots \Sigma_n}} \exp\left(-\frac{1}{2} \sum_{i=1}^n \frac{y_i^2}{\Sigma_i}\right) \quad (11)$$

$$\Phi_{(0, \Sigma)}(Y) = \frac{1}{\sqrt{(2\pi)^n \cdot |\Sigma|}} \exp\left(-\frac{1}{2} Y^T \begin{bmatrix} \frac{1}{\Sigma_1} & & \\ & \ddots & \\ & & \frac{1}{\Sigma_n} \end{bmatrix} Y\right) \quad (12)$$

$$\Phi_{(0, \Sigma)}(Y) = \frac{1}{\sqrt{(2\pi)^n \cdot |\Sigma|}} \exp\left(-\frac{1}{2} Y^T \Sigma^{-1} Y\right) \quad (13)$$

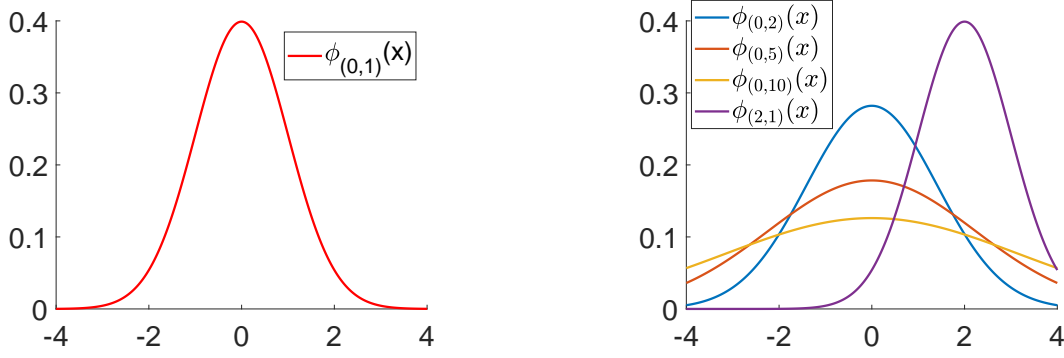


Figure 1: On the left the uni-variate Gaussian distribution having a 0 mean and unitary variance, while on the right examples Gaussians having different values for the variance and the mean.

Notice that:

$$\int_{\mathcal{Y}} \Phi_{(0,\Sigma)}(Y) dy_1, \dots, y_n = 1 \quad (14)$$

is ensured due to the fact that the following single integrals are all equal to 1 (equation (5)):

$$\Phi_{(0,\Sigma)}(Y) = \left(\int_{\mathcal{Y}} \frac{1}{\sqrt{2\pi\Sigma_1}} \exp\left(-\frac{1}{2} \frac{y_1^2}{\Sigma_1}\right) dy_1 \right) \cdots \left(\int_{\mathcal{Y}} \frac{1}{\sqrt{2\pi\Sigma_n}} \exp\left(-\frac{1}{2} \frac{y_n^2}{\Sigma_n}\right) dy_n \right) \quad (15)$$

Moreover it is true that the covariance matrix Σ can be obtained through the following expectation:

$$\Sigma = \mathbb{E}_{\Phi}[Y \cdot Y^T] = \int_{\mathcal{Y}} \Phi_{(0,\Sigma)}(Y) Y \cdot Y^T dy_1, \dots, y_n \quad (16)$$

In the domain of \mathcal{Y} , the multivariate Gaussian $\Phi(Y)$ is made of many independent Gaussians. Anyway, when performing a change of variables, the same distribution becomes a general multivariate Gaussian, having the variables in the new domain that are correlated. Consider to set $Y = R \cdot X'$, with R that is a rotation matrix, i.e. $R^{-1} = R^T$. X' is in turn described by a multivariate Gaussian having a density defined as follows:

$$\Phi_{(0,\Sigma)}(X') = \frac{1}{\sqrt{(2\pi)^n \cdot |\Sigma|}} \exp\left(-\frac{1}{2} X'^T R^T \Sigma^{-1} R X'\right) \quad (17)$$

$$\Phi_{(0,\Sigma')}(X') = \frac{1}{\sqrt{(2\pi)^n \cdot |\Sigma'|}} \exp\left(-\frac{1}{2} X'^T (\Sigma')^{-1} X'\right) \quad (18)$$

The above simplifications are due to the fact that:

$$\Sigma' = R^T \Sigma R \quad (19)$$

$$(\Sigma')^{-1} = R^{-1} \Sigma^{-1} R = R^T \Sigma^{-1} R \quad (20)$$

Moreover, since R is a rotation matrix, Σ' is symmetric positive definite (since $\Sigma_{1,\dots,n} \geq 0$). Additionally $|\Sigma| = |\Sigma'|$.

Applying an additional transformation of the form $X = X' + \mu$, a distribution with a non zero mean (in the new domain) can be obtained. X is the generic multivariate Gaussian, with a density equal to:

$$\Phi_{(\mu,\Sigma')}(X) = \frac{1}{\sqrt{(2\pi)^n \cdot |\Sigma'|}} \exp\left(-\frac{1}{2} (X - \mu)^T (\Sigma')^{-1} (X - \mu)\right) \quad (21)$$

Indeed, applying a rotation and a translation, every multivariate Gaussian can be obtained from a simple set of independent uni-variate Gaussians. Section 7 reports a way to compute in a precise way the change of variables to perform to obtain a Gaussian with a certain mean and covariance. Figure 2 summarizes the above considerations.

3 Mixture models

Mixture models are used to define a probability density function as the combination of a certain number of simpler ones. It is possible to define a generic mixture model, by combining N probability densities f_1, \dots, f_N satisfying equation

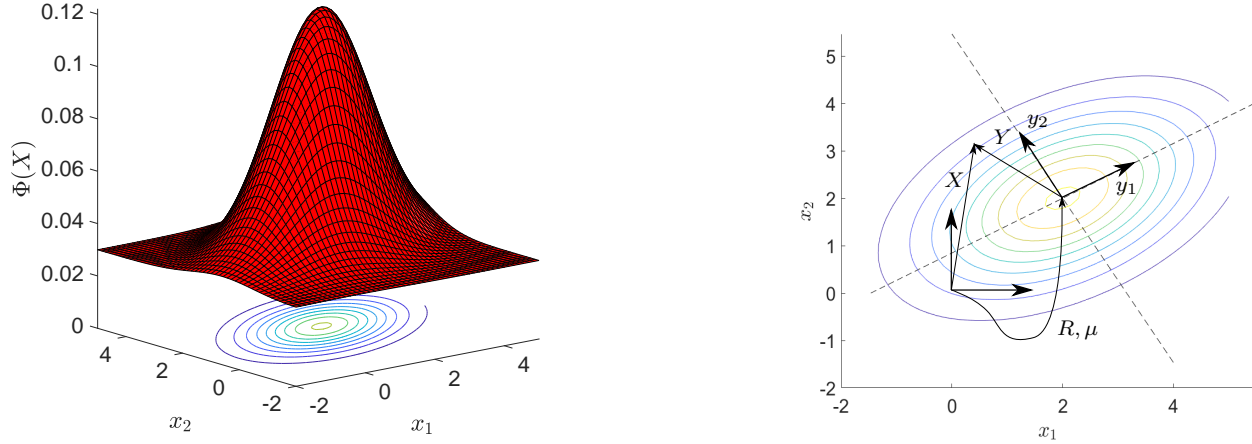


Figure 2: Example of a bi-variate Gaussian. When considering $X \in \mathcal{X}$, x_1 and x_2 are correlated, while when passing in the space of \mathcal{Y} , the bi-variate Gaussian is simply a union of two independent mono-variate Gaussians. On the left the distribution $\Phi(X)$ is represented, together with its level curve (bottom of the figure). On the right, the transformation to apply for obtaining the space \mathcal{Y} .

(2) and sharing the same domain \mathcal{X} . Indeed, considering N weights $\lambda_1, \dots, \lambda_N$, the density of the mixture f_{mix} is defined as follows:

$$f_{mix}(x) = \sum_{i=1}^N \lambda_i f_i(x) \quad (22)$$

To ensure that f_{mix} is in turn a valid probability density function satisfying equation 2, it is necessary to impose that the combination expressed in equation (22) should be convex, meaning that:

$$\sum_{i=1}^N \lambda_i = 1 \quad (23)$$

In fact, when the above specification holds, it is true that:

$$\begin{aligned} \int_X f_{mix}(x) dx &= \int_X [\lambda_1 f_1(x) + \dots + \lambda_N f_N(x)] dx \\ &= \int_X \lambda_1 f_1(x) dx + \dots + \int_X \lambda_N f_N(x) dx \\ &= \lambda_1 \int_X f_1(x) dx + \dots + \lambda_N \int_X f_N(x) dx \\ \int_X f_{mix}(x) dx &= \lambda_1 + \dots + \lambda_N = \sum_{i=1}^N \lambda_i = 1 \end{aligned} \quad (24)$$

The simplifications made in the above equation are justified by the fact that every function f_i is a probability distribution function satisfying equation (2):

$$\begin{cases} \int_X f_1(x) dx = 1 \\ \vdots \\ \int_X f_N(x) dx = 1 \end{cases} \quad (25)$$

It is worth noticing that no particular hypothesis were posed about the combining distributions f_1, \dots, f_N . We need only to require they are valid probability density functions defined over the same domain. The above considerations are true also when considering multivariate distributions.

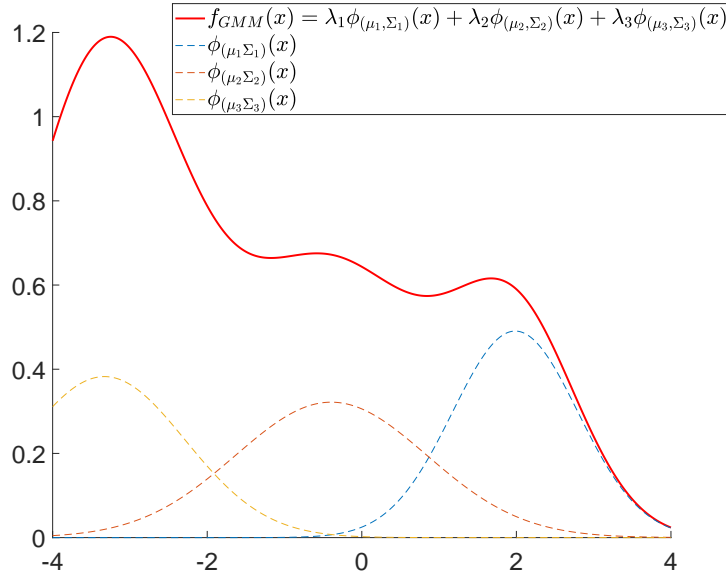


Figure 3: Example of Gaussian mixture involving univariate Gaussians. The parameters of the mixture have the same values: $\mu_1 = 1.9852$, $\mu_2 = -0.3957$, $\mu_3 = -3.3294$; $\Sigma_1 = 0.8131$, $\Sigma_2 = 1.24$, $\Sigma_3 = 1.0429$ and $\lambda_1 = 0.3$, $\lambda_2 = 0.6$, $\lambda_3 = 0.1$.

4 Gaussian Mixture models

Gaussian mixture models are particular mixture models combining N multivariate Gaussian distributions. The distribution of a GMM, f_{GMM} , is defined in this way:

$$f_{GMM}(x_1, \dots, x_n) = \sum_{i=1}^N \lambda_i \Phi_{(\mu_i, \Sigma_i)}(x_1, \dots, x_n) \quad (26)$$

where in the above equation Σ_i and μ_i are the variance and mean of the i^{th} Gaussian distribution in the mixture. An example of uni-variate Gaussian mixture is reported in Figure 3.

GMM can be also adopted for approximating complex unknown distributions. Indeed, considering the proper number of clusters, any kind of probability distribution can be approximated by a GMM.

As an example, consider to approximate the uniform distribution $U(0, 1)$ with a GMM, f_{approx} , made of n_{mix} components defined as follows:

$$f_{approx}(x|n_{mix}) = \sum_{i=1}^{n_{mix}} \Phi_{(\mu_i, \Sigma_i)} \quad \Sigma_i = \left(\frac{1}{n_{mix}} \right)^2 \quad \mu_i = \frac{1}{n_{mix}} (i - 0.5) \quad (27)$$

In Figure 4, the density of $U(0, 1)$ is compared with the one of f_{approx} , varying the number of approximating clusters. As can be visually appreciated, the approximation error decreases when increasing the number of clusters in the model. This phenomenon is not always verified when considering other kind of distributions.

5 Learning of GMM using Expectation maximization

The values of the weights $\lambda_1, \dots, \lambda_N$, as well the parameters characterizing every single density Φ_i in the mixture, i.e. the covariances and the means, are tuned through a learning process which considers as training set $S = \langle X^1, \dots, X^M \rangle$, i.e. M realizations of f_{GMM} . Such learning process is performed using the expectation maximization algorithm.

EM, is essentially an iterative algorithm that starts from an initial guess for the parameters of the clusters $\theta = \{\dots, \lambda_i, \Sigma_i, \mu_i, \dots\}$, and then adjusts the model values until the convergence to a maximum for the likelihood $L(\theta|S)$. EM is considered as an unsupervised algorithm, since only the number of clusters must be specified when performing learning, i.e. how many clusters consider for the mixture, omitting the labels¹ of the elements in the training set.

¹The Gaussian in the mixture that produced that sample

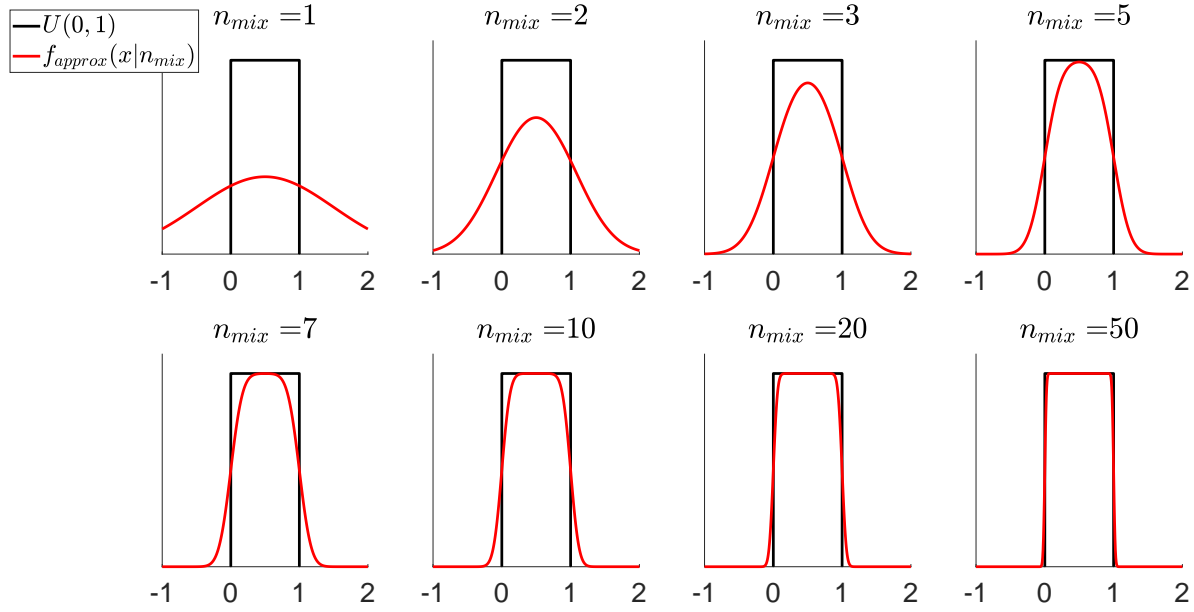


Figure 4: The uniform density between 0 and 1 is compared with an approximating mixture, varying the number of components.

The expectation-maximization algorithm can be adopted not only to learn GMM. It is in general used to learn the optimal parameter θ of a model having some observed variables X_1, \dots, X_n and also some latent ones Z_1, \dots, Z_n . Z are variables whose values are hidden, but are linked in a probabilistic way to those observed, i.e. X . Learning is done according to a training set $S = \langle X^1, \dots, X^M \rangle$, made of realizations of X : the corresponding values for Z^1, \dots, Z^M are not known. EM starts from an initial guess θ_0 and iteratively improves it. At every iteration, an Expectation and a Maximization are performed, explaining the name of the algorithm. The Expectation step is performed for taking into account the expectation of the likelihood w.r.t. Z , in order to maximise the likelihood of the training set, no matter the values for Z , which are, in a certain sense, eliminated.

As usually done, learning aims at maximizing a likelihood function involving the training set. In this case, we would like to find those θ maximising $L(X|\theta)$ ². $L(X|\theta)$ can be computed considering how the joint conditioned distribution of X, Z is factorizable:

$$\begin{aligned} \mathbb{P}(X, Z|\theta) &= \mathbb{P}(Z|X, \theta)\mathbb{P}(X|\theta) \\ \mathbb{P}(X|\theta) &= \frac{\mathbb{P}(X, Z|\theta)}{\mathbb{P}(Z|X, \theta)} \end{aligned} \quad (28)$$

Passing to the logarithms we obtain:

$$\log(\mathbb{P}(X|\theta)) = \log(\mathbb{P}(X, Z|\theta)) - \log(\mathbb{P}(Z|X, \theta)) \quad (29)$$

We are now in position to describe the Expectation step of EM algorithm. Right hand side of equation (29) is a function of Z , which is unfortunately unknown. For this reason, we want to marginalize Z , by passing to the expectations w.r.t to density $\mathbb{P}(Z|X, \theta_k)$, where θ_k are the values of the parameter at step k :

$$\begin{aligned} \sum_Z \mathbb{P}(Z|X, \theta_k) \log(\mathbb{P}(X|\theta)) &= \sum_Z \mathbb{P}(Z|X, \theta_k) \log(\mathbb{P}(X, Z|\theta)) + \dots \\ &\quad - \sum_Z \mathbb{P}(Z|X, \theta_k) \log(\mathbb{P}(Z|X, \theta)) \\ \log(\mathbb{P}(X|\theta)) \sum_Z \mathbb{P}(Z|X, \theta_k) &= \sum_Z \mathbb{P}(Z|X, \theta_k) \log(\mathbb{P}(X, Z|\theta)) + \dots \\ &\quad - \sum_Z \mathbb{P}(Z|X, \theta_k) \log(\mathbb{P}(Z|X, \theta)) \end{aligned} \quad (30)$$

²For the moment assume to have a single sample X in the training set, i.e. $S = \{X\}$ and $L(S|\theta) = L(X|\theta)$

Setting:

$$\begin{aligned} Q(\theta|\theta_k) &= \sum_Z \mathbb{P}(Z|X, \theta_k) \log(\mathbb{P}(X, Z|\theta)) \\ H(\theta_k|\theta_k) &= - \sum_Z \mathbb{P}(Z|X, \theta_k) \log(\mathbb{P}(Z|X, \theta)) \end{aligned} \quad (31)$$

leads to ³

$$\log(\mathbb{P}(X|\theta)) = Q(\theta|\theta_k) + H(\theta|\theta_k) \quad (32)$$

Considering the difference $\log(\mathbb{P}(X|\theta)) - \log(\mathbb{P}(X|\theta_k))$ and equation (32) leads to:

$$\log(\mathbb{P}(X|\theta)) - \log(\mathbb{P}(X|\theta_k)) = Q(\theta|\theta_k) - Q(\theta_k|\theta_k) + H(\theta|\theta_k) - H(\theta_k|\theta_k) \quad (33)$$

At this point we can apply the Gibbs inequality, prescribing that in case of two distributions $f_{1,2}$ defined over the same domain applies what follows:

$$- \sum_x f_1(x) \log(f_1(x)) \leq - \sum_x f_1(x) \log(f_2(x)) \quad (34)$$

Setting $f_1 = \mathbb{P}(Z|X, \theta_k)$ and $f_2 = \mathbb{P}(Z|X, \theta)$ the inequalities in equation (34) allows us to state that:

$$H(\theta|\theta_k) - H(\theta_k|\theta_k) \geq 0 \quad (35)$$

and consequently that:

$$\log(\mathbb{P}(X|\theta)) - \log(\mathbb{P}(X|\theta_k)) \geq Q(\theta|\theta_k) - Q(\theta_k|\theta_k) \quad (36)$$

For this reason, the Maximization step of the algorithm computes θ_{k+1} in order to increase Q , which leads indirectly (equation (36)) to an increase of the quantity of interest, i.e. $\log(\mathbb{P}(X|\theta))$. To be more precise, θ_{k+1} is computed as follows:

$$\theta_{k+1} = \operatorname{argmax}_{\theta} Q(\theta|\theta_k) \quad (37)$$

It is not difficult to prove that function Q , when considering a training set made of a certain number of independent samples, is a summation of terms:

$$Q(\theta|\theta_k) = \sum_i Q_i(\theta|\theta_k) = \sum_i \sum_j \mathbb{P}(Z^j|X^i, \theta_k) \log(\mathbb{P}(X^i, Z^j|\theta)) \quad (38)$$

5.0.1 Learning of Gaussian Mixture Models

In case of GMM, the latent variables Z^i are the labels specifying which cluster produced every sample X^i . Let be $\gamma_j^i = \mathbb{P}(X^i \in \text{Cluster}_j)$ (see Section 6). γ_j^i is a function of the model parameters and therefore varies along the iterations of the EM algorithm, i.e. γ_{jk}^i . Let n_{jk} be the sum of γ over samples in the training set, i.e. $n_{jk} = \sum_i \gamma_{jk}^i$. When considering GMM, it is true what follows:

$$\mathbb{P}(Z^i = j|X^i, \theta_k) = \gamma_{jk}^i \quad (39)$$

$$\begin{aligned} \mathbb{P}(X^i, Z = j|\theta) &= \mathbb{P}(X^i, |Z = j, \theta) \mathbb{P}(Z = j|\theta) \\ &= \lambda_j \left(\sqrt{2\pi |\Sigma_j|} \right)^{-1} \exp \left(-0.5 (X^i - \mu_j)^T \Sigma_j^{-1} (X^i - \mu_j) \right) \end{aligned} \quad (40)$$

In the second equation we exploited the fact that in a mixture model, weights λ expresses the a priori probability of a sample being generated by the corresponding cluster. Considering equations (39) (40), function Q in case of GMM is computable as follows:

$$\begin{aligned} Q(\theta|\theta_k) &= \sum_i \sum_j \gamma_{jk}^i \left(\log(\lambda_j) - 0.5 \log(|\Sigma_j|) - 0.5 (X^i - \mu_j)^T \Sigma_j^{-1} (X^i - \mu_j) \right) \\ &= \sum_j n_{jk} \left(\log(\lambda_j) - 0.5 \log(|\Sigma_j|) \right) + \dots \\ &\dots + \sum_i \sum_j -0.5 (X^i - \mu_j)^T \Sigma_j^{-1} (X^i - \mu_j) \end{aligned} \quad (41)$$

³Considering that $\sum_Z \mathbb{P}(Z|X, \theta_k) = 1$.

The maximization step described before, has to solve a constrained maximization problem, considering Q as objective function and $\sum_j \lambda_j = 1$ as a constraint, since GMM are mixture models. Since we deal with an equality constraints, we consider the Lagrangian function Q' :

$$Q'(\theta) = Q(\theta|\theta_k) + \xi(\sum_j \lambda_j - 1) \quad (42)$$

where ξ is the lagrangian multiplier. The maximum of Q' is obtained by finding those combinations of values for which the gradient is null.

Imposing the gradient of Q' w.r.t. the generic weight λ_j equal to 0 leads to:

$$\begin{aligned} \frac{\partial}{\partial \lambda_j} &= \frac{n_{jk}}{\lambda_j} + \xi = 0 \\ \lambda_j &= -\frac{n_{jk}}{\xi} \end{aligned} \quad (43)$$

In order to let the constraint $\sum_j \lambda_j = 1$ be satisfied, we have to prescribe that:

$$\sum_j \lambda_j = \sum_j \frac{n_{jk}}{\xi} \Rightarrow \xi = -\sum_j n_{jk} \quad (44)$$

Therefore, substituting into equation (43) leads to:

$$\lambda_{j \ k+1} = \frac{n_{jk}}{\sum_{j=1}^N n_{jk}} \quad (45)$$

The gradient of Q' w.r.t. the generic weight μ_j is equal to:

$$\begin{aligned} \frac{\partial}{\partial \mu_j} &= \sum_i \gamma_{jk}^i \frac{\partial}{\partial \mu_j} \left(-0.5(X^i - \mu_j)^T \Sigma_j^{-1} (X^i - \mu_j) \right) \\ &= \sum_i \gamma_{jk}^i \Sigma_j^{-1} (\mu_j - X^i) \\ &= \Sigma_j^{-1} \left(n_{jk} \mu_j - \sum_i \gamma_{jk}^i X^i \right) \end{aligned} \quad (46)$$

Imposing $n_{jk} \mu_j - \sum_i \gamma_{jk}^i X^i = 0$, ensures the entire gradient is null. Therefore, it applies what follows:

$$\mu_{j \ k+1} = \frac{\sum_i \gamma_{jk}^i X^i}{n_{jk}} \quad (47)$$

Finally, the gradient w.r.t. to Σ_k is equal to:

$$\begin{aligned} \frac{\partial}{\partial \Sigma_j} &= -0.5 n_{jk} \Sigma_j^{-1} + 0.5 \sum_i \gamma_{jk}^i \Sigma_j^{-1} (X^i - \mu_j)^T (X^i - \mu_j) \Sigma_j^{-1} \\ &= 0.5 \Sigma_j^{-1} \left(-n_{jk} I + \left(\sum_i \gamma_{jk}^i (X^i - \mu_j)^T (X^i - \mu_j) \right) \Sigma_j^{-1} \right) \end{aligned} \quad (48)$$

Imposing $-n_{jk} I + \left(\sum_i \gamma_{jk}^i (X^i - \mu_j)^T (X^i - \mu_j) \right) \Sigma_j^{-1} = 0$ leads to:

$$\Sigma_{j \ k+1} = \frac{\sum_i \gamma_{jk}^i (X^i - \mu_j)(X^i - \mu_j)^T}{n_{jk}} \quad (49)$$

At every step k of EM, every γ_{jk}^i is recomputed and equations (45), (47) and (49) are applied for updating the parameters of the mixture. After all the simplifications it turns out that the updating equations of EM, in case of training

a GMM, have an heuristic interpretation. Indeed, equation (45), the new value of the weight of a cluster simply consider the importance if that cluster w.r.t. to all the others, i.e. the summation of the probabilities that samples in the training set belongs to that cluster.

The new means of the clusters are computed as a weighted mean, equation (47), which gives more importance to those samples having an high probability to belongs to the cluster for which the mean is re evaluated. A similar consideration holds for the covariance of clusters, equation (49).

6 Classification

The functions characterizing the mixture, equation (27), can be interpreted as clusters. In such cases, the values of weights $\lambda_1, \dots, \lambda_N$ are priors for the probability that a certain value \bar{X} is part of a certain cluster. The classification of a value \bar{X} , according to the Bayes formula, can be done as follows:

$$\begin{aligned} \mathbb{P}(\bar{x} \in Cluster_i) &= \mathbb{P}(Cluster_i | \bar{X}) \propto L(\bar{X} | Cluster_i) \mathbb{P}_{prior}(Cluster_i) \\ &\propto f_i(\bar{X}) \lambda_i \\ &= \frac{f_i(\bar{X}) \lambda_i}{\sum_{j=1}^N f_j(\bar{X}) \lambda_j} = \frac{\Phi_{(\mu_i, \Sigma_i)}(\bar{X}) \lambda_i}{\sum_{j=1}^N \Phi_{(\mu_j, \Sigma_j)}(\bar{X}) \lambda_j} \end{aligned} \quad (50)$$

7 Drawing samples from a GMM

From what discussed in Section 2 (which leads to obtain equation (5)), it is always possible to find a linear transformation to obtain a multivariate Gaussian with a desired mean μ and covariance Σ , starting from an isotropic Gaussian $\Phi_{(0, I)}(Y)$, with I that is the identity matrix. Consider the change of variable $X = A \cdot Y + T$. If X is distributed as a multivariate Gaussian having a zero mean and a covariance equal to the identity matrix, then Y is a multivariate Gaussian with a covariance $\Sigma = A \cdot A^T$ and a mean $\mu = T$. The point here is how to compute A in order to obtain that desired Σ . Since Σ must be a symmetric positive definite matrix, it can be factorized as follows (due to the spectral theorem):

$$\Sigma = U \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix} U^T \quad (51)$$

$$= \left(U \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \right) \cdot \left(\begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} U^T \right) \quad (52)$$

$$= A \cdot A^T \quad (53)$$

where U are the eigenvectors of Σ and $\sigma_1^2, \dots, \sigma_n^2$ its eigenvalues. From the analysis equation (53), it is evident that:

$$A = U \cdot \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \quad (54)$$

In order to sample from a Gaussian distribution having Σ and μ as covariance and mean, it is sufficient to draw samples from the isotropic Gaussian (the one having a 0 mean and a covariance equal to the identity), i.e. take samples from n independent uni-variate Gaussians having a zero mean and a unitary variance and then transform such samples. Suppose to have sampled M values $\langle Y^1, \dots, Y^M \rangle$ from the isotropic Gaussian. Applying the following transformation to each sample:

$$X^j = U \cdot \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} Y^j + \mu \quad \forall j = \{1, \dots, M\} \quad (55)$$

M samples distributed as the desired Gaussian multivariate are obtained.

The above procedure is able to produce samples from a single Gaussian. When dealing with a GMM, a similar procedure can be followed. For each sample to draw, the cluster to consider for producing that sample must be firstly determined. The i^{th} cluster is selected with a probability equal to the corresponding weight λ_i ⁴ and then, a sample is draw from the i^{th} Gaussian by following the strategy described above.

8 Divergence of Kullback Leibler of GMMs

The divergence of two general pdfs f_1 and f_2 is a measure of how much that two distributions differ from each other. Indeed, in case that two distributions are defined over the same domain \mathcal{X} , the divergence of f_2 w.r.t f_1 is equal to:

$$D_{KL}(f_1||f_2) = \mathbb{E}_{f_1} \left[\log \left(\frac{f_1(X)}{f_2(X)} \right) \right] = \int_{\mathcal{X}} f_1(X) \log \left(\frac{f_1(X)}{f_2(X)} \right) dX_{1,\dots,n} \quad (56)$$

It is important to remark that the divergence is not symmetric, i.e. $D_{KL}(f_1||f_2) \neq D_{KL}(f_2||f_1)$. Such measure can be used to compare tow high dimensional multivariate distributions.

In case of GMM, the Kullback-Leibler divergence is not computable in a closed form. On the opposite, a Monte Carlo approach can be followed. Indeed, M samples $\langle X^1, \dots, X^M \rangle$ from f_1 can be drawn (as described in Section 7) and the following mean is assumed as an approximation of the expectation in equation (56):

$$D_{KL}(f_1||f_2) \cong \frac{1}{M} \sum_{i=1}^M \log \left(\frac{f_1(X^i)}{f_2(X^i)} \right) \quad (57)$$

Other approaches exist, providing a theoretical upper and lower bound for the real value of the divergence.

⁴The index of the cluster to consider is sampled from a discrete distribution having $\lambda_1, \dots, \lambda_N$ as values.