

Gaussian distribution by itself may not be enough in modelling real datasets

Sometimes a linear superposition of two or more Gaussians may capture data better. Such superpositions can be formulated as mixture distributions.

For a superposition of K Gaussians we have

$$p(x) = \sum_{k=1}^K w_k N(x | \mu_k, \Sigma_k)$$

Each $N(x | \mu_k, \Sigma_k)$ is called a component of the mixture. The parameters w_k are called mixing coefficients

We have $\sum_{k=1}^K w_k = 1$

Also $p(x) \geq 0$, $N(x | \mu_k, \Sigma_k) \geq 0$

$\Rightarrow w_k \geq 0 \quad \forall k$

$\therefore 0 \leq w_k \leq 1$ (can be thought of as probabilities)

$$p(x) = \sum_{k=1}^K p(k) p(x|k) \quad (w_k = p(k))$$

Think of $p(k)$ as prior of picking k^{th} component and $p(x|k) = N(x | \mu_k, \Sigma_k)$

i.e. probability of x conditioned on k

The posterior prob. $p(k|x)$ are given as

$$\begin{aligned} \gamma_k(x) &= p(k|x) \\ &= \frac{p(k) p(x|k)}{\sum_l p(l) p(x|l)} \\ &= \frac{w_k N(x | \mu_k, \Sigma_k)}{\sum_l w_l N(x | \mu_l, \Sigma_l)} \end{aligned}$$

Parameters of GMM are w , μ and Σ .

$$w = (w_1, \dots, w_K)$$

$$\mu = (\mu_1, \dots, \mu_K)$$

$$\Sigma = (\Sigma_1, \dots, \Sigma_K)$$

How to get the values?

Log of likelihood is given by

$$\ln p(\mathbf{x} | w, \mu, \Sigma) = \sum_{n=1}^N \ln \left[\sum_{k=1}^K w_k N(x_n | \mu_k, \Sigma_k) \right]$$

Another Interpretation:-

\mathbf{z} is K -dim binary random variable where at a time some z_k is 1 and all other coordinates are 0.

$$\therefore z_k \in \{0, 1\} \quad \text{and} \quad \sum_k z_k = 1$$

K possible states of the vector \mathbf{z}

Try to define joint dist. $p(\mathbf{x}, \mathbf{z})$

in terms of $p(\mathbf{z})$ and $p(\mathbf{x}|\mathbf{z})$

Now $p(z_k=1) = w_k \quad \left(\begin{array}{l} 0 \leq w_k \leq 1 \\ \sum_k w_k = 1 \end{array} \right)$

$$p(\mathbf{z}) = \prod_{k=1}^K w_k^{z_k}$$

Also $p(x|z_k=1) = N(x | \mu_k, \Sigma_k)$

$$\therefore p(x|\mathbf{z}) = \prod_{k=1}^K N(x | \mu_k, \Sigma_k)^{z_k}$$

Now

$$\begin{aligned} p(x) &= \sum_{\mathbf{z}} p(\mathbf{z}) p(x|\mathbf{z}) \\ &= \sum_{k=1}^K w_k N(x | \mu_k, \Sigma_k) \end{aligned}$$

\therefore Marginal of \mathbf{x} is GMM

If we have observations x_1, \dots, x_N

and $\approx p(x) = \sum_{\mathbf{z}} p(x, \mathbf{z})$

for every observed data point x_n there is corresponding latent variable \mathbf{z}_n

What is the advantage?

Now $\gamma(z_k)$ i.e. $p(z_k=1|x)$

$$\begin{aligned} &= \frac{p(z_k=1) p(x|z_k=1)}{\sum_l p(z_l=1) p(x|z_l=1)} \\ &= \frac{w_k N(x | \mu_k, \Sigma_k)}{\sum_l w_l N(x | \mu_l, \Sigma_l)} \end{aligned}$$

Maximizing log likelihood of GMM is more complex problem than the case of a single Gaussian. Because of summation over k inside log, log fn. does not act directly on Gaussian.

We don't get closed form soln by equating derivative to zero.