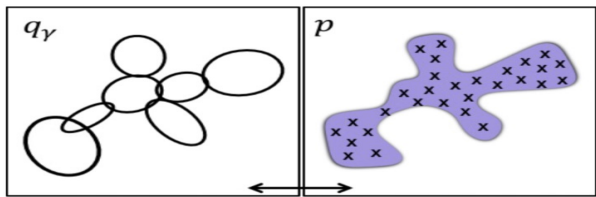


Gaussian Mixture Models

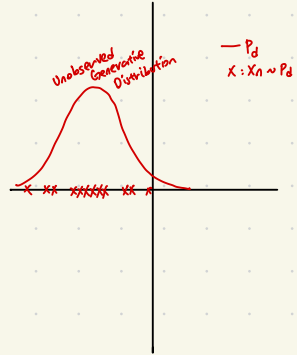


Reminder

$$\text{Gaussian Distribution: } p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Assume we have N i.i.d samples from Gaussian distribution, $\{x_n \sim p_\theta\}_{n=1}^N$

$$\mu = \frac{1}{N} \sum_{i=1}^N x^{(i)}, \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \mu)^2$$



1. Maximum Likelihood: \star [likelihood: Measure of how well a statistical model explains the observed data.]

$$\arg \max_{\mu, \sigma} p(x_1, x_2, \dots, x_n; \mu, \sigma) = \prod_{n=1}^N p(x_n; \mu, \sigma)$$

2. Maximum log-likelihood:

$$\arg \max_{\mu, \sigma} \log \left(\prod_{n=1}^N p(x_1, x_2, \dots, x_n) \right) = \sum_{n=1}^N \log(p(x_n; \mu, \sigma))$$

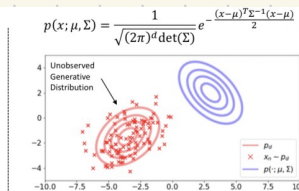
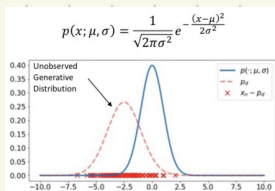
3. Minimizing Negative Log likelihood:

$$\arg \min_{\mu, \sigma} \sum_{n=1}^N \left(\frac{\log(2\pi\sigma^2)}{2} + \frac{(x_n - \mu)^2}{2\sigma^2} \right)$$

$$\frac{\partial l}{\partial \mu} = 0 \Rightarrow \mu^* = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\frac{\partial l}{\partial \sigma} = 0 \Rightarrow \sigma^* = \frac{1}{N} \sum_{n=1}^N (x_n - \mu^*)^2$$

Multivariate Gaussian



$$\frac{\partial l}{\partial \mu} = 0 \Rightarrow \mu^* = \frac{1}{N} \sum_{n=1}^N x_n$$

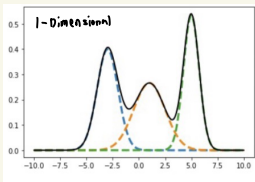
$$\frac{\partial l}{\partial \mu} = 0 \Rightarrow \mu^* = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\frac{\partial l}{\partial \sigma} = 0 \Rightarrow \sigma^* = \frac{1}{N} \sum_{n=1}^N (x_n - \mu^*)^2$$

$$\frac{\partial l}{\partial \sigma} = 0 \Rightarrow \sigma^* = \frac{1}{N} \sum_{n=1}^N (x_n - \mu^*)(x_n - \mu^*)^T$$

\star Maximizing log-likelihood is equivalent to minimizing KL divergence

Mixture of Gaussian



$$p(x_i | [a_k, \mu_k, \sigma_k]_{k=1}^K) = \sum_{k=1}^K a_k N(x_i | \mu_k, \sigma_k) = \sum_{k=1}^K \frac{a_k}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}$$

where $a_k \geq 0$, $\sum_{k=1}^K a_k = 1$, and $\sigma_k \geq 0$.

$$\begin{aligned} \mu &\in \mathbb{R}^{K \times \text{dim}}, \mu_k \in \mathbb{R}^{1 \times \text{dim}} \\ \sigma &\in \mathbb{R}^{K \times \text{dim}}, \sigma_k \in \mathbb{R}^{1 \times \text{dim}} \\ N(x_i | \mu_k, \sigma_k) &\in \mathbb{R}^{N \times \text{dim}} \end{aligned}$$

$$\text{Log-likelihood} : \log \left(\prod_{n=1}^N \left(\sum_{k=1}^K a_k N(x_{ni} | \mu_k, \sigma_k) \right) \right) = \sum_{n=1}^N \log \left(\sum_{k=1}^K a_k N(x_{ni} | \mu_k, \sigma_k) \right)$$

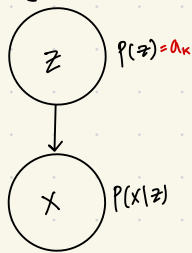
log(Sum(...)) doesn't have a closed form solution

↳ How do we deal with this? : Expectation Maximization

Expectation Maximization - Origins

$$Z = [z_1, \dots, z_K]$$

indicator of Gaussian



$$\begin{aligned} p(z_k=1) &= a_k \\ p(x | z_k=1) &= \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{(x - \mu_k)^2}{2\sigma_k^2}} \end{aligned}$$

$$p(x) = \sum_{k=1}^K p(x | z_k=1) p(z_k=1) = \sum_{k=1}^K \underbrace{\frac{a_k}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{(x - \mu_k)^2}{2\sigma_k^2}}}_{a_k \cdot N(x_{ni} | \mu_k, \sigma_k)}$$

$$p(x) = 0.3 \cdot N(x_{ni} | \mu_1, \sigma_1) + 0.7 \cdot N(x_{ni} | \mu_2, \sigma_2)$$

Since log-likelihood of GMM doesn't have a closed-form solution, we need gradient of likelihood (for optimizer, such as Gradient Descent)

$$\frac{d}{d\theta} \log p(x) = \frac{d}{d\theta} \log \left(\sum_z p(z|x) \right)$$

where $\theta = \{\mu_1, \dots, \mu_K, \sigma_1, \dots, \sigma_K\}$

$$\dots = \sum_z p(z|x) \frac{d}{d\theta} \log(p(x|z)p(z)) = \sum_z p(z|x) \frac{d}{d\theta} \log(p(x|z)) + \sum_z p(z|x) \frac{d}{d\theta} \log(p(z))$$

We can optimize μ, σ We can optimize a_k

$$\text{And, } p(z|x) = \frac{p(x|z)p(z)}{\sum_{z'} p(x|z')p(z')} : \text{Is soft assignment of data } x \text{ to each Gaussian}$$

Expectation Maximization Algorithm

• Expectation step: for fixed parameters $[(a_k, \mu_k, \sigma_k)]_{k=1}^K$ compute $r_n^k = p(z_k=1 | x_n)$ for each sample

$$r_n^k = \frac{a_k N(x_{ni} | \mu_k, \sigma_k)}{\sum_{j=1}^K a_j N(x_{ni} | \mu_j, \sigma_j)}$$

• Maximization Step: for fixed r_n^k solve the maximum log-likelihood to obtain optimal parameters:

- Mixture coefficients: $a_k = \frac{N_k}{N}$ for $N_k = \sum_n r_n^k$

- Means: $\mu_k = \frac{1}{N_k} \sum_n r_n^k x_n$

- Variances: $\sigma_k^2 = \frac{1}{N_k} \sum_n r_n^k (x_n - \mu_k)^2$

- Covariances: $\sum_k = \frac{1}{N_k} \sum_n r_n^k \underbrace{(x_n - \mu_k)}_{2.6 \times 1.2} \underbrace{(x_n - \mu_k)^T}_{2 \times 2.1 \times 2} : 2.100 \times 2$

EM-GMM Algorithm Pseudocode (Input: X , $k=3$)

- Initialize K Gaussian distributions and their mixture coefficients

$$p(x|z_{k=1}) = \mathcal{N}(M_k, \Sigma_k)$$

$$p(z_{k=1}) = a_k = \frac{1}{K}$$

- Iterate until convergence

- For fixed GMM parameters, calculate Soft Assignments

$$\star \text{Soft Assignment: } p(z|x) = \frac{p(x|z)p(z)}{\sum_z p(x|z)p(z)}$$

$$p(z_{k=1}) = \frac{p(x|z_{k=1})p(z_{k=1})}{\sum_{i=1}^K p(x|z_i)p(z_i)}$$

- For fixed soft assignments, calculate GMM parameters

$$\bullet M_k = \frac{\sum_{n=1}^N p(z_k=1|x_n) x_n}{\sum_{n=1}^N p(z_k=1|x_n)}$$

$$\bullet \Sigma_k = \frac{\sum_{n=1}^N p(z_k=1|x_n) (x_n - M_k)(x_n - M_k)^T}{\sum_{n=1}^N p(z_k=1|x_n)}$$

$$\bullet a_k = \frac{\sum_{n=1}^N p(z_k=1|x_n)}{N}$$