

## Gaussian Mixture Model. using EM method.

Why do we need gaussian model (Advantages)

- 1) It is very similar to the normal distribution
- 2) It is easier to do mathematical manipulations and calculations using a gaussian model as compared to the other ones.  
→ we can differentiate infinite times.

Sometimes the nature of the distributions is not gaussian; but we divide it into many models assuming each of the individual model is gaussian. This is called the gaussian mixture model or GMM.

### Univariate Gaussian Distribution (1D)

$$G(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

mean                      covariance

## Multi variate gaussian Distribution:

$$N(x|\mu, \Sigma) = \frac{1}{(2\pi|\Sigma|)^{1/2}} \exp\left[-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right]$$

↙   ↓   ↘

this is a mean covariance  
vector with  
dimension  
of  $x$

To apply GMM we need to estimate  
( $\Sigma, \mu$ ) of a distribution.

One method  $\rightarrow$  Maximum Likelihood  
estimation. (ML)

### \* ML method

1) consider the log of gaussian  
distribution

$$\ln P(x|\mu, \Sigma) = -\frac{1}{2} \ln 2\pi|\Sigma| - \frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)$$

2) we take its derivative and equate  
it to zero. ie.

$$\frac{\partial \ln p(x|\mu, \Sigma)}{\partial \mu} = 0, \quad \frac{\partial \ln p(x|\mu, \Sigma)}{\partial \Sigma} = 0$$

$\Downarrow$

$\Downarrow$

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n, \quad \Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})(x_n - \mu_{ML})^T$$



where  $N$  is the number of samples or data points.

## Gaussian Mixture.

\* it is the linear superposition of Gaussians !.

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

no. of gaussians  $\rightarrow$  normal multivariate gaussian distribution

mixing coeff: weightage for each gaussian distribution.

\* Normalization and positivity require:

$$0 \leq \pi_k \leq 1; \quad \sum_{k=1}^K \pi_k = 1$$

if we consider log-likelihood

$$\ln p(x | \mu, \Sigma, \pi) = \sum_{n=1}^N \ln P(x_n) = \sum_{n=1}^N \ln \left[ \pi_k N(x_n | \mu_k, \Sigma_k) \right]$$

Maximum likelihood doesn't work here as there is no closed form solution.

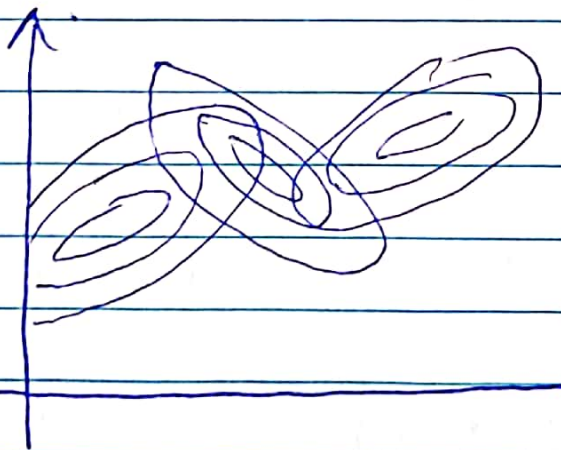
So for calculating parameters we can use the Expectation Maximization technique.

Or the EM technique.

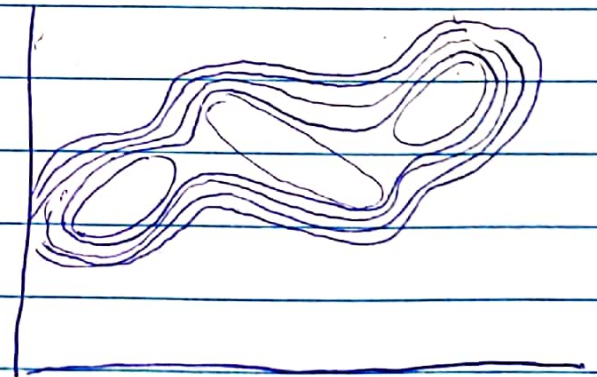
Extra

⇒ Mixture of 3 gaussians example:

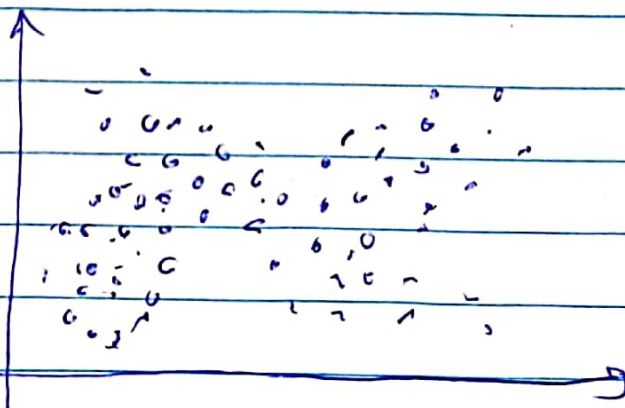
Let us represent the contour plot of 3 gaussians put together



This is the combined Contour plot of the same



Sampled data.





Latent variable: posterior prob.

\* we can think of the mixing coeff as prior probabilities for the components

\* For a given value of 'k', we can evaluate the corresponding posterior probabilities, called responsibilities.

From Bayes Rule.

$$\gamma_k(x) = P(k|x) = \frac{P(k)P(x|k)}{P(x)}$$

latent

$$\text{Variable} = \pi_k N(x|\mu_k, \Sigma_k) \quad \text{where} \quad \pi_k = \frac{N_k}{N}$$
$$\sum_{j=1}^K \pi_j N(x|\mu_j, \Sigma_j)$$

Interpret  $N_k$  as the effective no. of points assigned to cluster  $k$ .

What does EM algorithm do?

\* EM algorithm is an iterative optimization technique which is operated locally. to find set of parameters.

which all parameters do we need to estimate?

→ we need to estimate many parameters  
let's say we are planning to use  $K$  gaussian surfaces then we need to predict  $\mu_k$ ,  $\Sigma_k$  and  $\pi_k$

$\nwarrow \quad \downarrow \quad \swarrow$   
mean      covariance      this  
matrix    matrix      is a  
                         scalar.

for each gaussian.  
hence in total we need to guess  $K \times 3$  parameters.

\* There are two steps in EM

~~1) Estimation step~~

1) Estimation step!

for a given parameter values we can compute the expected values of the latent variable

2) Maximization step:

updates the parameters of our model based on the latent variable calculated using ML method.



## Algorithm for GMM;

Given a gaussian mixture model the goal is to maximize the likelihood function with respect to the parameters comprising the means and covariances of the components and the mixing coeffs.

- 1) Initialize the means  $\mu_j$ ; covariances  $\Sigma_j$  and mixing coeff  $\pi_j$ ; and evaluate the initial value of the log likelihood. The starting points can be completely random.

also we can initialize them at some known points such as the mean of the total dataset. etc.

- \* The idea is that closer we are to the final solution, faster it will be our solution.

- 2) E. Step, Evaluate the responsibilities using the current parameter values.

$$\gamma_k(x) = \frac{\pi_k N(x | \mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)}$$

summing up the gaussians with the estimated parameters.

3.) M. Step: Re estimate the parameters using the current responsibilities.

Till now we have calculated  $\gamma_j$ 's

$$\mu_j = \frac{\sum_{n=1}^N \gamma_j(x_n) x_n}{\sum_{n=1}^N \gamma_j(x_n)} \quad \left. \vphantom{\sum_{n=1}^N} \right\} \text{ same as ML}$$

$$\Sigma_j = \frac{\sum_{n=1}^N \gamma_j(x_n) (x_n - \mu_j)(x_n - \mu_j)^T}{\sum_{n=1}^N \gamma_j(x_n)}$$

$$\pi_j = \frac{1}{N} \sum_{n=1}^N \gamma_j(x_n)$$

4.) Evaluate log likelihood

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

if there is no convergence return to step 2.

otherwise we have reached our final answer.



We need to repeat the steps 2, 3, 4 until we reach convergence.

Again There can be more than one methods to check convergence

- 1) If the difference between the consecutive estimations of the log likelihood is below a certain threshold we can say that the estimation converged.

Another method may be to do the same but by using the parameters  $\pi_j, \xi_j, \mu_j$  instead of log likelihood.

\* EXTRA

What is  $\delta_k$  (Latent function)?

→ The mathematical definition of  $\delta_k$  is:

$$\delta_k = \begin{cases} 1 & \text{at gaussian } k \\ 0 & \text{otherwise} \end{cases}$$

It can be considered as a type of probability density function. It behaves like the dirac delta function of  $k$ -means.

\* EXTRA

How to pick the value of  $K$  in GMM

\* One method is to pick that value of  $K$  which minimizes the value of  $L = \log(P(x_1, x_2, \dots, x_n)) = \sum_{i=1}^n \log \sum_{k=1}^K P(x_i | k) P(k)$

using brute force.

→ but this will result in a value of  $K$  that is equal to ~~number~~ number of training data points.



\* Another option is to split the data points into training( $T$ ) and validation( $V$ ) set.

→ and then for each  $K$ : fit the parameters of  $T$  and measure the likelihood of  $V$ . based on the above results we can select a  $K$  which will give min. likelihood.

But again, in this method as well we may end up with a very large no. <sup>for</sup> ~~of~~  $K$ .

There are many methods of selecting a  $K$  which involve matrices and complex calculations.

According to Occam's Razor if we have two explanations to a certain problem in datasets we must always select the explanation that is simpler that way there are less chances of complications.

Hence: Occam's Razor: pick "simplest" of all models that fit. just

Now instead of using likelihood we penalize it as the complexity of the model increases.

→ Akaike info. criterion (AIC):  $-L + K$

→ Bayesian info. criterion (BIC):  $-L + \frac{1}{2} K \ln(N)$

here  $L \rightarrow$  is the likelihood;  
 $K \rightarrow$  is no. of gaussians.

These are two popular types of penalties.

Instead of maximizing  $L$ , minimize  $-L$  and add the defined penalty.

Finally now that we have defined this penalized model we can get  $K$  by the two methods discussed before.