

A Tutorial of Gaussian Mixture Model (GMM)

Nov. 2015

Qiuqiang Kong
q.kong@qmul.ac.uk

- Gaussian Distribution
- Gaussian Mixture Model (GMM)
- A General PGM View of EM algorithm
- GMM Revise
- Tricks

Gaussian Distribution

The probability density function of Gaussian distribution is

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right) \quad (1)$$

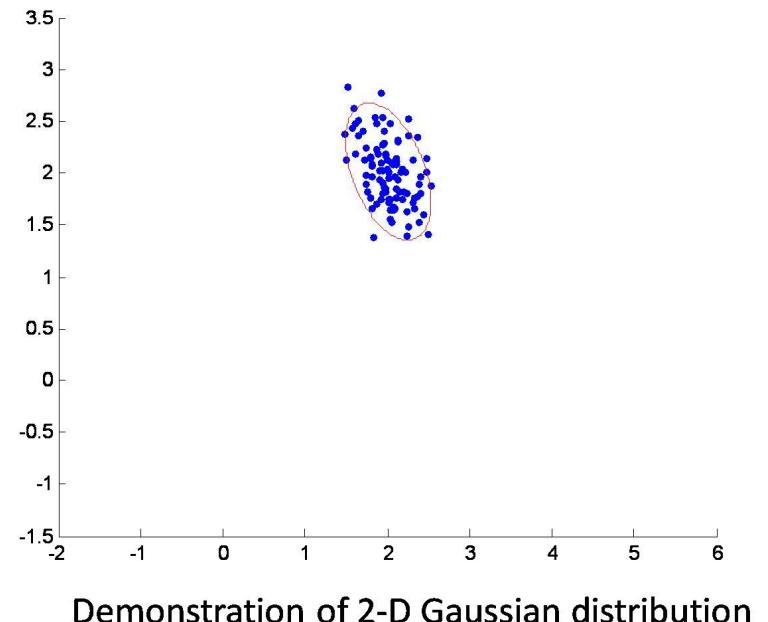
Gaussian distribution is denominated by parameters: $\theta = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$

where

$\boldsymbol{\mu}_{p \times 1}$ is mean $\boldsymbol{\Sigma}_{p \times p}$ is covariance

Q: How to estimate the parameters?

A: Using maximum likelihood estimation.



Applying maximum likelihood estimation method

$$\{\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}\} = \arg \max_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} L(\mathbf{x}_1, \dots, \mathbf{x}_N | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (2)$$

where $L(\mathbf{x}_1, \dots, \mathbf{x}_N | \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is likelihood function

$$\begin{aligned} L(\mathbf{x}_1, \dots, \mathbf{x}_N | \boldsymbol{\theta}) &= \ln \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\theta}) \\ &= \sum_{n=1}^N \left(-\frac{D}{2} \ln(2\pi) - \frac{1}{2} |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right) \end{aligned} \quad (3)$$

Optimize (3) with respect to $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}$ simultaneously is difficult.
However, sample mean is independent of sample variance. So the optimization can be done with respect to $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}$ separately.

Let the derivative of likelihood L with respect to μ, Σ be zero. We get estimation of μ, Σ .

$$\frac{\partial L}{\partial \mu} = \sum_{n=1}^N -\Sigma^{-1}(\mathbf{x}_n - \mu) = 0$$

$$\frac{\partial L}{\partial \Sigma} = \sum_{n=1}^N \left(-\frac{1}{2} \Sigma^{-T} + \frac{1}{2} \Sigma^{-T} (\mathbf{x}_n - \mu)(\mathbf{x}_n - \mu)^T \Sigma^{-T} \right) = 0$$

$$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \quad (4)$$

$$\hat{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \hat{\mu})(\mathbf{x}_n - \hat{\mu})^T$$

Because $L(\mathbf{x}_1, \dots, \mathbf{x}_N | \mu, \Sigma)$ is convex function,

The extreme point is global maximum point.

$\hat{\mu}$ is unbiased estimation

$\hat{\Sigma}$ is biased estimation, when $n > \infty$, it is asymptotic unbiased

Hint^[1]: $\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = 2\mathbf{A}\mathbf{x}$

$$\frac{\partial |\mathbf{X}|}{\mathbf{X}} = |\mathbf{X}| \mathbf{X}^{-T}$$

$$\frac{\partial \mathbf{a}^T \mathbf{X}^{-1} \mathbf{b}}{\mathbf{X}} = -\mathbf{X}^{-T} \mathbf{a} \mathbf{b}^T \mathbf{X}^{-T}$$

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

Gaussian Mixture Model (GMM)

Intuitively, GMM is the sum of several weighted Gaussian distributions.

Probability density function of GMM:

$$p(\mathbf{x}|\theta) = \sum_{k=1}^K \pi_k N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (5)$$

where $\sum_{k=1}^K \pi_k = 1$

Parameters: $\theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\} \quad k = 1, \dots, K$

Illustration

It is unknown which mixture produce \mathbf{x}_n .

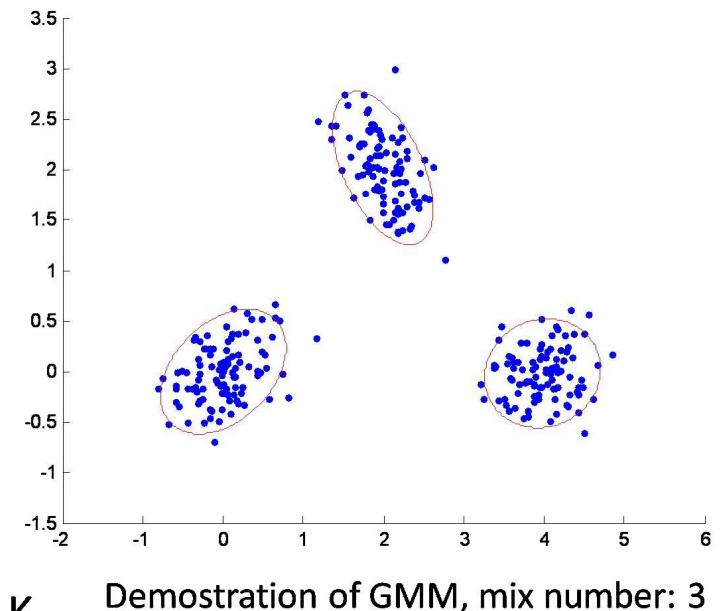
Define event A_k : "the k -th mixture occurs", $k=1, \dots, K$

with probability on A_k : $p(A_k) = \pi_k$

According to Total Probability Theorem

$$p(\mathbf{x}) = \sum_{k=1}^K p(\mathbf{x}|A_k)p(A_k) = \sum_{k=1}^K \pi_k N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (6)$$

which has same form as (5).



GMM parameter estimation using maximum likelihood

$$L(\mathbf{x}_1, \dots, \mathbf{x}_N | \theta) = \ln \prod_{n=1}^N p(\mathbf{x}_n | \theta) = \sum_{n=1}^N \ln \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (7)$$

Take constraint $\sum_{k=1}^K \pi_k = 1$ into account. Introduce Lagrange function

$$R(\theta, \lambda) = \sum_{n=1}^N \ln \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) \quad (8)$$

To optimize $R(\theta, \lambda)$ Let the derivative of $R(\theta, \lambda)$ with respect to

$\theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}$ $k = 1, \dots, K$ be zero.

$$\begin{aligned}
 \frac{\partial L}{\partial \boldsymbol{\mu}_k} &= \sum_{n=1}^N \frac{\pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \left(-\boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right) = 0 \quad k = 1, \dots, K \\
 \frac{\partial L}{\partial \boldsymbol{\Sigma}_k} &= \sum_{n=1}^N \frac{\pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \left(-\frac{1}{2} \boldsymbol{\Sigma}_k^{-T} + \frac{1}{2} \boldsymbol{\Sigma}_k^{-T} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-T} \right) = 0 \quad k = 1, \dots, K \\
 \frac{\partial L}{\partial \pi_k} &= \sum_{n=1}^N \frac{N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} + \lambda = 0 \quad k = 1, \dots, K
 \end{aligned} \tag{9}$$

This is a transcendental equation set, can not be solved analytically.

Solve: introduce iterative method (EM algorithm)
while (not converge)

**E step: use old parameters
 to calculate**

$$\frac{\partial L}{\partial \boldsymbol{\mu}_k} = \sum_{n=1}^N \frac{\pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \left(-\boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right) = 0 \quad k = 1, \dots, K \tag{10}$$

**M step: update
 parameters**

notice
$$\frac{\pi_k N(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} = \frac{p(\mathbf{x} | A_k) p(A_k)}{p(\mathbf{x})} = p(A_k | \mathbf{x}) \quad (11)$$

EM algorithm

1. Choose initial θ^{old}
2. E step: (Expectation)

$$\gamma_{nk} = p(A_k | \mathbf{x}_n \theta^{old}) \quad (12)$$

where $\theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\} \quad k = 1, \dots, K$

3. M step: (Maximization)

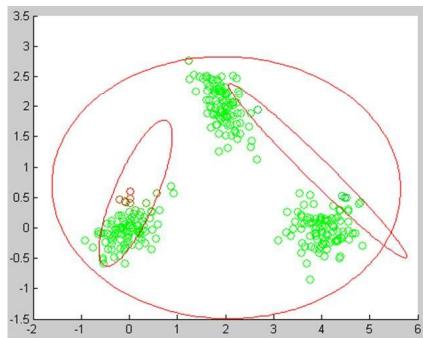
$$\begin{aligned} \pi_k &= \frac{1}{N} \sum_{n=1}^N \gamma_{nk} & \boldsymbol{\mu}_k &= \frac{\sum_{n=1}^N \gamma_{nk} \mathbf{x}_n}{\sum_{n=1}^N \gamma_{nk}} & \boldsymbol{\Sigma}_k &= \frac{\sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma_{nk}} \end{aligned} \quad (13)$$

4. If converge then stop. Otherwise goto 2.

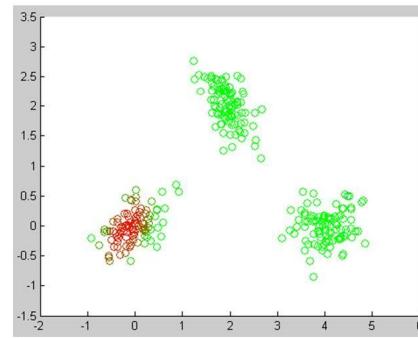
Q: Will likelihood decrease using EM algorithm?

Q: Will likelihood converge using EM algorithm?

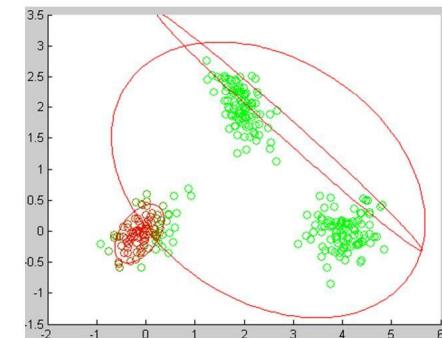
Example of EM algorithm



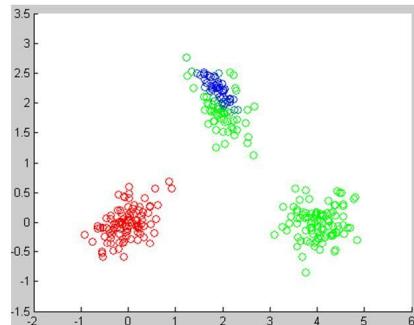
Initialize parameters



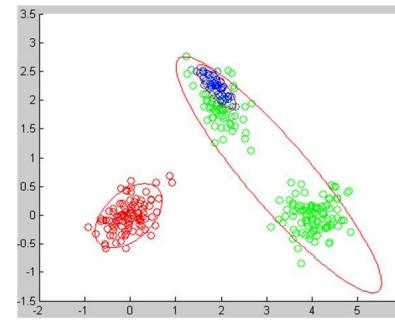
E step



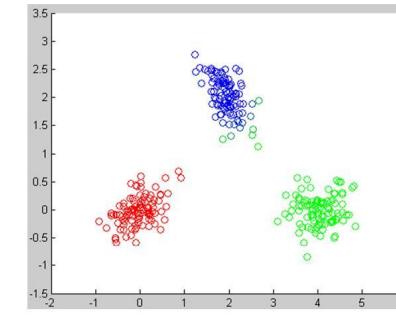
M step



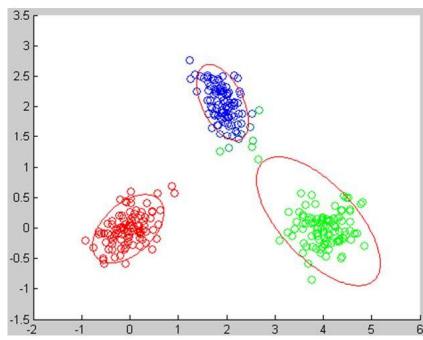
E step



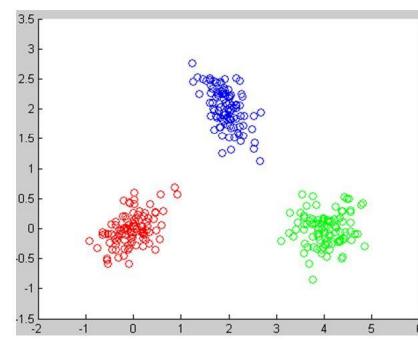
M step



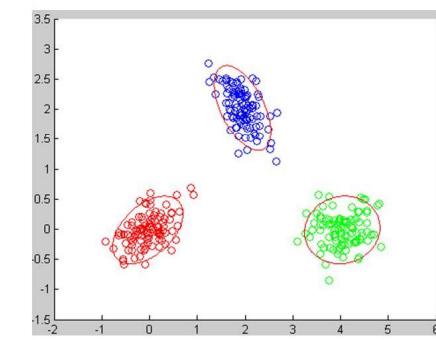
E step



M step



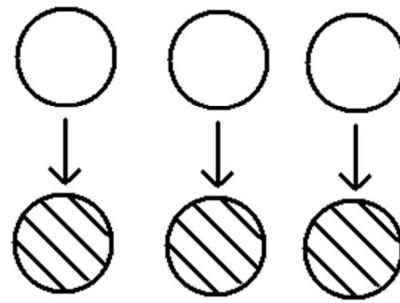
E step



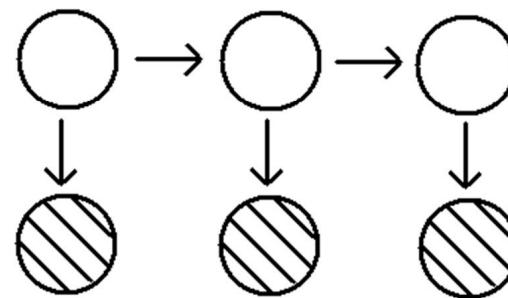
M step

A General View of EM Algorithm

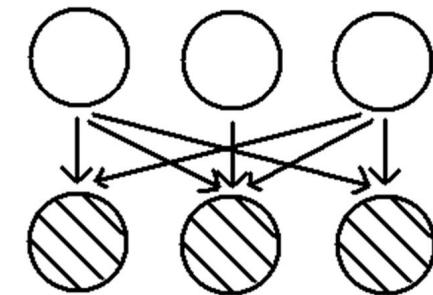
In general, EM algorithm is a method to estimate parameters in a probabilistic graphic model (PGM) with latent variables.



GMM



HMM



Others

Black circle is the observed variable. White circle is latent variable.

The estimation of θ will be $\hat{\theta} = \arg \max_{\theta} \ln p(\mathbf{X}|\theta)$

where $\ln p(\mathbf{X}|\theta) = \ln \sum_{\mathbf{Z}} p(\mathbf{XZ}|\theta)$

Obviously, it is needed to sum all \mathbf{Z} before applying log, which is difficult.
EM algorithm is aimed at solving this problem.

Introduce auxiliary function $q(\mathbf{Z})$ over \mathbf{Z}

$$\begin{aligned}
 \ln p(\mathbf{X}|\theta) &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln p(\mathbf{X}|\theta) \\
 &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{XZ}|\theta)}{p(\mathbf{Z}|\mathbf{X}\theta)} \\
 &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{XZ}|\theta)q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X}\theta)q(\mathbf{Z})} \\
 &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{XZ}|\theta)}{q(\mathbf{Z})} + \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X}\theta)} \\
 &= \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \frac{p(\mathbf{XZ}|\theta)}{q(\mathbf{Z})} + \text{KL}\left(q(\mathbf{Z}) \| p(\mathbf{Z}|\mathbf{X}\theta)\right)
 \end{aligned} \tag{14}$$

① ②

② is KL divergence between $q(\mathbf{Z})$ and $p(\mathbf{Z}|\mathbf{X}\theta)$, which is non negative.
It only equals 0 iff $q(\mathbf{Z})$ and $p(\mathbf{Z}|\mathbf{X}\theta)$.

E step: Let $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}\theta^{\text{old}})$

M step: optimize

$$\ln p(\mathbf{X}|\theta) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{XZ}|\theta) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{Z}|\mathbf{X}\theta^{\text{old}}) = Q(\theta, \theta^{\text{old}}) + \text{const} \tag{15}$$

with respect to θ

$$\text{where } Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{XZ}|\theta) \tag{16}$$

EM algorithm for PGM with latent variables

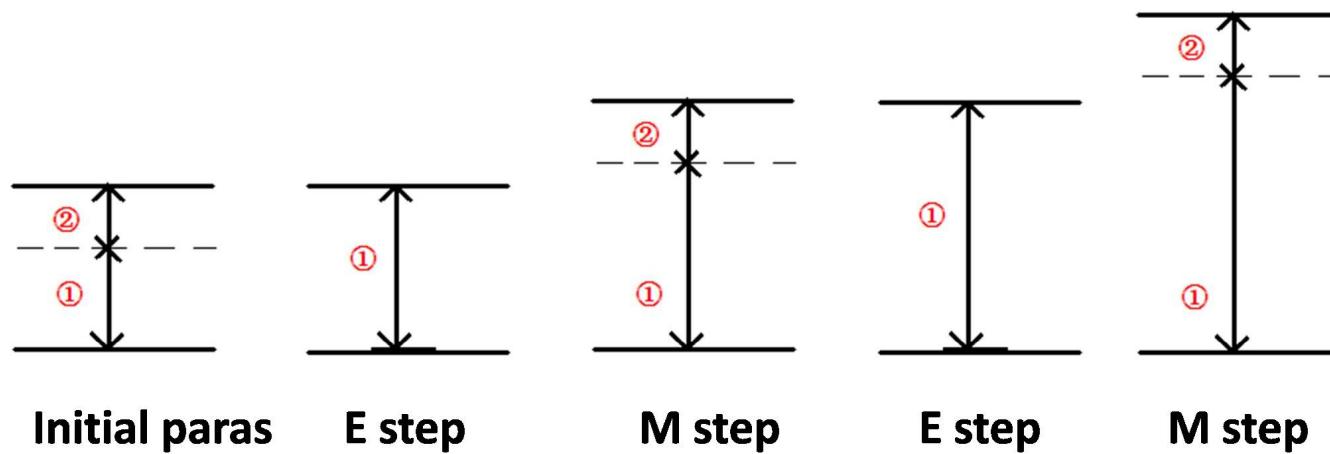
1. Init parameters

2. E step: $q(\mathbf{Z}) = p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}})$ (17)

3. M step: $\hat{\theta} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}})$ (18)

where $Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{XZ} | \theta)$

4. If converge then stop, otherwise goto 2



Illustration^[1]: For a function which is difficult to find extreme point (red curve), E step is find the complete data function (blue curve), M step is update θ .

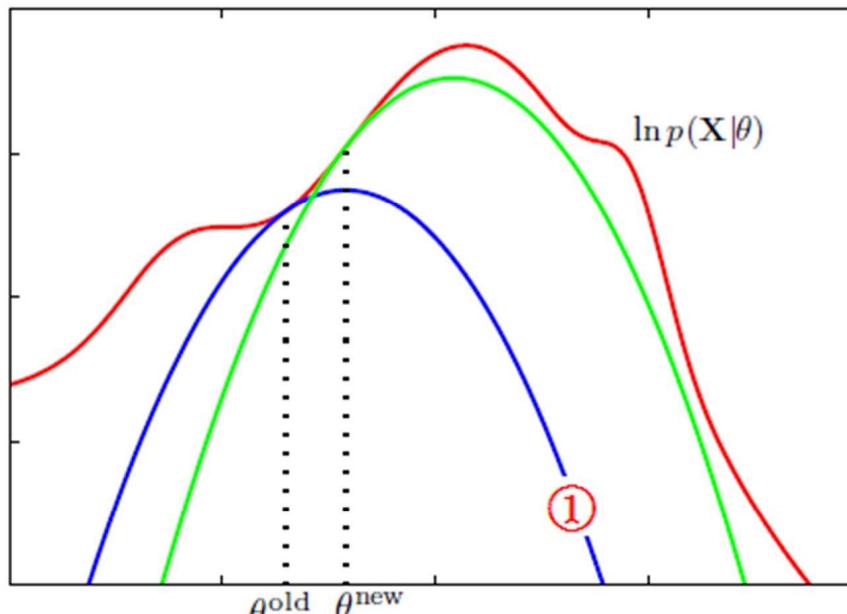


Illustration of EM algorithm

Q: Will EM algorithm converge?

A: Under some conditions, EM algorithm will converge to stationary point^[2].

[1] Bishop, Christopher M. *Pattern recognition and machine learning*. Springer, 2006. P453

[2] Wu, CF Jeff. "On the convergence properties of the EM algorithm." *The Annals of statistics* (1983): 95-103.

GMM Revise

Define event A_k : “the k-th mixture occurs”, $k=1, \dots, K$

with probability on A_k : $p(A_k) = \pi_k$

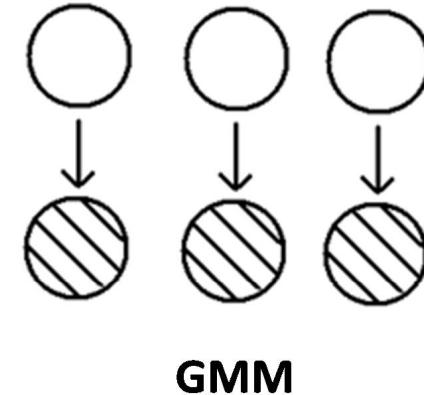
Let \mathbf{z} be K-dimensional binary random variable
having 1-of-K representation

$$\mathbf{z} : A_k \rightarrow (0, \dots, 0, 1, 0, \dots, 0) \quad k = 1, \dots, K \quad (19)$$

(the k-th position is 1)

The probability over \mathbf{z} is

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} \quad (20)$$



Thinking:

Q: Why choose random variable in this form?

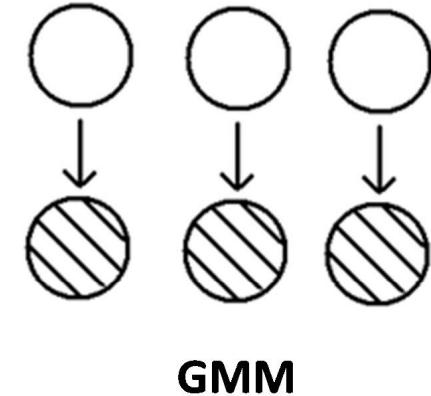
Can we choose others, eg.

$\mathbf{z} : A_k \rightarrow k, k=1, \dots, K$?

A: 1-of-K representation belongs to exponential family, which will make computation feasible

For GMM model, $Q(\theta, \theta^{\text{old}})$ is

$$\begin{aligned}
Q(\theta, \theta^{\text{old}}) &= \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{XZ} | \theta) \\
&= \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}}) \ln \prod_{n=1}^N p(\mathbf{x}_n \mathbf{z}_n | \theta) \\
&\quad \text{Independence of PGM} \\
&= \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}}) \sum_{n=1}^N \ln p(\mathbf{x}_n \mathbf{z}_n | \theta) \\
&= \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{x}_1 \mathbf{z}_1 | \theta) + \dots + \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}\theta^{\text{old}}) \ln p(\mathbf{x}_N \mathbf{z}_N | \theta) \\
&= \sum_{\mathbf{z}_1} p(\mathbf{z}_1 | \mathbf{x}_1 \theta^{\text{old}}) \ln p(\mathbf{x}_1 \mathbf{z}_1 | \theta) + \dots + \sum_{\mathbf{z}_N} p(\mathbf{z}_N | \mathbf{x}_N \theta^{\text{old}}) \ln p(\mathbf{x}_N \mathbf{z}_N | \theta) \\
&= \sum_{n=1}^N \sum_{\mathbf{z}_n} p(\mathbf{z}_n | \mathbf{x}_n \theta^{\text{old}}) \ln p(\mathbf{x}_n \mathbf{z}_n | \theta) \tag{21}
\end{aligned}$$



substitute

$$p(\mathbf{x}_n z_{nk} | \theta) = p(z_{nk} | \theta) p(\mathbf{x}_n | z_{nk} \theta) = \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \tag{22}$$

into (21), and denote

$$\gamma_{nk} = p(z_{nk} | \mathbf{x}_n \theta^{\text{old}}) \tag{23}$$

We obtain

$$Q(\theta, \theta^{\text{old}}) = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \left[\ln \pi_k - \frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right] \quad (24)$$

Take constraint $\sum_{k=1}^K \pi_k = 1$ into consideration, introduce Lagrange function

$$R(\theta, \lambda) = Q(\theta, \theta^{\text{old}}) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) \quad (25)$$

Let the derivative of $R(\theta, \lambda)$ be zero,
we get the estimation of $\theta = \{\pi_k, \boldsymbol{\mu}_k, \Sigma_k\} \quad k = 1, \dots, K$

$$\begin{aligned} \frac{\partial R(\theta, \lambda)}{\partial \pi_k} &= 0 \quad (k = 1, \dots, K) & \pi_k &= \frac{1}{N} \sum_{n=1}^N \gamma_{nk} \quad (k = 1, \dots, K) \\ \frac{\partial R(\theta, \lambda)}{\partial \boldsymbol{\mu}_k} &= 0 \quad (k = 1, \dots, K) & \boldsymbol{\mu}_k &= \frac{\sum_{n=1}^N \gamma_{nk} \mathbf{x}_n}{\sum_{n=1}^N \gamma_{nk}} \quad (k = 1, \dots, K) \\ \frac{\partial R(\theta, \lambda)}{\partial \Sigma_k} &= 0 \quad (k = 1, \dots, K) & \Sigma_k &= \frac{\sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma_{nk}} \quad (k = 1, \dots, K) \end{aligned} \quad (26)$$

EM algorithm for GMM

1. Choose initial θ^{old}
2. E step: (Expectation)

$$\eta_{nk} = p(A_k \mid \mathbf{x}_n, \theta^{old})$$

where $\theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\} \quad k = 1, \dots, K$

3. M step: (Maximization)

$$\pi_k = \frac{1}{N} \sum_{n=1}^N \gamma_{nk} \quad \boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma_{nk} \mathbf{x}_n}{\sum_{n=1}^N \gamma_{nk}} \quad \boldsymbol{\Sigma}_k = \frac{\sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma_{nk}}$$

4. If converge then stop. Otherwise goto 2.

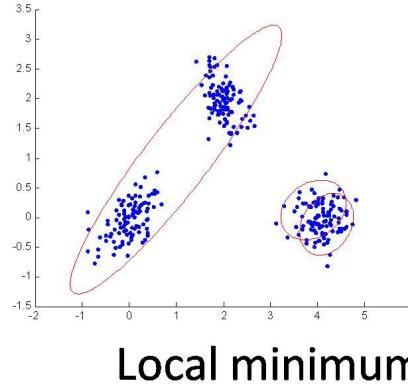
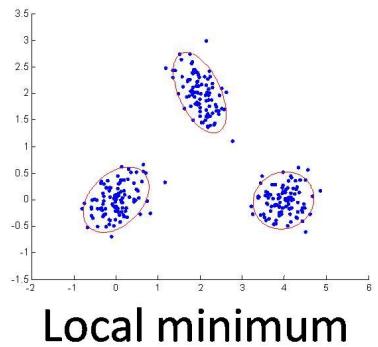
This is the same as previous result.

Tricks

1. GMM is a non convex problem. There are many local minimum.

Solve: Use kmeans to initialize parameter.

Initialize parameters and run EM algorithm for several times and choose the best model.

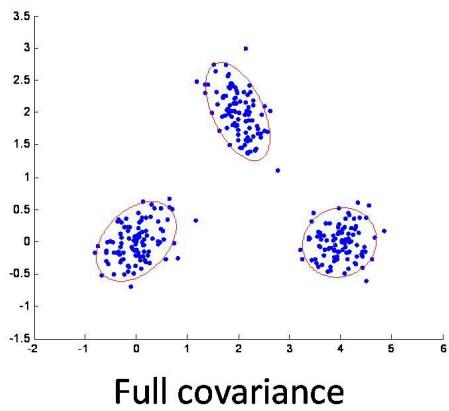


2. if $\text{eig}(\Sigma)$ is too small. Then Σ^{-1} will be unstable.

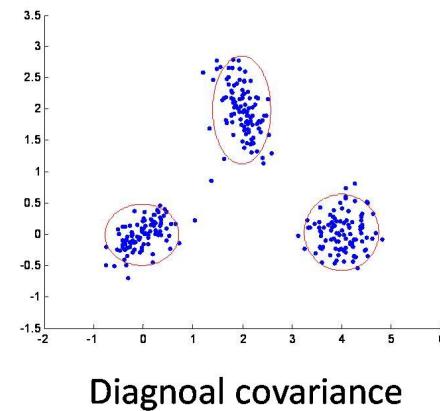
Solve: if $\text{eig}(\Sigma) < \varepsilon$ then $\Sigma = \Sigma + \sigma I$

3. The number of parameters of full GMM is $K+K*D+K*D*D$. If dataset is small then model will overfit.

Solve: Use diagonal Σ instead. The number of parameters will decrease to $K+K*D+K*D$.



Full covariance



Diagonal covariance

4. Gaussian probability density function can be underflow exponentially.

Solve: use log and normalization to solve this problem.

Matlab Code

<https://github.com/qiuqiangkong/matlab-gmm>

THANK YOU!