The maximization procedure involves alternating between two steps

$t$: step $\quad t = 0, 1, 2 \cdots$ $\quad$ Equivalent steps

**E-step**

I $\quad Q^{(t+1)} = \underset{Q}{\text{argmax}} \; G(Q, \underline{\theta}^{(t)})$ $\qquad$ Compute $Q^{(t+1)}$ with elements

closed form $\quad Q_{i,y}^{(t+1)} = P_{\underline{\theta}^{(t)}}\left[Y = y_i \mid X = \underline{x}_i\right]$

**M-step**

II $\quad \underline{\theta}^{(t+1)} = \underset{\underline{\theta}}{\text{argmax}} \; G(Q^{(t+1)}, \underline{\theta})$ $\qquad \underline{\theta}^{(t+1)} = \underset{\underline{\theta}}{\text{argmax}} \; F(Q^{t+1}, \underline{\theta})$

**EM algorithm** $\qquad \underline{\theta}_0, \underline{\theta}_1, \underline{\theta}_2 \cdots$

$Y \to X$

We know the property of $G$ function that $G(Q, \underline{\theta}) \le L(\underline{\theta})$

that means $\quad G(Q, \underline{\theta}^{(t)}) \le L(\underline{\theta}^{(t)})$

Now, if we substitute a $Q$ matrix with elements

$$Q_{i,y} = P_{\underline{\theta}}\left[Y = y_i \mid X = x_i\right] \qquad \text{in } G(Q, \underline{\theta})$$

$$G(Q, \underline{\theta}) = \sum_{i=1}^{m} \sum_{y_i=1}^{k} \left( \underbrace{Q_{i,y} \log P_{\underline{\theta}}\left[X = x_i, Y = y_i\right]}_{A} - \underbrace{Q_{i,y} \log Q_{i,y}}_{B} \right)$$

$$y_i \equiv y$$

$$= \sum_{i=1}^{m} \sum_{y_i=1}^{k} Q_{i,y} \log \frac{P_{\underline{\theta}}\{X = \underline{x}_i, Y = y_i\}}{Q_{i,y}} \qquad \frac{A}{B}$$

$$P(X, Y) = P(Y \mid X) P(X)$$

$$= \sum_{i=1}^{m} \sum_{y_i=1}^{k} \underbrace{P_{\underline{\theta}}\left[Y = y_i \mid X = \underline{x}_i\right]} \log \frac{\underbrace{P_{\underline{\theta}}\left[X = x_i, Y = y_i\right]}}{\underbrace{P_{\underline{\theta}}\left[Y = y_i \mid X = x_i\right]}}$$

$$= \sum_{i=1}^{m} \sum_{y_i=1}^{k} P_{\underline{\theta}}\left[Y = y_i \mid X = x_i\right] \log \underbrace{P_{\underline{\theta}}\left[X = \underline{x}_i\right]}$$

does not depend on $y$

$$= \sum_{i=1}^{m} \log P_{\underline{\theta}}\left[X = \underline{x}_i\right] \underbrace{\sum_{y_i=1}^{k} P_{\theta}\left[Y = y_i \mid X = x_i\right]}_{\text{label}} = 1$$

$$= \sum_{i=1}^{m} \log P_{\underline{\theta}}\left[X = \underline{x}_i\right] = L(\underline{\theta}).$$

Thus, substituting $Q_{i,y} = P_{\underline{\theta}}\left[Y = y_i \mid X = x_i\right]$ achieves the maximum value of the $G(Q, \underline{\theta})$ function when $\underline{\theta}$ is fixed.

$\therefore$ we write $Q^{(t+1)} = \underset{Q}{\text{argmax}}\ G\left(Q, \underline{\theta}^{(t)}\right)$

has elements $Q_{i,y}^{(t+1)} = P_{\underline{\theta}^{(t)}}\left[Y = y_i \mid X = x_i\right]$

The iterative procedure to maximize the $G$ function is called Expectation Maximization. The first step is the E step and the second step is the M step.

The EM procedure ensures that the observed data log likelihood never decreases.

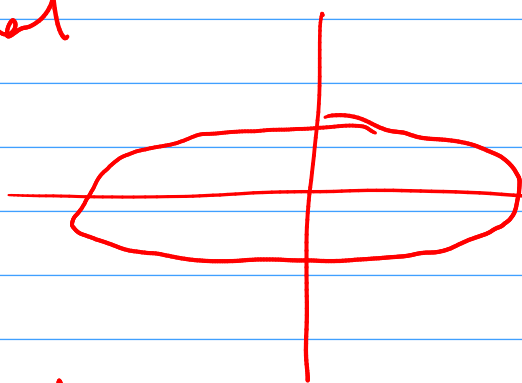$$L\left(\underline{\theta}^{(t+1)}\right) \geqslant L\left(\underline{\theta}^{(t)}\right).$$

$$L\left(\underline{\theta}^{t+1}\right) = G\left(Q^{(t+2)}, \underline{\theta}^{(t+1)}\right) \geqslant G\left(Q^{(t+1)}, \underline{\theta}^{(t)}\right) = L\left(\underline{\theta}^{(t)}\right)$$

$$\because\ G\left(Q^{(t+2)}, \underline{\theta}^{(t+1)}\right) \geqslant G\left(Q^{(t+1)}, \underline{\theta}^{(t+1)}\right)$$

$$\geqslant G\left(Q^{(t+1)}, \underline{\theta}^{(t)}\right)$$
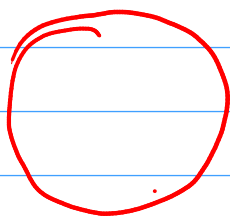
**Full Covariance**

**Diagonal**

variance values are same

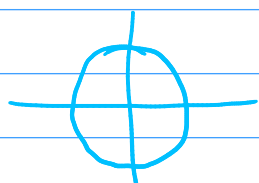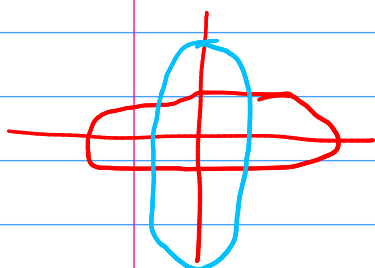$$C = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$\mu_1 \Sigma_1 \qquad \mu_2 \Sigma_2 \quad \cdots \quad \mu_k \Sigma_k$

$\Sigma :$ full / Diag / Identity

$$d_1 I \begin{bmatrix} d_1 & d_1 & 0 \\ 0 & d_1 & \ddots \\ & & d_1 \end{bmatrix}$$

hyper spheres

$$\begin{bmatrix} d_1 & & & \\ & d_2 & 0 & \\ 0 & & d_3 & \\ & & & \ddots \\ & & & d_k \end{bmatrix}$$

Axis aligned Gaussians
hyperellipses
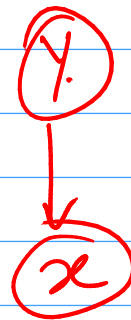
$$C = \sum_{i=1}^{m} \underline{x_i} \underline{x_i}^T$$

under: sum of outer product of vectors.

$\underline{x_i}$ $d \times 1$

$x_i^T x_i$
$1 \times d \times d \times 1$
$= 1 \times 1$
scalar
(dot prod)

Outer product $\underline{x_i} x_i^T$

$d \times 1 \times 1 \times d$
$= d \times d$ matrix

## GMM

Structure

Ⓨ
↓
Ⓧ

$P(X=x \mid Y=y)$ Gaussian
$\phi, \mu, \Sigma.$

Organize the data into
Gaussian (hyper ellipse) cluster

$\hat{y_i} = \underset{y}{argmax} \ P\left(Y=y \mid \underline{X} = \underline{x_i}\right)$

$\overline{x_i}$  $\overline{x_i}$ Mean



$P(X=x_i \mid Y=y)$
likelihood of
generating $x_i$
from Gaussian
$y$

## Mahalanobis distance

$$(\underline{x} - \underline{\mu})^T (\underline{x} - \underline{\mu})$$

Euclidean distance.
$$\| \underline{x} - \underline{\mu} \|^2$$

$$\rightarrow (\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu})$$

Minor

Major

5  10  15

$$\underbrace{(x_1 - \mu_1)^2} + \underbrace{(x_2 - \mu_2)^2}$$

$x_2$

5  $x_1$

1

2

3

4

$\underline{x}$

$$\underline{P(X = \underline{x})} = \sum_K P(X = \underline{x} \mid Y = y_K) P(Y = y)$$