

## EM Algorithm

- an iterative method for numerically calculating MLE
- increases the likelihood at each step

 $x[n] \rightarrow \text{real observed data}$ 

$y[n] \rightarrow$  latent / imaginary data  $\rightarrow$  used  
if  $y[n]$  makes MLE calculation easier

Example: Let  $x[n] = \sum_{i=1}^I \cos 2\pi f_i n + w[n]$   $n=0,1,\dots,N-1$   
 $\hookrightarrow N(0, \sigma^2)$

$$Q = [f_1 \dots f_p]^T \rightarrow \text{to be estimated.}$$

Normally MIE would req. minimization of  $\sum_{i=1}^N C_i = P_0 \sum_{i=1}^N \sqrt{2}$

$$J(f) = \sum_{n=0}^{N-1} \left( x[n] - \sum_{i=1}^P \cos(2\pi f_i n) \right)^2$$

$$\begin{bmatrix} \ell_1 \\ \vdots \\ \ell_p \end{bmatrix}$$

- min over  $p$  dimensions.
- $p$  coupled equations.

→  $p$  coupled equations.

On the other hand, if we had different data, say  $y_i[n] = \cos 2\pi f_i n + w_i[n]$

$$i = 1, 2, \dots, P$$

$$n = 0, 1, \dots, N-1$$
$$i = 1, 2, \dots, p$$
$$n = 0, 1, \dots, N-1$$

where  $w_i[n] \sim N(0, \sigma_i^2)$ ,  $w_i$  &  $w_j$  are ind.  
&  $\sum_{i=1}^P \sigma_i^2 = \sigma^2$

$$2 \sum_{i=1}^p \sigma_i^2 = \sigma^2$$

Then, we can calculate  $\hat{f}_i^o$  by minimizing

$$J(f_i^o) = \sum_{n=0}^{N-1} (y_i[n] - \cos 2\pi f_i^o n)^2$$

$$J(f_1^0) = \sum_{n=0}^{N-1} (y_1[n] - \cos 2\pi f_1 n)^2$$

separately for each  $i = 1, \dots, p$ .

$\Rightarrow$   $p$  - decoupled equations; one for each  $p$ .  
 $\Rightarrow$  Easier problem

⇒ Easier problem

$\{y_1[n], y_2[n], \dots, y_p[n]\} \rightarrow \text{Complete / latent Data}$   
 $(n=0, 1, \dots, N-1)$

$\{x[0], \dots, x[N-1]\} \rightarrow \text{Incomplete Data}$

Clearly  $x[n] = \sum_{i=1}^P y_i[n]$  if  $w[n] = \sum_{i=1}^P w_i[n]$

Note: This decomposition is not unique: e.g.

$$\left. \begin{aligned} y_1[n] &= \sum_{i=1}^P \cos 2\pi f_i n & y_2[n] &= w[n] \\ x[n] &= y_1[n] + y_2[n] \end{aligned} \right\} \text{ is also a possible decomp.}$$

In general:  $x = g(y_1, \dots, y_p) = g(y)$

Main problem:  $y_i[n]$  are not known in reality, whereas  $x[n]$  is observed.

Recall The mse estimate of  $\ln p(y; \theta)$  given  $x$  is:

$$E_{y/x} [\ln p_y(y; \theta)] = \int \ln p(y; \theta) p(y/x; \theta) dy$$

$\swarrow$  unknown       $\uparrow$  known       $\searrow$  unknown

Strategy: 1) Replace  $\theta$  with the current guess of  $\theta$  & then  $\max_{\theta} E_{y/x} [\ln p_y(y; \theta)]$ .  
 2) Iterate

E.M. Algo: (Expectation - Maximization)

E-step:  $U(\theta, \theta_k) = \int \ln p_y(y; \theta) p(y/x; \theta_k) dy$   
M-step:  $\theta_{k+1} = \arg \max_{\theta} U(\theta, \theta_k)$

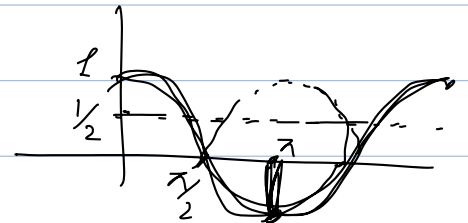
# Calculation of the conditional prob  $p(y/x; \theta)$  is often very difficult.  
For our example above:

$$\begin{aligned}
 \ln p(y; \theta) &= \sum_{i=1}^P \ln p(y_i; \theta_i) \\
 &= \sum_{i=1}^P \ln \left\{ \frac{1}{(2\pi\sigma_i^2)^{N/2}} \exp \left[ -\frac{1}{2\sigma_i^2} \sum_{n=0}^{N-1} (y_i[n] - \cos 2\pi f_i n)^2 \right] \right\} \\
 &\stackrel{\text{constant}}{=} C - \sum_{i=1}^P \frac{1}{2\sigma_i^2} \sum_{n=0}^{N-1} (y_i[n] - \cos 2\pi f_i n)^2 \\
 &= \underbrace{g(y)}_{\text{ind. of } f_i} + \sum_{i=1}^P \frac{1}{\sigma_i^2} \sum_{n=0}^{N-1} \left( y_i[n] \cos 2\pi f_i n - \frac{1}{2} \cos^2 2\pi f_i n \right)
 \end{aligned}$$

Approximate  $\sum_{n=0}^{N-1} \cos^2 2\pi f_i n \approx \frac{N}{2}$  for  $f_i^0$  not near 0 or  $\frac{1}{2}$

$$\begin{cases} \cos 0 = 1 \Rightarrow \sum_{n=0}^{N-1} 1 = N \\ \cos 2\pi \frac{1}{2} n = \cos \pi n = -1 \end{cases}$$

Otherwise  $\approx \sum_{n=0}^{N-1} \frac{1}{2} \approx \frac{N}{2}$



$$\text{Then } \ln p(y; \theta) = h(y) + \sum_{i=1}^P \frac{1}{\sigma_i^2} \sum_{n=0}^{N-1} y_i[n] \cos 2\pi f_i n$$

$$\text{Let } c_i^0 = [1 \quad \cos 2\pi f_i^0 \quad \dots \quad \cos 2\pi f_i^0 (N-1)]^T$$

$$\text{Then } \ln p(y; \theta) = h(y) + \sum_{i=1}^P \frac{1}{\sigma_i^2} c_i^{0T} y_i^0$$

$$y_i^0 = \begin{bmatrix} y_i[0] \\ \vdots \\ y_i[N-1] \end{bmatrix}$$

Further let  $C = \begin{bmatrix} \frac{1}{\sigma_1^2} c_1 \\ \vdots \\ \frac{1}{\sigma_p^2} c_p \end{bmatrix}$   $y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_p \end{bmatrix}$

Then  $\ln p(y; \theta) = h(y) + C^T y$

# The E-step :  $U(\theta; \theta_k) = E[\ln p(y; \theta) / x; \theta_k]$   
 $= E[h(y) / x; \theta_k] + C^T E(y / x; \theta_k)$

ind. of  $\theta$   
 hence ignored  
 in next M-step

Recall  $y$  &  $x$  are jointly Gaussian

since  $x = \sum y_i = \begin{bmatrix} I & I & \dots & I \end{bmatrix} y$   
 $N \times 1$   $N \times N$   $N \times N_p$   $N_p \times 1$

Since  $x$  &  $y$  are jointly Gaussian:  
 $E(y/x; \theta_k) = E(y) + C_{yx} C_{xx}^{-1} (x - E(x))$

Recall :  $E(y) = \begin{bmatrix} c_1 \\ \vdots \\ c_p \end{bmatrix}_{N_p \times 1}$

$E(x) = \sum_{i=1}^P c_i$

$C_{xx} = \sigma^2 I$

$C_{yx} = E \left( \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix} w^T \right)$

$= E \left( \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix} \left\{ [I \ I \ \dots \ I] \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix} \right\}^T \right)$

$$= \begin{bmatrix} \sigma_1^2 I & & \\ & \ddots & \\ & & \sigma_p^2 I \end{bmatrix} \begin{bmatrix} I \\ I \\ \vdots \\ I \end{bmatrix} = \begin{bmatrix} \sigma_1^2 I \\ \vdots \\ \sigma_p^2 I \end{bmatrix}_{N_p \times N}$$

$$E(y/x; \theta_k) = \begin{bmatrix} c_1 \\ \vdots \\ c_p \end{bmatrix} + \frac{1}{\sigma^2} \begin{bmatrix} \sigma_1^2 I \\ \vdots \\ \sigma_p^2 I \end{bmatrix} \left( x - \sum_{i=1}^p c_i \right)$$

$$= \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_p \end{bmatrix} + \begin{bmatrix} \frac{\sigma_1^2}{\sigma^2} \left( x - \sum_{i=1}^p c_i \right) \\ \vdots \\ \frac{\sigma_p^2}{\sigma^2} \left( x - \sum_{i=1}^p c_i \right) \end{bmatrix}$$

$$\Leftrightarrow E(y_i/x; \theta_k) = c_i^o + \frac{\sigma_i^2}{\sigma^2} \left( x - \sum_{i=1}^p c_i \right) \quad i=1, \dots, p$$

Recall this is the best estimate of  $y_i^o$

$$=: \hat{y}_i^o \quad \left| \quad \hat{y}_i^o[n] = \cos 2\pi f_i^o n + \frac{\sigma_i^2}{\sigma^2} \left( x[n] - \sum_{i=1}^p \cos 2\pi f_i^o n \right)$$

# Relevant part of the E-step:

$$u'(\theta, \theta_k) = \sum_{i=1}^p c_i^T \hat{y}_i^o$$

Note that this sum can be maximized by maximizing each term separately  
i.e.  $\underbrace{f_{i,k+1}^o}_{\text{iteration } n} = \arg \max_{f_i^o} c_i^T \hat{y}_i^o$

In effect: the E-M steps are:

E-step: For  $i=1, 2, \dots, P$  Est. of  $w_i[n]$   

$$\hat{y}_i[n] = \cos 2\pi f_k^0 n + \beta_i \left[ x[n] - \sum_{i=1}^P \cos 2\pi f_k^0 n \right]$$

M-step: For  $i=1, 2, \dots, P$   

$$f_{k+1}^0 = \arg \max_{f_k^0} \sum_{n=0}^{N-1} \hat{y}_i[n] \cos 2\pi f_k^0 n$$

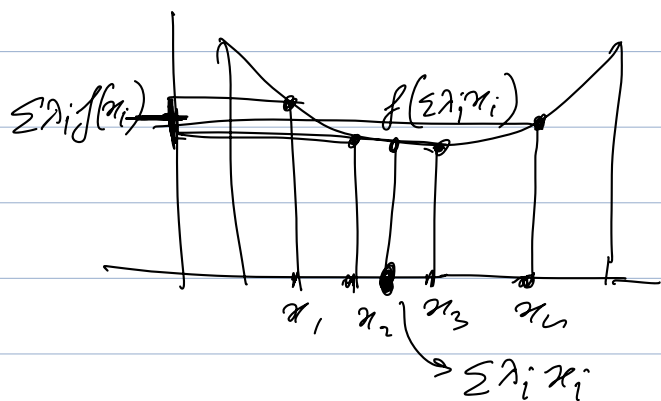
$\beta_i$ 's are arbitrary as long as  $\sum_{i=1}^P \beta_i = 1$

$\sigma_i^2$  are not unique, they can be chosen arbitrarily as long as  $\sum_{i=1}^P \sigma_i^2 = \sigma^2$  or  $\sum_{i=1}^P \beta_i = \sum_{i=1}^P \frac{\sigma_i^2}{\sigma^2} = 1$

Q. Why does it work?

Jensen's Inequality: Let  $f$  be a convex function defined on an interval  $I$ . If  $x_1, \dots, x_n \in I$  and  $\lambda_1, \lambda_2, \dots, \lambda_n \geq 0$  with  $\sum_{i=1}^n \lambda_i = 1$ , then

$$f\left(\sum_{i=1}^n \lambda_i x_i\right) \leq \sum_{i=1}^n \lambda_i f(x_i)$$



Proof:  $n=1$  trivial  
 $n=2$  def. of convexity

Induction: assume true for  $n$ .

$$f\left(\sum_{i=1}^{n+1} \lambda_i x_i\right) = f\left(\lambda_{n+1} x_{n+1} + \sum_{i=1}^n \lambda_i x_i\right)$$

$$\begin{aligned}
&= f\left(\lambda_{n+1} x_{n+1} + (1-\lambda_{n+1}) \left(\frac{1}{1-\lambda_{n+1}}\right) \leq \lambda_i x_i\right) \\
&\leq \lambda_{n+1} f(x_{n+1}) + (1-\lambda_{n+1}) f\left(\frac{1}{1-\lambda_{n+1}} \leq \lambda_i x_i\right) \\
&= \lambda_{n+1} f(x_{n+1}) + (1-\lambda_{n+1}) f\left(\sum_{i=1}^n \frac{\lambda_i}{1-\lambda_{n+1}} x_i\right) \\
&\leq \lambda_{n+1} f(x_{n+1}) + (1-\lambda_{n+1}) \sum_{i=1}^n \frac{\lambda_i}{1-\lambda_{n+1}} f(x_i) \\
&= \lambda_{n+1} f(x_{n+1}) + \sum_{i=1}^n \lambda_i f(x_i) \quad \left[ \because \lambda_1 + \dots + \lambda_{n+1} = 1 \right] \\
&= \sum_{i=1}^{n+1} \lambda_i f(x_i)
\end{aligned}$$


---

# Since  $\ln(x)$  is concave,

$$\ln\left(\sum_{i=1}^n \lambda_i x_i\right) \geq \sum_{i=1}^n \lambda_i \ln(x_i)$$


---

# Technical Assumption on  $y$  (the latent variable)

$\theta \rightarrow y \rightarrow x$  is a Markov Chain

$\Leftrightarrow P(x/y)$  is independent of  $\theta$ .

---

Define  $L(\theta) = \ln p(x; \theta)$  i.e.  $L(\theta_k) = \ln p(x; \theta_k)$

Then,  $L(\theta) - L(\theta_k) = \ln \left( \frac{p(x; \theta)}{p(x; \theta_k)} \right)$

$$= \ln \int \frac{p(y, x; \theta)}{p(x; \theta_k)} dy = \ln \int \frac{p(y, x; \theta)}{\frac{p(y, x; \theta_k)}{p(y/x; \theta_k)}} dy$$

$$= \ln \int \frac{P(y, x; Q)}{P(y, x; Q_k)} P(y/x; Q_k) dy$$

$$= \ln \int \frac{P(y; Q) \cancel{P(x/y)}}{P(y; Q_k) \cancel{P(x/y)}} P(y/x; Q_k) dy$$

due to Markov assumption.

$$\geq \int \left\{ \ln \left( \frac{P(y; Q)}{P(y; Q_k)} \right) \right\} P(y/x; Q_k) dy \quad \left[ \begin{array}{l} \text{By Jensen's} \\ \text{inequality} \end{array} \right]$$

(2)

$$= \int \ln P(y; Q) P(y/x; Q_k) dy - \underbrace{\int \ln P(y; Q_k) P(y/x; Q_k) dy}_{\text{Ind of } Q \rightarrow \text{Ignore}}$$

Recall that this is the quantity we maximize in the E-M step.

Define  $\Delta(Q/Q_k) := \int \left\{ \ln \left( \frac{P(y; Q)}{P(y; Q_k)} \right) \right\} P(y/x; Q_k) dy$

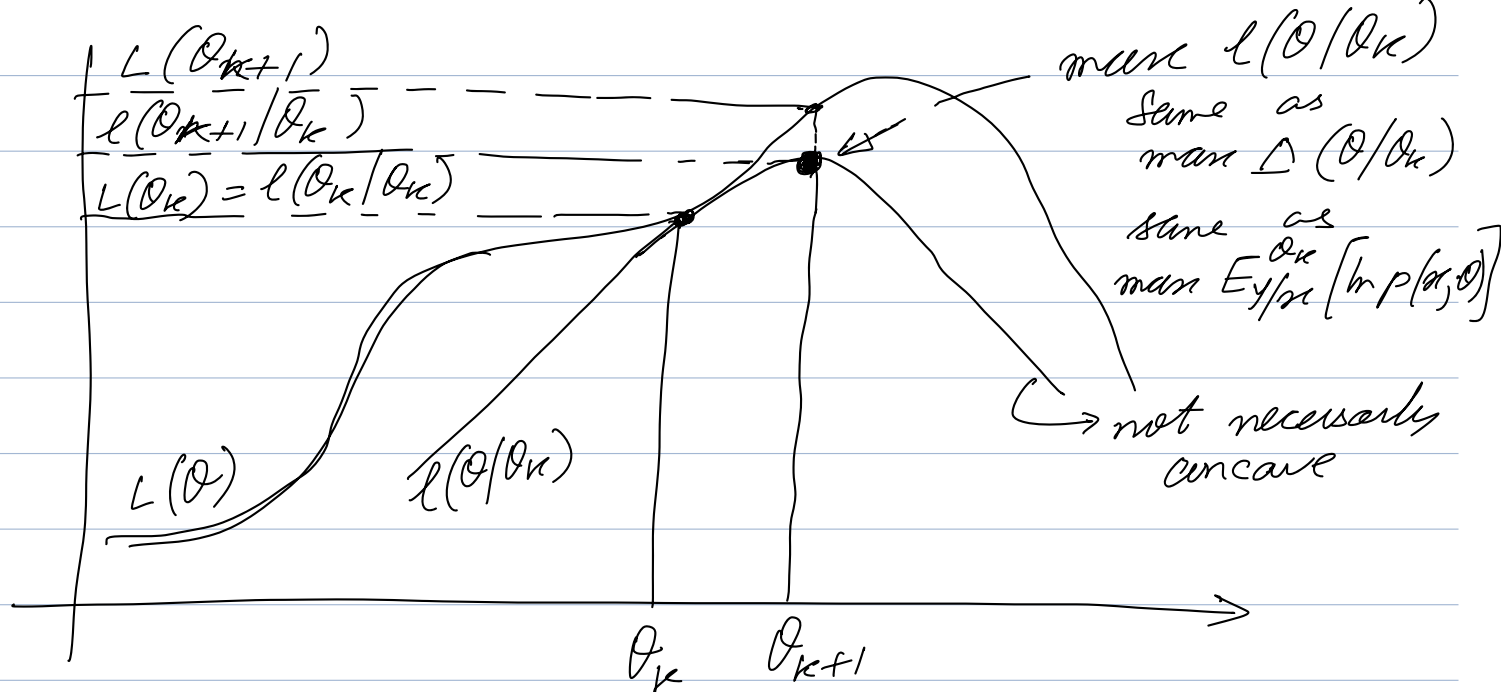
From (2),  $L(Q) - L(Q_k) \geq \Delta(Q/Q_k)$

$$\Leftrightarrow L(Q) \geq \underbrace{L(Q_k) + \Delta(Q/Q_k)}_{=: l(Q/Q_k)}$$

Clearly,  $\Delta(Q_k/Q_k) = \int \ln 1 P(y/x; Q_k) dy = 0$

$$\Rightarrow l(Q_k/Q_k) = L(Q_k)$$





# From fig. it is 'clear' that value of  $L(\theta_{k+1}) \geq L(\theta_k) \quad \forall k$ .

# This also follows from the fact that we maximize  $l(\theta/\theta_k)$  at each iteration

$$\boxed{L(\theta_{k+1}) \geq L(\theta_k) \quad \forall k}$$

Monotonicity Property of EM.