## 1.1 Convergence of the EM algorithm

We will show that $\ell\left(\hat{\boldsymbol{\theta}}^{(k+1)}\right) \geq \ell\left(\hat{\boldsymbol{\theta}}^{(k)}\right)$.

In other words, each step of the EM algorithm leads to an improvement of the log-likelihood value.

Thus, if the likelihood is <u>well behaved</u>, it will achieve the MLE, otherwise the EM will achieve a local maxima (if there is one).
$\quad\hookrightarrow$ bounded, unimodal.

$y$ = observed data.
$z$ = hidden data.

We know $f_{Z|Y}(z|y;\boldsymbol{\theta}) = \frac{f_{YZ}(y,z;\boldsymbol{\theta})}{f_Y(y|\boldsymbol{\theta})}$.  def'n of conditional density
true for any $y, z$

$\Rightarrow f_Y(y;\theta) = \dfrac{f_{YZ}(y,z;\theta)}{f_{Z|Y}(z|y;\theta)}$  just rewritten (not clear why).

Assume we observe $\boldsymbol{y} = (y_1, \ldots, y_n)$, then

$L(\theta|\underline{Y}) = f_{\underline{Y}}(\underline{Y};\theta) = \dfrac{f_{\underline{YZ}}(\underline{Y},\underline{Z};\theta)}{f_{\underline{Z}|\underline{Y}}(\underline{Z}|\underline{Y};\theta)}$ (if iid, product of univariate densities)

holds for any $\underline{Z}$!

$\Rightarrow \ell(\theta) = \log f_{\underline{Y}}(\underline{Y};\theta) = \log f_{\underline{YZ}}(\underline{Y},\underline{Z};\theta) - \log f_{\underline{Z}|\underline{Y}}(\underline{Z}|\underline{Y};\theta) = \ell_c(\theta|\underline{Y},\underline{Z}) - \ell(\theta|\{\underline{Z}|\underline{Y}\})$

log likelihood of data $\underline{Y}$ want to optimize.

log likelihood of "completedata" $\underline{Y}, \underline{Z}$    "conditional likelihood"

So, in order to show that $\ell\left(\hat{\boldsymbol{\theta}}^{(k+1)}\right) \geq \ell\left(\hat{\boldsymbol{\theta}}^{(k)}\right)$, this is the same as  $\Rightarrow$ take expected value wrt $\underline{Z}|\underline{Y};\hat{\theta}^{(k)}$

$Q(\hat{\theta}^{(k+1)}, \hat{\theta}^{(k)}) - H(\hat{\theta}^{(k+1)}, \hat{\theta}^{(k)}) \geq Q(\hat{\theta}^{(k)}, \hat{\theta}^{(k)}) - H(\hat{\theta}^{(k)}, \hat{\theta}^{(k)})$.

$\int \ell(\theta|\underline{Y}) f_{\underline{Z}|\underline{Y}}(\underline{Z}|\underline{Y};\hat{\theta}^{(k)}) d\underline{Z} = \int \log f_{\underline{YZ}}(y,z;\theta) f_{\underline{Z}|\underline{Y}}(\underline{Z}|\underline{Y};\hat{\theta}^{(k)}) d\underline{Z}$
$\qquad - \int \log f_{\underline{Z}|\underline{Y}}(\underline{Z}|\underline{Y};\theta) f_{\underline{Z}|\underline{Y}}(\underline{Z}|\underline{Y};\hat{\theta}^{(k)}) d\underline{Z}$

$\ell(\theta|\underline{Y}) \underbrace{\int f_{\underline{Z}|\underline{Y}}(\underline{Z}|\underline{Y};\hat{\theta}^{(k)}) d\underline{Z}}_{=1} = \ldots$

$\Rightarrow \ell(\theta|\underline{Y}) = Q(\theta, \hat{\theta}^{(k)}) - H(\theta, \hat{\theta}^{(k)})$

(function of $\theta$!)

**Step 1:** Show that $H(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(k)})$ is maximized when $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}^{(k)}$.

i.e. $H\left(\hat{\underline{\theta}}^{(k)}, \hat{\underline{\theta}}^{(k)}\right) \geq H(\theta, \hat{\theta}^{(k)})$ for any $\theta \in |\Theta|$.

Recall: Jensen's Inequality. A function $\Phi$ is convex if $\Phi(\frac{x_1+x_2}{2}) \leq \frac{1}{2}\Phi(x_1) + \frac{1}{2}\Phi(x_2)$. Then
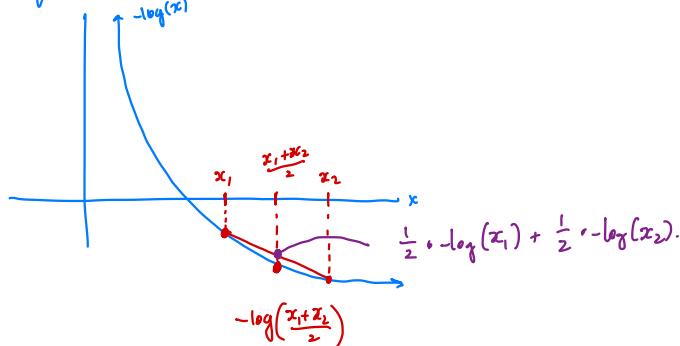
$$\Phi(\mathrm{E}[g(X)]) \leq \mathrm{E}[\Phi(g(X))],$$

where $g$ is a real-valued integrable function. $\quad \Longleftrightarrow$

Fact: $-\log$ is convex

$\Phi\left(\int g(x) f(x) dx\right) \leq \int \Phi(g(x)) f(x) dx$ where $f(x)$ is density of $X$.



$-\log(x)$

$x_1 \quad \frac{x_1+x_2}{2} \quad x_2$

$x$

$\frac{1}{2} \cdot -\log(x_1) + \frac{1}{2} \cdot -\log(x_2).$

$-\log\left(\frac{x_1+x_2}{2}\right)$

(non negative).

Consider $H\left(\hat{\underline{\theta}}^{(k)}, \hat{\underline{\theta}}^{(k)}\right) - H(\theta, \hat{\theta}^{(k)})$. WTS this is positive $\forall \theta$

$$H\left(\theta, \hat{\theta}^{(k)}\right) = \int \log\left(f_{Z|Y}(z|y;\theta)\right) f_{Z|Y}(z|y;\hat{\theta}^{(k)}) dz$$

$$\Longrightarrow H\left(\hat{\theta}^{(k)}, \hat{\theta}^{(k)}\right) - H\left(\theta, \hat{\theta}^{(k)}\right) = \int \left(\log\left(f_{Z|Y}(z|y;\hat{\theta}^{(k)})\right) - \log\left(f_{Z|Y}(z|y;\theta)\right)\right) f_{Z|Y}(z|y;\hat{\theta}^{(k)}) dz$$

$$= \int -\log\left(\frac{f_{Z|Y}(z|y;\theta)}{f_{Z|Y}(z|y;\hat{\theta}^{(k)})}\right) f_{Z|Y}(z|y;\hat{\theta}^{(k)}) dz$$

$$\geq -\log \int \frac{f_{Z|Y}(z|y;\theta)}{f_{Z|Y}(z|y;\hat{\theta}^{(k)})} f_{Z|Y}(z|y;\hat{\theta}^{(k)}) dz$$

$$= -\log \underbrace{\int f_{Z|Y}(z|y;\theta) dz}_{1}$$

$$= 0$$

$$\Longrightarrow H\left(\hat{\theta}^{(k)}, \hat{\theta}^{(k)}\right) \geq H(\theta, \hat{\theta}^{(k)}) \; \forall \theta. \; /\!/$$

**Step 2:** Find a $\hat{\boldsymbol{\theta}}^{k+1}$ that will optimize $Q$.

Recall goal is to find $\hat{\theta}^{(k+1)}$ s.t. $\ell\left(\hat{\underline{\theta}}^{(k+1)}\right) \geq \ell\left(\hat{\underline{\theta}}^{(k)}\right) + \ell\left(\underline{\theta}\right) = Q\left(\underline{\theta}, \hat{\underline{\theta}}^{(k)}\right) - H\left(\underline{\theta}, \hat{\underline{\theta}}^{(k)}\right)$

Let $\hat{\theta}^{(k+1)} = \underset{\underline{\theta}}{\text{argmax}} \ Q\left(\underline{\theta}, \hat{\theta}^{(k)}\right).$

$\nearrow$
This is the
EM algorithm.

We know $H\left(\hat{\underline{\theta}}^{(k+1)}, \hat{\underline{\theta}}^{(k)}\right) \leq H\left(\hat{\underline{\theta}}^{(k)}, \hat{\underline{\theta}}^{(k)}\right)$ because true for all $\underline{\theta}$

$+ \quad Q\left(\hat{\underline{\theta}}^{(k+1)}, \hat{\underline{\theta}}^{(k)}\right) \geq Q\left(\hat{\underline{\theta}}^{(k)}, \hat{\underline{\theta}}^{(k)}\right)$ by optimization.

So, $\ell\left(\hat{\underline{\theta}}^{(k)}\right) = Q\left(\hat{\underline{\theta}}^{(k)}, \hat{\underline{\theta}}^{(k)}\right) - H\left(\hat{\underline{\theta}}^{(k)}, \hat{\underline{\theta}}^{(k)}\right)$

$\leq Q\left(\hat{\underline{\theta}}^{(k+1)}, \hat{\underline{\theta}}^{(k)}\right) - H\left(\hat{\underline{\theta}}^{(k+1)}, \hat{\underline{\theta}}^{(k)}\right) = \ell\left(\hat{\theta}^{(k+1)}\right) \checkmark$

**Example (Two-Component Mixture, Cont'd):**

$$Q(\theta, \hat{\theta}^{(k)}) = \int \log f_{YZ}(y, z; \theta) f_{Z|Y}(z|y; \hat{\theta}^{(k)}) dz$$

For the Gaussian mixture, the complete log-likelihood:

$$\log f_{YZ}(Y, Z; \theta) = \sum_{i=1}^{n} \left\{ z_i \log f_1(Y_i; \mu_1, \Sigma_1) + (1-z_i) \log f_2(Y_i; \mu_2, \Sigma_2) + z_i \log p + (1-z_i) \log(1-p) \right\}.$$

To get the conditional density, $\quad f_{Z|Y}(Z|y; \hat{\theta}^{(k)}) = \prod_{i=1}^{n} f_{Z|Y}(z_i|y_i; \hat{\theta}^{(k)})$

$$f_{Z|Y}(z_i|y_i; \hat{\theta}^{(k)}) = \frac{f_{YZ}(y_i, z_i; \hat{\theta}^{(k)})}{f_Y(y_i; \hat{\theta}^{(k)})} \qquad \text{complete density contribution}$$
$$\text{observed density.}$$

$$= \frac{\left[ \hat{p}^{(k)} f_1(y_i; \hat{\mu}_1^{(k)}, \hat{\Sigma}_1^{(k)}) \right]^{z_i} \left[ (1-\hat{p}^{(k)}) f_2(y_i; \hat{\mu}_2^{(k)}, \hat{\Sigma}_2^{(k)}) \right]^{1-z_i}}{\hat{p}^{(k)} f_1(y_i; \hat{\mu}_1^{(k)}, \hat{\Sigma}_1^{(k)}) + (1-\hat{p}^{(k)}) f_2(y_i; \hat{\mu}_2^{(k)}, \hat{\Sigma}_2^{(k)})}$$

$z_i$ can only take values 0 or 1.
$\Rightarrow$ Bernoulli!

$$P(z_i = 1 | Y_i = y_i, \hat{\theta}^{(k)}) = \frac{\hat{p}^{(k)} f_1(y_i; \hat{\mu}_1^{(k)}, \hat{\Sigma}_1^{(k)})}{\hat{p}^{(k)} f_1(y_i; \hat{\mu}_1^{(k)}, \hat{\Sigma}_1^{(k)}) + (1-\hat{p}^{(k)}) f_2(y_i; \hat{\mu}_2^{(k)}, \hat{\Sigma}_2^{(k)})} = \hat{w}_i^{(k)} \quad \overset{\text{define.}}{} \quad \text{Then}$$

$$z_i | Y_i = y, \hat{\theta}^{(k)} \sim \text{Bern}(\hat{w}_i^{(k)}).$$

$$\Rightarrow Q(\theta, \hat{\theta}^{(k)}) = \sum_{i=1}^{n} E_{z_i|y_i}\Big|_{\theta = \hat{\theta}^{(k)}} \left[ \log f_{YZ}(y_i, z_i; \theta) \right] \quad \text{and} \quad E_{z_i|y_i}\Big|_{\theta=\hat{\theta}^{(k)}} [g(z_i)] = g(1)\hat{w}_i^{(k)} + g(0)(1-\hat{w}_i^{(k)}) \Rightarrow$$

$$Q(\theta, \hat{\theta}^{(k)}) = \sum_{i=1}^{n} \left\{ \hat{w}_i^{(k)} \left[ \overset{1}{z_i}\log f_1(y_i; \mu_1, \Sigma_1) + \overset{0}{(1-z_i)} \log f_2(y_i; \mu_2, \Sigma_2) + \overset{1}{z_i} \log p + \overset{0}{(1-z_i)} \log(1-p) \right]\Big|_{z_i=1} + \right.$$
$$\left. (1-\hat{w}_i^{(k)}) \left[ \underset{0}{z_i} \log f_1(y_i; \mu_1, \Sigma_1) + (1-z_i) \log f_2(y_i; \mu_2, \Sigma_2) + \underset{0}{z_i} \log p + \underset{1}{(1-z_i)} \log(1-p) \right]_{z_i=0} \right.$$

which yields the intuitive expression from before!

So "plugging in the weights" makes sense from an optimization standpoint <u>in this example.</u>

In general can't always separate E + M in this way for Q.

The EM algorithm allows us to obtain $\hat{\boldsymbol{\theta}}_{\text{EM}}$, the parameter estimate which optimizes the algorithm.

Which, if the likelihood is "nice" can $= \hat{\theta}_{\text{MLE}}.$