

09 Finding MLE

(1) EM algorithm.

when do we use EM?

- when we have missing values or assume hidden (latent) parameter for problem simplification
- when maximizing incomplete-data log likelihood directly is hard.

Algorithm:

$L(\theta|y, z)$ → observation: known and fixed
 $L(\theta|y, z) = f(y, z|\theta)$ → complete-data likelihood.
 DGP: unknown fixed

So $L(\theta|y, z)$ is a function of z : $h_\theta(z)$.

(1) Start with $\theta^{(0)}$

(2) Given the current $\theta^{(r)}$, $r=1, 2, \dots$

E-step:

$$Q(\theta, \theta^{(r)}) = E[\ln f(y, z|\theta) | y, \theta^{(r)}] \rightarrow \text{expected value of complete-data likelihood}$$

$$= \int_y \ln f(y, z|\theta) \underbrace{h(z|y, \theta^{(r)}) dz}_{\text{conditional pdf of latent data depends on current value of } \theta^{(r)} \text{ and observed data.}} \quad \text{recall: } E[h(Y|X=x)] = \int_y h(y)f(y|x) dy.$$

M-step:

$$\theta^{(r+1)} = \arg \max_{\theta} Q(\theta, \theta^{(r)}) \quad [\text{if } \theta^{(r+1)} \text{ s.t. } Q(\theta^{(r+1)}, \theta^{(r)}) > Q(\theta^{(r)}, \theta^{(r)}) \Rightarrow \text{it's called GEM}].$$

(3) The iteration continues until $\|\theta^{(r+1)} - \theta^{(r)}\|$ or $|Q(\theta^{(r+1)}, \theta^{(r)}) - Q(\theta^{(r)}, \theta^{(r)})|$ is smaller than tolerance level $\varepsilon > 0$, e.g. 10^{-6} .

Why does it work?

- Each iteration increases log-likelihood i.e. $\ln L(\theta^{(r+1)}) \geq \ln L(\theta^{(r)})$
- and the EM algorithm increases the incomplete-loglikelihood.
 i.e. if θ' is found s.t. $Q(\theta', \theta') = \max_{\theta} Q(\theta, \theta')$ $\Rightarrow \ln L(\theta'|y) \geq \ln L(\theta|y)$

Example:

$Y \sim \text{multinomial}(P_1, P_2, P_3, P_4)$ with $(P_1, P_2, P_3, P_4) = (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$ with $0 \leq \theta \leq 1$.
 and $y = (y_1, y_2, y_3, y_4) = (125, 18, 20, 34)$.

$X \sim \text{multinomial}(P_1, P_2, P_3, P_4, P_5)$ with $(P_1, P_2, P_3, P_4, P_5) = (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$

and $x = (x_1, x_2, x_3, x_4, x_5) \Rightarrow x_1 + x_2 = y_1, x_3 = y_2, x_4 = y_3, x_5 = y_4$.

$$\Rightarrow L(\theta|x) = \frac{(x_1+x_2+\dots+x_5)!}{x_1! x_2! \dots x_5!} \left(\frac{1}{2}\right)^{x_1} \left(\frac{1}{4}\right)^{x_2} \left(\frac{1}{4}\right)^{x_3} \left(\frac{1}{4}\right)^{x_4}$$

$$\Rightarrow l(\theta|x) = C + \underbrace{x_1 \ln \frac{1}{2} + (x_2+x_5) \ln \frac{1}{4}}_{x_1 + x_5} + (x_3+x_4) \ln \left(\frac{1}{4}\right) = C + (x_2+x_5) \ln \theta + (x_3+x_4) \ln (1-\theta).$$

$$\Rightarrow Q(\theta, \theta^{(r)}) = E[l(\theta|x) | y, \theta^{(r)}] = E[(x_2+x_5) \ln \theta + (x_3+x_4) \ln (1-\theta) | \theta^{(r)}, y].$$

$$= E[x_2 \ln \theta | y, \theta^{(r)}] + y_4 \ln (1-\theta) + (y_2+y_3) \ln (1-\theta)$$

$$= \frac{\theta^{(r)}}{\frac{1}{2} + \frac{\theta^{(r)}}{4}} y_1 \ln \theta + y_4 \ln (1-\theta) + (y_2+y_3) \ln (1-\theta)$$

$$= \frac{\theta^{(r)}}{\frac{1}{2} + \frac{\theta^{(r)}}{4}} [y_1 + y_4] \ln \theta + (y_2+y_3) \ln (1-\theta)$$

To maximise $Q(\theta, \theta^{(r)})$:

$$\text{let } \frac{\partial}{\partial \theta} Q(\theta, \theta^{(r)}) = \frac{1}{\theta} \left[\frac{\theta^{(r)}}{\frac{1}{2} + \frac{\theta^{(r)}}{4}} y_1 + y_4 \right] - \frac{1}{1-\theta} (y_2+y_3) = 0 \Rightarrow (1-\theta) \left(\frac{\theta^{(r)}}{\frac{1}{2} + \frac{\theta^{(r)}}{4}} y_1 + y_4 \right) = \theta (y_2+y_3).$$

$$\Rightarrow \theta^{(r+1)} = \frac{y_1 \theta^{(r)} / (\frac{1}{2} + \frac{\theta^{(r)}}{4}) + y_4}{y_1 \theta^{(r)} / (\frac{1}{2} + \frac{\theta^{(r)}}{4}) + y_4 + y_2 + y_3}.$$

(2) Newton-Raphson (faster but less stable numerically and less tractable analytically). When will it fail?

formulas:

$$\text{for 1-dimension: } x^{(n+1)} = x^{(n)} - \frac{f(x^{(n)})}{f'(x^{(n)})} \text{ to solve } f(x) = 0$$

for higher dimension:

$$\text{let } \boldsymbol{x}_n = (x_1, \dots, x_n)^T$$

$\theta = (\theta_1, \dots, \theta_q)^T$ for q parameters.

$$l(\theta) = \ln L(\theta) = \sum_{i=1}^n \ln f(x_i | \theta).$$

$U(\theta) = \frac{\partial}{\partial \theta} l(\theta) \rightarrow$ score function, is a $q \times 1$ vector.

$$H(\theta) = \frac{\partial^2}{\partial \theta^2} U(\theta) = \frac{\partial^2}{\partial \theta^2} l(\theta) \rightarrow$$
 Hessian function, is a $q \times q$ matrix. $\begin{pmatrix} \frac{\partial^2}{\partial \theta_1^2} l(\theta) & \cdots & \frac{\partial^2}{\partial \theta_1 \partial \theta_q} l(\theta) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial \theta_q \partial \theta_1} l(\theta) & \cdots & \frac{\partial^2}{\partial \theta_q^2} l(\theta) \end{pmatrix}$

$J(\theta) = -H(\theta) \rightarrow$ observed information function.

$$I(\theta) = E[J(\theta)] = E[-H(\theta)] \rightarrow$$
 Fisher information function.

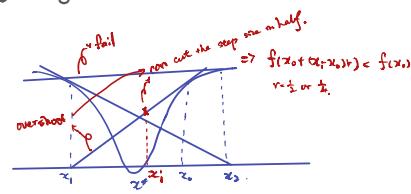
(1) Then start with $\theta^{(0)}$,

(2) for $k=0, 1, 2, \dots$

$$\theta^{(k+1)} = \theta^{(k)} - [H(\theta^{(k)})]^{-1} U(\theta^{(k)}) \quad \text{if } \theta^{(k+1)} = \theta^{(k)} + [I(\theta^{(k)})]^{-1} U(\theta^{(k)}) \Rightarrow \text{Fisher-scoring}$$

$$\stackrel{?}{=} \theta^{(k)} + [J(\theta^{(k)})]^{-1} U(\theta^{(k)}) \text{ to solve } Q(\theta) = 0$$

(3) The iteration continues until $\|\theta^{(m+1)} - \theta^{(m)}\|$ or $|Q(\theta^{(m+1)}, \theta^{(m)}) - Q(\theta^{(m)}, \theta^{(m)})|$ is smaller than tolerance level $\varepsilon > 0$, e.g. 10^{-6} .



$$|f'(x)|^{-1} \leq 1 \text{ for all } x - d < x < x + d.$$

note:

• need to check $H(\theta) < 0$ for θ s.t. $U(\theta) = 0$.

• NR can be extended to BFGS, BHHH, Nelder-Mead and simulated annealing etc.

Example continued:

The observed data log-likelihood is:

$$l_Y(\theta) = \ln C_1 + y_1 \ln(2+\theta) + (y_2 + y_3) \ln(1-\theta) + y_4 \ln \theta.$$

and

$$\frac{\partial}{\partial \theta} l_Y(\theta) = \frac{y_1}{2+\theta} - \frac{y_2+y_3}{1-\theta} + \frac{y_4}{\theta}$$

$$\text{let } U(\theta) = \frac{\partial}{\partial \theta} l_Y(\theta) = \frac{y_1}{2+\theta} - \frac{y_2+y_3}{1-\theta} + \frac{y_4}{\theta}$$

$$\Rightarrow H(\theta) = \frac{\partial^2}{\partial \theta^2} l_Y(\theta) = \frac{-y_1}{(2+\theta)^2} - \frac{y_2+y_3}{(1-\theta)^2} - \frac{y_4}{\theta^2}$$

$$\Rightarrow I(\theta) = E[-H(\theta)] = E\left[\frac{y_1}{(2+\theta)^2} + \frac{y_2+y_3}{(1-\theta)^2} + \frac{y_4}{\theta^2}\right] = \frac{\frac{(2+\theta)}{4} \cdot n}{(2+\theta)^2} + \frac{\frac{1-\theta}{4} \cdot 2 \cdot n}{(1-\theta)^2} + \frac{\frac{\theta}{4} n}{\theta^2}$$

$$= n\left(\frac{1}{4(2+\theta)} + \frac{1}{2(1-\theta)} + \frac{1}{4\theta}\right)$$

$$\Rightarrow \theta^{(m+1)} = \theta^{(m)} - \frac{U(\theta^{(m)})}{H(\theta^{(m)})} \text{ or } \theta^{(m+1)} = \theta^{(m)} + \frac{U(\theta^{(m)})}{I(\theta^{(m)})} \text{ of Fisher scoring.}$$

(3) maxLik() in R.

(4) Estimation of variance of MLE.

(i) Using asymptotic efficiency of MLE,

$$\bullet \text{Var}(\hat{\theta}) \approx (J(\hat{\theta}))^{-1} = \left[-\frac{\partial^2}{\partial \theta \partial \theta} l(\hat{\theta})\right]^{-1} \text{ or } (I(\hat{\theta}))^{-1} = E\left[-\frac{\partial^2}{\partial \theta \partial \theta} l(\hat{\theta})\right]^{-1}$$

$$\text{s.e.}(\hat{\theta}_i) = \sqrt{\text{Var}(\hat{\theta})}_{ii}, i=1, \dots, q \text{ i.e. the } i\text{th diagonal of } \text{Var}(\hat{\theta}).$$

(ii) Using Louis method,

$$\bullet \text{Var}(\hat{\theta}) \approx (J(\hat{\theta}))^{-1} \text{ where } J(\phi) = -H(\phi) = -\frac{\partial^2 Q(\phi, \phi)}{\partial \phi \partial \phi} \Big|_{\phi=\hat{\theta}} - \text{Var}\left[\frac{\partial \ln l(y, \phi)}{\partial \phi} \Big|_{\phi=\hat{\theta}}\right]$$

example continued:

V

$$\begin{aligned}
 -\frac{\partial^2 Q(\theta, \hat{\theta})}{\partial \theta^2} \Big|_{\theta=\hat{\theta}} &= -\frac{\partial^2}{\partial \theta^2} \left[\frac{\hat{\theta}}{2+\hat{\theta}} (y_1 + y_4) \ln \theta + (y_2 + y_3) \ln(1-\theta) \right] \\
 &= \frac{1}{\hat{\theta}^2} \left(\frac{\hat{\theta}}{2+\hat{\theta}} y_1 + y_4 \right) + \frac{1}{(1-\hat{\theta})^2} (y_2 + y_3). \quad (\text{replace } \theta^{(n)} \text{ with } \hat{\theta} \text{ in } Q(\theta, \theta^{(n)}); \text{ and then} \\
 &\quad \text{differentiate with respect to } \theta; \text{ then replace } \theta \text{ with } \hat{\theta}). \\
 \dots (1)
 \end{aligned}$$

$$\frac{\partial}{\partial \theta} l(\theta | y, z) = \frac{\partial}{\partial \theta} \left[C + (x_2 + x_5) \ln \theta + (x_3 + x_4) \ln(1-\theta) \right] = \frac{x_2 + x_5}{\theta} - \frac{x_3 + x_4}{1-\theta}.$$

$$\frac{\partial^2}{\partial \theta^2} l(\theta | y, z) = -\frac{x_2 + x_5}{\theta^2} - \frac{x_3 + x_4}{(1-\theta)^2}$$

$$\begin{aligned}
 \Rightarrow -\text{Var} \left[\frac{\partial l(\theta | y, z)}{\partial \theta} \Big| y, \hat{\theta} \right] \Big|_{\theta=\hat{\theta}} &= -\text{Var} \left[-\frac{x_2 + x_5}{\theta^2} - \frac{x_3 + x_4}{(1-\theta)^2} \Big| y, \hat{\theta} \right] = -\text{Var} \left(\frac{x_2}{\theta^2} \Big| y, \hat{\theta} \right) = -\frac{y_1 \frac{\hat{\theta}}{\hat{\theta}+2} (1-\frac{\hat{\theta}}{\hat{\theta}+2})}{\hat{\theta}^2} \\
 \Rightarrow \text{Var}(\hat{\theta}) &\approx [(1) + (2)]^{-1} \\
 &= \frac{-2y_1}{\hat{\theta}(\hat{\theta}+2)^2} \quad \dots (2)
 \end{aligned}$$