

Solution to Homework #10

Question 1 (Q4 of page 127)

(a) Using Taylor Expansion, we get

$$\sum_{i=1}^n \ddot{\ell}_{\hat{\theta}_n}(x_i) = \sum_{i=1}^n \ddot{\ell}_{\theta}(x_i) + \left\{ \sum_{i=1}^n \ddot{\ell}_{\tilde{\theta}_n}(x_i) \right\} (\hat{\theta}_n - \theta) \quad \text{where } \tilde{\theta}_n \text{ lies between } \theta \text{ \& } \hat{\theta}_n.$$

$$\left| \frac{1}{n} \sum_{i=1}^n \ddot{\ell}_{\hat{\theta}_n}(x_i) - \frac{1}{n} \sum_{i=1}^n \ddot{\ell}_{\theta}(x_i) \right| \leq \left\{ \frac{1}{n} \sum_{i=1}^n |M_{\tilde{\theta}_n}(x_i)| + o_p(1) \right\} |\hat{\theta}_n - \theta|$$

$$= \mathbb{P}_n [M_{\tilde{\theta}_n}(x)] |\hat{\theta}_n - \theta| + o_p(1).$$

$$\xrightarrow{P} 0 \quad \left[\begin{array}{l} |\hat{\theta}_n - \theta| \xrightarrow{P} 0 \text{ as } \sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, I_{\theta}^{-1}) \\ \& \sup_{\theta} \mathbb{E}_{\theta} [M(x)] = M < \infty \Rightarrow \mathbb{P} [M_{\tilde{\theta}_n}(x)] \text{ is bounded} \end{array} \right]$$

Now, by WLLN, $-\frac{1}{n} \sum_{i=1}^n \ddot{\ell}_{\theta}(x_i) \rightarrow I_{\theta}$

$$\begin{aligned} \hat{I}_n &= -\frac{1}{n} \sum_{i=1}^n \ddot{\ell}_{\hat{\theta}_n}(x_i) \\ &= -\left[\frac{1}{n} \sum_{i=1}^n \ddot{\ell}_{\hat{\theta}_n}(x_i) - \frac{1}{n} \sum_{i=1}^n \ddot{\ell}_{\theta}(x_i) \right] - \frac{1}{n} \sum_{i=1}^n \ddot{\ell}_{\theta}(x_i) \end{aligned}$$

$$\begin{aligned} \mathbb{E}(\hat{I}_n) &= 0 - \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \ddot{\ell}_{\theta}(x_i) \right] \\ &= I_{\theta} \end{aligned}$$

\hat{I}_n is a consistent estimator of I_{θ} .

by the continuous mapping theorem, \hat{I}_n^{-1} is a consistent estimator of I_{θ}^{-1}

[Proved]

(b) $\sqrt{n} I_{\theta}^{-1/2} (\hat{\theta}_n - \theta) \xrightarrow{d} N(0, I_{k \times k})$ [all the regularity conditions are satisfied]

Now, $\hat{I}_n \xrightarrow{P} I_{\theta} \Rightarrow \frac{\hat{I}_n}{I_{\theta}} \xrightarrow{P} 1 \Rightarrow \frac{\hat{I}_n^{-1/2}}{I_{\theta}^{-1/2}} \xrightarrow{P} 1$ [by Cont. Mapping Thm]

by Slutsky's Theorem, $\sqrt{n} \hat{I}_n^{-1/2} (\hat{\theta}_n - \theta) \xrightarrow{d} N(0, I_{k \times k})$

$\therefore n(\hat{\theta}_n - \theta)' \hat{I}_n (\hat{\theta}_n - \theta) \xrightarrow{d} \chi_k^2$

So, the required $(1-\alpha)$ CI is given by:-

$\{ \theta : n(\hat{\theta}_n - \theta)' \hat{I}_n (\hat{\theta}_n - \theta) < \chi_{\alpha, k}^2 \}$ where $\chi_{\alpha, k}^2$ is the $(1-\alpha)^{th}$ percentile of χ_k^2 distribution.

(c) $\sqrt{n} I_{\theta}^{-1/2} (\hat{\theta}_n - \theta) \xrightarrow{d} N(0, I_{k \times k})$

$\Rightarrow \sqrt{n} I_{\xi_n}^{-1/2} (\hat{\theta}_n - \theta) \xrightarrow{d} N(0, I_{k \times k})$ [if $\xi_n \rightarrow \theta$, by Slutsky's Thm]

$\Rightarrow n(\hat{\theta}_n - \theta)' I_{\xi_n} (\hat{\theta}_n - \theta) \xrightarrow{d} \chi_k^2 \rightarrow \text{eq (i)}$

Now, $\ln(\theta) = \sum_{i=1}^n \ln_{\theta}(x_i)$

$= \sum_{i=1}^n \ln_{\hat{\theta}_n}(x_i) + \cancel{\frac{1}{n} \sum_{i=1}^n \ln_{\hat{\theta}_n}(x_i)} (\hat{\theta}_n - \theta)' \left(\sum_{i=1}^n \dot{\ln}_{\hat{\theta}_n}(x_i) \right) + \frac{1}{2} (\hat{\theta}_n - \theta)' \left[\sum_{i=1}^n \ddot{\ln}_{\hat{\theta}_n}(x_i) \right] (\hat{\theta}_n - \theta)$, for some $\tilde{\theta}_n$

$\Rightarrow -2(\ln(\theta) - \ln(\hat{\theta}_n)) = (\hat{\theta}_n - \theta)' \left(\sum_{i=1}^n \dot{\ln}_{\hat{\theta}_n}(x_i) \right)$

$+ n(\hat{\theta}_n - \theta)' I_{\tilde{\theta}_n} (\hat{\theta}_n - \theta)$

$= 0 + n(\hat{\theta}_n - \theta)' I_{\tilde{\theta}_n} (\hat{\theta}_n - \theta)$

$\xrightarrow{d} \chi_k^2$ (by eq (i), as $\hat{\theta}_n \xrightarrow{as} \theta \Rightarrow \tilde{\theta}_n \xrightarrow{as} \theta$)

$\Rightarrow -2(\ln(\theta) - \ln(\hat{\theta}_n)) \xrightarrow{d} \chi_k^2$

\therefore the required $(1-\alpha)$ CI is given by:-

$\{ \theta : -2(\ln(\theta) - \ln(\hat{\theta}_n)) < \chi_{\alpha, k}^2 \}$

Question 2 (No 5 of Class Notes)

(a) We assume that the complete data is the genotype data which is a multinomial model with 6 possible outcomes (OO, AA, AO, BB, BO, AB) occurring with prob $(r^2, p^2, 2rp, q^2, 2rq, 2pq)$.

Let us assume that the complete data has n_1 individuals with genotype AA & n_2 individuals with genotype BB

the complete log-likelihood is given by

$$\begin{aligned} \ell/n_1, n_2 &= \text{constant} + N_O \log(r^2) + n_1 \log(p^2) + (N_A - n_1) \log(2rp) \\ &\quad + n_2 \log(q^2) + (N_B - n_2) \log(2rq) + N_{AB} \log(2pq) \\ &= k + (2N_O + N_A + N_B - n_1 - n_2) \log(r) + (n_1 + N_A + N_{AB}) \log(p) \\ &\quad + (n_2 + N_B + N_{AB}) \log(q) \end{aligned}$$

$$= k + a \log r + b \log p + c \log q \quad [\text{where } a, b \text{ \& } c \text{ depend on } n_1 \text{ \& } n_2]$$

$$\frac{\partial \ell}{\partial p} = -\frac{a}{1-p-q} + \frac{b}{p} \quad [\because r=1-p-q]$$

$$\frac{\partial \ell}{\partial q} = -\frac{a}{1-p-q} + \frac{c}{q}$$

E-step

For the E-step $\hat{\lambda}_k$ at the k -th step, we must have

$$E \left[\frac{\partial \ell}{\partial p} \mid p^{(k)}, q^{(k)}, r^{(k)} \right] = 0$$

$$\& \quad E \left[\frac{\partial \ell}{\partial q} \mid p^{(k)}, q^{(k)}, r^{(k)} \right] = 0$$

Let $a_k = E(a | P^{(k)}, q^{(k)}, r^{(k)})$, $b_k = E(b | P^{(k)}, q^{(k)}, r^{(k)})$

& $c_k = E(c | P^{(k)}, q^{(k)}, r^{(k)})$

E-step yields :-

$$-\frac{a_k}{1-p-q} + \frac{b_k}{p} = 0 \quad \longrightarrow \text{eq (i)}$$

$$-\frac{a_k}{1-p-q} + \frac{c_k}{q} = 0 \quad \longrightarrow \text{eq (ii)}$$

$$a_k = E \left[2N_0 + N_A + N_B - n_1 - n_2 \mid P^{(k)}, q^{(k)}, r^{(k)} \right]$$

$$= 2N_0 + N_A + N_B - \frac{N_A p^{(k)}}{p^{(k)} + 2r^{(k)}} - \frac{N_B q^{(k)}}{q^{(k)} + 2r^{(k)}}$$

$\therefore n_1 / N_A, P^{(k)}, q^{(k)}, r^{(k)}$
 $\sim \text{Bin} \left(N_A, \frac{p^{(k)}}{p^{(k)} + 2r^{(k)}} \right)$

& $n_2 / N_B, P^{(k)}, q^{(k)}, r^{(k)}$
 $\sim \text{Bin} \left(N_B, \frac{q^{(k)}}{q^{(k)} + 2r^{(k)}} \right)$

$$= 2 \left[N_0 + \frac{N_A r^{(k)}}{p^{(k)} + 2r^{(k)}} + \frac{N_B r^{(k)}}{q^{(k)} + 2r^{(k)}} \right]$$

Similarly, $b_k = \frac{N_A p^{(k)}}{p^{(k)} + 2r^{(k)}} + N_A + N_B$

& $c_k = \frac{N_B q^{(k)}}{q^{(k)} + 2r^{(k)}} + N_B + N_B$

M-step

Solving eq (i) & (ii) we get,

$$\hat{p} = \frac{b_k}{a_k + b_k + c_k} \quad \hat{q} = \frac{c_k}{a_k + b_k + c_k}, \quad \hat{r} = \frac{a_k}{a_k + b_k + c_k}$$

So, to apply the E-M algorithm, we start with an initial choice $P^{(0)}, q^{(0)}, r^{(0)}$ & compute P, q, r iteratively as follows:-

$$p^{(k+1)} = \frac{N_A p^{(k)} / \{p^{(k)} + 2r^{(k)}\} + N_A + N_{AB}}{2(N_C + N_A + N_B + N_{AB})}$$

$$q^{(k+1)} = \frac{N_B q^{(k)} / \{q^{(k)} + 2r^{(k)}\} + N_B + N_{AB}}{2(N_C + N_A + N_B + N_{AB})}$$

$$r^{(k+1)} = \frac{N_C + N_A r^{(k)} / \{p^{(k)} + 2r^{(k)}\} + N_B r^{(k)} / \{q^{(k)} + 2r^{(k)}\}}{N_C + N_A + N_B + N_{AB}}$$

Stopping Condition

Fix some small $\epsilon > 0$ & stop when $\max\{|p^{(k+1)} - p^{(k)}|, |q^{(k+1)} - q^{(k)}|, |r^{(k+1)} - r^{(k)}|\} < \epsilon$
 & conclude that $(p^{(k+1)}, q^{(k+1)}, r^{(k+1)})$ are the required EM-algorithm estimates of (p, q, r)

(b) We take $N_C = 176$, $N_A = 182$, $N_B = 60$ & $N_{AB} = 17$. We start with $p^{(0)} = q^{(0)} = r^{(0)} = 1/3$ & take $\epsilon = 10^{-4}$. The iterations are shown below:-

<u>k</u>	<u>p^(k)</u>	<u>q^(k)</u>	<u>r^(k)</u>	<u>Max Diff</u>
				0.2567
1	0.2985	0.1115	0.5900	
2	0.2710	0.0945	0.6346	0.0445
3	0.2655	0.0933	0.6412	0.0066
4	0.2646	0.0932	0.6422	0.0010
5	0.2645	0.0932	0.6424	0.0002
6	0.2644	0.0932	0.6424	$< 0.0001 = \epsilon$

the required estimates of (p, q, r) are $(0.2644, 0.0932, 0.6424)$

Question 3 (No 6 of Class Notes)

(a) X has density $f(x)$ & $Y|X=x \sim N(\beta x, \sigma^2)$.
 $(x_1, y_1), \dots, (x_n, y_n)$ are iid obs from (X, Y) , but x_{m+1}, \dots, x_n are missing
 $1 < m < n$ & the missingness satisfies MAR condition. The observed x_i 's are
 distinct & we assume that X has point mass $P_i > 0$ at the observed
 data $x_i = x_i$, for $i=1(1)m$, $\sum_{i=1}^m P_i = 1$

the likelihood function is given by

$$L = \prod_{i=1}^m \left[P_i \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_i - \beta x_i)^2}{2\sigma^2}\right\} \right] \times \prod_{j=m+1}^n \left[\sum_{i=1}^m P_j \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_j - \beta x_j)^2}{2\sigma^2}\right\} \right]$$

(b) $l = \log(L) =$ ~~$\log(L)$~~

$$\frac{\partial l}{\partial \beta} = 0 \Rightarrow \sum_{i=1}^m \frac{(y_i - \beta x_i) x_i}{\sigma^2} + \sum_{j=m+1}^n \frac{\sum_{i=1}^m P_j \exp\left\{-\frac{(y_j - \beta x_j)^2}{2\sigma^2}\right\} x_j (y_j - \beta x_j)}{\sum_{j=1}^m P_j \exp\left\{-\frac{(y_j - \beta x_j)^2}{2\sigma^2}\right\}} = 0$$

$$\frac{\partial l}{\partial \sigma^2} = 0 \Rightarrow -\frac{n}{2\sigma^2} + \sum_{i=1}^m \frac{(y_i - \beta x_i)^2}{2\sigma^4} + \sum_{j=m+1}^n \frac{\sum_{i=1}^m P_j \exp\left\{-\frac{(y_j - \beta x_j)^2}{2\sigma^2}\right\} \frac{(y_j - \beta x_j)^2}{\sigma^4}}{\sum_{j=1}^m P_j \exp\left\{-\frac{(y_j - \beta x_j)^2}{2\sigma^2}\right\}} = 0$$

$$\frac{\partial l}{\partial P_a} = 0 \Rightarrow \frac{1}{P_a} - \frac{1}{1 - \sum_{j=1}^{m-1} P_j} + \sum_{j=m+1}^n \frac{\left[\exp\left\{-\frac{(y_j - \beta x_a)^2}{2\sigma^2}\right\} - \exp\left\{-\frac{(y_j - \beta x_m)^2}{2\sigma^2}\right\} \right]}{\sum_{j=1}^m P_j \exp\left\{-\frac{(y_j - \beta x_j)^2}{2\sigma^2}\right\}} = 0$$

$a = 1(1)m-1 \quad \left[\sum_{j=1}^m P_j = 1 \right]$

(c) For the EM algorithm, we assume that the missing data x_{m+1}, \dots, x_n is known. Then the complete likelihood is given by -

$$L_c = \prod_{i=1}^m \left[P_i \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_i - \beta x_i)^2}{2\sigma^2}\right\} \right] \times \prod_{i=m+1}^n \left[f_i \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_i - \beta x_i)^2}{2\sigma^2}\right\} \right]$$

$$\log L_c = \sum_{i=1}^m \log(P_i) - n \log(\sigma^2) - \sum_{i=1}^m \frac{(y_i - \beta x_i)^2}{2\sigma^2} - \sum_{i=m+1}^n \frac{(y_i - \beta x_i)^2}{2\sigma^2} + \sum_{i=m+1}^n \log(f_i) + \text{constant}$$

$$L_2 = \prod_{i=1}^m \left[\frac{P_i}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y_i - \beta x_i)^2}{\sigma^2} \right\} \right] \prod_{i=m+1}^n \prod_{j=1}^m \left[P_j \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y_i - \beta x_j)^2}{\sigma^2} \right\} \right]^{I_{ij}}$$

$$I_{ij} = I \{ x_i = x_j \} \quad i=(m+1)(1)m, j=1(1)m \quad \sum_{j=1}^m I_{ij} = 1$$

$$L_2 = \sum_{i=1}^m \log P_i - \frac{n}{2} \log(\sigma^2) - \sum_{i=1}^m \frac{(y_i - \beta x_i)^2}{\sigma^2} + \sum_{i=m+1}^n \sum_{j=1}^m I_{ij} \log(P_j) - \sum_{i=m+1}^n \sum_{j=1}^m I_{ij} \frac{(y_i - \beta x_j)^2}{\sigma^2} + \text{const (indep of other parameters)}$$

$$\frac{\partial L_2}{\partial P_j} = \frac{1}{P_j} + \frac{1}{P_j} \sum_{i=m+1}^n I_{ij} - \lambda \quad \left[\text{using Lagrange multipliers since } \sum_{i=1}^m P_i = 1 \right]$$

$$\frac{\partial L_2}{\partial \beta} = \frac{\sum_{i=1}^m (y_i - \beta x_i) x_i + \sum_{i=m+1}^n \sum_{j=1}^m (y_i - \beta x_j) x_j}{2\sigma^2}$$

$$\frac{\partial L_2}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \left[\sum_{i=1}^m (y_i - \beta x_i)^2 + \sum_{i=m+1}^n \sum_{j=1}^m I_{ij} (y_i - \beta x_j)^2 \right]$$

E-step

$$I_{ij}^{(k)} = E \left[I_{ij} / P_i^{(k)}, \dots, P_m^{(k)}, y, \beta^{(k)}, \sigma^{2(k)} \right]$$

$$= P \left(x_i = x_j / P_i^{(k)}, \dots, P_m^{(k)}, y, \beta^{(k)}, \sigma^{2(k)} \right)$$

$$= \frac{P_j^{(k)} \frac{1}{\sqrt{2\pi\sigma^{2(k)}}} \exp \left\{ -\frac{(y_i - \beta^{(k)} x_j)^2}{\sigma^{2(k)}} \right\}}{\sum_{l=1}^m P_l^{(k)} \frac{1}{\sqrt{2\pi\sigma^{2(k)}}} \exp \left\{ -\frac{(y_i - \beta^{(k)} x_l)^2}{\sigma^{2(k)}} \right\}}$$

$$= \frac{P_j^{(k)} \exp \left\{ -\frac{(y_i - \beta^{(k)} x_j)^2}{\sigma^{2(k)}} \right\}}{\sum_{l=1}^m P_l^{(k)} \exp \left\{ -\frac{(y_i - \beta^{(k)} x_l)^2}{\sigma^{2(k)}} \right\}}$$

So, the score equations at the $(k+1)$ -th step are: \rightarrow

$$E \left[\frac{\partial l_2}{\partial P_j} \right] = 0$$

$$\Rightarrow E \left[\frac{\partial l_2}{\partial P_j} \mid P_1^{(k)}, \dots, P_m^{(k)}, y, \beta^{(k)}, \sigma^{2(k)} \right] = 0$$

$$\Rightarrow \frac{1}{P_j} + \frac{1}{P_j} \sum_{i=m+1}^n I_{ij}^{(k)} - \lambda = 0 \quad \rightarrow (i)$$

$$E \left[\frac{\partial l_2}{\partial \beta} \mid P_1^{(k)}, \dots, P_m^{(k)}, y, \beta^{(k)}, \sigma^{2(k)} \right] = 0$$

$$\Rightarrow \sum_{i=1}^m (y_i - \beta x_i) x_i + \sum_{i=m+1}^n \sum_{j=1}^m I_{ij}^{(k)} (y_i - \beta x_j) x_j = 0 \quad \rightarrow (ii)$$

$$E \left[\frac{\partial l_2}{\partial \sigma^2} \mid P_1^{(k)}, \dots, P_m^{(k)}, y, \beta^{(k)}, \sigma^{2(k)} \right] = 0$$

$$\Rightarrow -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \left[\sum_{i=1}^m (y_i - \beta x_i)^2 + \sum_{i=m+1}^n \sum_{j=1}^m I_{ij}^{(k)} (y_i - \beta x_j)^2 \right] = 0 \quad \rightarrow (iii)$$

M-step

Solving eq (i), (ii) & (iii) we get,

$$P_j^{(k+1)} = \frac{1 + \sum_{i=m+1}^n I_{ij}^{(k)}}{\sum_{i=1}^m x_i y_i + \sum_{i=m+1}^n \sum_{j=1}^m I_{ij}^{(k)} y_i x_j}$$

$$\beta^{(k+1)} = \frac{\sum_{i=1}^m x_i^2 y_i + \sum_{i=m+1}^n \sum_{j=1}^m I_{ij}^{(k)} y_i x_j}{\sum_{i=1}^m x_i^2 (1 + \sum_{j=m+1}^n I_{ij}^{(k)})}$$

$$\hat{\sigma}^{2(k+1)} = \frac{1}{n} \left[\sum_{i=1}^m (y_i - \beta^{(k+1)} x_i)^2 + \sum_{i=m+1}^n \sum_{j=1}^m I_{ij}^{(k)} (y_i - \beta^{(k+1)} x_j)^2 \right]$$

Stopping Condition

We fix some small $\epsilon > 0$ & stop when

$$\max \left\{ \left| \hat{\beta}^{(k+1)} - \hat{\beta}^{(k)} \right|, \left| \hat{\sigma}^{2(k+1)} - \hat{\sigma}^{2(k)} \right|, \left| \hat{P}_1^{(k+1)} - \hat{P}_1^{(k)} \right|, \dots, \left| \hat{P}_m^{(k+1)} - \hat{P}_m^{(k)} \right| \right\} < \epsilon$$

& conclude that $(\hat{\beta}^{(k+1)}, \hat{\sigma}^{2(k+1)}, \hat{P}_1^{(k+1)}, \dots, \hat{P}_m^{(k+1)})$ are the required estimates using the EM algorithm.