

On estimating multivariate discrete probabilities by pooled incomplete samples and related topics

関東学院大学 経済学部 布能 英一郎

1. Introduction

Proposition 1.1 k, m は $m < k$ なる自然数。確率変数 \mathbf{X}, \mathbf{Y} は互いに独立で

$$\begin{aligned} \mathbf{X} &= (X_1, \dots, X_m, \dots, X_k) \sim \text{Multinomial}(N_1; p_0, \dots, p_m, \dots, p_k), \\ \mathbf{Y} &= (Y_1, \dots, Y_m) \sim \text{Multinomial}(N_2; \frac{p_0}{\sum_{j=0}^m p_j}, \frac{p_1}{\sum_{j=0}^m p_j}, \dots, \frac{p_m}{\sum_{j=0}^m p_j}) \end{aligned}$$

とする。この状況下で、Asano(1965) は p_i の MLE \hat{p}_i が

$$\frac{x_i + y_i}{N_1 \left(1 + \frac{N_2}{\sum_{j=0}^m x_j}\right)}, \quad \text{if } i \leq m, \quad \frac{x_i}{N_1}, \quad \text{if } i > m \tag{1}$$

であることを示し、更に、 \hat{p}_i に関する諸性質を研究した。これと同様な現象は、他の分布の下でもいくつか生じている。

Proposition 1.2 確率変数 \mathbf{X}, \mathbf{Y} は互いに独立で

$$\begin{aligned} \mathbf{X} &= (X_1, \dots, X_m, \dots, X_k) \sim \text{NegativeMultinomial}(r_1; p_1, \dots, p_m, \dots, p_k), \\ \mathbf{Y} &= (Y_1, \dots, Y_m) \sim \text{NegativeMultinomial}(r_2; \frac{(1-p_0)p_1}{\sum_{j=1}^m p_j}, \dots, \frac{(1-p_0)p_m}{\sum_{j=1}^m p_j}) \end{aligned}$$

ならば

$$\hat{p}_i = \begin{cases} \frac{T_x + T_y}{T_x + T_y + r_1 + r_2} \frac{x_i + y_i}{T_x \left(1 + \frac{T_y}{\sum_{j=1}^m x_j}\right)}, & 1 \leq i \leq m, \\ \frac{T_x + T_y}{T_x + T_y + r_1 + r_2} \frac{x_i}{T_x}, & i > m. \end{cases} \tag{2}$$

但し、 $T_x = \sum_{j=1}^k x_j$, $T_y = \sum_{j=1}^m y_j$.

Proposition 1.3 $X_1, \dots, X_m, \dots, X_k, Y_1, \dots, Y_m$ はすべて独立で、

$$X_i \sim \text{Poisson}(\lambda_i), \quad Y_i \sim \text{Poisson}\left(\frac{\lambda_i}{\lambda_1 + \dots + \lambda_m}\right)$$

と仮定。このとき、 λ_i の MLE $\hat{\lambda}_i$ は

$$\hat{\lambda}_i = \begin{cases} \frac{T_x + T_y}{2} \frac{x_i + y_i}{T_x \left(1 + \frac{T_y}{\sum_{j=1}^m x_j}\right)} & \text{if } i \leq m, \\ \frac{T_x + T_y}{2} \frac{x_i}{T_x} & \text{if } i > m. \end{cases} \quad (3)$$

Proposition 1.2, 1.3 にて、MLE が Asano の形と同じようになることを、次表で見るとわかりやすい。

	$i \leq m$ に対して	$i \geq m + 1$ に対して
Proposition 1.1 (多項分布)(Asano)	$\frac{x_i + y_i}{N_1 \left(1 + \frac{N_2}{\sum_{j=0}^m x_j}\right)}$	$\frac{x_i}{N_1}$
Proposition 1.2 (負の多項分布)	$\frac{T_x + T_y}{T_x + T_y + r_1 + r_2} \frac{x_i + y_i}{T_x \left(1 + \frac{T_y}{\sum_{j=1}^m x_j}\right)}$	$\frac{T_x + T_y}{T_x + T_y + r_1 + r_2} \frac{x_i}{T_x}$
Proposition 1.3 (ポアソン分布)	$\frac{T_x + T_y}{2} \frac{x_i + y_i}{T_x \left(1 + \frac{T_y}{\sum_{j=1}^m x_j}\right)}$	$\frac{T_x + T_y}{2} \frac{x_i}{T_x}$

このように、MLE が Asano の形と同じようになる例が、多項分布以外にもある。しかしながら、どのような場合に、MLE が Asano が示した形と同じになるのか？について、はつきとした見解を見出せていなかった。これが本研究に着手した動機である。

2. パラメーター変換によって MLE を求める

Proposition 1.1, 1.2, 1.3 いずれの場合においても、MLE を以下のようなパラメーター変換によって求めるのが容易である。

Proposition 1.1 (Asano) の場合 $u_i = \frac{p_i}{\sum_{j=0}^m p_j}$, ($i \leq m$), $t_0 = \sum_{j=0}^m p_j$, $t_i = p_i$, ($i > m$) なるパラメーター変換を用いる。逆変換が $p_i = t_0 u_i$, ($i \leq m$) ゆえ

$$L \propto t_0^{x_0 + \dots + x_m} t_{m+1}^{x_{m+1}} \dots t_k^{x_k} u_0^{y_0 + x_0} \dots u_m^{y_m + x_m}$$

である。よって $t_0 = \frac{\sum_{j=0}^m x_j}{N_1}$, $t_j = \frac{x_j}{N_1}$, ($j > m$), $u_i = \frac{x_i + y_i}{\sum_{j=0}^m (x_j + y_j)}$. これより直ちに (1) を得る。

Proposition 1.2 の場合 $p_0 = 1 - (p_1 + \dots + p_k)$ と定める。そして、パラメーター変換

$$t = p_0, \quad v = \frac{p_1 + \dots + p_m}{p_1 + \dots + p_m + p_{m+1} + \dots + p_k}, \quad \xi_i = \frac{p_i}{p_1 + \dots + p_m}, \quad (i = 1, 2, \dots, m)$$

$$\eta_j = \frac{p_j}{p_{m+1} + \dots + p_k}, \quad (j = m + 1, \dots, k)$$

を用いる。逆変換が $p_0 = t, p_i = (1-t)v\xi_i$ ($i \leq m$), $p_i = (1-t)(1-v)\eta_i$ ($i > m$) であるから $L \propto t^{r_1+r_2}(1-t)^{T_x+T_y} v^{\sum_{j=1}^m x_j} (1-v)^{\sum_{j=m+1}^k x_j} \xi_1^{x_1+y_1} \dots \xi_m^{x_m+y_m} \eta_{m+1}^{x_{m+1}} \dots \eta_k^{x_k}$. これより

$$\hat{t} = \frac{r_1+r_2}{T_x+T_y+r_1+r_2}, \quad \hat{v} = \frac{\sum_{j=1}^m x_j}{T_x}, \quad \hat{\xi}_i = \frac{x_i+y_i}{\sum_{j=1}^m (x_j+y_j)}, \quad (i \leq m),$$

$$\hat{\eta}_i = \frac{x_i}{\sum_{j=m+1}^k x_j}, \quad (i > m)$$

を得る。 $\hat{p}_0 = \hat{t}$, $\hat{p}_i = (1-\hat{t})\hat{v}\hat{\xi}_i$, ($i \leq m$), $\hat{p}_i = (1-\hat{t})(1-\hat{v})\hat{\eta}_i$, ($i > m$) に代入することで (2) が得られる。

Proposition 1.3 の場合 便宜的に

$$\xi_i = (\lambda_1 + \dots + \lambda_m + \dots + \lambda_k) \frac{\lambda_i}{\lambda_1 + \dots + \lambda_m}$$

と置くと、 $X_i \sim \text{Poisson}(\lambda_i)$, $Y_j \sim \text{Poisson}(\xi_j)$ である。変数変換

$$s = \lambda_1 + \dots + \lambda_m + \lambda_{m+1} + \dots + \lambda_k, \quad t_0 = \frac{\lambda_1 + \dots + \lambda_m}{\lambda_1 + \dots + \lambda_m + \lambda_{m+1} + \dots + \lambda_k},$$

$$u_i = \frac{\lambda_i}{\lambda_1 + \dots + \lambda_m}, \quad (i = 1, 2, \dots, m), \quad t_j = \frac{\lambda_j}{\lambda_{m+1} + \dots + \lambda_k}, \quad (j = m+1, \dots, k)$$

を用いると、逆変換が $\lambda_i = st_0 u_i$, ($i \leq m$), $\lambda_i = s(1-t_0)t_i$, ($i > m$), $\xi_i = st_i$ であるから

$$L \propto s^{\sum_{i=1}^k x_i + \sum_{i=1}^m y_i} \exp(-2s) u_1^{x_1+y_1} \dots u_m^{x_m+y_m} t_0^{\sum_{i=1}^m x_i} (1-t_0)^{\sum_{i=m+1}^k x_i} t_{m+1}^{x_{m+1}} \dots t_k^{x_k}.$$

これより、

$$\hat{t}_0 = \frac{\sum_{j=1}^m x_j}{\sum_{j=1}^k x_j} = \frac{\sum_{j=1}^m x_j}{T_x}, \quad \hat{t}_j = \frac{x_j}{\sum_{j=m+1}^k x_j}, \quad j = m+1, m+2, \dots, k,$$

$$\hat{u}_i = \frac{x_i+y_i}{\sum_{l=1}^m (x_l+y_l)}, \quad i = 1, 2, \dots, m, \quad \hat{s} = \frac{T_x+T_y}{2}$$

が得られる。よって、 $\hat{\lambda}_i = \hat{s}\hat{t}_0\hat{u}_i$, ($i \leq m$), $\hat{\lambda}_i = \hat{s}(1-\hat{t}_0)\hat{t}_i$, ($i > m$) により、(3) が得られる。

3. Several extensions

本節では、Proposition 1.3 をいくつかの状況に拡張することを考察する。

Proposition 3.1 l, m, k は $l < m < k$ なる自然数。

$$\xi_i = (\lambda_1 + \dots + \lambda_l + \dots + \lambda_m) \frac{\lambda_i}{\lambda_1 + \dots + \lambda_l}, \quad (i = 1, 2, \dots, l) \quad (4)$$

と定める。確率変数 $X_1, \dots, X_k, Y_1, \dots, Y_l, Y_{m+1}, \dots, Y_k$ は互いに独立で $X_i \sim \text{Poisson}(\lambda_i)$ for $i = 1, \dots, k$, $Y_i \sim \text{Poisson}(\xi_i)$ for $i = 1, \dots, l$, $Y_i \sim \text{Poisson}(\lambda_i)$ for $i = m+1, \dots, k$, ならば、 λ_i の MLE $\hat{\lambda}_i$ は

$$\hat{\lambda}_i = \begin{cases} \frac{T'_x + T'_y}{2} \frac{x_i + y_i}{T'_x \left(1 + \frac{T'_y}{\sum_{i=1}^l x_i}\right)} & \text{if } i \leq l, \\ \frac{T'_x + T'_y}{2} \frac{x_i}{T'_x} & \text{if } l < i \leq m, \\ \frac{x_i + y_i}{2} & \text{if } i > m \end{cases} \quad (5)$$

である。但し $T'_x = \sum_{i=1}^m x_i$, $T'_y = \sum_{i=1}^l y_i$.

これは、パラメータ変換

$$\begin{aligned} s &= \lambda_1 + \dots + \lambda_m + \lambda_{m+1} + \dots + \lambda_s, & t &= \frac{\lambda_1 + \dots + \lambda_m}{\lambda_1 + \dots + \lambda_m + \dots + \lambda_k}, \\ u &= \frac{\lambda_1 + \dots + \lambda_l}{\lambda_1 + \dots + \lambda_m}, & \theta_i &= \frac{\lambda_i}{\lambda_1 + \dots + \lambda_l}, \quad (i = 1, 2, \dots, l), \\ v_i &= \frac{\lambda_i}{\lambda_{l+1} + \dots + \lambda_m}, \quad (j = l+1, \dots, m), \\ w_i &= \frac{\lambda_i}{\lambda_{m+1} + \dots + \lambda_k}, \quad (j = m+1, \dots, k) \end{aligned}$$

を用いれば容易に求められる。

Proposition 3.2 $l < m < k$ および ξ_i は、Proposition 2.1 と同じ。 $X_1, \dots, X_k, Y_1, \dots, Y_l$, および Y_* は互いに独立で、 $X_i \sim \text{Poisson}(\lambda_i)$, for all $i = 1, \dots, k$, $Y_i \sim \text{Poisson}(\xi_i)$, for $i = 1, \dots, l$, $Y_* \sim \text{Poisson}(\sum_{i=m+1}^k \lambda_i)$, とする。このとき λ_i の MLE は

$$\hat{\lambda}_i = \begin{cases} \frac{T'_x + T'_y}{2} \frac{x_i + y_i}{T'_x \left(1 + \frac{T'_y}{\sum_{i=1}^l x_i}\right)}, & \text{if } i \leq l, \\ \frac{T'_x + T'_y}{2} \frac{x_i}{T'_x}, & \text{if } l < i \leq m, \\ \frac{1}{2} \frac{\sum_{i=m+1}^k x_i + y_*}{\sum_{i=m+1}^k x_i} x_i, & \text{if } i > m \end{cases} \quad (6)$$

である。

これも、Proposition 3.1 と同じ変数変換で求められる。

Proposition 3.3 $l < m < k$ は、Proposition 2.1 と同じ。 $\{X_i\}_{i=1, \dots, k}$ および $\{Y_i\}_{i=1, \dots, l}$ は互いに独立で $X_i \sim \text{Poisson}(\lambda_i)$, $Y_i \sim \text{Poisson}(\xi_i)$, ここで

$$\xi_i = (\lambda_1 + \dots + \lambda_l + \dots + \lambda_m) \frac{\lambda_i}{\lambda_1 + \dots + \lambda_l}, \quad i = 1, 2, \dots, l. \quad (7)$$

このとき、 λ_i の MLE は

$$\hat{\lambda}_i = \begin{cases} \frac{T'_x + T_y}{2} \frac{x_i + y_i}{T'_x \left(1 + \frac{T_y}{\sum_{i=1}^l x_i}\right)}, & \text{if } i \leq l, \\ \frac{T'_x + T_y}{2} \frac{x_i}{T'_x}, & \text{if } l < i \leq m, \\ x_i, & \text{if } i > m \end{cases} \quad (8)$$

である。

Proposition 3.1 で用いたパラメータ変換を用いると、

$$L \propto s^{T_x + T_y} t^{(\sum_{i=1}^m x_i) + T_y} (1-t)^{\sum_{i=m+1}^k x_i} u^{\sum_{i=1}^l x_i} (1-u)^{\sum_{i=l+1}^m x_i} \\ \times \prod_{i=1}^l \theta_i^{x_i + y_i} \prod_{i=l+1}^m v_i^{x_i} \prod_{i=m+1}^k w_i^{x_i} \times \exp(-s(1+t)).$$

これより直ちに

$$\hat{\theta}_i = \frac{x_i + y_i}{\sum_{i=1}^l (x_i + y_i)}, \quad (i = 1, 2, \dots, l), \quad \hat{v}_i = \frac{x_i}{\sum_{i=l+1}^m x_i}, \quad (i = l+1, l+2, \dots, m), \\ \hat{w}_i = \frac{x_i}{\sum_{i=m+1}^k x_i}, \quad (i = m+1, m+2, \dots, k), \quad \hat{u} = \frac{\sum_{i=1}^l x_i}{\sum_{i=1}^m x_i}$$

を得る。 s および t の MLE を、対数尤度

$$\log L = C + (T_x + T_y) \log s + \left(\sum_{i=1}^m x_i + T_y \right) \log t + \left(\sum_{i=m+1}^k x_i \right) \log(1-t) - s(1+t)$$

より求めてみる。尤度方程式は

$$0 = \frac{\partial L}{\partial s} = \frac{T_x + T_y}{s} - (1+t), \quad (9)$$

$$0 = \frac{\partial L}{\partial t} = \frac{\sum_{i=1}^m x_i + T_y}{t} - \frac{\sum_{i=m+1}^k x_i}{1-t} - s \quad (10)$$

であるから、(9) より $s = \frac{T_x + T_y}{1+t}$ を得る。これを (10) に代入することで

$$0 = \frac{\sum_{i=1}^m x_i + T_y}{t} - \frac{\sum_{i=m+1}^k x_i}{1-t} - \frac{T_x + T_y}{1+t} \quad (11)$$

が得られる。この式は

$$0 = \left(-\sum_{i=1}^m x_i - T_y - \sum_{i=m+1}^k x_i + T_x + T_y \right) t^2 - \left(\sum_{i=m+1}^k x_i + T_x + T_y \right) t + \sum_{i=1}^m x_i + T_y \\ = \left(\sum_{i=m+1}^k x_i + T_x + T_y \right) t + \sum_{i=1}^m x_i + T_y$$

と同値なので、 $\hat{t} = \frac{\sum_{i=1}^m x_i + T_y}{\sum_{i=m+1}^k x_i + T_x + T_y}$ が得られる。これより

$$\begin{aligned}\hat{s} &= \frac{T_x + T_y}{1 + \hat{t}} = \frac{T_x + T_y}{1 + \frac{\sum_{i=1}^m x_i + T_y}{\sum_{i=m+1}^k x_i + T_x + T_y}} = \frac{(T_x + T_y)(\sum_{i=m+1}^k x_i + T_x + T_y)}{T_x + T_y + \sum_{i=m+1}^k x_i + \sum_{i=1}^m x_i + T_y} \\ &= \frac{T_x + T_y + \sum_{i=m+1}^k x_i}{2}\end{aligned}$$

が得られ、あとは MLE の invariant property によって (8) が求められる。

Proposition 3.4 $\{X_i\}_{i=1,\dots,k}$ および $\{Y_i\}_{i=1,\dots,m}$ ($m < k$) は互いに独立で、 $X_i \sim \text{Poisson}(\lambda_i)$, ($i = 1, \dots, k$), $Y_i \sim \text{Poisson}(\xi_i)$, ($i = 1, \dots, l$), $Y_i \sim \text{Poisson}(\eta_i)$, ($i = l+1, \dots, m$) と仮定する。ここで

$$\begin{aligned}\xi_i &= (\lambda_1 + \dots + \lambda_l + \dots + \lambda_m) \frac{\lambda_i}{\lambda_1 + \dots + \lambda_l}, \quad (i = 1, 2, \dots, l), \\ \eta_i &= (\lambda_{m+1} + \dots + \lambda_k) \frac{\lambda_i}{\lambda_{l+1} + \dots + \lambda_m}, \quad (i = l+1, l+2, \dots, m).\end{aligned}$$

そうすると、 λ_i の MLE は

$$\begin{aligned}\frac{T'_x + T'_y}{2} \frac{x_i + y_i}{T'_x \left(1 + \frac{T'_y}{\sum_{i=1}^l x_i}\right)}, \quad (i \leq l), \quad \frac{T'_x + T'_y}{2} \frac{x_i + y_i}{T'_x \left(1 + \frac{\sum_{i=l+1}^m y_i}{\sum_{i=l+1}^m x_i}\right)}, \quad (l < i \leq m), \\ \left\{1 + \frac{\sum_{i=l+1}^m y_i}{\sum_{i=m+1}^k x_i}\right\} \frac{x_i}{2}, \quad (i > m)\end{aligned}$$

である。なお、 $T'_y = \sum_{i=1}^l y_i$.

これも、Proposition 3.1 と同じ変数変換で求められる。

4. Related multinomial models

前章の結果は、多項分布の場合においてもほぼ同様に成り立つ。つまり、Prop. 1.3 の拡張として Prop 3.1 から Prop 3.4 を得たのであるが、ほぼ同様に Prop 1.1 を拡張することができる。

Proposition 4.1 l, k, m は $m < k$ なる自然数。確率変数 \mathbf{X}, \mathbf{Y} は互いに独立で

$$\mathbf{X} = (X_0, \dots, X_l, \dots, X_m, \dots, X_k) \sim \text{Multinomial}(N_1; p_0, \dots, p_k),$$

$$\mathbf{Y} = (Y_0, \dots, Y_l, Y_{m+1}, \dots, Y_k)$$

$$\sim \text{Multinomial}(N_2; \sum_{j=0}^m p_j \frac{p_0}{\sum_{j=0}^l p_j}, \dots, \sum_{j=0}^m p_j \frac{p_l}{\sum_{j=0}^l p_j}, p_{m+1}, \dots, p_k),$$

とすると、 $\hat{p}_i(\mathbf{x}, \mathbf{y})$ は

$$\hat{p}_i(\mathbf{x}, \mathbf{y}) = \begin{cases} \frac{N'_1 + N'_2}{N_1 + N_2} \frac{x_i + y_i}{N'_1 \left(1 + \frac{N'_2}{\sum_{i=0}^l x_i}\right)}, & \text{if } i \leq l, \\ \frac{N'_1 + N'_2}{N_1 + N_2} \frac{x_i}{N'_1}, & \text{if } l < i \leq m, \\ \frac{x_i + y_i}{N_1 + N_2}, & \text{if } i > m \end{cases} \quad (12)$$

である。但し、 $N'_1 = \sum_{i=0}^m x_i$, $N'_2 = \sum_{i=0}^l y_i$.

Proposition 4.2 l, k, m は $m < k$ なる自然数。確率変数 \mathbf{X}, \mathbf{Y} は互いに独立で

$$\mathbf{X} = (X_0, \dots, X_l, \dots, X_m, \dots, X_k) \sim \text{Multinomial}(N_1; p_0, \dots, p_k),$$

$$\mathbf{Y} = (Y_0, \dots, Y_l, Y_*)$$

$$\sim \text{Multinomial}(N_2; \sum_{j=0}^m p_j \frac{p_0}{\sum_{j=0}^l p_j}, \dots, \sum_{j=0}^m p_j \frac{p_l}{\sum_{j=0}^l p_j}, \sum_{i=m+1}^k p_i)$$

ならば、 $\hat{p}_i(\mathbf{x}, \mathbf{y})$ は

$$\hat{p}_i(\mathbf{x}, \mathbf{y}) = \begin{cases} \frac{N'_1 + N'_2}{N_1 + N_2} \frac{x_i + y_i}{N'_1 \left(1 + \frac{N'_2}{\sum_{i=0}^l x_i}\right)}, & \text{if } i \leq l, \\ \frac{N'_1 + N'_2}{N_1 + N_2} \frac{x_i}{N'_1}, & \text{if } l < i \leq m, \\ \frac{\sum_{i=m+1}^k x_i + y_*}{N_1 + N_2} \frac{x_i}{\sum_{i=m+1}^k x_i}, & \text{if } i > m \end{cases} \quad (13)$$

と書ける。

Proposition 4.3 l, k, m は $m < k$ なる自然数。確率変数 \mathbf{X}, \mathbf{Y} は互いに独立で

$$\mathbf{X} = (X_0, \dots, X_l, \dots, X_m, \dots, X_k) \sim \text{Multinomial}(N_1; p_0, \dots, p_k),$$

$$\mathbf{Y} = (Y_0, \dots, Y_l, Y_{l+1}, \dots, Y_m)$$

$$\sim \text{Multinomial}(N_2; \sum_{j=0}^m p_j \frac{p_0}{\sum_{j=0}^l p_j}, \dots, \sum_{j=0}^m p_j \frac{p_l}{\sum_{j=0}^l p_j}, \sum_{j=m+1}^k p_j \frac{p_{l+1}}{\sum_{j=l+1}^m p_j}, \dots, \sum_{j=m+1}^k p_j \frac{p_m}{\sum_{j=l+1}^m p_j})$$

ならば、 p_i は

$$\hat{p}_i(\mathbf{x}, \mathbf{y}) = \begin{cases} \frac{N'_1 + N'_2}{N_1 + N_2} \frac{x_i + y_i}{N'_1 \left(1 + \frac{N'_2}{\sum_{i=0}^l x_i}\right)}, & \text{if } i \leq l, \\ \frac{N'_1 + N'_2}{N_1 + N_2} \frac{x_i + y_i}{N'_1 \left(1 + \frac{\sum_{i=l+1}^m y_i}{\sum_{i=l+1}^m x_i}\right)}, & \text{if } l < i \leq m, \\ \left\{1 + \frac{\sum_{i=l+1}^m y_i}{\sum_{i=m+1}^k x_i}\right\} \frac{x_i}{N_1 + N_2}, & \text{if } i > m \end{cases} \quad (14)$$

である。但し、 $N'_1 = \sum_{i=0}^m x_i$, $N'_2 = \sum_{i=0}^l y_i$.

Proposition 4.1, 4.2 および 4.3 いずれも、パラメータ変数変換

$$t = p_0 + \cdots + p_m, \quad u = \frac{p_0 + \cdots + p_l}{p_0 + \cdots + p_l + \cdots + p_m}, \quad \theta_i = \frac{p_i}{p_0 + \cdots + p_l} \quad (i \leq l), \\ v_i = \frac{p_i}{p_{l+1} + \cdots + p_m} \quad (l < i \leq m), \quad w_i = \frac{p_i}{p_{m+1} + \cdots + p_k} \quad (m < i).$$

を用いれば容易。

5. パラメータが樹木構造でない場合

Proposition 1.1 ~ 1.3 は、パラメータが樹木構造で書けている。これに対して、Proposition 2.1 ~ 2.5 は、パラメータが「完全に樹木構造」とは言えないものの、MLE が exact に求められたものである。そこで、パラメータが樹木構造ではないために MLE が exact には求められないが、EM アルゴリズムを用いて MLE が求められるようなモデルについて議論する。

Problem $\{X_i\}_{i=1, \dots, k}$, $\{Y_i\}_{i=1, \dots, m}$, $\{Z_i\}_{i=l, \dots, k}$, および $\{W_i\}_{i=1, \dots, l, m+1, \dots, k}$ は、互いに独立で、 $X_i \sim \text{Poisson}(\lambda_i)$, $Y_i \sim \text{Poisson}(\mu_i)$, $Z_i \sim \text{Poisson}(\eta_i)$, $W_i \sim \text{Poisson}(\xi_i)$, とする。但し、

$$\mu_i = (\lambda_1 + \cdots + \lambda_m) \frac{\lambda_i}{\lambda_1 + \cdots + \lambda_l} \quad i = 1, \dots, l, \\ \eta_i = (\lambda_l + \cdots + \lambda_k) \frac{\lambda_i}{\lambda_{l+1} + \cdots + \lambda_m} \quad i = l+1, \dots, m, \\ \xi_i = (\lambda_1 + \cdots + \lambda_l + \lambda_{m+1} + \cdots + \lambda_k) \frac{\lambda_i}{\lambda_{m+1} + \cdots + \lambda_k} \quad i = m+1, \dots, k.$$

このとき、 λ_i の MLE を求めたい。

$$s_1 = \sum_{i=1}^l \lambda_i, \quad s_2 = \sum_{i=l+1}^m \lambda_i, \quad s_3 = \sum_{i=m+1}^k \lambda_i, \quad t_i = \frac{\lambda_i}{\sum_{i=1}^l \lambda_i} \quad (i \leq l) \\ u_i = \frac{\lambda_i}{\sum_{i=l+1}^m \lambda_i} \quad (l < i \leq m), \quad v_i = \frac{\lambda_i}{\sum_{i=m+1}^k \lambda_i} \quad (i \geq m)$$

なるパラメータ変換を用いる。上記の逆変換は

$$\begin{aligned} \lambda_i &= s_1 t_i, \quad (i \leq l), \quad \lambda_i = s_2 u_i, \quad (l < i \leq m), \quad \lambda_i = s_3 v_i, \quad (i \geq m), \\ \mu_i &= (s_1 + s_2) t_i, \quad (i \leq l), \quad \eta_i = (s_2 + s_3) u_i, \quad (l < i \leq m), \quad \xi_i = (s_1 + s_3) v_i, \quad (i \geq m). \end{aligned}$$

これを用いて

$$\begin{aligned} L &= \left(\prod_{i=1}^k \frac{\lambda_i^{x_i}}{x_i!} \exp(-\lambda_i) \right) \left(\prod_{i=1}^l \frac{\mu_i^{y_i}}{y_i!} \exp(-\mu_i) \right) \left(\prod_{i=l+1}^m \frac{\eta_i^{z_i}}{z_i!} \exp(-\eta_i) \right) \left(\prod_{i=m+1}^k \frac{\xi_i^{w_i}}{w_i!} \exp(-\xi_i) \right) \\ &\propto s_1^{\sum_{i=1}^l x_i} s_2^{\sum_{i=l+1}^m x_i} s_3^{\sum_{i=m+1}^k x_i} (s_1 + s_2)^{\sum_{i=1}^l y_i} (s_2 + s_3)^{\sum_{i=l+1}^m z_i} (s_1 + s_3)^{\sum_{i=m+1}^k w_i} \\ &\quad \times \left(\prod_{i=1}^l t_i^{x_i+y_i} \right) \left(\prod_{i=l+1}^m u_i^{x_i+z_i} \right) \left(\prod_{i=m+1}^k v_i^{x_i+w_i} \right) \exp(-3(s_1 + s_2 + s_3)). \end{aligned}$$

これより直ちに

$$\hat{t}_i = \frac{x_i + y_i}{\sum_{i=1}^l (x_i + y_i)}, \quad \hat{u}_i = \frac{x_i + z_i}{\sum_{i=l+1}^m (x_i + z_i)}, \quad \hat{v}_i = \frac{x_i + w_i}{\sum_{i=m+1}^k (x_i + w_i)}$$

を得る。さて、 $\hat{s}_1, \hat{s}_2, \hat{s}_3$ は、次のようにして求められる：

$$s = s_1 + s_2 + s_3, \quad \theta_i = \frac{s_i}{s_1 + s_2 + s_3} \quad \text{と変数変換すると、} \quad s_1 = s\theta_1, \quad s_2 = s\theta_2, \quad s_3 = s\theta_3.$$

そして、表記の簡略化のため

$$T(x) = \sum_{i=1}^k x_i, \quad T(y) = \sum_{i=1}^l y_i, \quad T(z) = \sum_{i=l+1}^m z_i, \quad T(w) = \sum_{i=m+1}^k w_i$$

と定めると、

$$\begin{aligned} L &\propto (s\theta_1)^{\sum_{i=1}^l x_i} (s\theta_2)^{\sum_{i=l+1}^m x_i} (s\theta_3)^{\sum_{i=m+1}^k x_i} (s(\theta_1 + \theta_2))^{T(y)} (s(\theta_2 + \theta_3))^{T(z)} \\ &\quad \times (s(\theta_1 + \theta_3))^{T(w)} \exp(-3s) \\ &= s^{T(x)} \exp(-3s) \times \theta_1^{\sum_{i=1}^l x_i} \theta_2^{\sum_{i=l+1}^m x_i} \theta_3^{\sum_{i=m+1}^k x_i} (\theta_1 + \theta_2)^{T(y)} (\theta_2 + \theta_3)^{T(z)} (\theta_1 + \theta_3)^{T(w)}. \end{aligned}$$

これより $\hat{s} = \frac{T(x) + T(y) + T(z) + T(w)}{3}$. さて、

$$L(\theta_1, \theta_2, \theta_3) \propto \theta_1^{\sum_{i=1}^l x_i} \theta_2^{\sum_{i=l+1}^m x_i} \theta_3^{\sum_{i=m+1}^k x_i} \times (\theta_1 + \theta_2)^{T(y)} (\theta_2 + \theta_3)^{T(z)} (\theta_1 + \theta_3)^{T(w)}$$

であるから、 $\hat{\theta}_i$, ($i = 1, 2, 3$) を EM アルゴリズムによって次のように求めることができる。最初に、初期値 $\theta_i^{(0)}$ ($i = 1, 2, 3$) を

$$\begin{aligned} \theta_1^{(0)} &= \frac{\sum_{i=1}^l x_i + T(y)/2 + T(w)/2}{T_x + T_y + T_z + T_w}, \quad \theta_2^{(0)} = \frac{\sum_{i=l+1}^m x_i + T(y)/2 + T(z)/2}{T_x + T_y + T_z + T_w}, \\ \theta_3^{(0)} &= \frac{\sum_{i=m+1}^k x_i + T(z)/2 + T(w)/2}{T_x + T_y + T_z + T_w} \end{aligned}$$

に選ぶ。各 $k = 0, 1, \dots$, に対して、 $T^{(k)}(y)[1]$, $T^{(k)}(y)[2]$, $T^{(k)}(z)[2]$, $T^{(k)}(z)[3]$, $T^{(k)}(w)[1]$, $T^{(k)}(w)[3]$, (E-step) および $\theta_1^{(k+1)}$, $\theta_2^{(k+1)}$, $\theta_3^{(k+1)}$, (M-step) を次のように構築する。

E-step : 完全データの条件付期待値を計算するステップ

$$T^{(k)}(y)[i] = E(T(y)[i] | T(y); \theta_1^{(k)}, \theta_2^{(k)}) = T(y) \frac{\theta_i^{(k)}}{\theta_1^{(k)} + \theta_2^{(k)}}, \quad i = 1, 2,$$

$$T^{(k)}(z)[i] = E(T(z)[i] | T(z); \theta_2^{(k)}, \theta_3^{(k)}) = T(z) \frac{\theta_i^{(k)}}{\theta_2^{(k)} + \theta_3^{(k)}}, \quad i = 2, 3,$$

$$T^{(k)}(w)[i] = E(T(w)[1] | T(w); \theta_1^{(k)}, \theta_3^{(k)}) = T(w) \frac{\theta_i^{(k)}}{\theta_1^{(k)} + \theta_3^{(k)}}, \quad i = 1, 3.$$

M-step : 完全データ (の近似値) を使って、パラメータの値を更新するステップ。すなわち、データ $\sum_{i=1}^l x_i + T^{(k)}(y)[1] + T^{(k)}(w)[1]$, $\sum_{i=l+1}^m x_i + T^{(k)}(y)[2] + T^{(k)}(z)[2]$, $\sum_{i=m+1}^k x_i + T^{(k)}(z)[3] + T^{(k)}(w)[3]$ が与えられたときの尤度を最大にする $\theta_1^{(k+1)}$, $\theta_2^{(k+1)}$, $\theta_3^{(k+1)}$ の値は

$$\theta_1^{(k+1)} = \frac{\sum_{i=1}^l x_i + T^{(k)}(y)[1] + T^{(k)}(w)[1]}{T(x) + T(y) + T(z) + T(w)},$$

$$\theta_2^{(k+1)} = \frac{\sum_{i=l+1}^m x_i + T^{(k)}(y)[2] + T^{(k)}(z)[2]}{T(x) + T(y) + T(z) + T(w)},$$

$$\theta_3^{(k+1)} = \frac{\sum_{i=m+1}^k x_i + T^{(k)}(z)[3] + T^{(k)}(w)[3]}{T(x) + T(y) + T(z) + T(w)}$$

である。各 $i = 1, 2, 3$ に対して、EM アルゴリズムにより $\lim_{k \rightarrow \infty} \theta_i^{(k)} = \hat{\theta}_i$ である。MLE の invariance property により、

$$\hat{\lambda}_i = \hat{s}_1 \hat{\theta}_i \hat{t}_i \quad (i \leq l), \quad \hat{\lambda}_i = \hat{s}_2 \hat{\theta}_i \hat{u}_i \quad (l < i \leq m), \quad \hat{\lambda}_i = \hat{s}_3 \hat{\theta}_i \hat{v}_i \quad (i > m).$$

を得る。下記の表は、E-step, M-step における 観測された値および完全データと確率との対応を示したものである

確率	観測された値
θ_1	$\sum_{i=1}^l x_i$
$\theta_1 + \theta_2$	$T(y)$
θ_2	$\sum_{i=l+1}^m x_i$
$\theta_2 + \theta_3$	$T(z)$
θ_3	$\sum_{i=m+1}^k x_i$
$\theta_3 + \theta_1$	$T(w)$

セル確率	完全データ
θ_1	$\sum_{i=1}^l x_i + T(y)[1] + T(w)[1]$
θ_2	$\sum_{i=l+1}^m x_i + T(y)[2] + T(z)[2]$
θ_3	$\sum_{i=m+1}^k x_i + T(z)[3] + T(w)[3]$

Acknowledgements

本研究は、科学研究費 基盤研究 (C) 課題番号 20500261 の研究成果の一部である。この科学研究費の援助を、記して感謝申し上げます。

References

- [1] Asano, C. (1965) On estimating multinomial probabilities by pooling incomplete samples, *Annals of the Institute of Statistical Mathematics* **17**, 1-13.
- [2] Dempster, A.P., Laird, N.M., and Rubin, D.B. (1977). Maximum likelihood from incomplete data via the EM algorithm (with discussion). *Journal of the Royal Statistical Society* **B 39**, 1-38.
- [3] Johnson, N.L., Kotz, S. and Balakrishnan, N. (1997). *Discrete Multivariate Distributions*. Wiley.
- [4] McLachlan, G. J., and Krishnan, T. ((1977) *The EM Algorithm and Extensions*. Wiley.