

R-project 3

EM algorithm

Let X_1, \dots, X_n and Y_1, \dots, Y_n are all mutually independent random variables, where $Y_i \sim \text{Poisson}(\beta\tau_i)$ and $X_i \sim \text{Poisson}(\tau_i)$. Suppose that Y_i models the incidence of a disease, where the underlying rate is a function of an overall effect β and an additional factor τ_i , which measures population density in area i . We do not see τ_i but get information on it through X_i .

The complete data likelihood is

$$\mathcal{L}(\beta, \tau_1, \dots, \tau_n; \mathbf{x}^{obs}, \mathbf{y}) \propto f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}; \beta, \tau_1, \dots, \tau_n) = \prod_{i=1}^n \frac{e^{-\beta\tau_i} (\beta\tau_i)^{y_i}}{y_i!} \frac{e^{-\tau_i} (\tau_i)^{x_i}}{x_i!}$$

The likelihood estimators which can be found by straightforward differentiation, are

$$\hat{\beta} = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i} \quad \text{and} \quad \hat{\tau}_i = \frac{x_i + y_i}{\hat{\beta} + 1} \quad i = 1, \dots, n$$

Suppose that value x_1 was missing. The observed-data likelihood with x_1 missing is

$$\mathcal{L}_{obs}(\beta, \tau_1, \dots, \tau_n; \mathbf{x}^{obs}, \mathbf{y}) = \sum_{x_1=0}^{\infty} f_{\mathbf{X}, \mathbf{Y}}(\mathbf{x}, \mathbf{y}; \beta, \tau_1, \dots, \tau_n)$$

with $(\mathbf{x}^{obs}, \mathbf{y}) = (y_1, (x_2, y_2), \dots, (x_n, y_n))$. Differentiation leads to the MLE equations:

$$\hat{\beta} = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n \hat{\tau}_i}$$

$$y_1 = \hat{\tau}_1 \hat{\beta}$$

$$x_i + y_i = \hat{\tau}_i (\hat{\beta} + 1) \quad i = 2, \dots, n$$

which we can solve with the EM algorithm using the expected complete data log likelihood:

$$\begin{aligned} Q(\boldsymbol{\theta} | \tilde{\boldsymbol{\theta}}) = & \sum_{i=1}^n [-\beta\tau_i + y_i (\log(\beta) + \log(\tau_i))] + \sum_{i=2}^n [-\tau_i + x_i \log(\tau_i)] + \sum_{x_1=0}^{\infty} [-\tau_1 + x_1 \log(\tau_1)] \frac{e^{-\tilde{\tau}_1} (\tilde{\tau}_1)^{x_1}}{x_1!} \\ & - \left(\sum_{i=1}^n \log(y_i!) + \sum_{i=2}^n \log(x_i!) + \sum_{x_1=0}^{\infty} \log(x_1!) \frac{e^{-\tilde{\tau}_1} (\tilde{\tau}_1)^{x_1}}{x_1!} \right) \end{aligned}$$

where $\boldsymbol{\theta} = (\beta, \boldsymbol{\tau}) = (\beta, \tau_1, \dots, \tau_n)$, and

$$\arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}^{(j)}) = \begin{cases} \hat{\beta}^{j+1} &= \frac{\sum_{i=1}^n y_i}{\hat{\tau}_1^{(j)} + \sum_{i=2}^n x_i} \\ \hat{\tau}_1^{j+1} &= \frac{\hat{\tau}_1^{(j)} + y_1}{\hat{\beta}^{(j+1)} + 1} \\ \hat{\tau}_i^{j+1} &= \frac{x_i + y_i}{\hat{\beta}^{(j+1)} + 1} \quad i = 2, \dots, n \end{cases}$$

Consider the following data on $n = 18$ leukemia counts and the associated populations for a number of areas in New York State (Lange et al. 1994)

Population	Number of cases	Population	Number of cases
?	3	948	0
3560	4	1172	1
3739	1	1047	3
2784	1	3138	5
2571	3	5485	4
2729	1	5554	6
3952	2	2943	2
993	0	4969	5
1908	2	4828	4

1. Develop an EM algorithm to fit the Poisson model to these data.
2. A direct solution of the original incomplete-data likelihood equations is possible and is given by

$$\hat{\beta} = \frac{\sum_{i=1}^n y_i}{\sum_{i=2}^n x_i}$$

$$\hat{\tau}_1 = \frac{y_1}{\hat{\beta}}$$

$$\hat{\tau}_i = \frac{x_i + y_i}{\hat{\beta} + 1} \quad i = 2, \dots, n$$

Compare the direct solution with the EM solution.