

Integrating the likelihood function of the somatic mutation latent variable model

Louis Dijkstra*

November 10, 2014

The goal is to compute/estimate the integral ¹

$$\int_a^b L(\theta_h, \theta_c \mid Z^h, Z^c) d\theta_c \quad (\theta_h \text{ is fixed}) \quad (0.1)$$

where

$$L(\theta_h, \theta_c \mid Z^h, Z^c) = \prod_{i=1}^k g_i^h(Z_i^h \mid \theta_h) \times \prod_{j=1}^l g_j^c(Z_j^c \mid \theta_h, \theta_c; \alpha). \quad (0.2)$$

Note that we can rewrite the integral to

$$\underbrace{\prod_{i=1}^k g_i^h(Z_i^h \mid \theta_h)}_{\text{constant}} \times \int_a^b \prod_{j=1}^l g_j^c(Z_j^c \mid \theta_h, \theta_c; \alpha) d\theta_c. \quad (0.3)$$

The conditional probability distribution for observation $g_j^c(Z_j^c \mid \theta_h, \theta_c; \alpha)$ can be written as:

$$g_j^c(Z_j^c \mid \theta_h, \theta_c; \alpha) = s_j \theta_c + t_j \quad (0.4)$$

where

$$s_j = (1 - \alpha) [p_j^c(Z_j) - a_j^c(Z_j)] \quad (0.5)$$

and

$$t_j = \pi_j^c \alpha [\theta_h p_j^c(Z_j) + (1 - \theta_h) a_j^c(Z_j)] + (1 - \alpha) a_j^c(Z_j) + (1 - \pi_j^c) u_j^c(Z_j). \quad (0.6)$$

The product $\prod_{j=1}^l g_j^c(Z_j^c \mid \theta_h, \theta_c; \alpha)$ can, therefore, be written as m -order polynomial

$$\prod_{j=1}^l (s_j \theta_c + t_j) = \sum_{m=0}^k \beta_m \theta_c^m \quad (0.7)$$

with coefficients

$$\beta_m = \sum_{\substack{I \subset \{1, \dots, l\} \\ |I|=m}} \prod_{i \in I} s_i \prod_{j \notin I} t_j \quad (0.8)$$

*E-mail: dijkstra@cwi.nl

¹The prior $h(\theta_h, \theta_c)$ is initially assumed uniform over the parameter space.

The integral in (0.3) can be written as

$$\int_a^b \sum_{m=0}^k \beta_m \theta_c^m d\theta_h = \sum_{m=0}^k \frac{\beta_m}{m+1} (b^{m+1} - a^{m+1}). \quad (0.9)$$

What remains is determining the coefficients β_0, \dots, β_m efficiently. Computing them directly from eq. (0.8) is infeasible, since it requires us to sum over all possible subsets of observations. We propose the following iterative process:

Initialization $\beta_0^{(0)} := 1$

Update steps add every observation $i = 1, 2, \dots, l$ one by one and ‘update’ the coefficients with the following rules:

$$\begin{aligned} \beta_0^{(i)} &= \beta_0^{(i-1)} t_i \\ \beta_m^{(i)} &= \beta_m^{(i-1)} t_i + \beta_{m-1}^{(i-1)} s_i \quad (m = 1, 2, \dots, i-1) \\ \beta_i^{(i)} &= \beta_{i-1}^{(i-1)} s_i. \end{aligned} \quad (0.10)$$

Termination the coefficients β_1, \dots, β_m are

$$\beta_m = \beta_m^{(l)} \text{ for } m = 1, 2, \dots, l. \quad (0.11)$$

Arbitrary prior distribution In case of an arbitrary prior, the integral

$$\int_a^b L(\theta_h, \theta_c \mid Z^h, Z^c) h(\theta_h, \theta_c) d\theta_c \quad (\theta_h \text{ is fixed}) \quad (0.12)$$

is equal to

$$\prod_{i=1}^k g_i^h(Z_i^h \mid \theta_h) \times \sum_{m=0}^l \beta_m \int_a^b \theta_c^m h(\theta_h, \theta_c) d\theta_c \quad (0.13)$$

where the coefficient β_m can be computed as before.