

## Appendix 1 – Calculation of $P_{ik} = E(\mathbf{b}_{ik}|\mathbf{y}, \widehat{\mathbf{Pr}}_k)$

In the expectation step of the EM algorithm we require the  $E_{b|y}$  of equation 6b. This requires the  $E(\mathbf{b}_{ik}|\mathbf{y}, \widehat{\mathbf{Pr}}_k)$  which is derived in this appendix.

The model is  $\mathbf{y} = \mathbf{1}_n\mu + \mathbf{Z}_i\mathbf{g}_i + \mathbf{u} + \mathbf{e}$ ,

Then,

$$\begin{aligned} E_{\mathbf{u}}(\mathbf{b}_{ik}|\mathbf{y}, \widehat{\mathbf{Pr}}_k) &= p(\mathbf{b}_{ik} = 1|\mathbf{y}, \widehat{\mathbf{Pr}}_k) \\ &= \frac{p(\mathbf{y}|\mathbf{b}_{ik}=1) \times p(\mathbf{b}_{ik}=1|\widehat{\mathbf{Pr}}_k)}{p(\mathbf{y})} \\ &\propto p(\mathbf{y}|\mathbf{b}_{ik} = 1) \times p(\mathbf{b}_{ik} = 1|\widehat{\mathbf{Pr}}_k) \end{aligned} \quad (A1)$$

where,

$p(\mathbf{b}_{ik} = 1|\widehat{\mathbf{Pr}}_k) = \widehat{\mathbf{Pr}}_k$ , and

$p(\mathbf{y}|\mathbf{b}_{ik} = 1) = \frac{1}{\sqrt{|\mathbf{W}_k|}} \exp(-\frac{1}{2}(\mathbf{y} - \mathbf{1}_n\mu - \mathbf{u})'\mathbf{W}_k^{-1}(\mathbf{y} - \mathbf{1}_n\mu - \mathbf{u}))$ , so

$\log p(\mathbf{y}|\mathbf{b}_{ik} = 1) = -0.5 (\log|\mathbf{W}_k| + (\mathbf{y} - \mathbf{1}_n\mu - \mathbf{u})'\mathbf{W}_k^{-1}(\mathbf{y} - \mathbf{1}_n\mu - \mathbf{u}))$

based on  $(\mathbf{y} - \mathbf{1}_n\mu - \mathbf{u})|(\mathbf{b}_{ik} = 1) \sim N(0, \mathbf{W}_k)$ , and  $\mathbf{W}_k = \mathbf{Z}_i\mathbf{Z}_i'\sigma_k^2 + \mathbf{I}\sigma_e^2$ .

Therefore,

$$\log l_{ik} = \log p(\mathbf{b}_{ik} = 1|\mathbf{y}, \widehat{\mathbf{Pr}}_k) = \log p(\mathbf{y}|\mathbf{b}_{ik} = 1) + \log p(\mathbf{b}_{ik} = 1|\widehat{\mathbf{Pr}}_k) + \text{constant}$$

*constant* appear on both denominator term and numerator term of equation (A7), and therefore could be ignored.

The expression above for  $\log p(\mathbf{y}, |\mathbf{b}_{ik} = 1)$  involves the unknown  $\mathbf{u}$ . Therefore, we take the expectation over  $\mathbf{u}|\mathbf{y}$ . That is,

$$\log p(\mathbf{y}|\mathbf{b}_{ik} = 1) = -0.5 E_{\mathbf{u}|\mathbf{y}}\{\log|\mathbf{W}_k| + (\mathbf{y} - \mathbf{1}_n\mu - \mathbf{u})'\mathbf{W}_k^{-1}(\mathbf{y} - \mathbf{1}_n\mu - \mathbf{u})\}$$

Only the quadratic form  $Q = (\mathbf{y} - \mathbf{1}_n\mu - \mathbf{u})'\mathbf{W}_k^{-1}(\mathbf{y} - \mathbf{1}_n\mu - \mathbf{u})$  of  $\log p(\mathbf{y}|\mathbf{b}_{ik} = 1)$  involves  $\mathbf{u}$ . Therefore, apply Searle's expectation rule[32] for  $Q$  as follows:

$$E_{\hat{\mathbf{u}}}Q = (\mathbf{y} - \mathbf{1}_n\mu - \hat{\mathbf{u}})'\mathbf{W}_k^{-1}(\mathbf{y} - \mathbf{1}_n\mu - \hat{\mathbf{u}}) + \text{tr}(\mathbf{W}_k^{-1}PEV((\hat{\mathbf{u}})))$$

Hence,

$$\begin{aligned} \log p(\mathbf{y}, |\mathbf{b}_{ik} = 1) &= -0.5\{\log |\mathbf{W}_k| + E_{\hat{\mathbf{u}}}Q\} \\ &= -0.5\{\log |\mathbf{W}_k| + \mathbf{y}^\dagger \mathbf{W}_k^{-1} \mathbf{y}^\dagger + \text{tr}(\mathbf{W}_k^{-1}PEV((\hat{\mathbf{u}})))\} \end{aligned}$$

where,  $\mathbf{y}^\dagger = (\mathbf{y} - \mathbf{1}_n\mu - \hat{\mathbf{u}})$ .

Although  $\mathbf{W}_k$  is an  $n \times n$  matrix. the calculation of  $\log |\mathbf{W}_k|$  and  $\mathbf{W}_k^{-1}$  can be simplified by using the Woodbury identity so that

$$\mathbf{W}_k^{-1} = (\mathbf{Z}_i\mathbf{Z}_i'\sigma_k^2 + \mathbf{I}\sigma_e^2)^{-1} = \sigma_e^{-2} \left( \mathbf{I} - \frac{\mathbf{Z}_i\mathbf{Z}_i'\sigma_k^2}{\sigma_k^2\mathbf{Z}_i'\mathbf{Z}_i + \sigma_e^2} \right) \quad (\text{A2})$$

$$|\mathbf{W}_k| = \sigma_e^{(2n-2)} (\sigma_k^2\mathbf{Z}_i'\mathbf{Z}_i + \sigma_e^2), \text{ so}$$

$$\log |\mathbf{W}_k| = (2n - 2)\log \sigma_e^2 + \log(\sigma_k^2\mathbf{Z}_i'\mathbf{Z}_i + \sigma_e^2) \quad (\text{A3})$$

Such transformation could transfer the inverse calculation of a large matrix  $\mathbf{W}_k$  to the multiplication of the vectors, which could reduce the cost for matrix calculation.

Therefore, substitute (A3) and (A4) into  $\log p(\mathbf{y}|\mathbf{b}_{ik}, \hat{\mathbf{u}})$  as follow:

$$\begin{aligned} \log p(\mathbf{y}|\mathbf{b}_{ik} = 1) &= -0.5\{(n - 1)\log \sigma_e^2 + \log(\sigma_k^2\mathbf{Z}_i'\mathbf{Z}_i + \sigma_e^2)\} \\ &\quad -0.5\{(\mathbf{y}^*'\mathbf{y}^*)\sigma_e^{-2} - (\mathbf{y}^*\mathbf{Z}_i)'\sigma_k^2\sigma_e^{-2}/(\sigma_k^2\mathbf{Z}_i'\mathbf{Z}_i + \sigma_e^2)\} \\ &\quad -0.5\{\text{tr}(PEV(\hat{\mathbf{u}}))\sigma_e^{-2} - \text{tr}(\mathbf{Z}_i\mathbf{Z}_i'PEV(\hat{\mathbf{u}}))\sigma_k^2\sigma_e^{-2}/(\sigma_k^2\mathbf{Z}_i'\mathbf{Z}_i + \sigma_e^2)\} \end{aligned} \quad (\text{A4})$$

Then,

$$\begin{aligned} \log l_{ik} &= \log p(\mathbf{y}|\mathbf{b}_{ik} = 1) + \log p(\mathbf{b}_{ik} = 1|\widehat{\mathbf{P}}\mathbf{r}_k) \\ &= -0.5\{2(n - 1)\log \sigma_e^2 + \log V\} \end{aligned}$$

$$\begin{aligned}
& -0.5 \left\{ (\mathbf{y}^\dagger \mathbf{y}^\dagger) \sigma_e^{-2} - (\mathbf{y}^\dagger \mathbf{Z}_i)^2 \sigma_k^2 \sigma_e^{-2} / V \right\} \\
& -0.5 \left\{ \text{tr}(\text{PEV}(\hat{\mathbf{u}})) \sigma_e^{-2} - \text{tr}(\mathbf{Z}_i \mathbf{Z}_i' \text{PEV}(\hat{\mathbf{u}})) \sigma_k^2 \sigma_e^{-2} / V \right\} \\
& + \log \text{Pr}_k
\end{aligned} \tag{A5}$$

where,  $\mathbf{y}^\dagger = \mathbf{y} - \mathbf{1}_n \mu - \hat{\mathbf{u}}$ ,  $V = \sigma_k^2 \mathbf{Z}_i' \mathbf{Z}_i + \sigma_e^2$  and  $n$  is the number of animals.  $\text{PEV}(\hat{\mathbf{u}})$  ( $n \times n$  symmetric matrix) could be approximated by  $\text{PEV}(\hat{\mathbf{u}}^*)$  as derived in appendix 2 and could be calculated based on GBLUP, outside the iterations of EM algorithm.  $\text{tr}(\mathbf{Z}_i \mathbf{Z}_i' \text{PEV}(\hat{\mathbf{u}}))$  means to add up the diagonal elements of symmetric matrix. In other words, we just need to calculate and then add up the diagonal elements of the multiplication of  $\mathbf{Z}_i \mathbf{Z}_i'$  (also a  $n \times n$  symmetric matrix) and  $\text{PEV}(\hat{\mathbf{u}})$ . Because  $\text{tr}(\mathbf{Z}_i \mathbf{Z}_i' \text{PEV}(\hat{\mathbf{u}}))$  and  $\text{tr}(\text{PEV}(\hat{\mathbf{u}}))$  does not change each iterations, they could be calculated once and stored in front of the EM steps.

With the expression for  $\log l_{ik} = \log p(b_{ik} = 1 | \mathbf{y}, \hat{\mathbf{P}}_{r_k})$ , we can now calculate the probability that each SNP is in one of four normal distributions:

$$P_{ik} = E_{\mathbf{u}}(b_{ik} | \mathbf{y}, \hat{\mathbf{P}}_{r_k}) = \frac{\exp(\log l_{ik})}{\sum_{k=1}^4 \exp(\log l_{ik})} \tag{A6}$$