

Supplemental material

Equivalence between Standard Selection Index and BLUP

Consider a standard single-trait model of the form

$$\mathbf{y} = \mathbf{u} + \boldsymbol{\varepsilon}$$

where $\mathbf{y} = (y_1, \dots, y_n)'$, $\mathbf{u} = (u_1, \dots, u_n)'$, and $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$ are vectors of phenotypes, genetic, and environmental effects, respectively. Here, for simplicity we assume that all these vectors have zero-mean.

In a standard G-BLUP model, \mathbf{u} and $\boldsymbol{\varepsilon}$ are assumed to be independent (i.e., $\text{cov}(\mathbf{u}, \boldsymbol{\varepsilon}') = \mathbf{0}$), both have null means (i.e., $\mathbb{E}(\mathbf{u}) = \mathbb{E}(\boldsymbol{\varepsilon}) = \mathbf{0}$), and (co)variance matrices $\text{var}(\mathbf{u}) = \sigma_u^2 \mathbf{G}$ and $\text{var}(\boldsymbol{\varepsilon}) = \sigma_\varepsilon^2 \mathbf{I}$, respectively; here \mathbf{G} is a relationship matrix that could be derived from a pedigree or from DNA sequences.

Consider now a partition of each of the data in into a training (trn) and a testing (tst) set. The objective is to predict the genetic values of the individuals in the testing set (\mathbf{u}_{tst}) using the phenotype data available from the training set (\mathbf{y}_{trn}). The (co)variance matrix of the vector of breeding values can be partitioned as follows

$$\text{var} \left(\begin{bmatrix} \mathbf{u}_{trn} \\ \mathbf{u}_{tst} \end{bmatrix} \right) = \sigma_u^2 \begin{bmatrix} \mathbf{G}_{trn} & \mathbf{G}_{trn,tst} \\ \mathbf{G}_{trn,tst}' & \mathbf{G}_{tst} \end{bmatrix}$$

where \mathbf{G}_{trn} and \mathbf{G}_{tst} are the genetic relationship submatrices for the training and testing data points, respectively, and $\mathbf{G}_{trn,tst}$ is the genetic relationship submatrix between training and testing subjects. The Best Linear Predictor (BLP) of \mathbf{u}_{tst} ($\hat{\mathbf{u}}_{tst}$) takes the form (e.g., Searle *et al.* 1992):

$$\begin{aligned} \mathbb{E}(\mathbf{u}_{tst} | \mathbf{y}_{trn}) &= \mathbb{E}(\mathbf{u}_{tst}) + \text{cov}(\mathbf{u}_{tst}, \mathbf{y}_{trn}') [\text{var}(\mathbf{y}_{trn})]^{-1} (\mathbf{y}_{trn} - \mathbb{E}(\mathbf{y}_{trn})) \\ &= \mathbf{G}_{trn,tst}' (\mathbf{G}_{trn} + \lambda_0 \mathbf{I})^{-1} \mathbf{y}_{trn}. \end{aligned}$$

Alternatively, one can write $\hat{\mathbf{u}}_{tst} = \mathbf{H} \cdot \mathbf{y}_{trn}$, where $\mathbf{H} = \mathbf{G}_{trn,tst}' (\mathbf{G}_{trn} + \lambda_0 \mathbf{I})^{-1}$ is a ‘‘Hat’’ matrix. Thus, the BLUP of the genetic value of the i^{th} testing individual is $\hat{u}_{tst(i)} = \mathbf{H}_i' \mathbf{y}_{trn}$ where \mathbf{H}_i' is the i^{th} row of \mathbf{H} , that is $\mathbf{H}_i' = \mathbf{G}_i' (\mathbf{G}_{trn} + \lambda_0 \mathbf{I})^{-1}$ which is equal to the weights of the standard selection index, $\hat{\beta}_i' = \mathbf{G}_i' (\mathbf{G}_{trn} + \lambda_0 \mathbf{I})^{-1}$ (see Equation 2 in the manuscript).

References

Searle S. R., G. Casella, and C. E. McCulloch, 1992 *Variance components*. John Wiley & Sons, Inc.