

**S1 Text. Supplementary material for the manuscript
“On the accuracy of genomic selection”**

1. Proof of the main result

In what follows, the quantities \mathbf{X} and \mathbf{Q} are supposed to be known. In other words, all the results will be conditional on \mathbf{X} and \mathbf{Q} . Recall that $\boldsymbol{\theta}$ is fixed, and also that $\mathbf{q}_{n_{\text{TRN}}+1}$ and $\mathbf{x}_{n_{\text{TRN}}+1}$ are considered random.

Using the causal model, we have

$$\begin{aligned} \text{Cov}\left(\hat{Y}_{n_{\text{TRN}}+1}, Y_{n_{\text{TRN}}+1}\right) &= \text{Cov}\left(\mathbf{x}'_{n_{\text{TRN}}+1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{Y}, \mathbf{q}'_{n_{\text{TRN}}+1}\boldsymbol{\theta} + e_{n_{\text{TRN}}+1}\right) \\ &= \text{Cov}\left(\mathbf{x}'_{n_{\text{TRN}}+1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{Q}\boldsymbol{\theta}, \mathbf{q}'_{n_{\text{TRN}}+1}\boldsymbol{\theta}\right). \end{aligned}$$

Besides, since $\mathbf{x}_{n_{\text{TRN}}+1}$ and $\mathbf{q}_{n_{\text{TRN}}+1}$ are centered, we have

$$\text{Cov}\left(\hat{Y}_{n_{\text{TRN}}+1}, Y_{n_{\text{TRN}}+1}\right) = \mathbb{E}\left(\mathbf{x}'_{n_{\text{TRN}}+1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{Q}\boldsymbol{\theta}\mathbf{q}'_{n_{\text{TRN}}+1}\boldsymbol{\theta}\right). \quad (1)$$

Let us now focus on the terms present at the denominator, that is to say $\text{Var}(Y_{n_{\text{TRN}}+1})$ and $\text{Var}\left(\hat{Y}_{n_{\text{TRN}}+1}\right)$. By definition,

$$\text{Var}(Y_{n_{\text{TRN}}+1}) = \sigma_G^2 + \sigma_e^2 \quad \text{where} \quad \sigma_G^2 = \boldsymbol{\theta}'\text{Var}(\mathbf{q}_{n_{\text{TRN}}+1})\boldsymbol{\theta}. \quad (2)$$

Besides, we have the relationship

$$\text{Var}\left(\hat{Y}_{n_{\text{TRN}}+1}\right) = \mathbb{E}\left(\text{Var}\left(\hat{Y}_{n_{\text{TRN}}+1} \mid \mathbf{x}_{n_{\text{TRN}}+1}\right)\right) + \text{Var}\left(\mathbb{E}\left(\hat{Y}_{n_{\text{TRN}}+1} \mid \mathbf{x}_{n_{\text{TRN}}+1}\right)\right). \quad (3)$$

To begin with, let us compute the quantity $\mathbb{E}\left(\text{Var}\left(\hat{Y}_{n_{\text{TRN}}+1} \mid \mathbf{x}_{n_{\text{TRN}}+1}\right)\right)$. We have

$$\text{Var}\left(\hat{Y}_{n_{\text{TRN}}+1} \mid \mathbf{x}_{n_{\text{TRN}}+1}\right) = \mathbb{E}\left(\hat{Y}_{n_{\text{TRN}}+1}^2 \mid \mathbf{x}_{n_{\text{TRN}}+1}\right) - \left(\mathbf{x}'_{n_{\text{TRN}}+1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{Q}\boldsymbol{\theta}\right)^2$$

and

$$\begin{aligned} \mathbb{E}\left(\hat{Y}_{n_{\text{TRN}}+1}^2 \mid \mathbf{x}_{n_{\text{TRN}}+1}\right) &= \mathbb{E}\left(\left(\mathbf{x}'_{n_{\text{TRN}}+1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{Q}\boldsymbol{\theta} + \mathbf{x}'_{n_{\text{TRN}}+1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{e}\right)^2 \mid \mathbf{x}_{n_{\text{TRN}}+1}\right) \\ &= \mathbb{E}\left(\left(\mathbf{x}'_{n_{\text{TRN}}+1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{Q}\boldsymbol{\theta}\right)^2 \mid \mathbf{x}_{n_{\text{TRN}}+1}\right) + \mathbb{E}\left(\left(\mathbf{x}'_{n_{\text{TRN}}+1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{e}\right)^2 \mid \mathbf{x}_{n_{\text{TRN}}+1}\right) \\ &= \left(\mathbf{x}'_{n_{\text{TRN}}+1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{Q}\boldsymbol{\theta}\right)^2 + \sigma_e^2 \left\|\mathbf{x}'_{n_{\text{TRN}}+1}\mathbf{X}'\mathbf{V}^{-1}\right\|^2 \end{aligned}$$

where $\|\cdot\|$ is the L^2 norm. As a result,

$$\mathbb{E}\left(\text{Var}\left(\hat{Y}_{n_{\text{TRN}}+1} \mid \mathbf{x}_{n_{\text{TRN}}+1}\right)\right) = \sigma_e^2 \mathbb{E}\left(\left\|\mathbf{x}'_{n_{\text{TRN}}+1}\mathbf{X}'\mathbf{V}^{-1}\right\|^2\right).$$

On the other hand, the second term in formula (3) is

$$\begin{aligned}\text{Var}\left(\mathbb{E}\left(\hat{Y}_{n_{\text{TRN}}+1} \mid \mathbf{x}_{n_{\text{TRN}}+1}\right)\right) &= \text{Var}\left(\mathbf{x}'_{n_{\text{TRN}}+1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Q} \boldsymbol{\theta}\right) \\ &= \boldsymbol{\theta}' \mathbf{Q}' \mathbf{V}^{-1} \mathbf{X} \text{Var}\left(\mathbf{x}_{n_{\text{TRN}}+1}\right) \mathbf{X}' \mathbf{V}^{-1} \mathbf{Q} \boldsymbol{\theta}\end{aligned}$$

where $\text{Var}\left(\mathbf{x}_{n_{\text{TRN}}+1}\right)$ is the covariance matrix of size $p \times p$. Then,

$$\text{Var}\left(\hat{Y}_{n_{\text{TRN}}+1}\right) = \sigma_e^2 \mathbb{E}\left(\left\|\mathbf{x}'_{n_{\text{TRN}}+1} \mathbf{X}' \mathbf{V}^{-1}\right\|^2\right) + \boldsymbol{\theta}' \mathbf{Q}' \mathbf{V}^{-1} \mathbf{X} \text{Var}\left(\mathbf{x}_{n_{\text{TRN}}+1}\right) \mathbf{X}' \mathbf{V}^{-1} \mathbf{Q} \boldsymbol{\theta} . \quad (4)$$

To conclude, according to formulae (1), (2), (4), the accuracy ρ_{RR} satisfies the following expression

$$\rho_{\text{RR}} = \frac{\mathbb{E}\left(\mathbf{x}'_{n_{\text{TRN}}+1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Q} \boldsymbol{\theta} \boldsymbol{\theta}' \mathbf{Q}' \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\theta}\right)}{\left(\sigma_e^2 \mathbb{E}\left(\left\|\mathbf{x}'_{n_{\text{TRN}}+1} \mathbf{X}' \mathbf{V}^{-1}\right\|^2\right) + \boldsymbol{\theta}' \mathbf{Q}' \mathbf{V}^{-1} \mathbf{X} \text{Var}\left(\mathbf{x}_{n_{\text{TRN}}+1}\right) \mathbf{X}' \mathbf{V}^{-1} \mathbf{Q} \boldsymbol{\theta}\right)^{1/2} (\sigma_G^2 + \sigma_e^2)^{1/2}} . \quad (5)$$

Note that the formula can be rewritten:

$$\rho_{\text{RR}} = \frac{\boldsymbol{\theta}' \mathbb{E}\left(\mathbf{q}_{n_{\text{TRN}}+1} \mathbf{x}'_{n_{\text{TRN}}+1}\right) \mathbf{X}' \mathbf{V}^{-1} \mathbf{Q} \boldsymbol{\theta}}{\left(\sigma_e^2 \mathbb{E}\left(\left\|\mathbf{x}'_{n_{\text{TRN}}+1} \mathbf{X}' \mathbf{V}^{-1}\right\|^2\right) + \boldsymbol{\theta}' \mathbf{Q}' \mathbf{V}^{-1} \mathbf{X} \text{Var}\left(\mathbf{x}_{n_{\text{TRN}}+1}\right) \mathbf{X}' \mathbf{V}^{-1} \mathbf{Q} \boldsymbol{\theta}\right)^{1/2} (\sigma_G^2 + \sigma_e^2)^{1/2}} .$$

2. A new proxy (QTLs in perfect LD with some markers)

Let us assume that we have the relationship

$$\mathbf{x}'_{n_{\text{TRN}}+1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Q} \boldsymbol{\theta} = \mathbf{q}'_{n_{\text{TRN}}+1} \boldsymbol{\theta} .$$

Let us consider the different terms present in the general formula (5). First, we have

$$\mathbb{E}\left(\mathbf{x}'_{n_{\text{TRN}}+1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Q} \boldsymbol{\theta} \boldsymbol{\theta}' \mathbf{Q}' \mathbf{V}^{-1} \mathbf{X} \boldsymbol{\theta}\right) = \mathbb{E}\left(\mathbf{q}'_{n_{\text{TRN}}+1} \boldsymbol{\theta} \boldsymbol{\theta}' \mathbf{q}_{n_{\text{TRN}}+1}\right) = \sigma_G^2 .$$

Besides, since

$$\text{Var}\left(\mathbb{E}\left(\hat{Y}_{n_{\text{TRN}}+1} \mid \mathbf{x}_{n_{\text{TRN}}+1}\right)\right) = \text{Var}\left(\mathbf{x}'_{n_{\text{TRN}}+1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Q} \boldsymbol{\theta}\right) = \text{Var}\left(\mathbf{q}'_{n_{\text{TRN}}+1} \boldsymbol{\theta}\right) = \sigma_G^2 ,$$

we have

$$\text{Var}\left(\hat{Y}_{n_{\text{TRN}}+1}\right) = \sigma_e^2 \mathbb{E}\left(\left\|\mathbf{x}'_{n_{\text{TRN}}+1} \mathbf{X}' \mathbf{V}^{-1}\right\|^2\right) + \sigma_G^2 .$$

Then, the accuracy becomes,

$$\rho_{\text{PLD}} = \frac{\sigma_G^2}{\left(\sigma_e^2 \mathbb{E}\left(\left\|\mathbf{x}'_{n_{\text{TRN}}+1} \mathbf{X}' \mathbf{V}^{-1}\right\|^2\right) + \sigma_G^2\right)^{1/2} (\sigma_G^2 + \sigma_e^2)^{1/2}} .$$

To conclude, since $\frac{\sigma_G^2}{\sigma_e^2} = \frac{h^2}{1-h^2}$, we obtain the final result

$$\rho_{\text{PLD}} = h \sqrt{\frac{h^2/(1-h^2)}{\mathbb{E}\left(\left\|\mathbf{x}'_{n_{\text{TRN}+1}}\mathbf{X}'\mathbf{V}^{-1}\right\|^2\right) + \frac{h^2}{1-h^2}}}.$$

3. Link with the previous work of Daetwyler et al. [2008]

Estimators computed from Ridge Regression are equal to OLS estimators if λ is set to zero (see for instance Fan and Lv [2008]). So, by setting $\lambda = 0$, we obtain the prediction

$$\hat{\mathbf{Y}}_{n_{\text{TRN}+1}} = \mathbf{q}'_{n_{\text{TRN}+1}} (\mathbf{Q}'\mathbf{Q})^{-1} \mathbf{Q}'\mathbf{Y}.$$

Having replaced the terms $\mathbf{X}'\mathbf{V}^{-1}$ by $(\mathbf{Q}'\mathbf{Q})^{-1}\mathbf{Q}'$ and $\mathbf{x}_{n_{\text{TRN}+1}}$ by $\mathbf{q}_{n_{\text{TRN}+1}}$ in our general formula (5), the accuracy becomes

$$\begin{aligned} \rho &= \frac{\mathbb{E}\left(\mathbf{q}'_{n_{\text{TRN}+1}}(\mathbf{Q}'\mathbf{Q})^{-1}\mathbf{Q}'\mathbf{Q}\boldsymbol{\theta}\mathbf{q}'_{n_{\text{TRN}+1}}\boldsymbol{\theta}\right)}{\left(\sigma_e^2\mathbb{E}\left(\left\|\mathbf{q}'_{n_{\text{TRN}+1}}(\mathbf{Q}'\mathbf{Q})^{-1}\mathbf{Q}'\right\|^2\right) + \boldsymbol{\theta}'\mathbf{Q}'\mathbf{Q}(\mathbf{Q}'\mathbf{Q})^{-1}\text{Var}(\mathbf{q}_{n_{\text{TRN}+1}})(\mathbf{Q}'\mathbf{Q})^{-1}\mathbf{Q}'\mathbf{Q}\boldsymbol{\theta}\right)^{1/2} (\sigma_G^2 + \sigma_e^2)^{1/2}} \\ &= \frac{\mathbb{E}\left(\mathbf{q}'_{n_{\text{TRN}+1}}\boldsymbol{\theta}\mathbf{q}'_{n_{\text{TRN}+1}}\boldsymbol{\theta}\right)}{\left(\sigma_e^2\mathbb{E}\left(\left\|\mathbf{q}'_{n_{\text{TRN}+1}}(\mathbf{Q}'\mathbf{Q})^{-1}\mathbf{Q}'\right\|^2\right) + \boldsymbol{\theta}'\text{Var}(\mathbf{q}_{n_{\text{TRN}+1}})\boldsymbol{\theta}\right)^{1/2} (\sigma_G^2 + \sigma_e^2)^{1/2}} \\ &= \frac{\sigma_G^2}{\left(\sigma_e^2\mathbb{E}\left(\left\|\mathbf{q}'_{n_{\text{TRN}+1}}(\mathbf{Q}'\mathbf{Q})^{-1}\mathbf{Q}'\right\|^2\right) + \sigma_G^2\right)^{1/2} (\sigma_G^2 + \sigma_e^2)^{1/2}} \\ &= \frac{h\sigma_G}{\left(\sigma_e^2\mathbb{E}\left(\left\|\mathbf{q}'_{n_{\text{TRN}+1}}(\mathbf{Q}'\mathbf{Q})^{-1}\mathbf{Q}'\right\|^2\right) + \sigma_G^2\right)^{1/2}}. \end{aligned}$$

To finish, we use that (proof given in next section)

$$\mathbb{E}\left(\left\|\mathbf{q}'_{n_{\text{TRN}+1}}(\mathbf{Q}'\mathbf{Q})^{-1}\mathbf{Q}'\right\|^2\right) \approx \frac{C}{n_{\text{TRN}}},$$

so the phenotypic accuracy and the genotypic accuracy are the following

$$\rho = h \sqrt{\frac{h^2/(1-h^2)}{\frac{C}{n_{\text{TRN}}} + \frac{h^2}{1-h^2}}}, \quad \tilde{\rho} = \sqrt{\frac{h^2/(1-h^2)}{\frac{C}{n_{\text{TRN}}} + \frac{h^2}{1-h^2}}}.$$

In Daetwyler et al. [2008], the authors consider the case $\sigma_G^2 + \sigma_e^2 = 1$. As a result,

$$\rho = \frac{h^2}{\sqrt{\sigma_e^2 \frac{C}{n_{\text{TRN}}} + h^2}}, \quad \tilde{\rho} = \frac{h}{\sqrt{\sigma_e^2 \frac{C}{n_{\text{TRN}}} + h^2}}. \quad (6)$$

Besides, they use the approximation $\sigma_e^2 \approx 1$. Using this approximation in (6) and simplifying by h , we obtain

$$\rho = \frac{h}{\sqrt{\frac{1}{\eta h^2} + 1}} \quad , \quad \tilde{\rho} = \sqrt{\frac{h^2 \eta}{1 + h^2 \eta}} \quad \text{where } \eta = n_{\text{TRN}}/C .$$

We can notice that this expression of $\tilde{\rho}$ is the same as the one presented in formula (1) of Daetwyler et al. [2008]. In the same way, the expression for ρ is the same as the one given at the end of Appendix A of Visscher et al. (2010), except that the focus was on the quantity ρ^2 .

Later, in their paper, Daetwyler et al. [2008] relaxed the assumption $\sigma_e^2 \approx 1$, and studied another approximation: $\sigma_e^2 \approx (1 - h^2) + h^2(1 - \tilde{\rho}^2)$. Using this new approximation in formula (6), we obtain

$$\tilde{\rho}^2 = \frac{h^2 \eta}{(1 - h^2 \tilde{\rho}^2) + h^2 \eta}$$

which is the same quantity as presented in formula (1) of Appendix S1 of Daetwyler et al. [2008].

4. Proof of $\mathbb{E} \left(\left\| \mathbf{q}'_{n_{\text{TRN}}+1} (\mathbf{Q}'\mathbf{Q})^{-1} \mathbf{Q}' \right\|^2 \right) \approx C/n_{\text{TRN}}$

Recall that Daetwyler et al. [2008] suppose that the matrix $(\mathbf{Q}'\mathbf{Q})^{-1}$ is diagonal and then,

$$(\mathbf{Q}'\mathbf{Q})_{j,j}^{-1} = \left(\sum_{i=1}^{n_{\text{TRN}}} Q_{i,j}^2 \right)^{-1} .$$

Let d_1, \dots, d_C denote the quantities such as

$$(\mathbf{Q}'\mathbf{Q})_{j,j} = d_j \quad (j = 1, \dots, C) .$$

Let $q_{n_{\text{TRN}}+1,j}$ denote the genotype of the TST individual at the j th QTL. We have the relationship

$$\begin{aligned} \left\| \mathbf{q}'_{n_{\text{TRN}}+1} (\mathbf{Q}'\mathbf{Q})^{-1} \mathbf{Q}' \right\|^2 &= \sum_{j=1}^C \frac{(q_{n_{\text{TRN}}+1,j})^2}{d_j} \\ &+ 2 \sum_{j=1}^{C-1} \sum_{j'=j+1}^C \frac{q_{n_{\text{TRN}}+1,j} q_{n_{\text{TRN}}+1,j'} \sum_{i=1}^{n_{\text{TRN}}} (Q_{i,j} Q_{i,j'})}{d_j d_{j'}} . \end{aligned}$$

Since the QTLs are assumed to be in linkage equilibrium, we have

$$\forall j \neq j', \quad \mathbb{E}(q_{n_{\text{TRN}}+1,j} q_{n_{\text{TRN}}+1,j'}) = 0 .$$

Recall that computations are conditional on \mathbf{Q} . As a consequence,

$$\mathbb{E} \left(2 \sum_{j=1}^{C-1} \sum_{j'=j+1}^C \frac{q_{n_{\text{TRN}}+1,j} q_{n_{\text{TRN}}+1,j'} \sum_{i=1}^{n_{\text{TRN}}} (Q_{i,j} Q_{i,j'})}{d_j d_{j'}} \right) = 0 .$$

Then,

$$\mathbb{E} \left(\left\| \mathbf{q}'_{n_{\text{TRN}}+1} (\mathbf{Q}'\mathbf{Q})^{-1} \mathbf{Q}' \right\|^2 \right) = \sum_{j=1}^C \frac{\mathbb{E} \left((q_{n_{\text{TRN}}+1,j})^2 \right)}{d_j} .$$

Finally, the authors use the approximation, $d_j \approx n_{\text{TRN}} \mathbb{E} \left((Q_{n_{\text{TRN}},j})^2 \right)$ suitable when n_{TRN} is large. Besides, they assume that the TST and TRN samples come from the same population. In this context, $Q_{1,j}, \dots, Q_{n_{\text{TRN}},j}, q_{n_{\text{TRN}}+1,j}$ are i.i.d. variables, and we have the relationship

$$\begin{aligned} \mathbb{E} \left(\sum_{j=1}^C \frac{(q_{n_{\text{TRN}}+1,j})^2}{d_j} \right) &\approx \sum_{j=1}^C \frac{\mathbb{E} \left((q_{n_{\text{TRN}}+1,j})^2 \right)}{n_{\text{TRN}} \mathbb{E} \left((Q_{n_{\text{TRN}},j})^2 \right)} \\ &\approx C/n_{\text{TRN}} . \end{aligned}$$

As a result,

$$\mathbb{E} \left(\left\| \mathbf{q}'_{n_{\text{TRN}}+1} (\mathbf{Q}'\mathbf{Q})^{-1} \mathbf{Q}' \right\|^2 \right) \approx C/n_{\text{TRN}} .$$

It concludes the proof.

References

- Daetwyler, H.D., Villanueva, B., Woolliams, J.A. (2008) Accuracy of predicting the genetic risk of disease using a genome-wide approach. *PLoS One*, **3(10)**, e3395.
- Daetwyler, H.D., Villanueva, B., Woolliams, J.A. (2010) The impact of genetic architecture on genome-wide evaluation methods. *Genetics*, **185(3)**, 1021-1031.
- Visscher, P.M., Yang, J., Goddard, M.E. (2010) A commentary on “common SNPs explain a large proportion of the heritability for human height” by Yang et al.(2010). *Twin Research and Human Genetics*, **13(06)**, 517-524.
- Fan, J., Lv, J. (2008) Sure independence screening for ultrahigh dimensional feature space. *Journal of the Royal Statistical Society Series B*, **70(5)**, 849-911.
- Meuwissen, T.H., Hayes, B., Goddard, M.E. (2001) Prediction of total genetic value using genome-wide dense marker maps. *Genetics*, **157(4)**, 1819-1829.