

Supplementary information for “The benefits of selecting phenotype-specific variants for applications of mixed models in genomics”

Christoph Lippert, Gerald Quon, Eun Yong Kang, Carl M. Kadie, Jennifer Listgarten, and David Heckerman

In the following we demonstrate the well-known fact that, in expectation, misspecification of the genetic similarity matrix leads to reduced predictive accuracy, both in terms of an inflated squared error as well as a lower predictive log-likelihood.

Let \mathbf{y} and \mathbf{y}^* denote training data and a test point, respectively. For the linear mixed model, \mathbf{y} and \mathbf{y}^* are jointly distributed as

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{y}^* \end{bmatrix} \sim N\left(\mathbf{m}; \sigma_g^2 \begin{bmatrix} \mathbf{Z} \\ \mathbf{Z}^* \end{bmatrix} \begin{bmatrix} \mathbf{Z} \\ \mathbf{Z}^* \end{bmatrix}^T + \sigma_e^2 \mathbf{I}\right),$$

where, for the sake of simplicity and without loss of generality, we suppress the dependence on the number S of variants contained in \mathbf{Z} .

It follows that

$$\mathbf{y}^* | \mathbf{y} \sim N\left(\underbrace{\mathbf{m} + \sigma_g^2 \mathbf{Z}^* \mathbf{Z}^T (\sigma_g^2 \mathbf{Z} \mathbf{Z}^T + \sigma_e^2 \mathbf{I})^{-1} (\mathbf{y} - \mathbf{m})}_{\hat{\boldsymbol{\mu}}_{\mathbf{Z}}^*}; \underbrace{\sigma_g^2 \mathbf{Z}^* \mathbf{Z}^{*T} + \sigma_e^2 \mathbf{I} - \sigma_g^4 \mathbf{Z}^* \mathbf{Z}^T (\sigma_g^2 \mathbf{Z} \mathbf{Z}^T + \sigma_e^2 \mathbf{I})^{-1} \mathbf{Z} \mathbf{Z}^{*T}}_{\boldsymbol{\Sigma}^*}\right).$$

The expected value of the squared error under the correct model equals

$$\begin{aligned} & E\left(\sum_{i=1}^N (y_i^* - \mu_{\mathbf{Z}_i}^*)^2\right) \\ &= E((\mathbf{y}^* - \boldsymbol{\mu}_{\mathbf{Z}}^*)^T (\mathbf{y}^* - \boldsymbol{\mu}_{\mathbf{Z}}^*)) \\ &= \text{Trace}[\boldsymbol{\Sigma}^*]. \end{aligned}$$

Let the misspecified joint genomic covariance matrix be $\begin{bmatrix} \mathbf{U} \\ \mathbf{U}^* \end{bmatrix} \begin{bmatrix} \mathbf{U} \\ \mathbf{U}^* \end{bmatrix}^T$. Then the misspecified predictor is

$$\boldsymbol{\mu}_{\mathbf{U}}^* = \mathbf{m} + \sigma_g^2 \mathbf{U}^* \mathbf{U}^T (\sigma_g^2 \mathbf{U} \mathbf{U}^T + \sigma_e^2 \mathbf{I})^{-1} (\mathbf{y} - \mathbf{m}).$$

It follows that the squared error of the misspecified model equals

$$\begin{aligned} & E\left(\sum_{i=1}^N (y_i^* - \mu_{\mathbf{U}_i}^*)^2\right) \\ &= E((\mathbf{y}^* - \boldsymbol{\mu}_{\mathbf{U}}^*)^T (\mathbf{y}^* - \boldsymbol{\mu}_{\mathbf{U}}^*)) \end{aligned}$$

$$= \frac{\text{Trace}[\boldsymbol{\Sigma}^*]}{E\left(\sum_{i=1}^N (y_i^* - \mu_{z_i}^*)^2\right)} + \underbrace{(\boldsymbol{\mu}_U^* - \boldsymbol{\mu}_Z^*)^T (\boldsymbol{\mu}_U^* - \boldsymbol{\mu}_Z^*)}_{\geq 0}.$$

Here, we used the well-known identity regarding the expectation of a squared form of a normally distributed variable (see, e.g. Section 0.5 in Roweis¹).

Misspecification also leads to a lower predictive log likelihood. The expected difference between the correctly specified model and the incorrectly specified model is

$$\begin{aligned} & E\left(\underbrace{\ln P(\mathbf{y}^*|\mathbf{y})}_{\text{correct model}} - \underbrace{\ln P_U(\mathbf{y}^*|\mathbf{y})}_{\text{incorrect model}}\right) \\ &= \int P(\mathbf{y}^*|\mathbf{y}) \ln \frac{P(\mathbf{y}^*|\mathbf{y})}{P_U(\mathbf{y}^*|\mathbf{y})} d\mathbf{y}^* \\ & \quad KL(P(\mathbf{y}^*|\mathbf{y}); P_U(\mathbf{y}^*|\mathbf{y})) \\ & \quad \geq 0, \end{aligned}$$

where the KL divergence is zero if and only if the distributions are equal, namely $P_U(\mathbf{y}^*|\mathbf{y}) = P(\mathbf{y}^*|\mathbf{y})$. It follows that, in expectation, the correctly specified model achieves a higher predictive log likelihood than the misspecified model.

References

1. Roweis, S. *Gaussian identities*. (1999, <http://www.cs.nyu.edu/~roweis/notes/gaussid.pdf>).