

# Prediction in high dimensional linear models and application to genomic prediction with a sparse genetic map

Charles-Elie Rabier\*, Simona Grusea#

\* IMAG, UMR 5149, CNRS, Université de Montpellier

# Institut de Mathématiques de Toulouse, Université de Toulouse, INSA



## Summary

Genomic selection (GS) consists in selecting individuals on the basis of genomic predictions, using a large number of genetic markers. An important question in GS is to determine the number of markers required for a good prediction. In order to answer this question, we present here new statistical results regarding Ridge regression. We analyzed rice data from the Philippines and focused on the flowering time collected during the dry season 2012. Using different densities of markers, we show that at least 1553 markers are required to implement GS.

## Causal model vs. Prediction model

Learning sample of size  $n$

**Causal model\*** ( $Q$  true regressors, with  $Q$  bounded)

$\theta^*$  vector of effects,  $M^*$  matrix of measures,  $Y$  vector of phenotypes

$$Y = M^* \theta^* + e$$

where  $Y = (Y_1, \dots, Y_n)'$ ,  $\theta^* = (\theta_1^*, \dots, \theta_Q^*)'$ ,  $e \sim N(0, \sigma_e^2 I_n)$

**Bayesian prediction model** ( $K$  regressors, with  $K \gg n$ )

$\theta$  vector of effects,  $M$  matrix of design

$$Y = M \theta + \epsilon$$

where  $Y = (Y_1, \dots, Y_n)'$ ,  $\theta = (\theta_1, \dots, \theta_K)'$ ,  $\theta \sim N(0, \sigma_\theta^2 I_K)$ ,  $\epsilon \sim N(0, \sigma_\epsilon^2 I_n)$ ,  $\epsilon_j \perp \theta_k$

We assume that the prediction model does not necessarily contain the true regressors

In other words, each column of  $M^*$  does not necessarily match a column of  $M$

## Validation sample + Accuracy criteria

— Let  $\text{new}$  denote an individual from the validation set

$$Y_{\text{new}} = m_{\text{new}}' \theta^* + e_{\text{new}} \quad \text{where } e_{\text{new}} \sim N(0, \sigma_e^2) \text{ and } m_{\text{new}}' \text{ vector of measures for the individual new}$$

— Prediction of the continuous variable  $Y_{\text{new}}$

$$\hat{Y}_{\text{new}} = m_{\text{new}}' \hat{\theta} = m_{\text{new}}' M' (M M' + \lambda I_n)^{-1} Y = m_{\text{new}}' (M' M + \lambda I_K)^{-1} M' Y$$

⇒ Accuracy criteria

$$\rho = \frac{\text{Cov}(\hat{Y}_{\text{new}}, Y_{\text{new}})}{\sqrt{\text{V}(\hat{Y}_{\text{new}}) \text{V}(Y_{\text{new}})}} \quad \text{with } m_{\text{new}} \text{ and } m_{\text{new}}^* \text{ random, } M \text{ is known}$$

Component present in the breeder's equation (cf. Lynch and Walsh, 1998)

## About the accuracy in the Ridge regression framework

The predictor is  $\hat{Y}_{\text{new}} = m_{\text{new}}' (M' M + \lambda I_K)^{-1} M' Y$

Let us define :

$$A_1 := \theta^{*\prime} \mathbb{E} (m_{\text{new}}^* m_{\text{new}}') M' V^{-1} M^* \theta^*, \quad A_2 := \sigma_e^2 \mathbb{E} \left( \left\| m_{\text{new}}' M' V^{-1} \right\|^2 \right) \\ A_3 := \theta^{*\prime} M^* V^{-1} M V (m_{\text{new}}) M' V^{-1} M^* \theta^*, \quad A_4 := \theta^{*\prime} \text{V} (m_{\text{new}}^*) \theta^* + \sigma_e^2$$

For the Ridge regression, we have

$$\rho = \frac{A_1}{(A_2 + A_3)^{1/2} (A_4)^{1/2}}$$

Under the oracle situation, we have  $\hat{Y}_{\text{new}} = m_{\text{new}}^* \theta^*$  and

$$\rho = \text{Cor} (m_{\text{new}}^* \theta^*, Y_{\text{new}}) = \sqrt{\frac{\text{V} (m_{\text{new}}^* \theta^*)}{\text{V} (Y_{\text{new}})}} = h$$

In genetics,  $h^2$  is called the heritability of the trait

## Results

Using the SVD decomposition of  $M = U D W'$  and  $M^* = U^* D^* W'^*$ , an estimation of the accuracy for the Ridge regression is (under several conditions)

$$\hat{\rho} = \frac{\hat{A}_1}{(\hat{A}_2 + \hat{A}_3)^{1/2} (\hat{A}_4)^{1/2}}$$

where

$$\hat{A}_1 = \frac{1}{n} \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \left\| U^{(s)} U^{(s)\prime} M^* \theta^* \right\|^2, \quad \hat{A}_2 = \frac{\sigma_e^2}{n} \sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} \\ \hat{A}_3 = \frac{1}{n} \sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} \left\| U^{(s)} U^{(s)\prime} M^* \theta^* \right\|^2, \quad \hat{A}_4 = \frac{1}{n} \sum_{s=1}^r d_s^{*2} \left\| W^{*(s)} W^{*(s)\prime} \theta^* \right\|^2 + \sigma_e^2$$

We discuss the link between

- the singular values  $d_1^* \geq d_2^* \geq \dots \geq d_r^* > 0$  and  $d_1 \geq d_2 \geq \dots \geq d_r > 0$
- the regularization parameter  $\lambda$
- the projection of the signal  $M^* \theta^*$  on the different subspaces

For that, we consider  $\lambda \rightarrow +\infty$  with  $\lambda = o(d_\ell^{*2})$ , and also the partitions  $\Omega_1^*$ ,  $\Omega_2^*$ ,  $\Omega_3^*$  of  $\{1, \dots, r^*\}$  and  $\Omega_1$ ,  $\Omega_2$ ,  $\Omega_3$  of  $\{1, \dots, r\}$ , such as

$$\Omega_1^* := \left\{ \ell \mid \lambda = o(d_\ell^{*2}) \right\}, \quad \Omega_1 := \left\{ s \mid \lambda = o(d_s^2) \right\} \\ \Omega_2^* := \left\{ \ell \mid d_\ell^{*2} \sim \frac{1}{C_\ell} \lambda \text{ with } C_\ell > 0 \right\}, \quad \Omega_2 := \left\{ s \mid d_s^2 \sim \frac{1}{C_s} \lambda \text{ avec } C_s > 0 \right\} \\ \Omega_3^* := \left\{ \ell \mid d_\ell^{*2} = o(\lambda) \right\}, \quad \Omega_3 := \left\{ s \mid d_s^2 = o(\lambda) \right\}$$

and several technical conditions

- (C1\*)  $\frac{n^{2\tau}}{r^*} \sum_{\ell \in \Omega_1^*} d_\ell^{*2} \rightarrow +\infty$
- (C2)  $\sum_{s \in \Omega_3} d_s^2 = o(\lambda)$
- (C3)  $\sum_{s \in \Omega_3} d_s^4 = o(\lambda^2)$
- (C4\*)  $\frac{n^{2\tau}}{r^*} = o(1/\lambda)$
- (C5)  $\#\Omega_1 = O(1)$
- (C6)  $\#\Omega_2 = O(1)$
- (C7\*)  $\frac{n^{2\tau}}{r^*} \sum_{\ell \in \Omega_1^*} \xi_2^{(\ell)} d_\ell^{*2} = o(1)$
- (C8\*)  $\frac{n^{2\tau}}{r^*} \sum_{\ell \in \Omega_1^*} \xi_3^{(\ell)} d_\ell^{*2} = o(1)$

to show the convergence to the oracle accuracy. For large  $n$ ,  $\hat{\rho} \sim \sqrt{\xi(n)} h$ .  $1 - \xi(n)$  can be viewed as a loss coefficient : it is the percentage of the  $L^2$  norm of  $U^{*(\ell)}$  that is unable to be captured

## Applications on simulated data and on rice real data

Is  $\hat{\rho}$  a good proxy for the accuracy  $\rho$ ?

We need to estimate the nuisance parameters  $\theta^*$  and  $M^*$   
⇒ we consider more markers for TRN than for TST

Comparison between :

- new proxy with more markers for TRN than for TST (imperfect linkage disequilibrium)
- old proxy with the same number of markers for TRN and TST (complete linkage)

**Simulated data** (panmictic population)

A small example with 1,000 markers for TRN and 500 markers for TST

Method	50 generations	70 generations	100 generations	MSE
Emp. Acc.	0.3909	0.3772	0.3217	
$\hat{\rho}(M^*, \hat{\theta}_{\text{LASSO}}^*)$	0.3397 (0.0112)	0.3436 (0.0132)	0.2629 (0.0146)	<b>0.0130</b>
$\hat{\rho}(M^*, \hat{\theta}_{\text{GPLASSO}}^*)$	0.2413 (0.0334)	0.3059 (0.0179)	0.2178 (0.0228)	0.0247
$\hat{\rho}(M^*, \hat{\theta}_{\text{ADLASSO}}^*)$	0.4677 (0.01293)	0.4821 (0.0222)	0.4093 (0.0164)	0.0172
$\hat{\rho}^{\text{pLD}}(\hat{\theta}_{\text{ADLASSO}})$	0.2970 (0.0336)	0.3182 (0.0306)	0.0986 (0.0693)	0.0445

**Rice flowering time** (data from Spindel et al, Plos Genetics 2015)

- $K = 73,147$  for TRN
- 4 densities of markers for TST (448, 781, 1553 and 3076)
- 252 TRN, 63 TST (i.e. 80% and 20%) + 100 samplings

**Number of markers (SNPs) required for predicting the TST individuals**

Method	448 SNPs	781 SNPs	1553 SNPs	3076 SNPs	MSE
Emp. Acc.	0.4789	0.4919	0.5275	0.5242	
$\hat{\rho}(M^*, \hat{\theta}_{\text{LASSO}}^*)$	0.4621 (0.0244)	0.4653 (0.0226)	0.4737 (0.0254)	0.4728 (0.0263)	<b>0.0247</b>
$\hat{\rho}(M^*, \hat{\theta}_{\text{ADLASSO}}^*)$	0.4269 (0.0355)	0.4379 (0.0376)	0.4520 (0.0419)	0.4461 (0.0430)	0.0395
$\hat{\rho}^{\text{pLD}}(\hat{\theta}_{\text{ADLASSO}})$	0.3662 (0.0454)	0.4202 (0.0281)	0.4919 (0.0215)	0.4952 (0.0342)	0.0323

## Reference

— Rabier CE and Grusea S. Prediction in high dimensional linear models and application to genomic selection ..., *Journal of the Royal Statistical Society : Series C (Applied Statistics)*, Vol 70(4), 2021

