

On the accuracy in high dimensional linear models and its application to genomic selection

Charles-Elie Rabier, Brigitte Mangin, Simona Grusea

► **To cite this version:**

Charles-Elie Rabier, Brigitte Mangin, Simona Grusea. On the accuracy in high dimensional linear models and its application to genomic selection. 2018. <hal-01456310v2>

HAL Id: hal-01456310

<https://hal.archives-ouvertes.fr/hal-01456310v2>

Submitted on 11 Mar 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Running Headline: On the accuracy in high dimension

Title: On the accuracy in high dimensional linear models and its application to genomic selection

C.E. Rabier^{a,b,c,d}, B. Mangin^e, S. Grusea^c

^a*ISEM, Université de Montpellier, CNRS, France*

^b*LIRMM, Université de Montpellier, CNRS, France*

^c*Institut de Mathématiques de Toulouse, Université de Toulouse, INSA de Toulouse, France*

^d*MIAT, Université de Toulouse, INRA, Castanet-Tolosan, France*

^e*LIPM, Université de Toulouse, INRA, CNRS, Castanet-Tolosan, France*

Abstract. Genomic selection is today a hot topic in genetics. It consists in predicting breeding values of selection candidates, using the large number of genetic markers now available thanks to the recent progress in molecular biology. One of the most popular method chosen by geneticists is Ridge regression. In this context, we focus on some predictive aspects of Ridge regression and present theoretical results regarding the accuracy criteria, i.e., the correlation between predicted value and true value. We show the influence of the singular values, the regularization parameter, and the projection of the signal on the space spanned by the rows of the design matrix. Asymptotic results in a high dimensional framework are also given; in particular, we prove that the convergence to an optimal accuracy highly depends on a weighted projection of the signal on each subspace. We discuss also on how to improve the prediction. Last, illustrations on simulated and real data are proposed.

Keywords: Accuracy, Genomic Selection, High Dimension, Linear Model, Prediction, Ridge Regression, Singular Value Decomposition, Sparsity.

1. Introduction and background

This year 2016, professor Michael Goddard and professor Theodorus Meuwissen were awarded The John J. Carty Award for the Advancement of Science by the National Academy of Science. They are considered as pioneers in the development of genomic selection (GS), because of their stimulating paper Meuwissen et al. (2001). In this context, our manuscript is devoted to methodological aspect of GS, a hot topic in genomics.

1.1. Preliminaries

For many years, geneticists focused on linkage analysis (LA) in order to detect on a given chromosome a Quantitative Trait Locus, so-called QTL: a

QTL is a section of the DNA that contains one or more genes influencing a quantitative trait which is able to be measured.

In this context, the most popular statistical method was Interval Mapping (Lander and Botstein (1989)). It consists in performing statistical tests along the genome. Using the information brought by genetic markers, the presence of a QTL is tested at every location in the genome.

Later, geneticists moved on to genome-wide association studies (GWAS). In contrast to LA, GWAS are based on unrelated individuals and as a result, larger sample sizes can be considered. GWAS enabled the discovery of many SNP-trait associations in humans (e.g. age-related macular degeneration, Fritsche et al. (2016), autism spectrum disorder, Connolly et al. (2017)).

However, both approaches (LA and GWAS) suffered from the fact that they were unable to detect QTLs with very small effects. Recall that most traits of interest are governed by a large number of small-effect QTLs (Goddard and Hayes (2008); Buckler et al. (2009)). It turns out that predictions based on selected SNPs could not be considered as reliable.

Today, Genomic Selection, motivated by the seminal paper of Meuwissen et al. (2001), is an extremely popular technique in genetics. It consists in predicting breeding values of selection candidates using a large number of genetic markers, thanks to the recent progress in molecular biology. The goal is not to detect QTLs anymore, but to predict the future phenotype of young candidates as soon as their DNA has been collected. GS relies on the expectation that each QTL will be highly correlated with at least one marker (Schulz-Streeck et al. (2012)). In genetics, this correlation is named Linkage Disequilibrium (LD): it refers to the non independence of alleles at 2 different loci (see Duret (2008) for more details).

GS was first applied to animal breeding (see Hayes et al. (2009)) and later to plant breeding (Jannink et al. (2010)): it was recently investigated on apple (Kumar et al. (2012)), sugar beet (Wurschum et al. (2013)), pea (Burstin et al. (2015)), and on inbred lines of rice (Spindel et al. (2015)).

1.2. A linear model

Let us introduce the statistical model associated to GS. The quantitative trait is observed on n training (TRN) individuals and we denote by Y_1, \dots, Y_n the observations. p markers lie on the genome, and β_j refers to the fixed marker effect of the j -th marker. In what follows, X is a matrix of size $n \times p$, and $'$ denotes transposition. The i -th row of X , written as $x'_i = (X_{i,1}, \dots, X_{i,p})$, represents the genome information at each marker available for the i -th individual.

A fixed number of QTLs lie on the genome, having an effect on the quantitative trait. For $1 \leq j \leq p$, $\beta_j = 0$ means that the corresponding marker is not a QTL, whereas $\beta_j \neq 0$ refers to a QTL. In genetics, this setting is named complete LD. In what follows, $\|\beta\|_0^0 := \sum_{j=1}^p |\beta_j|^0$ (with $0^0 = 0$) will denote the number of QTLs (i.e. non null marker effects).

We assume the following causal linear model for the quantitative trait:

$$Y = X\beta + \varepsilon, \tag{1}$$

where $Y = (Y_1, \dots, Y_n)'$, $\beta = (\beta_1, \dots, \beta_p)'$, $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma_e^2 I_n)$, I_n is the identity matrix of size n and σ_e^2 refers to the environmental variance.

In this manuscript we propose an analysis conditional on the observed values x_1, \dots, x_n . However, before imposing this conditioning, we have to precise that the matrix X is independent of ε . Simulated data will be generated accordingly. In what follows, r will denote the rank of the matrix X , and $\mathcal{R}(X)$ will refer to the linear space generated by the rows of X .

1.3. Introducing a test individual

A supplementary individual, so-called test (TST) individual (denoted *new*) is genotyped but not phenotyped. Using same notations as those used for the TRN population, x_{new} denotes the column vector containing the genome information at the p markers of the individual *new*. As a result, the quantitative trait Y_{new} can be written

$$Y_{new} = x'_{new} \beta + \varepsilon_{new},$$

where $\varepsilon_{new} \sim \mathcal{N}(0, \sigma_e^2)$.

We suppose that x'_{new} , ε_{new} and ε are all independent.

1.4. Introducing the accuracy

In GS, we are interested in predicting either the genotypic value $x'_{new} \beta$, or the phenotypic value Y_{new} . In both cases, an estimator \hat{Y}_{new} is constructed from a prediction model learned on n TRN individuals. \hat{Y}_{new} is a function of the random variables x_{new} and ε . Then, the quality of the prediction is evaluated according to some accuracy criteria, i.e. the correlation between predicted and true values. This criteria is a key element in genetics: it plays a role in the rate of genetic gain. Indeed, the accuracy is one component present in the breeders equation (see for instance Lynch and Walsh (1998)). The *phenotypic accuracy* ρ_{ph} , also called predictive ability, is defined in the following way

$$\rho_{ph} := \frac{\text{Cov}(\hat{Y}_{new}, Y_{new})}{\sqrt{\text{Var}(\hat{Y}_{new}) \text{Var}(Y_{new})}}, \quad (2)$$

whereas the *genotypic accuracy* ρ_g is defined as

$$\rho_g := \frac{\text{Cov}(\hat{Y}_{new}, x'_{new} \beta)}{\sqrt{\text{Var}(\hat{Y}_{new}) \text{Var}(x'_{new} \beta)}}. \quad (3)$$

Note that, when x_{new} , ε_{new} and ε are all independent, these two accuracies are linked by the relationship $\rho_{ph}/\rho_g = h$, where h is defined as the squared root of the heritability of the trait:

$$h^2 := \frac{\text{Var}(x'_{new} \beta)}{\text{Var}(Y_{new})} = \frac{\text{Var}(x'_{new} \beta)}{\text{Var}(x'_{new} \beta) + \text{Var}(\varepsilon_{new})} = \frac{\beta' \text{Var}(x_{new}) \beta}{\beta' \text{Var}(x_{new}) \beta + \text{Var}(\varepsilon_{new})}. \quad (4)$$

In what follows, we set $\sigma_G^2 := \text{Var}(x'_{new} \beta)$ and, as a consequence, we have the relationship $h^2 = \sigma_G^2 / (\sigma_G^2 + \sigma_e^2)$.

Depending on the authors, one focuses either on the phenotypic accuracy (e.g. Visscher et al. (2010)) or on the genotypic accuracy (e.g. Daetwyler et al. (2008, 2010)).

In what follows, the *oracle situation* will denote the settings where the QTLs locations and their effects are known. Then, under the oracle situation, the natural predictor is $\hat{Y}_{new} = x'_{new} \beta$. As a result, according to formula (2), the oracle accuracies are the following

$$\rho_g^{oracle} = 1, \quad \rho_{ph}^{oracle} = h.$$

1.5. Some background on Ridge regression

In the present study we propose to focus on Ridge regression, one of the most popular methods for prediction of breeding values. Ridge regression (Tihonov (1963); Hoerl et al. (1970)) has been studied for many years. In genetics, this regression model, initially proposed by Meuwissen et al. (2001) and Whittaker et al. (2000), is called random regression best linear unbiased predictor (RRBLUP) or genomic best linear unbiased predictor (GBLUP).

The Ridge estimator, suitable in a high dimensional setting (i.e. $p > n$), is the following:

$$\hat{\beta} := (X'X + \lambda I_p)^{-1} X'Y, \tag{5}$$

where λ refers to a regularization (or tuning) parameter.

Although Ridge regression is approximately 60 years old, statisticians keep studying this topic and excellent papers have been published recently (e.g. Shao and Deng (2012); Bühlmann (2013); Dicker (2016)).

1.6. Our contributions and roadmap

Our study starts, in Section 2, by recalling recent results on the accuracy. After a quick reminder on the singular value decomposition we introduce our main result, Theorem 1, that presents a general formula for the genotypic accuracy ρ_g . This is a key formula for the rest of the manuscript, since the other theorems and lemmas are built on it. According to Theorem 1, ρ_g depends on the projection of the signal β on $\mathcal{R}(X)$. This projection can be named “weighted projection” since some weights depending on singular values and on the tuning parameter λ act as multiplying factors.

Section 3 focuses on the case where TRN and TST samples come from the same probability distribution. In this context, Theorem 2 gives an estimation $\hat{\rho}_g$ of ρ_g which does not require the genome information of TST individuals. In other words, before genotyping TST individuals, it is possible to evaluate the accuracy of the future predictions on TST individuals. This estimation can help geneticists to figure out whether or not their population is appropriate for GS. Lemma 1 introduces a lower bound for $\hat{\rho}_g$: as in ii) of Theorem 1 of Shao

and Deng (2012), it only takes into account a global projection of the signal on $\mathcal{R}(X)$, with a global weight (i.e. same weights on each subspace).

Lemmas 2 and 3 propose a sharper analysis. In particular, Lemma 2 deals with the case where the projected signal is spread out uniformly on each vector of an orthonormal basis of $\mathcal{R}(X)$. It shows that under six given conditions, the estimation $\hat{\rho}_g$ tends to the oracle genotypic accuracy. These conditions are basically imposed on the singular values and on the ratio between the rank r of the matrix X and the projected signal on $\mathcal{R}(X)$.

Lemma 3 investigates certain extreme cases where the projected signal belongs either to the subspace spanned by the vector of an orthonormal basis of $\mathcal{R}(X)$ associated to the largest singular value of X , or to the subspace spanned by the vector associated to the smallest singular value. This setting is particularly interesting, since Ridge regression imposes shrinkage without taking into account the signal.

In Section 4 we tackle the problem of TRN and TST samples not coming from the same probability distribution. Theorem 3 introduces an estimator $\check{\rho}_g$ of ρ_g . In contrast to $\hat{\rho}_g$, $\check{\rho}_g$ requires genome informations on TST individuals. From a theoretical point of view, $\check{\rho}_g$ relies on the scalar product between a random projection of the signal and the usual projection of the signal on $\mathcal{R}(X)$. Lemma 1 presented in Supplementary material is the analogue of Lemma 1 under this new configuration.

Last, in Section 5 we propose a “modified” predictor for the genotypic accuracy; it is still derived from Ridge regression, but it may present better performances (cf. Theorem 4). We propose to project the vector Y on a well chosen subspace of the space spanned by the columns of X . In Lemma 6 we will give conditions for having an increase in terms of accuracy.

Our paper ends in Section 6 with an illustration on simulated data, mimicking the evolution of a population over time. We show the impact of the different probability distributions (TRN and TST) on the quality of the estimated accuracy. Furthermore, we highlight the fact that proxies built on our theoretical results outperform existing proxies in GS. Note that most of the existing proxies are built on Daetwyler et al. (2008)’s seminal formula: it consists in substituting an estimation of the effective number of independent loci M_e into this formula (see paragraph 6.2.3 for more details). Performances of the “modified” Ridge estimator are also illustrated. Finally, a real data analysis is proposed; it relies on the recent paper of Spindel et al. (2015) dealing with GS in rice.

2. General expression for the accuracy

2.1. Introducing Ridge regression and the corresponding accuracy

Recall the expression of the Ridge estimator:

$$\hat{\beta} = (X'X + \lambda I_p)^{-1} X'Y.$$

Since we have the well-known relationship

$$(X'X + \lambda I_p)^{-1} X' = X' (XX' + \lambda I_n)^{-1}, \quad (6)$$

the computation of $\hat{\beta}$ only requires the inversion of a $n \times n$ matrix.

In this context, the prediction for the so-called *new* individual is the following:

$$\hat{Y}_{new} := x'_{new}\hat{\beta} = x'_{new}X'V^{-1}Y, \quad \text{where } V = XX' + \lambda I_n.$$

In what follows we will assume that Y , the columns of X , Y_{new} and x_{new} are all centered.

According to formula (5) of Rabier et al. (2016), assuming that x_1, \dots, x_n are known and that ε , x_{new} and ε_{new} are random, the genotypic accuracy has the following expression:

$$\rho_g = \frac{\beta' \text{Var}(x_{new}) X'V^{-1}X\beta}{\left(\sigma_\varepsilon^2 \mathbb{E}\left(\|x'_{new}X'V^{-1}\|^2\right) + \beta'X'V^{-1}X\text{Var}(x_{new})X'V^{-1}X\beta\right)^{1/2} \sigma_G}, \quad (7)$$

where $\|\cdot\|$ is the L^2 norm and $\text{Var}(x_{new})$ is the covariance matrix of size $p \times p$. Note that this accuracy can be viewed as a conditional accuracy, since this expression was obtained conditionally on the TRN design matrix X .

2.2. SVD decomposition

Following Shao and Deng (2012) and Bühlmann (2013), let us consider the singular value decomposition of X :

$$X = PDQ', \quad (8)$$

where P is an $n \times r$ matrix satisfying $P'P = I_r$, Q is a $p \times r$ matrix satisfying $Q'Q = I_r$, and $D = \text{Diag}(d_1, \dots, d_r)$ with $d_1 \geq \dots \geq d_r > 0$. The columns of Q (resp. P) constitute an orthogonal basis of the space spanned by the rows (resp. columns) of X . In what follows, $Q^{(s)}$ will denote the s -th column of Q , and as a consequence $\mathcal{R}(X) = \text{Span}\{Q^{(1)}, \dots, Q^{(r)}\}$. By construction QQ' is an idempotent matrix, and $QQ'\beta$ is the projection of β onto $\mathcal{R}(X)$. We set

$$\theta := QQ'\beta$$

and, as mentioned in Shao and Deng (2012), we have the relationship

$$\hat{\theta} := QQ'\hat{\beta} = \hat{\beta}.$$

The Ridge estimator $\hat{\beta}$ presents therefore the particularity that it belongs to $\mathcal{R}(X)$.

2.3. Results

We introduce the following notations

$$A_1 := \beta' \text{Var}(x_{new}) X'V^{-1}X\beta, \quad A_2 := \sigma_\varepsilon^2 \mathbb{E}\left(\|x'_{new}X'V^{-1}\|^2\right)$$

$$A_3 := \beta' X'V^{-1}X\text{Var}(x_{new}) X'V^{-1}X\beta, \quad A_4 := \sigma_G^2.$$

Our main result is the following.

Theorem 1. Let $\Sigma = \text{Var}(x_{new})$ be the covariance matrix of size $p \times p$. Furthermore, let us assume that X is known and that ε , x_{new} and ε_{new} are random. Then, the genotypic accuracy has the following expression

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

, where

$$A_1 = \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \beta' \Sigma Q^{(s)} Q^{(s)'} \beta, \quad A_2 = \sigma_e^2 \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \mathbb{E} \left(\left\| Q^{(s)} Q^{(s)'} x_{new} \right\|^2 \right)$$

$$A_3 = \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)} Q^{(s)'} \beta \right)' \Sigma \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)} Q^{(s)'} \beta \right), \quad A_4 = \beta' \Sigma \beta.$$

The proof is given in Section 1 of the Supplementary material.

The phenotypic accuracy ρ_{ph} is obtained by replacing the term A_4 at the denominator by $A_4 + \sigma_e^2$. Note that $Q^{(s)} Q^{(s)'} v$ is the projection of a column vector v of size p on the vector space spanned by $Q^{(s)}$. In view of Theorem 1, ρ_g depends on the projections $Q^{(s)} Q^{(s)'} \beta$ of the signal and also on the projections $Q^{(s)} Q^{(s)'} x_{new}$ of the genome information for the individual *new*.

In what follows we are interested in estimating the genotypic accuracy ρ_g . A consistent estimator of A_2 is easily derived from the Law of large numbers. Besides, by Slutsky's lemma in the matrix case, consistent estimators of A_1 , A_3 and A_4 can be obtained provided that a consistent estimator of the covariance matrix Σ is used. This finally leads to a consistent estimator of ρ_g .

However, finding a consistent estimator for Σ is very challenging in the high dimensional setting; it is nowadays a hot topic in statistics. Some recent results (see e.g. Cai et al. (2010)) address this question, but the authors make quite restrictive assumptions on the covariance matrix Σ .

In our present work we have chosen the empirical covariance estimator, since it is the classical estimator used by geneticists in practice. We will show on simulated data that our estimators perform in a very satisfactory manner.

3. Estimation when TRN and TST samples come from the same probability distribution

In this section, let us consider the case where the TRN and TST samples come from the same probability distribution. In this context, using the empirical covariance matrix $X'X/n$ as an estimation of the covariance matrix Σ from Theorem 1, we obtain the following theorem.

Theorem 2. Let us assume that x_1, \dots, x_n and x_{new} are independent and identically distributed (i.i.d.). Besides, let us consider that x_1, \dots, x_n have been observed (i.e. X is known), and that ε , x_{new} and ε_{new} are random. Then,

an estimation of the genotypic accuracy is

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

where

$$\begin{aligned} \hat{A}_1 &= \frac{1}{n} \sum_{s=1}^r \frac{d_s^4}{d_s^2 + \lambda} \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2, \quad \hat{A}_2 = \frac{\sigma_e^2}{n} \sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} \\ \hat{A}_3 &= \frac{1}{n} \sum_{s=1}^r \frac{d_s^6}{(d_s^2 + \lambda)^2} \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2, \quad \hat{A}_4 = \frac{1}{n} \sum_{s=1}^r d_s^2 \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2. \end{aligned}$$

In contrast to Theorem 1, the projections $Q^{(s)} Q^{(s)'} x_{new}$ are not present in this new expression. Theoretical developments rely on the following estimation \hat{A}_2 of A_2 :

$$\hat{A}_2 := \frac{\sigma_e^2}{n} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \text{Tr} \left(X Q^{(s)} Q^{(s)'} Q^{(s)} Q^{(s)'} X' \right).$$

The proof is given in Section 2 of the Supplementary material.

This estimation $\hat{\rho}_g$ relies only on information collected on TRN (phenotypes and markers). As a consequence, this accuracy estimation can be used to evaluate GS accuracy before genotyping of the TST individuals. Note that the unknown quantity β present in Theorem 2 can be estimated for instance by LASSO (Tibshirani (1996)), Adaptive LASSO (Zou (2006)) or Group LASSO (Yuan and Lin (2006)). We refer to our applications in Section 6.

Let us now give bounds for the quantity $\hat{\rho}_g$.

Lemma 1 (Bounds on $\hat{\rho}_g$). *Using same assumptions as in Theorem 2, we always have*

$$\frac{\|QQ'\beta\|^2 \min_s \frac{d_s^4}{d_s^2 + \lambda}}{\sqrt{\sigma_e^2 r + \|QQ'\beta\|^2 \max_s \frac{d_s^6}{(d_s^2 + \lambda)^2}} \sqrt{\|QQ'\beta\|^2 \max_s d_s^2}} \leq \hat{\rho}_g \leq \rho_g^{oracle}.$$

The proof is given in Section 3 of the Supplementary material.

According to this lemma, the smaller the ratio $\frac{r}{\|QQ'\beta\|^2}$ is, the larger the lower bound is. Furthermore, the quantity $\min_s \frac{d_s^4}{d_s^2 + \lambda}$ should be large enough, and the term $\max_s \frac{d_s^6}{(d_s^2 + \lambda)^2}$ not too large.

Although the above lower bound can give a first indication on the quality of the prediction, a sharper analysis is needed (see below). Indeed, until now, as in ii) of Theorem 1 of Shao and Deng (2012), we have only taken into account a global projection $QQ'\beta$ of the signal on $\mathcal{R}(X)$, with a global weight.

3.1. *Convergence of $\hat{\rho}_g$ to ρ_g^{oracle} when $n \rightarrow +\infty$ and $p \rightarrow +\infty$*

Recall that $d_1 \geq d_2 \geq \dots \geq d_r > 0$ are the singular values of X . To study asymptotic properties of $\hat{\rho}_g$, we consider that

$$\begin{aligned} d_1^2 &\sim n^\psi \text{ with } 0 < \psi \leq 1, \\ d_r^2 &\sim n^\eta \text{ with } \eta \leq \psi \leq 1 \text{ and } \eta \text{ and } \psi \text{ not depending on } n. \end{aligned}$$

Recall that the notation $u_n \sim v_n$ means that $\frac{u_n}{v_n} \rightarrow 1$ when $n \rightarrow \infty$.

Moreover, we will assume that

$$\|QQ'\beta\|^2 \sim n^{2\tau} \text{ with } \tau < \eta \text{ and } \tau \text{ not depending on } n.$$

These conditions are largely inspired from Shao and Deng (2012) and Fan and Lv (2008).

Let us consider a regularization parameter λ such as $\lambda \rightarrow \infty$ and $\lambda = o(d_1^2)$. Besides, Ω_1 , Ω_2 and Ω_3 will denote the three following sets:

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

Note that Ω_1 contains at least the index 1. On simulated data, after having chosen λ by Restricted Maximum Likelihood (Corbeil and Searle (1976)), so-called REML (cf. Section 6), these different sets were not empty.

In what follows, we will call respectively “largest singular values” the ones whose index s belong to the set Ω_1 . In the same way, “intermediate singular values” and “smallest singular values” refers to the sets Ω_2 and Ω_3 respectively.

Let us introduce a few extra conditions:

- (C1) $\frac{n^{2\tau}}{r} \sum_{s \in \Omega_1} d_s^2 \rightarrow +\infty$
- (C2) $\sum_{s \in \Omega_3} d_s^2 = o(\lambda)$
- (C3) $\sum_{s \in \Omega_3} d_s^4 = o(\lambda^2)$
- (C4) $n^{2\tau}/r = o(1/\lambda)$, i.e. $\lambda = o(r/n^{2\tau})$
- (C5) $\#\Omega_1 = O(1)$
- (C6) $\#\Omega_2 = O(1)$,

where $\#\Omega$ refers to the cardinal of the set Ω . We refer to the Supplementary material for some explanations on these technical assumptions.

The following lemma assumes that the signal is spread out uniformly on each subspace.

Lemma 2 (Convergence to the oracle accuracy). *Let us consider same assumptions as in Theorem 2. Besides, let us suppose that the projected signal is spread out uniformly on each subspace $\text{Span}\{Q^{(s)}\}$, i.e.*

$$\left\|Q^{(s)}Q^{(s)'}\beta\right\|^2 \sim \frac{n^{2\tau}}{r}, \quad s = 1, \dots, r \quad (9)$$

and let us assume that conditions (C1-C2-C3-C4-C5-C6) hold. Then we have $\hat{\rho}_g \rightarrow \rho_g^{oracle}$.

The proof is given in Section 4 of the Supplementary material.

If we set $r = n^\gamma$ with $0 < \gamma \leq 1$, then the condition (C4) implies that $\tau < \gamma/2$. In other words, when trying to recover the oracle accuracy, the lower the rank r is, the weaker the signal can be.

Recall that the tuning parameter λ is such as $\lambda \rightarrow \infty$, $\lambda = o(d_1^2)$. Let us now introduce the following lemma which deals with some extreme cases.

Lemma 3 (Extreme cases). *Let us consider same assumptions as in Theorem 2.*

1. *If the projected signal belongs only to $\text{Span}\{Q^{(1)}\}$, that is to say*

$$\left\| Q^{(1)} Q^{(1)'} \beta \right\|^2 \sim n^{2\tau}, \quad \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2 = 0, \text{ for } 1 < s \leq r, \text{ then}$$

- *if $2\tau + \psi > 1$, then $\hat{\rho}_g \rightarrow \rho_g^{oracle}$.*
- *if $2\tau + \psi < 1$*
 - *if $\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} = o(n^{2\tau + \psi})$, then $\hat{\rho}_g \rightarrow \rho_g^{oracle}$*
 - *if $n^{2\tau + \psi} = o\left(\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2}\right)$, then $\hat{\rho}_g \rightarrow 0$.*

2. *If the projected signal belongs only to $\text{Span}\{Q^{(r)}\}$, that is to say*

$$\left\| Q^{(r)} Q^{(r)'} \beta \right\|^2 \sim n^{2\tau}, \quad \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2 = 0, \text{ for } 1 \leq s < r, \text{ and}$$

moreover $\lambda \sim Cn^{\eta + \kappa}$ with $\kappa > \max(0, -\eta)$, $C > 0$, then,

- *if $\tau + \eta/2 - \kappa < 0$, then $\hat{\rho}_g \rightarrow 0$.*
- *if $\tau + \eta/2 - \kappa > 0$*
 - *if $\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} = o(n^{2\tau + \eta - 2\kappa})$, then $\hat{\rho}_g \rightarrow \rho_g^{oracle}$*
 - *if $n^{2\tau + \eta - 2\kappa} = o\left(\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2}\right)$, then $\hat{\rho}_g \rightarrow 0$.*

The proof is given in Section 5 of the Supplementary material.

Since $\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} \leq r$, we have $r = o(n^{2\tau + \psi})$ and thus the condition $\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} = o(n^{2\tau + \psi})$ can be replaced by $r = o(n^{2\tau + \psi})$. In the same way, condition $\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} = o(n^{2\tau + \eta - 2\kappa})$ can be replaced by $r = o(n^{2\tau + \eta - 2\kappa})$.

According to this lemma, when the projected signal belongs only to $\text{Span}\{Q^{(r)}\}$, κ should be not too large in order to ensure that $\tau + \eta/2 - \kappa > 0$ and also that $r = o(n^{2\tau + \eta - 2\kappa})$. As a consequence, the tuning parameter λ should be chosen appropriately.

Summary of the results in Section 3: We present an estimation $\hat{\rho}_g$ of the genotypic accuracy which is suitable when TRN and TST individuals are sampled

from the same population. This estimation relies only on information collected on TRN (phenotypes and markers). So, it is possible to evaluate the accuracy of the future prediction of TST individuals before genotyping them. In other words, our formula can help geneticists to figure out whether or not their population is appropriate for GS.

4. Estimation when TRN and TST samples do not come from the same probability distribution

In this section we consider the general case where the TRN and TST samples do not come necessarily from the same probability distribution. Furthermore, let us assume that genome informations for n_{new} new individuals are available and that we are willing to predict the phenotypes of those individuals. Let X_{new} be a random matrix of size $n_{new} \times p$ containing the genomic markers of the new individuals. The singular value decomposition of X_{new} is the following:

$$X_{new} = WFZ',$$

where W is a $n_{new} \times r_{new}$ matrix satisfying $W'W = I_{r_{new}}$, Z is a $p \times r_{new}$ matrix satisfying $Z'Z = I_{r_{new}}$, and F is $r_{new} \times r_{new}$ diagonal matrix of full rank.

In what follows, $\langle \cdot, \cdot \rangle$ denotes the usual scalar product. Using $X'_{new}X_{new}/n_{new}$ as estimator of the covariance matrix Σ , we obtain the following Theorem 3, a random version of Theorem 2.

Theorem 3. *Let us assume that X is given and that X_{new} is random, with its rows being i.i.d. Then, an estimator of the genotypic accuracy is*

$$\check{\rho}_g = \frac{\check{A}_1}{(\check{A}_2 + \check{A}_3)^{1/2} (\check{A}_4)^{1/2}}, \quad (10)$$

where

$$\begin{aligned} \check{A}_1 &= \frac{1}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \left(\sum_{\alpha=1}^{r_{new}} f_\alpha^2 \langle Z^{(\alpha)} Z^{(\alpha)'} \beta, Q^{(s)} Q^{(s)'} \beta \rangle \right), \\ \check{A}_2 &= \frac{\sigma_e^2}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \sum_{i=1}^{n_{new}} \left(\sum_{\alpha=1}^{r_{new}} f_\alpha Q^{(s)'} Z^{(\alpha)} W_i^{(\alpha)} \right)^2, \\ \check{A}_3 &= \frac{1}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)'} \beta \sum_{\ell=1}^r \frac{d_\ell^2}{d_\ell^2 + \lambda} Q^{(\ell)'} \beta \left(\sum_{\alpha=1}^{r_{new}} f_\alpha^2 \langle Z^{(\alpha)} Z^{(\alpha)'} Q^{(s)}, Z^{(\alpha)} Z^{(\alpha)'} Q^{(\ell)} \rangle \right), \\ \check{A}_4 &= \frac{1}{n_{new}} \sum_{\alpha=1}^{r_{new}} f_\alpha^2 \left\| Z^{(\alpha)} Z^{(\alpha)'} \beta \right\|^2. \end{aligned}$$

Note that the expression in Equation (10) was obtained with the help of the estimator \check{A}_2 defined in the following way

$$\check{A}_2 := \frac{\sigma_e^2}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \text{Tr} \left(X_{new} Q^{(s)} Q^{(s)'} Q^{(s)} Q^{(s)'} X'_{new} \right).$$

The proof is given in Section 6 of the Supplementary material. In contrast to Section 3, we need to collect the genome information on TST in order to compute the quantity $\check{\rho}_g$. Lemma 1 of the Supplementary material is the analogue of Lemma 1. Contrary to the lower bound introduced in Lemma 1, the new lower bound can take negative values, since the scalar product $\langle ZZ'\beta, QQ'\beta \rangle$ is present in the numerator. This may happen when the rows of X are not i.i.d. or when the probability distributions of TRN and TST are very different.

In Section 6 we will illustrate performances of $\check{\rho}_{ph}$ and $\hat{\rho}_{ph}$ on simulated and real data.

Summary of the results in Section 4: We present an estimator $\check{\rho}_g$ which can be used to estimate the genotypic accuracy when TRN and TST individuals are not sampled from the same population. An example of application is plant breeding: since a large number of generations can be obtained easily, the fitted model is not readjusted at each generation, in order to save time or costs due to genotyping. In contrast to Section 3, our estimator $\check{\rho}_g$ relies on informations collected on TRN (phenotypes and markers) and on TST (markers).

5. How to improve the quality of the prediction

In this section, we introduce another estimator of marker effects β derived from Ridge regression, which may present, in some cases, better performances than previously studied estimators. We propose to project the vector Y on a well chosen subspace of the space spanned by the columns of X . Let $1 \leq \tilde{r} \leq r$ and $\sigma(\cdot)$ a one-to-one map $\sigma : \{1, \dots, \tilde{r}\} \rightarrow \{1, \dots, r\}$. We thus have $\sigma(k) \neq \sigma(k')$ for $k \neq k'$.

Let us consider the estimator

$$\tilde{\beta} = X'V^{-1}\tilde{P}\tilde{P}'Y, \text{ where } \tilde{P} = (P^{\sigma(1)}, \dots, P^{\sigma(\tilde{r})}).$$

Note that $\tilde{P}\tilde{P}'Y$ is the projection of Y on $Span\{P^{\sigma(1)}, \dots, P^{\sigma(\tilde{r})}\}$.

Besides, we set $\tilde{Q} = (Q^{\sigma(1)}, \dots, Q^{\sigma(\tilde{r})})$. Then, the corresponding prediction for the so-called *new* individual is the following:

$$\tilde{Y}_{new} = x'_{new}\tilde{\beta} = x'_{new}X'V^{-1}\tilde{P}\tilde{P}'Y.$$

We refer to Subsection 6.2.4, where we describe a procedure for choosing $\sigma(\cdot)$ and \tilde{r} .

Let $\tilde{\rho}_g$ be the analogue of ρ_g , with \hat{Y}_{new} replaced by \tilde{Y}_{new} (cf. formula (3)):

$$\tilde{\rho}_g := \frac{\text{Cov}(\tilde{Y}_{new}, x'_{new}\beta)}{\sqrt{\text{Var}(\tilde{Y}_{new}) \text{Var}(x'_{new}\beta)}}. \quad (11)$$

A more explicit formula for $\tilde{\rho}_g$ is given in Lemma 2 of the Supplementary material. This lemma can be viewed as a version of Theorem 1 based on this new estimator. Let us now present a lemma which is the analogue of Theorem 2.

Lemma 4. *Let us consider same assumptions as in Theorem 2. Then an estimation of the quantity $\tilde{\rho}_g$ is*

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

where

$$\begin{aligned} \hat{A}_1 &:= \frac{1}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{d_{\sigma(s)}^2 + \lambda} \left\| Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right\|^2, \quad \hat{A}_2 := \frac{\sigma_e^2}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{(d_{\sigma(s)}^2 + \lambda)^2} \\ \hat{A}_3 &:= \frac{1}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^6}{(d_{\sigma(s)}^2 + \lambda)^2} \left\| Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right\|^2, \quad \hat{A}_4 := \hat{A}_4. \end{aligned}$$

The proof is given in Section 8 of the Supplementary material.

Note that the quantities $\tilde{A}_1, \dots, \tilde{A}_4$ are the analogues of A_1, \dots, A_4 in this new setting.

Let us next introduce our Lemma 5, which is the analogue of Lemma 1 regarding bounds for the genotypic accuracy.

Lemma 5 (Bounds on $\hat{\rho}_g$). *Let us consider same assumptions as in Theorem 2. Then we always have*

$$\frac{\left\| \tilde{Q} \tilde{Q}' \beta \right\|^2 \min_{1 \leq s \leq \tilde{r}} \frac{d_{\sigma(s)}^4}{d_{\sigma(s)}^2 + \lambda}}{\sqrt{\sigma_e^2 \tilde{r} + \left\| \tilde{Q} \tilde{Q}' \beta \right\|^2} \max_{1 \leq s \leq \tilde{r}} \frac{d_{\sigma(s)}^6}{(d_{\sigma(s)}^2 + \lambda)^2} \sqrt{\left\| Q Q' \beta \right\|^2} \max_{1 \leq s \leq r} d_s^2} \leq \hat{\rho}_g \leq \rho_g^{oracle}.$$

The proof relies heavily on the proof of Lemma 1, using the expressions of \hat{A}_1, \hat{A}_2 and \hat{A}_3 given in Lemma 4. We can notice that at the denominator, the quantities \tilde{r} and $\left\| \tilde{Q} \tilde{Q}' \beta \right\|^2$ replace now the quantities r and $\left\| Q Q' \beta \right\|^2$ of Lemma 1. This decrease at the denominator will be profitable provided that the numerator does not decrease too much.

Remark : Note that if $\hat{A}_1 - \tilde{A}_1 = 0$, then $\hat{\rho}_g \geq \tilde{\rho}_g$.

This is the case for example if $\left\| Q^{(\ell)} Q^{(\ell)'} \beta \right\|^2 = 0$, for all $\ell \notin \{\sigma(1), \dots, \sigma(\tilde{r})\}$. Indeed, in this case we have $\hat{A}_1 = \tilde{A}_1$ and thus $\hat{\rho}_g \geq \tilde{\rho}_g$.

For fixed n , we can obtain the following comparison between $\hat{\rho}_g$ and $\tilde{\rho}_g$.

Lemma 6. *Let us suppose that $\hat{A}_1 - \tilde{A}_1 \neq 0$. Then we have $\hat{\rho}_g \geq \tilde{\rho}_g$ if and only if the following relation holds:*

$$\frac{\hat{A}_1}{\hat{A}_1 - \tilde{A}_1} \geq \frac{(\hat{A}_2 + \hat{A}_3)}{\hat{A}_2 + \hat{A}_3 - (\hat{A}_2 + \hat{A}_3)} \left(1 + \sqrt{\frac{\hat{A}_2 + \hat{A}_3}{\hat{A}_2 + \hat{A}_3}} \right).$$

The proof is given in Section 9 of the Supplementary material. In other words, under this condition, the accuracy is improved and as a consequence, we should choose the estimator \hat{Y}_{new} instead of the classical estimator \hat{Y}_{new} .

In what follows, we propose to give extra explanations on the results in this subsection. Let us introduce the estimator $\vec{\beta}$ such as

$$\vec{\beta} := X'V^{-1}\vec{P}\vec{P}'Y$$

where \vec{P} denotes the matrix obtained from P by removing the columns vectors $P^{(1)}, \dots, P^{(r)}$. \vec{Y}_{new} will denote the corresponding prediction:

$$\vec{Y}_{new} = x'_{new}\vec{\beta}.$$

Let P^\perp be a matrix of size $n \times n - r$ whose columns form an orthogonal basis of $(Span\{P^{(1)}, \dots, P^{(r)}\})^\perp$. We then have the relationship:

$$Y = PP'Y + P^\perp P^{\perp'}Y = \tilde{P}\tilde{P}'Y + \vec{P}\vec{P}'Y + P^\perp P^{\perp'}Y.$$

By definition we have $X'P^\perp P^{\perp'}Y = QDP'P^\perp P^{\perp'}Y = 0_p$, where 0_p denotes the column vector of size p with all components equal to 0. Then, we have

$$\begin{aligned} \hat{\beta} &= (X'X + \lambda I_p)^{-1}X'Y = (X'X + \lambda I_p)^{-1}X'(\tilde{P}\tilde{P}'Y + \vec{P}\vec{P}'Y) \\ &= X'V^{-1}\tilde{P}\tilde{P}'Y + X'V^{-1}\vec{P}\vec{P}'Y. \end{aligned}$$

As a result, $\hat{\beta} = \tilde{\beta} + \vec{\beta}$ and $\hat{Y}_{new} = \tilde{Y}_{new} + \vec{Y}_{new}$. We further have

$$\begin{aligned} \text{Cov}(\tilde{\beta}, \vec{\beta}) &= \text{Cov}(X'V^{-1}\tilde{P}\tilde{P}'Y, X'V^{-1}\vec{P}\vec{P}'Y) \\ &= X'V^{-1}\tilde{P}\tilde{P}'\text{Var}(Y)\vec{P}\vec{P}'V^{-1}X \\ &= \sigma_e^2 X'V^{-1}\tilde{P}\tilde{P}'I_n\vec{P}\vec{P}'V^{-1}X \\ &= O_p, \end{aligned}$$

where O_p denotes the zero matrix of size $p \times p$.

As a consequence, $\text{Cov}(\tilde{Y}_{new}, \vec{Y}_{new}) = O_p$. We deduce that the variance of \hat{Y}_{new} can be decomposed in the following way:

$$\text{Var}(\hat{Y}_{new}) = \text{Var}(\tilde{Y}_{new}) + \text{Var}(\vec{Y}_{new}).$$

By definition, according to Rabier et al. (2016), we have

$$A_1 = \text{Cov}(\hat{Y}_{new}, Y_{new}), \quad A_2 + A_3 = \text{Var}(\hat{Y}_{new})$$

As a consequence, we can obtain the following estimations

$$\widehat{\text{Cov}}(\hat{Y}_{new}, Y_{new}) = \hat{A}_1, \quad \widehat{\text{Var}}(\hat{Y}_{new}) = \hat{A}_2 + \hat{A}_3.$$

We have analogue estimates for \tilde{Y}_{new} and \vec{Y}_{new} .

In what follows, $\vec{\rho}_g$ is the analogue of $\hat{\rho}_g$, with \tilde{Y}_{new} replaced by \vec{Y}_{new} . $\hat{\rho}_g$ will refer to an estimation of $\vec{\rho}_g$. With these notations, we obtain the following corollary of Lemma 6.

Corollary 1. *Suppose that $\widehat{\text{Cov}}(\tilde{Y}_{new}, Y_{new}) \neq 0$ and $\widehat{\text{Cov}}(\vec{Y}_{new}, Y_{new}) \neq 0$. Then we have the following three possible situations:*

1. *We have $\hat{\rho}_g \geq \hat{\rho}_g$ if and only if*

$$\frac{\widehat{\text{Cov}}(\tilde{Y}_{new}, Y_{new})}{\widehat{\text{Cov}}(\vec{Y}_{new}, Y_{new})} \geq \frac{\widehat{\text{Var}}(\tilde{Y}_{new})}{\widehat{\text{Var}}(\vec{Y}_{new})} \left(1 + \sqrt{1 + \frac{\widehat{\text{Var}}(\tilde{Y}_{new})}{\widehat{\text{Var}}(\vec{Y}_{new})}} \right).$$

In this case, we also have $\hat{\rho}_g \geq \hat{\rho}_g$.

2. *We have $\hat{\rho}_g \geq \hat{\rho}_g$ if and only if*

$$\frac{\widehat{\text{Cov}}(\tilde{Y}_{new}, Y_{new})}{\widehat{\text{Cov}}(\vec{Y}_{new}, Y_{new})} \leq \sqrt{1 + \frac{\widehat{\text{Var}}(\tilde{Y}_{new})}{\widehat{\text{Var}}(\vec{Y}_{new})}} - 1.$$

In this case, we also have $\hat{\rho}_g \geq \hat{\rho}_g$.

3. *We have $\hat{\rho}_g \geq \hat{\rho}_g$ and $\hat{\rho}_g \geq \hat{\rho}_g$ if and only if*

$$\sqrt{1 + \frac{\widehat{\text{Var}}(\tilde{Y}_{new})}{\widehat{\text{Var}}(\vec{Y}_{new})}} - 1 \leq \frac{\widehat{\text{Cov}}(\tilde{Y}_{new}, Y_{new})}{\widehat{\text{Cov}}(\vec{Y}_{new}, Y_{new})} \leq \frac{\widehat{\text{Var}}(\tilde{Y}_{new})}{\widehat{\text{Var}}(\vec{Y}_{new})} \left(1 + \sqrt{1 + \frac{\widehat{\text{Var}}(\tilde{Y}_{new})}{\widehat{\text{Var}}(\vec{Y}_{new})}} \right).$$

The proof is deferred to the Section 10 of the Supplementary material.

For a given \tilde{r} and a given partition function $\sigma(\cdot)$, the above results allow to choose between the estimators $\hat{\rho}_g$, $\hat{\rho}_g$ and $\hat{\rho}_g$. These conditions are expressed in terms of the ratios of the respective covariances and variances. An opened question is how to choose the best partition among all possible partitions.

We further introduce the following notations : for $i = 1, \dots, 3$,

$$\tilde{\Omega}_i := \Omega_i \cap \{\sigma(1), \dots, \sigma(\tilde{r})\}.$$

We then have the following analogue of Lemma 2 which treats the case when the signal is spread out uniformly among the different subspaces.

Lemma 7. *Let us consider the same assumptions as in Lemma 2. Moreover, we suppose that we have the relation*

$$\sum_{s \in \tilde{\Omega}_1} d_s^2 \sim \sum_{s \in \Omega_1} d_s^2.$$

Then we have $\hat{\rho}_g \rightarrow \rho_g^{oracle}$ and $\hat{\rho}_g \rightarrow \rho_g^{oracle}$.

The proof is given in Section 11 of the Supplementary material. In other words, we have to impose that the L^2 norm of the singular values belonging to $\tilde{\Omega}_1$, respectively to Ω_1 , are equivalent.

In the same way as for the classical Ridge estimator, let us focus on a few extreme cases.

Lemma 8 (Extreme cases). *Let us consider same assumptions as in Theorem 2.*

1. If $1 \in \{\sigma(1), \dots, \sigma(\tilde{r})\}$ and the projected signal belongs only to $\text{Span}\{Q^{(1)}\}$, that is to say

$$\|Q^{(1)}Q^{(1)'}\beta\|^2 \sim n^{2\tau}, \quad \|Q^{(s)}Q^{(s)'}\beta\|^2 = 0, \text{ for } 1 < s \leq r,$$

and moreover $2\tau + \psi < 1$ and the following two conditions hold

- $\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{(d_{\sigma(s)}^2 + \lambda)^2} = o(n^{2\tau + \psi});$
- $n^{2\tau + \psi} = o\left(\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2}\right),$

then $\hat{\rho}_g \rightarrow \rho_g^{\text{oracle}}$, whereas $\hat{\rho}_g \rightarrow 0$.

2. If $r \in \{\sigma(1), \dots, \sigma(\tilde{r})\}$ and the projected signal belongs only to $\text{Span}\{Q^{(r)}\}$, that is to say

$$\|Q^{(r)}Q^{(r)'}\beta\|^2 \sim n^{2\tau}, \quad \|Q^{(s)}Q^{(s)'}\beta\|^2 = 0, \text{ for } 1 \leq s < r,$$

and moreover $\lambda \sim Cn^{\eta + \kappa}$ with $\kappa > \max(0, -\eta)$, $C > 0$, $\tau + \eta/2 - \kappa > 0$ and the following two conditions hold

- $\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{(d_{\sigma(s)}^2 + \lambda)^2} = o(n^{2\tau + \eta - 2\kappa});$
- $n^{2\tau + \eta - 2\kappa} = o\left(\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2}\right),$

then $\hat{\rho}_g \rightarrow \rho_g^{\text{oracle}}$, whereas $\hat{\rho}_g \rightarrow 0$.

The proof is largely inspired from the proof of Lemma 3. According to this lemma, there are a few cases where at the same time, the new accuracy $\hat{\rho}_g$ is optimal and the classical accuracy $\hat{\rho}_g$ is null.

Note that the condition $\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{(d_{\sigma(s)}^2 + \lambda)^2} = o(n^{2\tau + \psi})$ can be replaced by the condition $\tilde{r} = o(n^{2\tau + \psi})$. In the same way, the condition $\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{(d_{\sigma(s)}^2 + \lambda)^2} = o(n^{2\tau + \eta - 2\kappa})$ can be replaced by the condition $\tilde{r} = o(n^{2\tau + \eta - 2\kappa})$.

In Supplementary material we also investigate the same setting as in Theorem 3, when X_{new} is random. Lemmas 3 and 4 in Supplementary material are

the analogues of the previous Theorem 3 and of Lemma 1 in Supplementary material, respectively.

Summary of the results in Section 5: We introduce a new estimator $\tilde{\beta}$, derived from Ridge regression, which may lead to an increase in term of GS accuracy. It consists in projecting the vector Y of phenotypes on a well chosen subspace of the space spanned by the columns of X (i.e. the marker genotypes of the TRN individuals). \tilde{r} denotes the rank of this new subspace. The accuracy will particularly increase when the rank r of X is large and most of the projection of the signal β belongs to the space spanned by the axis associated to the largest singular values of X . In order to find an appropriate subspace, the quantity \hat{A}_1/\hat{A}_1 (resp. \check{A}_1/\hat{A}_1) has to be computed for different values of \tilde{r} , when TRN and TST are (resp. not) sampled from the same population; a procedure for choosing \tilde{r} is described in Section 6.2.4. To sum up, we propose a new genomic predictor of TST individuals that outperforms the classical RRBLUP or equivalently GBLUP. Besides, we derive its accuracy estimates, $\hat{\rho}_g$ or $\check{\rho}_g$ corresponding to TRN and TST individuals sampled from the same population or not sampled in the same population, respectively.

6. Applications

In this section we propose to illustrate our theoretical results with the help of simulated data.

6.1. Simulation framework

Genomic data were generated by means of the R package *hypred* of Technov (2014) and according to the same process as in Rabier et al. (2016). In particular, populations were simulated by random mating between haploid individuals (i.e. with only one copy of each chromosome), during (a) 30, (b) 50, or (c) 70 generations. Recombination was modeled according to Haldane (1919). Recall that Haldane modeling assumes that the number of recombinations follow a standard Poisson process.

In generation zero, two haploid founder lines were crossed. These two lines were completely different genetically, the first (resp. the other) line having allele +1 (resp. -1) at each marker. Generation 1 consisted of (a) 400 or (b) 500 haploid offsprings of these two founders. After that, the population kept evolving by random mating with a constant size at each generation. In the final generation, either 500 individuals or 400 individuals were sampled. Under the 400 offsprings scenario, 2 individuals were randomly selected, and 100 full sibs were generated in order to obtain some closely related individuals (as in classical genomic studies). So, this procedure allows to deal with two kinds of TRN populations, both based on 500 individuals: one containing 100 full sibs, and the other not containing any full sib. The prediction model was evaluated on 100 TST (in all cases) produced in the last generation. Note that our simulated

population is a simplified population, since our individuals are haploid and our two founders present a complete Linkage Disequilibrium.

We focused on one chromosome of length 1 Morgan. We considered 3 different densities of genetic markers equally spaced on the chromosome: (a) 100, (b) 1,000, or (c) 2,000 SNPs. We studied two configurations for the phenotypic model: (a) 2 QTLs located at 3cM and 80cM with effects +1 and -2 , respectively, and (b) 100 QTLs located every centimorgan, with the same effect +0.15. The environmental variance σ_e^2 was set to 1.

In what follows, we focus on the phenotypic accuracy criteria. $\hat{\rho}_{ph}$ and $\check{\rho}_{ph}$ denote the analogue of the quantities $\hat{\rho}_g$ and $\check{\rho}_g$ for the phenotypic accuracy. As in Rabier et al. (2016), we set the value of σ_e^2 to 1 and we consider this true value in the expressions of $\hat{\rho}_{ph}$ and $\check{\rho}_{ph}$. Recall that ρ_{ph} is obtained by replacing the term A_4 by $A_4 + \sigma_e^2$ in our Theorem 1. Indeed, in what follows, since we consider h unknown, we cannot use the relation $\rho_{ph} = h\rho_g$.

The empirical accuracy was computed with the R software, using the empirical correlation between the predicted values and the true values. Note also that all the quantities presented in the different tables are averages based on 100 simulations. Since we analyze the case where X does not vary across replicates, one simulation consists (a) in regenerating 100 TST individuals by random mating between individuals from the penultimate generation, and (b) in regenerating new phenotypes (TRN+TST).

The regularization parameter λ was estimated by REML. The R package *rrBLUP*, and in particular its function *kin.blup* were used in order to compute the variance components.

6.2. Illustrations on simulated data

6.2.1. Different probability distributions

To begin with, we propose to investigate the long-term behavior of GS, i.e. the reliability of the predicted model as a function of time (Habier et al. (2007); Goddard et al. (2009)). For instance, in plants, since a large number of generations can be obtained easily, the fitted model is usually not readjusted at each generation, in order to save time or costs due to genotyping.

In this context, Figure 1 compares different estimators of the phenotypic accuracy as functions of the number of generations during which the TST sample evolved. The TRN sample was always based on 30 generations. We can notice that $\check{\rho}_{ph}(\beta)$ matches the empirical accuracy whatever the number of generations for TST. In contrast, $\hat{\rho}_{ph}(\beta)$ deteriorates overtime. This is as expected, since $\check{\rho}_{ph}(\beta)$ handles explicitly the TST matrix X_{new} , which is not the case of $\hat{\rho}_{ph}(\beta)$ that relies on the TRN matrix X .

Table 1 in Supplementary material considers the same number of generations for TRN and TST and focuses on the case where a few siblings (100 or none) are included in the TRN sample. Recall that when full sibs are incorporated, the TRN and TST samples do not come from the same probability distribution. According to this table, even in the presence of 100 full sibs, we observe a good agreement between the empirical accuracy and estimations based on $\hat{\rho}_{ph}(\beta)$. In

view of Figure 1 and Table 1 in Supplementary material, it seems that the fact of not readjusting the model overtime has more impact on prediction than the presence of full sibs in the TRN set. To sum up, $\check{\rho}_{ph}(\beta)$ appears to be a reliable estimator whatever the simulation framework.

6.2.2. Behavior of the accuracy when β is estimated

In practice, the vector β containing the marker effects is an unknown quantity. Therefore, we propose to consider here different estimators of β suitable in a high-dimensional setting. We concentrate here on the LASSO (Tibshirani (1996)), the Adaptative LASSO (Zou (2006)) and on the Group LASSO (Yuan and Lin (2006)) estimators. Note that other estimators could have been chosen. Recall that the LASSO is a L^1 penalization method, and that the Adaptative LASSO replaces the L^1 penalty by a weighted penalty. Zou (2006) proved that Adaptative LASSO enjoys oracle properties. Last, the Group LASSO estimator differs from his cousins, since it allows to handle a group structure for β . We used the R packages *glmnet*, *parcor* and *gglasso* for computing the LASSO, the Adaptative LASSO and the Group LASSO estimators, respectively.

Tables 1 and 2 focus on the scenario with 2 large QTLs and 100 small QTLs, respectively. According to Table 1, the Adaptative LASSO estimator presents better performances than his cousins, whatever the density of markers and the number of generations. As expected, the best estimators are the ones assuming known β . Note that since the TRN and TST are based on the same number of generations, we did not observe significative differences between $\hat{\rho}_{ph}$ and $\check{\rho}_{ph}$.

Table 2 shows that the accuracy based on LASSO and cousin methods deteriorates slightly for a high density of markers (1,000 or 2,000). This accuracy also decreases when the number of generations increases. In view of the two tables, the Adaptative LASSO estimator is closer to the empirical accuracy under the 2 QTLs scenario. Indeed, when two large QTLs well separated lied on the genome, the Adaptative LASSO method was able to recover perfectly those QTLs, whereas the 100 QTLs scenario makes the signal recovery less trivial.

To complete our simulation study, it is worth to consider the case of a mixture between major genes and multiple small QTLs, which mimics probably better the common architecture for a lot of traits. We thus generated two large QTLs located at 3cM and 80cM, and 98 small QTLs located every centimorgan (except at 3cM and 80cM). We considered three scenarios: (a) large QTLs with effects +0.5 and -0.6, small QTLs with the same effect +0.07, (b) large QTLs with effects +1 and -0.7, small QTLs with the same effect +0.1, (c) large QTLs with effects +2 and -2, small QTLs with the same effect +0.1.

According to Table 2 in Supplementary material, under these new configurations the performances are still fair, even if they deteriorate slightly in presence of a high density of markers.

To conclude, in view of all our results presented in this section, the Adaptative LASSO seems to be the most appropriate method for substituting $\hat{\beta}$ into the expressions of $\hat{\rho}_{ph}$ and $\check{\rho}_{ph}$.

6.2.3. Comparison with existing methods

A large number of formulas for accuracy are available in the literature. One of the most popular was proposed in Daetwyler et al. (2008): the authors assumed that the gene locations are known (i.e. the indices of the non null coefficients of β are perfectly known) and focused on an orthogonal design. A general version of their formula regarding the genotypic accuracy is

$$\sqrt{\frac{h^2/(1-h^2)}{\frac{\|\beta\|_0^0}{n} + \frac{h^2}{1-h^2}}}. \quad (12)$$

Recall that $\|\beta\|_0^0 = \sum_{j=1}^p |\beta_j|^0$, with $(0^0 = 0)$. Later, the authors allowed for the presence in the genome of a large number of loci that cannot be considered independent due to linkage and a fixed genome size. In particular, they proposed to substitute the effective number of independent loci M_e for $\|\beta\|_0^0$ into their original formula. Subsequently, a large number of research groups built on this concept and proposed different ways of estimating M_e .

Table 3 compares the performances of seven different proxies in terms of the phenotypic accuracy. Three of these proxies, the ones based on M_{e1} , M_{e2} and M_{e3} , rely on the effective population size (e.g. Goddard et al. (2011)), whereas the M_{LJ} -based proxy comes from association studies (see Li and Ji (2005)). The expressions of M_{e1} , M_{e2} and M_{e3} are the following:

$$M_{e1} = \frac{2N_e L}{\log(4N_e l)}, M_{e2} = \frac{2N_e L}{\log(2N_e l)}, M_{e3} = \frac{2N_e L}{\log(N_e l)},$$

where L , l and N_e denote the genome length, the average chromosome length and the effective population size, respectively. M_{e1} was proposed by Goddard et al. (2009), whereas M_{e2} and M_{e3} by Goddard et al. (2011). We refer to Rabier et al. (2016) for more details on the estimation of N_e based on Hill and Weir (1998). The fifth proxy is the one introduced in Rabier et al. (2016). Note that the heritability h^2 was estimated with the help of variance components obtained by the R package *rrBLUP*. Last, the remaining proxies are those suggested in our present paper: $\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$ and $\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$.

Table 3 reports the Mean Squared Error (MSE) associated to each method, based on 15 architectures. An architecture refers to a fixed number of: (a) SNPs; (b) QTL numbers, effects, and locations. There are 15 architectures as we considered (a) either 100, 1,000 or 2,000 SNPs, and (b) either 2 large QTLs, 100 small QTLs, or the 3 scenarios of Table 2 in Supplementary material.

According to Table 3, $\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$ and $\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$ are the most competitive proxies. They outperform our recent proxy of Rabier et al. (2016), and other classical proxies used by geneticists. As expected, $\check{\rho}_{ph}(\beta)$ yields the best performances. However, it cannot be computed in practice, since it depends on the unknown β .

6.2.4. The quality of the prediction can be improved

In this subsection we propose to illustrate the quality of the predictions based on $\tilde{\beta}$. Recall that this estimator is built after having projected the vector

Y on a well chosen subspace of the space spanned by the columns of X .

In order to find an appropriate subspace, we used the following procedure. We decided that $\frac{d_{\sigma^{(k)}}^4}{d_{\sigma^{(k)}}^2 + \lambda} \|Q^{(\sigma^{(k)})} Q^{(\sigma^{(k)})'} \beta\|^2$ was the k -th largest term of the sequence $\left(\frac{d_s^4}{d_s^2 + \lambda} \|Q^{(s)} Q^{(s)'} \beta\|^2\right)_{s=1, \dots, r}$. The value of \tilde{r} was chosen as the largest

value satisfying the condition $\hat{A}_1 / \hat{A}_1 \leq v$, where v denotes a tuning parameter. The corresponding accuracy was then computed for a given value of v .

Since v was unknown, we performed an optimization over the grid $\{0.7, 0.8, 0.9, 0.925, 0.95, 0.975, 0.99\}$ and kept the highest accuracy.

Tables 4 and 5 focus on the cases $n = 500$ and $n = 800$, respectively. According to these tables, in all the cases we studied, the quantity $\hat{\rho}_{ph}(\beta)$ was greater than $\hat{\rho}_{ph}(\beta)$. In the same way, the empirical accuracy associated to the new estimator (i.e. estimated correlation $cor(\tilde{Y}_{new}, Y_{new})$), was always greater than the classical empirical accuracy based on the Ridge estimator (i.e. $cor(\hat{Y}_{new}, Y_{new})$).

Last, Table 6 focuses on the case where the vector β belongs to $\mathcal{R}(X)$. In particular, we considered $\beta = 0.3Q^{(1)} + 0.3Q^{(2)} + 0.3Q^{(3)}$. As expected (cf. remark below Lemma 5), $\hat{\rho}_{ph}(\beta)$ takes greater values than $\hat{\rho}_{ph}(\beta)$.

6.3. Real data: GS in rice

To conclude this article, we propose to analyze some data from the recent paper of Spindel et al. (2015) dealing with GS in rice.

We considered the dataset of 13,101 SNPs, randomly chosen by the authors from their 73,147 collected SNPs. We decided to focus on two rice traits: flowering and yield. Besides, our analysis relies on the dry season 2012. 80% of the observations were chosen for the TRN set, and the remaining 20% were affected to the TST set. According to the data, the number of TRN individuals was $n = 252$ for flowering and $n = 248$ for yield. In both cases we considered $n_{new} = 63$. Table 7 shows a comparison of the performances of seven different proxies in terms of phenotypic accuracy. The computed MSE relies on 100 data sets (with random individuals in the TRN and TST sets).

In order to compute proxies based on M_{e1} , M_{e2} and M_{e3} , we used the value $L = 13.101$ for the genome length (from Section ‘‘GS using marker subsets’’ of Spindel et al. (2015)), and $l = 1.09175$ for the average chromosome length. Recall that the rice presents 12 chromosomes. In order to make calculations easier, the effective population size N_e was obtained by using only 1,007 SNPs spread out in the genome (a SNP every 0.012 Morgan). Furthermore, we used the Adaptive LASSO method to compute our suggested proxies, $\hat{\rho}_{ph}(\beta)$ and $\check{\rho}_{ph}(\beta)$. Note that since σ_e^2 was unknown, we considered the estimation of σ_e^2 given by REML.

According to Table 7, $\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$ is the most interesting proxy. Indeed, for flowering and yield, the associated MSE was the smallest among all proxies, and the associated mean accuracies were pretty close to the empirical accuracies (0.5485 vs. 0.5576 for flowering, and 0.2650 vs. 0.3361 for yield). As a consequence, the results presented in this manuscript should be of interest for geneticists.

7. Conclusion

In this manuscript we propose somewhat complementary estimators for the accuracy in GS.

With the help of the quantities $\hat{\rho}_g$ and $\hat{\check{\rho}}_g$, geneticists can now figure out whether or not their population is appropriate for GS. Indeed, before genotyping TST individuals, they can have an idea of the reliability of the future predictions. In contrast, the estimators $\check{\rho}_g$ and $\check{\check{\rho}}_g$ seem to be more appropriate to give answers to a challenging question in GS: the choice of the TRN individuals with respect to the TST set (e.g. Rincent et al. (2012)). Given a set of young candidates recently genotyped, breeders want to know the most informative individuals to be phenotyped among the ones collected over years. In that sense, it is always profitable to handle explicitly the link between TRN and TST populations. From a technical point of view, our study focuses exclusively on the asymptotic properties of $\hat{\rho}_g$ and $\hat{\check{\rho}}_g$. Indeed, we did not investigate the asymptotic properties of the estimators $\check{\rho}_g$ and $\check{\check{\rho}}_g$. Since X_{new} is a random matrix, it would have required to consider explicitly the limiting distribution of the eigenvalues of $X'_{new}X_{new}/n_{new}$ (e.g. Dicker (2016)). This could be investigated in future research. Another topic of interest is the behavior of the estimators when a few individuals are removed from the TRN set. Removing one row of X leads to a compression of the singular values interval: the difference between the largest and the smallest singular values decreases (see for instance Chafaï et al. (2009)). Last, since the Adaptive LASSO estimator was found to be the most appropriate substitute for β , it should be interesting to quantify theoretically the loss of information due to an estimated β .

To conclude, GS is a very fruitful topic for geneticists and statisticians and a large amount of methodological questions remain open.

Supporting information. Additional information for this article is available online

Text S1 : Supplementary material containing several proofs and tables.

References

- Buckler, E.S., Holland, J.B., Bradbury, P.J., Acharya, C.B., Brown, P.J., Browne, C., et al (2009). The genetic architecture of maize flowering time. *Science*. **325**, (5941), 714-718.
- Bühlmann, P. (2013). Statistical significance in high-dimensional linear models. *Bernoulli*. **19**, (4), 1212-1242.
- Burstin, J., Salloignon, P., Martinello, M., Magnin-Robert, J.B., Siol, M., Jacquin, F., et al (2015). Genetic diversity and trait genomic prediction in a pea diversity panel. *BMC genomics*. **16**, (1), 105.

- Cai, T.T., Zhang, C., & Zhou, H.H. (2010). Optimal rates of convergence for covariance matrix estimation. *Ann. Statist.* **38**, (4), 2118-2144.
- Chafaï, D. (2009). Singular values of random matrices. *Lecture Notes*.
- Connolly, S., Anney, R., Gallagher, L., et al (2017). A genome-wide investigation into parent-of-origin effects in autism spectrum disorder identifies previously associated genes including SHANK3. *European Journal of Human Genetics.* **25**, (2), 234-239.
- Corbeil, R.R., & Searle, S.R. (1976). Restricted maximum likelihood (REML) estimation of variance components in the mixed model. *Technometrics.* **18**, (1), 31-38.
- 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.
- Dicker, L. H. (2016). Ridge regression and asymptotic minimax estimation over spheres of growing dimension. *Bernoulli.* **22**, (1), 1-37.
- Durrett, R. (2008). *Probability models for DNA sequence evolution*. Springer Science & Business Media.
- 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.
- Friedman, J., Hastie, T. & Tibshirani, R. (2001). *The elements of statistical learning*, Springer series in statistics Springer, Berlin.
- Fritsche, L. G., Igl, W., Bailey, J. N. C., Grassmann, F., Sengupta, S., Bragg-Gresham, J. L., ... & Kim, I. K. (2016). A large genome-wide association study of age-related macular degeneration highlights contributions of rare and common variants. *Nature genetics*, **48**, (2), 134.
- Goddard, M.E. & Hayes, B.J. (2009). Mapping genes for complex traits in domestic animals and their use in breeding programmes. *Nature Reviews Genetics.* **10**, (6), 381-391.
- Goddard, M. (2009). Genomic selection: prediction of accuracy and maximisation of long term response. *Genetica.* **136**, (2), 245-257.
- Goddard, M.E., Hayes, B.J., & Meuwissen, T.H.E. (2011). Using the genomic relationship matrix to predict the accuracy of genomic selection. *Journal of Animal Breeding and Genetics.* **128**, (6), 409-421.

- Habier, D., Fernando, R. & Dekkers, J. (2007). The impact of genetic relationship information on genome-assisted breeding values. *Genetics*. **177**, (4), 2389-2397.
- Haldane J. (1919). The combination of linkage values and the calculation of distances between the loci of linked factors. *J Genet.* **8**, (29), 299-309.
- Hayes, B., Bowman, P., Chamberlain, A. & Goddard, M. (2009). Invited review: Genomic selection in dairy cattle: Progress and challenges. *Journal of dairy science*. **92**, (2), 433-443.
- Hill, W. & Weir, B. (1998). Variances and covariances of squared linkage disequilibria in finite populations. *Theoretical population biology*. **33**, (1), 54-78.
- Hoerl, A.E. & Kennard, R. W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*. **12**, (1), 55-67.
- Jannink, J.L., Lorenz, A.J. & Iwata, H. (2010). Genomic selection in plant breeding: from theory to practice. *Briefings in functional genomics*. **9**, (2), 166-177.
- Kumar, S., Chagné, D., Bink, M.C., Volz, R.K., Whitworth, C. & Carlisle, C. (2012). Genomic selection for fruit quality traits in apple (*Malus × domestica* Borkh.). *PloS One*. **7**, (5), e36674.
- Lander, E.S. & Botstein, D. (1989). Mapping mendelian factors underlying quantitative traits using RFLP linkage maps. *Genetics*. **121**, (1), 185-199.
- Li, J. & Ji, L. (2005). Adjusting multiple testing in multilocus analyses using the eigenvalues of a correlation matrix. *Heredity*. **95**, (3), 221-227.
- Lynch, M. & Walsh, B. (1998). *Genetics and analysis of quantitative traits*, Sinauer Sunderland, MA.
- 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.
- Rabier, C.E., Barre, P., Asp, T., Charmet, G. & Mangin, B. (2016). On the Accuracy of Genomic Selection. *PloS One*. **11**, (6), e0156086. doi:10.1371/journal.pone.0156086.
- Rincent R, Laloë D, Nicolas S, Altmann T, Brunel D, Revilla P, et al. (2012). Maximizing the Reliability of Genomic Selection by Optimizing the Calibration Set of Reference Individuals: Comparison of Methods in Two Diverse Groups of Maize Inbreds (*Zea mays* L.). *Genetics*. **192**, (2), 715-728.
- Schulz-Streeck, T., Ogutu, J., Karaman, Z., Knaak, C. & Piepho, H. (2012). Genomic selection using multiple populations. *Crop Science*. **52**, (6), 2453-2461.

- Shao, J. & Deng, X. (2012). Estimation in high-dimensional linear models with deterministic design matrices. *The Annals of Statistics*. **40**, (2), 812-831.
- Spindel, J., Begum, H., Akdemir, D., Virk, P., Collard, B., Redoña, E., et al (2015). Genomic Selection and Association Mapping in rice (*Oryza sativa*): Effect of trait genetic architecture, training population composition, marker number and statistical model on accuracy of rice genomic selection in elite, tropical rice breeding lines. *PLoS Genetics*. **11**, (2), e1004982.
- Technow, F. (2014). *R Package hypred: Simulation of Genomic Data in Applied Genetics*. University of Hohenheim.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*. 267-288.
- Tikhonov, A.N. (1963). On the solution of ill-posed problems and the method of regularization. *Dokl. Akad. Nauk.. SSSR* **151**, 501-504.
- 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.
- Whittaker, J.C., Thompson, R., & Denham, M.C. (2000). Marker-assisted selection using ridge regression. *Genetics Research*. **75**, (2), 249-252.
- Würschum, T., Reif, J.C., Kraft, T., Janssen, G. & Zhao, Y. (2013). Genomic selection in sugar beet breeding populations. *BMC genetics*. **14**, (1), 85.
- Yuan, M. & Lin, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society Series B*. **68**, (1), 49-67.
- Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American statistical association*. **101**, (476), 1418-1429.

Charles-Elie Rabier (ce.rabier@gmail.com)

ISEM, Université de Montpellier, CNRS, France.

Brigitte Mangin (brigitte.mangin@inra.fr)

LIPM, Université de Toulouse, INRA, CNRS, Castanet-Tolosan, France

Simona Grusea (grusea@insa-toulouse.fr)

Institut de Mathématiques de Toulouse, Université de Toulouse, INSA de Toulouse, France.

Figure 1: Comparison among different estimators of the phenotypic accuracy as a function of the number of generations during which the TST sample evolved (TRN sample is always based on 30 generations). The chromosome is of length 1M and 2 QTLs are located at 3cM and 80cM with effects +1 and -2, respectively ($n = 500$, $n_{new} = 100$, $\sigma_e^2 = 1$). Emp. Acc. refers to the empirical *phenotypic accuracy*.

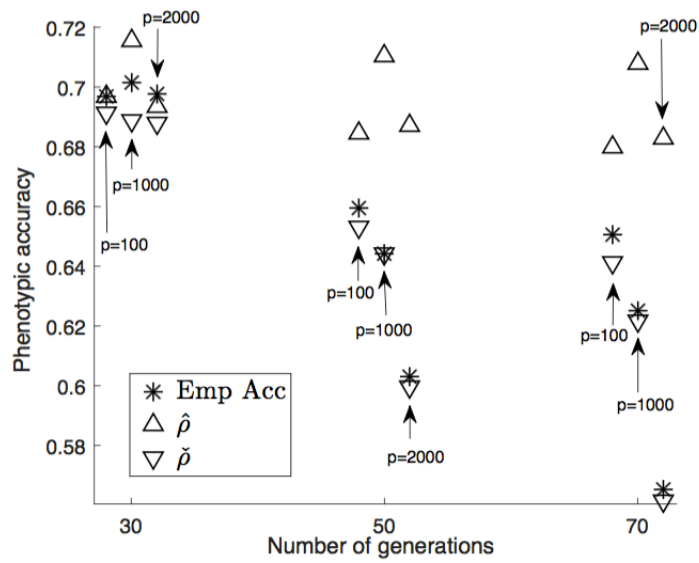


Table 1: Comparison among different estimators of the phenotypic accuracy, in presence of a few major genes ($n = 500$, $n_{new} = 100$, $\sigma_e^2 = 1$). The chromosome is of length 1M and the 2 QTLs are located at 3cM and 80cM with effects +1 and -2, respectively. $\hat{\beta}_{LASSO}$, $\hat{\beta}_{ADLASSO}$, $\hat{\beta}_{GPLASSO}$ refer to the LASSO, Adaptive LASSO and Group LASSO estimators of β , respectively. Emp. Acc. refers to the empirical phenotypic accuracy. The standard errors for the estimates are given in brackets.

Nb markers	Method	30 generations	50 generations	70 generations
100	Emp. Acc.	0.6967 (0.0048)	0.6804 (0.0049)	0.6708 (0.0055)
	$\hat{\rho}_{ph}(\beta)$	0.6969 (0.0001)	0.6765 (0.0001)	0.6717 (0.0001)
	$\hat{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.5962 (0.0028)	0.5767 (0.0029)	0.5735 (0.0030)
	$\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.6927 (0.0017)	0.6675 (0.0018)	0.6676 (0.0020)
	$\hat{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.5934 (0.0028)	0.5595 (0.0029)	0.5484 (0.0031)
	$\check{\rho}_{ph}(\beta)$	0.6915 (0.0015)	0.6731 (0.0015)	0.6654 (0.0017)
	$\check{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.5907 (0.0029)	0.5742 (0.0029)	0.5677 (0.0036)
	$\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.6872 (0.0020)	0.6712 (0.0022)	0.6614 (0.0027)
	$\check{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.5857 (0.0030)	0.5580 (0.0031)	0.5411 (0.0037)
1,000	Emp. Acc.	0.7015 (0.0047)	0.6683 (0.0053)	0.6713 (0.0043)
	$\hat{\rho}_{ph}(\beta)$	0.7155 (0.0001)	0.6597 (0.0001)	0.685 (0.0001)
	$\hat{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.6197 (0.0027)	0.5354 (0.0031)	0.5720 (0.0026)
	$\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.7066 (0.0015)	0.6488 (0.0017)	0.675 (0.0016)
	$\hat{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.6244 (0.0025)	0.5471 (0.0029)	0.586 (0.0023)
	$\check{\rho}_{ph}(\beta)$	0.6889 (0.0017)	0.6576 (0.0021)	0.6642 (0.0023)
	$\check{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.5965 (0.0030)	0.5347 (0.0037)	0.5544 (0.0038)
	$\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.6812 (0.0021)	0.6454 (0.0027)	0.6548 (0.0030)
	$\check{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.6022 (0.0028)	0.5495 (0.0034)	0.5708 (0.0035)
2,000	Emp. Acc.	0.6977 (0.0055)	0.6316 (0.0060)	0.4174 (0.0077)
	$\hat{\rho}_{ph}(\beta)$	0.6933 (0.0001)	0.6254 (0.0001)	0.4600 (0.0004)
	$\hat{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.5872 (0.0030)	0.4794 (0.0042)	0.2790 (0.0044)
	$\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.6783 (0.0017)	0.6134 (0.0020)	0.4399 (0.0033)
	$\hat{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.5904 (0.0027)	0.4814 (0.0035)	0.2801 (0.0045)
	$\check{\rho}_{ph}(\beta)$	0.6881 (0.0016)	0.6264 (0.0025)	0.4095 (0.0041)
	$\check{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.5842 (0.0034)	0.4831 (0.0042)	0.2522 (0.0050)
	$\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.6830 (0.0029)	0.6138 (0.0031)	0.3890 (0.0058)
	$\check{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.5902 (0.0032)	0.4878 (0.0042)	0.2601 (0.0050)

Table 2: Comparison among different estimators of the phenotypic accuracy, in presence of multiple small QTLs ($n = 500$, $n_{new} = 100$, $\sigma_e^2 = 1$). The chromosome is of length 1M and 100 QTLs with the same effect +0.15, are located every centimorgan. Same notations as in Table 1.

Nb markers	Method	30 generations	50 generations	70 generations
100	Emp. Acc.	0.8504 (0.0027)	0.8055 (0.0035)	0.7056 (0.0049)
	$\hat{\rho}_{ph}(\beta)$	0.8346 (0.0001)	0.8007 (0.0001)	0.6938 (0.0001)
	$\hat{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.7990 (0.0013)	0.7010 (0.0013)	0.6043 (0.0034)
	$\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.8366 (0.0007)	0.8036 (0.0008)	0.6998 (0.0017)
	$\hat{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.7813 (0.0012)	0.7370 (0.0014)	0.5611 (0.0026)
	$\check{\rho}_{ph}(\beta)$	0.8434 (0.0015)	0.7941 (0.0020)	0.6981 (0.0026)
	$\check{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.8020 (0.0024)	0.7471 (0.0025)	0.6006 (0.0049)
	$\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.8426 (0.0018)	0.7959 (0.0019)	0.7029 (0.0033)
	$\check{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.7889 (0.0023)	0.7250 (0.0025)	0.5611 (0.0037)
1,000	Emp. Acc.	0.8700 (0.0022)	0.8143 (0.0036)	0.7233 (0.0048)
	$\hat{\rho}_{ph}(\beta)$	0.8781 (0.0000)	0.8086 (0.0001)	0.7308 (0.0001)
	$\hat{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.8558 (0.0007)	0.7635 (0.0017)	0.6532 (0.0027)
	$\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.8495 (0.0006)	0.7627 (0.0012)	0.6718 (0.0016)
	$\hat{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.8508 (0.0006)	0.7581 (0.0014)	0.6466 (0.0023)
	$\check{\rho}_{ph}(\beta)$	0.8604 (0.0013)	0.8045 (0.0019)	0.7162 (0.0027)
	$\check{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.8299 (0.0019)	0.7502 (0.0026)	0.6233 (0.0039)
	$\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.8226 (0.0018)	0.7489 (0.0023)	0.6452 (0.0034)
	$\check{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.8273 (0.0018)	0.7479 (0.0024)	0.6224 (0.0036)
2,000	Emp. Acc.	0.8590 (0.0024)	0.8045 (0.0043)	0.7387 (0.0047)
	$\hat{\rho}_{ph}(\beta)$	0.8464 (0.0001)	0.8113 (0.0001)	0.7319 (0.0001)
	$\hat{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.8116 (0.0011)	0.7662 (0.0014)	0.6503 (0.0030)
	$\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.8102 (0.0007)	0.7641 (0.0009)	0.6697 (0.0016)
	$\hat{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.8062 (0.0009)	0.7607 (0.0012)	0.6495 (0.0024)
	$\check{\rho}_{ph}(\beta)$	0.8510 (0.0015)	0.7936 (0.0023)	0.7300 (0.0026)
	$\check{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.8096 (0.0023)	0.7339 (0.0035)	0.6317 (0.0041)
	$\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.8093 (0.0020)	0.7358 (0.0033)	0.6542 (0.0031)
	$\check{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.8074 (0.0020)	0.7322 (0.0033)	0.6364 (0.0039)

Table 3: Mean squared error (with respect to the Empirical accuracy) corresponding to 7 proxies. The MSE corresponding to $\check{\rho}_{ph}(\beta)$ is also shown. $MSE = \sum_{a=1}^{15} (\text{AccP}_a - \text{AccE}_a)^2 / 15$ where 15 is the number of studied architectures. AccE_a and AccP_a are averages on 100 replicates, and denote respectively, for architecture a , the Empirical Accuracy and the Accuracy based on the chosen proxy (30 generations for TRN).

MSE based on	50 generations for TST	70 generations for TST
$\check{\rho}_{ph}(\beta)$	5.9685×10^{-5}	3.8455×10^{-5}
$\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	1.2108×10^{-3}	1.2118×10^{-3}
$\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	2.2677×10^{-3}	1.5168×10^{-3}
Plos One (2016)	3.3056×10^{-3}	1.007×10^{-2}
M_{e1}	3.7936×10^{-3}	1.3779×10^{-2}
M_{e2}	3.7508×10^{-3}	1.3518×10^{-2}
M_{e3}	3.6970×10^{-3}	1.3165×10^{-2}
M_{LJ}	5.5578×10^{-3}	6.1021×10^{-3}

Table 4: Illustration of the predictions based on $\tilde{\beta}$. $cor(\hat{Y}_{new}, Y_{new})$ (resp. $cor(\tilde{Y}_{new}, Y_{new})$) refers to the empirical correlation between \hat{Y}_{new} (resp. \tilde{Y}_{new}) and Y_{new} . The chromosome is of length 4M and 100 QTLs with the same effect +0.15, are located every centimorgan on [0, 1M]. 4,000 markers (p=4,000) are equally spaced on [0, 4M] ($n = 500$, $n_{new} = 100$). The standard errors for the estimates are given in brackets.

σ_e^2	Method	50 generations	100 generations
1	$cor(\hat{Y}_{new}, Y_{new})$	0.7478 (0.0049)	0.5959 (0.0074)
	$cor(\tilde{Y}_{new}, Y_{new})$	0.7682 (0.0048)	0.6132 (0.0068)
	$\hat{\rho}_{ph}(\beta)$	0.7399 (0.0002)	0.6352 (0.0002)
	$\hat{\rho}_{ph}(\beta)$	0.7570 (0.0001)	0.6541 (0.0001)
9	$cor(\hat{Y}_{new}, Y_{new})$	0.2874 (0.0086)	0.1949 (0.0098)
	$cor(\tilde{Y}_{new}, Y_{new})$	0.3152 (0.0087)	0.2163 (0.0099)
	$\hat{\rho}_{ph}(\beta)$	0.3023 (0.0003)	0.2320 (0.0003)
	$\hat{\rho}_{ph}(\beta)$	0.3306 (0.0002)	0.2604 (0.0002)

Table 5: Same as Table 4 except that $n = 800$.

σ_e^2	Method	50 generations	100 generations
1	$cor(\hat{Y}_{new}, Y_{new})$	0.7908 (0.0041)	0.6101 (0.0063)
	$cor(\tilde{Y}_{new}, Y_{new})$	0.8087 (0.0036)	0.6289 (0.0053)
	$\hat{\rho}_{ph}(\beta)$	0.7965 (0.0000)	0.6508 (0.0001)
	$\hat{\rho}_{ph}(\beta)$	0.8095 (0.0005)	0.6663 (0.0001)
9	$cor(\hat{Y}_{new}, Y_{new})$	0.3725 (0.0097)	0.1981 (0.0106)
	$cor(\tilde{Y}_{new}, Y_{new})$	0.4044 (0.0093)	0.2302 (0.0108)
	$\hat{\rho}_{ph}(\beta)$	0.3766 (0.0003)	0.2248 (0.0032)
	$\hat{\rho}_{ph}(\beta)$	0.4041 (0.0001)	0.2494 (0.0035)

Table 6: Comparison among the quantities $\hat{\rho}_{ph}(\beta)$ and $\tilde{\hat{\rho}}_{ph}(\beta)$, when the vector β belongs to $\mathcal{R}(X)$. The chromosome is of length 1M, $\beta = 0.3Q^{(1)} + 0.3Q^{(2)} + 0.3Q^{(3)}$ and $n_{new} = 100$. The standard errors for the estimates are given in brackets.

σ_e^2	n	Method	200 generations	400 generations
1	500	$\hat{\rho}_{ph}(\beta)$	0.7550 (0.0001)	0.6721 (0.0002)
		$\tilde{\hat{\rho}}_{ph}(\beta)$	0.7810 (0.0001)	0.7041 (0.0001)
	800	$\hat{\rho}_{ph}(\beta)$	0.7487 (0.0001)	0.7037 (0.0001)
		$\tilde{\hat{\rho}}_{ph}(\beta)$	0.7728 (0.0001)	0.7312 (0.0001)
9	500	$\hat{\rho}_{ph}(\beta)$	0.3370 (0.0004)	0.2623 (0.0004)
		$\tilde{\hat{\rho}}_{ph}(\beta)$	0.3809 (0.0001)	0.3079 (0.0001)
	800	$\hat{\rho}_{ph}(\beta)$	0.3317 (0.0003)	0.2904 (0.0004)
		$\tilde{\hat{\rho}}_{ph}(\beta)$	0.3734 (0.0001)	0.3330 (0.0000)

Table 7: Mean squared error (with respect to the Empirical accuracy) corresponding to 7 proxies, and based on rice data from Spindel et al. (2015) (dry season 2012). The computed MSE rely on 100 data sets (with random individuals in TRN and TST sets). The average, for each proxy, is given in brackets. The Empirical accuracy was 0.5576 for flowering, and 0.3361 for yield, in average ($n = 252$ for flowering, $n = 248$ for yield, $n_{new} = 63$ in both cases).

MSE based on	Flowering	Yield
$\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	1.6248×10^{-2} (0.5485)	2.807×10^{-2} (0.2650)
$\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	2.41×10^{-2} (0.6201)	4.85×10^{-2} (0.4571)
Plos One (2016)	7.08×10^{-2} (0.7903)	1.25×10^{-1} (0.6647)
M_{e1}	4.49×10^{-2} (0.7055)	5.70×10^{-2} (0.5234)
M_{e2}	4.18×10^{-2} (0.6917)	5.10×10^{-2} (0.5064)
M_{e3}	3.83×10^{-2} (0.6741)	4.43×10^{-2} (0.4854)
M_{LJ}	4.71×10^{-2} (0.7142)	6.27×10^{-2} (0.5383)

Text S1: Supplementary material of “On the accuracy
in high dimensional linear models and its application to
genomic selection”

C.E. Rabier^{a,b,c,d}, B. Mangin^e, S. Grusea^c

^a*ISEM, Université de Montpellier, CNRS, France*

^b*LIRMM, Université de Montpellier, CNRS, France*

^c*INSA de Toulouse, Institut de Mathématiques de Toulouse, Université de Toulouse, France*

^d*MIAT, Université de Toulouse, INRA, Castanet-Tolosan, France*

^e*LIPM, Université de Toulouse, INRA, CNRS, Castanet-Tolosan, France*

1. Proof of Theorem 1 of the main manuscript

By definition,

$$A_1 = \beta' \text{Var}(x_{new}) X' V^{-1} X \beta.$$

We set $\bar{D} = \text{Diag}\left(\frac{d_1}{d_1^2 + \lambda}, \dots, \frac{d_r}{d_r^2 + \lambda}\right)$. With this notation, we have the relation:

$$X' V^{-1} = Q \bar{D} P'. \quad (1)$$

Using formula (8) of the main manuscript, we easily have

$$X' V^{-1} X \beta = \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)} Q^{(s)'} \beta. \quad (2)$$

As a consequence, since $\Sigma = \mathbb{E}(x_{new} x_{new}')$, we have the relationship

$$A_1 = \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \beta' \Sigma Q^{(s)} Q^{(s)'} \beta. \quad (3)$$

By definition,

$$A_2 = \sigma_e^2 \mathbb{E}\left(\|x'_{new} X' V^{-1}\|^2\right).$$

According to formula (1), we have

$$\begin{aligned} \|x'_{new} X' V^{-1}\|^2 &= x'_{new} X' V^{-1} (X' V^{-1})' x_{new} \\ &= x'_{new} Q \bar{D} P' P \bar{D} Q' x_{new} \\ &= x'_{new} Q \bar{D}^2 Q' x_{new}. \end{aligned}$$

Furthermore, we have

$$Q\bar{D}^2Q' = \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} Q^{(s)}Q^{(s)'}$$

Since $Q^{(s)}Q^{(s)'}$ is an idempotent matrix, we obtain

$$\begin{aligned} \|x'_{new}X'V^{-1}\|^2 &= \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} x'_{new}Q^{(s)}Q^{(s)'}x_{new} \\ &= \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} x'_{new}Q^{(s)}Q^{(s)'}Q^{(s)}Q^{(s)'}x_{new} \\ &= \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \|Q^{(s)}Q^{(s)'}x_{new}\|^2. \end{aligned}$$

Finally,

$$A_2 = \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \mathbb{E} \left(\|Q^{(s)}Q^{(s)'}x_{new}\|^2 \right).$$

By definition,

$$A_3 = \beta'X'V^{-1}X\text{Var}(x_{new})X'V^{-1}X\beta.$$

Then, according to formula (2),

$$A_3 = \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)}Q^{(s)'}\beta \right)' \Sigma \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)}Q^{(s)'}\beta \right).$$

2. Proof of Theorem 2 of the main manuscript

Let us define \hat{A}_1 in the following way:

$$\hat{A}_1 = \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \beta' \hat{\Sigma} Q^{(s)}Q^{(s)'}\beta,$$

where $\hat{\Sigma} := X'X/n$ is the empirical covariance matrix.

Then, using the SVD decomposition $X = PDQ'$, we obtain

$$\begin{aligned} \hat{A}_1 &= \frac{1}{n} \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \beta' X'X Q^{(s)}Q^{(s)'}\beta \\ &= \frac{1}{n} \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \beta' QD^2Q' Q^{(s)}Q^{(s)'}\beta \\ &= \frac{1}{n} \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \beta' \left(\sum_{u=1}^r d_u^2 Q^{(u)}Q^{(u)'} \right) Q^{(s)}Q^{(s)'}\beta. \end{aligned}$$

Since $Q'Q = I_r$, we further deduce

$$\begin{aligned}\hat{A}_1 &= \frac{1}{n} \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \beta' d_s^2 Q^{(s)} Q^{(s)'} Q^{(s)} Q^{(s)'} \beta \\ &= \frac{1}{n} \sum_{s=1}^r \frac{d_s^4}{d_s^2 + \lambda} \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2.\end{aligned}$$

A natural estimation of A_2 is

$$\begin{aligned}\hat{A}_2 &= \frac{\sigma_e^2}{n} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \sum_{i=1}^n \left\| Q^{(s)} Q^{(s)'} x_i \right\|^2 \\ &= \frac{\sigma_e^2}{n} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \text{Tr} \left(X Q^{(s)} Q^{(s)'} Q^{(s)} Q^{(s)'} X' \right) \\ &= \frac{\sigma_e^2}{n} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \text{Tr} \left(X Q^{(s)} Q^{(s)'} X' \right) \\ &= \frac{\sigma_e^2}{n} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \text{Tr} \left(P D Q' Q^{(s)} Q^{(s)'} Q D P' \right).\end{aligned}$$

Note that

$$D Q' Q^{(s)} = d_s e_s,$$

where e_s denotes the s -th vector of the canonical basis of \mathbb{R}^r .

$$\begin{aligned}\hat{A}_2 &= \frac{\sigma_e^2}{n} \sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} \text{Tr} (P e_s e_s' P') \\ &= \frac{\sigma_e^2}{n} \sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} \text{Tr} (P' P e_s e_s') \\ &= \frac{1}{n} \sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2}.\end{aligned}$$

Let us consider the following estimation of A_3 :

$$\begin{aligned}\hat{A}_3 &= \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)} Q^{(s)'} \beta \right)' \hat{\Sigma} \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)} Q^{(s)'} \beta \right) \\ &= \frac{1}{n} \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)} Q^{(s)'} \beta \right)' X' X \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)} Q^{(s)'} \beta \right) \\ &= \frac{1}{n} \left(X \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)} Q^{(s)'} \beta \right)' \left(X \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)} Q^{(s)'} \beta \right).\end{aligned}$$

Note that

$$XQ^{(s)}Q^{(s)'}\beta = PDQ'Q^{(s)}Q^{(s)'}\beta = d_sPe_sQ^{(s)'}\beta = d_sP^{(s)}Q^{(s)'}\beta.$$

As a consequence,

$$\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} XQ^{(s)}Q^{(s)'}\beta = \sum_{s=1}^r \frac{d_s^3}{d_s^2 + \lambda} P^{(s)}Q^{(s)'}\beta.$$

Last, we obtain

$$\begin{aligned} \hat{A}_3 &= \frac{1}{n} \left(\sum_{\ell=1}^r \frac{d_\ell^3}{d_\ell^2 + \lambda} \beta' Q^{(\ell)} P^{(\ell)'} \right) \left(\sum_{s=1}^r \frac{d_s^3}{d_s^2 + \lambda} P^{(s)} Q^{(s)'} \beta \right) \\ &= \frac{1}{n} \sum_{\ell=1}^r \frac{d_\ell^3}{d_\ell^2 + \lambda} \sum_{s=1}^r \frac{d_s^3}{d_s^2 + \lambda} \beta' Q^{(\ell)} P^{(\ell)'} P^{(s)} Q^{(s)'} \beta \\ &= \frac{1}{n} \sum_{\ell=1}^r \frac{d_\ell^6}{(d_\ell^2 + \lambda)^2} \beta' Q^{(\ell)} Q^{(\ell)'} \beta \\ &= \frac{1}{n} \sum_{\ell=1}^r \frac{d_\ell^6}{(d_\ell^2 + \lambda)^2} \left\| Q^{(\ell)} Q^{(\ell)'} \beta \right\|^2. \end{aligned}$$

Finally, let us consider the following estimation of A_4 :

$$\hat{A}_4 = \beta' \hat{\Sigma} \beta = \frac{1}{n} \beta' X' X \beta.$$

We have

$$\begin{aligned} \hat{A}_4 &= \frac{1}{n} \beta' Q D^2 Q' \beta = \frac{1}{n} \sum_{s=1}^r d_s^2 \beta' Q^{(s)} Q^{(s)'} \beta \\ &= \frac{1}{n} \sum_{s=1}^r d_s^2 \beta' Q^{(s)} Q^{(s)'} Q^{(s)} Q^{(s)'} \beta = \frac{1}{n} \sum_{s=1}^r d_s^2 \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2. \end{aligned}$$

3. Proof of Lemma 1 of the main manuscript

$$\begin{aligned} \hat{A}_1 &= \frac{1}{n} \sum_{s=1}^r \frac{d_s^4}{d_s^2 + \lambda} \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2 \\ &= \frac{1}{n} \sum_{s=1}^r \left(\frac{d_s^3}{d_s^2 + \lambda} \left\| Q^{(s)} Q^{(s)'} \beta \right\| \right) \left(d_s \left\| Q^{(s)} Q^{(s)'} \beta \right\| \right) \\ &\leq \frac{1}{n} \left(\sum_{s=1}^r \frac{d_s^6}{(d_s^2 + \lambda)^2} \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2 \right)^{1/2} \left(\sum_{s=1}^r d_s^2 \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2 \right)^{1/2} \\ &= \hat{A}_3^{1/2} \hat{A}_4^{1/2}, \end{aligned}$$

using the Cauchy-Schwartz inequality. Since $\hat{A}_2 \geq 0$, we obtain

$$\hat{\rho}_g \leq \frac{\hat{A}_1}{\hat{A}_3^{1/2} \hat{A}_4^{1/2}} \leq 1.$$

In order to obtain the lower bound, we just have to notice that

$$\|QQ'\beta\|^2 = \sum_{s=1}^r \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2.$$

Then,

$$\begin{aligned} n\hat{A}_1 &= \sum_{s=1}^r \frac{d_s^4}{d_s^2 + \lambda} \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2 \geq \|QQ'\beta\|^2 \min_s \frac{d_s^4}{d_s^2 + \lambda} \\ n\hat{A}_3 &= \sum_{s=1}^r \frac{d_s^6}{(d_s^2 + \lambda)^2} \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2 \leq \|QQ'\beta\|^2 \max_s \frac{d_s^6}{(d_s^2 + \lambda)^2} \\ n\hat{A}_4 &= \sum_{s=1}^r d_s^2 \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2 \leq \|QQ'\beta\|^2 \max_s d_s^2. \end{aligned}$$

Since $\frac{d_s^4}{(d_s^2 + \lambda)^2}$ is bounded by one, we have $n\hat{A}_2 = \sigma_e^2 \sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} \leq \sigma_e^2 r$. As a consequence,

$$\hat{\rho}_g \geq \frac{\|QQ'\beta\|^2 \min_s \frac{d_s^4}{d_s^2 + \lambda}}{\sqrt{\sigma_e^2 r + \|QQ'\beta\|^2 \max_s \frac{d_s^6}{(d_s^2 + \lambda)^2}} \sqrt{\|QQ'\beta\|^2 \max_s d_s^2}}.$$

4. Proof of Lemma 2 of the main manuscript

Using Theorem 2 of the main manuscript, we have:

$$n\hat{A}_1 \sim \sum_{s \in \Omega_1} d_s^2 \frac{n^{2\tau}}{r} + \sum_{s \in \Omega_2} \frac{d_s^4}{d_s^2 + C_s d_s^2} \frac{n^{2\tau}}{r} + \sum_{s \in \Omega_3} \frac{d_s^4}{\lambda} \frac{n^{2\tau}}{r}.$$

According to our conditions (C3) and (C4),

$$\sum_{s \in \Omega_3} \frac{d_s^4}{\lambda} \frac{n^{2\tau}}{r} = o(1).$$

Then,

$$n\hat{A}_1 \sim \sum_{s \in \Omega_1} d_s^2 \frac{n^{2\tau}}{r} + \sum_{s \in \Omega_2} \frac{d_s^2}{1 + C_s} \frac{n^{2\tau}}{r}. \quad (4)$$

We have

$$\sum_{s \in \Omega_2} \frac{d_s^2}{1+C_s} \frac{n^{2\tau}}{r} \leq \sum_{s \in \Omega_2} d_s^2 \frac{n^{2\tau}}{r} \leq \frac{n^{2\tau}}{r} \#\Omega_2 \tilde{C} \lambda,$$

with $\tilde{C} > 0$.

Since $\#\Omega_2 = O(1)$ by (C6) and $\lambda \frac{n^{2\tau}}{r} = o(1)$, we have

$$\sum_{s \in \Omega_2} d_s^2 \frac{n^{2\tau}}{r} = o(1) \quad (5)$$

and thus $\sum_{s \in \Omega_2} \frac{d_s^2}{1+C_s} \frac{n^{2\tau}}{r} = o(1)$. Therefore

$$n\hat{A}_1 \sim \sum_{s \in \Omega_1} d_s^2 \frac{n^{2\tau}}{r}. \quad (6)$$

In the same way, using condition (C3), we have

$$n\hat{A}_2 \sim \sigma_e^2 \#\Omega_1 + \sigma_e^2 \sum_{s \in \Omega_2} \frac{1}{(1+C_s)^2}.$$

Let us now focus on the quantity \hat{A}_3 .

$$n\hat{A}_3 \sim \sum_{s \in \Omega_1} d_s^2 \frac{n^{2\tau}}{r} + \sum_{s \in \Omega_2} \frac{d_s^2}{(1+C_s)^2} \frac{n^{2\tau}}{r} + \sum_{s \in \Omega_3} \frac{d_s^6}{\lambda^2} \frac{n^{2\tau}}{r}.$$

Since $\sum_{s \in \Omega_3} d_s^6 \leq \sum_{s \in \Omega_3} d_s^2 \sum_{s \in \Omega_3} d_s^4$, we have $\sum_{s \in \Omega_3} d_s^6 = o(\lambda^3)$ (cf. (C2) and (C3)). Then, according to (C4), $\sum_{s \in \Omega_3} \frac{d_s^6}{\lambda^2} \frac{n^{2\tau}}{r} = o(1)$. This yields

$$n\hat{A}_2 + n\hat{A}_3 \sim \sigma_e^2 \#\Omega_1 + \sigma_e^2 \sum_{s \in \Omega_2} \frac{1}{(1+C_s)^2} + \sum_{s \in \Omega_1} d_s^2 \frac{n^{2\tau}}{r} + \sum_{s \in \Omega_2} \frac{d_s^2}{(1+C_s)^2} \frac{n^{2\tau}}{r}.$$

We further have

$$\sum_{s \in \Omega_2} \frac{d_s^2}{(1+C_s)^2} \frac{n^{2\tau}}{r} \leq \sum_{s \in \Omega_2} d_s^2 \frac{n^{2\tau}}{r}.$$

Using the previous relation (5), we have $\sum_{s \in \Omega_2} \frac{d_s^2}{(1+C_s)^2} \frac{n^{2\tau}}{r} = o(1)$. As a result,

$$n\hat{A}_2 + n\hat{A}_3 \sim \sigma_e^2 \#\Omega_1 + \sigma_e^2 \sum_{s \in \Omega_2} \frac{1}{(1+C_s)^2} + \sum_{s \in \Omega_1} d_s^2 \frac{n^{2\tau}}{r}.$$

Then, conditions (C1), (C5) and (C6) ensure that

$$n\hat{A}_2 + n\hat{A}_3 \sim \sum_{s \in \Omega_1} d_s^2 \frac{n^{2\tau}}{r}. \quad (7)$$

Last,

$$n\hat{A}_4 \sim \sum_{s \in \Omega_1} d_s^2 \frac{n^{2\tau}}{r} + \sum_{s \in \Omega_2} d_s^2 \frac{n^{2\tau}}{r} + \sum_{s \in \Omega_3} d_s^2 \frac{n^{2\tau}}{r}.$$

According to conditions (C4) and (C2), $\sum_{s \in \Omega_3} d_s^2 \frac{n^{2\tau}}{r} = o(1)$. Using again the relation (5) we deduce

$$n\hat{A}_4 \sim \sum_{s \in \Omega_1} d_s^2 \frac{n^{2\tau}}{r}. \quad (8)$$

To conclude, using formulae (6), (7) and (8), we have $\hat{\rho}_g \rightarrow 1$.

5. Proof of Lemma 3 of the main manuscript

5.1. The projected signal belongs only to $\text{Span}\{Q^{(1)}\}$

Using Theorem 2, we have:

$$\hat{\rho}_g = \frac{\frac{d_1^3}{d_1^2 + \lambda} \|Q^{(1)}Q^{(1)'}\beta\|}{\left(\sigma_e^2 \sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} + \frac{d_1^6}{(d_1^2 + \lambda)^2} \|Q^{(1)}Q^{(1)'}\beta\|^2\right)^{1/2}}. \quad (9)$$

From Lemma 1 and the fact that $\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} \leq r \leq n$, we deduce that

$$1 \geq \hat{\rho}_g \geq \frac{\frac{d_1^3}{d_1^2 + \lambda} \|Q^{(1)}Q^{(1)'}\beta\|}{\left(\sigma_e^2 n + \frac{d_1^6}{(d_1^2 + \lambda)^2} \|Q^{(1)}Q^{(1)'}\beta\|^2\right)^{1/2}}. \quad (10)$$

Using further the fact that $d_1^2 \sim n^\psi$ and $\lambda = o(d_1^2)$, we obtain

$$\frac{d_1^6}{(d_1^2 + \lambda)^2} \|Q^{(1)}Q^{(1)'}\beta\|^2 \sim n^{2\tau + \psi}, \quad \frac{d_1^3}{d_1^2 + \lambda} \|Q^{(1)}Q^{(1)'}\beta\| \sim n^{\tau + \psi/2}.$$

If $2\tau + \psi > 1$, then

$$\frac{\frac{d_1^3}{d_1^2 + \lambda} \|Q^{(1)}Q^{(1)'}\beta\|}{\left(\sigma_e^2 n + \frac{d_1^6}{(d_1^2 + \lambda)^2} \|Q^{(1)}Q^{(1)'}\beta\|^2\right)^{1/2}} \rightarrow 1.$$

Finally, according to formula (10), $\hat{\rho}_g \rightarrow 1$.

Let us now consider the case $2\tau + \psi < 1$. Then, it is obvious from expression (9) that $\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} = o(n^{2\tau + \psi})$ entails $\hat{\rho}_g \rightarrow 1$.

In contrast, if $n^{2\tau + \psi} = o\left(\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2}\right)$ then $\hat{\rho}_g \rightarrow 0$.

5.2. *The projected signal belongs only to Span* $\{Q^{(r)}\}$

Using again Theorem 2, we have:

$$\hat{\rho}_g = \frac{\frac{d_r^3}{d_r^2 + \lambda} \|Q^{(r)} Q^{(r)'} \beta\|}{\left(\sigma_e^2 \sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} + \frac{d_r^6}{(d_r^2 + \lambda)^2} \|Q^{(r)} Q^{(r)'} \beta\|^2 \right)^{1/2}}. \quad (11)$$

Recall that $d_r^2 \sim n^\eta$ with $\eta < \psi \leq 1$. If we suppose moreover that $\lambda \sim Cn^{\kappa+\eta}$ with $\kappa > \max(0, -\eta)$ and $C > 0$, then we have

$$\begin{aligned} \frac{d_r^3}{d_r^2 + \lambda} &= \frac{d_r}{1 + \lambda/d_r^2} \sim \frac{1}{C} n^{\eta/2 - \kappa} \\ \frac{d_r^3}{d_r^2 + \lambda} \|Q^{(r)} Q^{(r)'} \beta\| &\sim \frac{1}{C} n^{\tau + \eta/2 - \kappa}. \end{aligned}$$

It is obvious that $\hat{\rho}_g \rightarrow 0$ when $\tau + \eta/2 - \kappa < 0$. Indeed, since $d_1^2 = o(n)$, in the denominator we have asymptotically $\sigma_e^2 \frac{d_1^4}{(d_1^2 + \lambda)^2} \sim \sigma_e^2$, which is bounded away from 0.

If $\tau + \eta/2 - \kappa > 0$, then we have to separate two different cases. If $\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2} = o(n^{2\tau + \eta - 2\kappa})$, then $\hat{\rho}_g \rightarrow 1$.

In contrast, if $n^{2\tau + \eta - 2\kappa} = o\left(\sum_{s=1}^r \frac{d_s^4}{(d_s^2 + \lambda)^2}\right)$, then $\hat{\rho}_g \rightarrow 0$.

6. Proof of Theorem 3 of the main manuscript

Let us consider the following natural estimator of A_1 :

$$\check{A}_1 = \frac{1}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \beta' X'_{new} X_{new} Q^{(s)} Q^{(s)'} \beta.$$

We have

$$\begin{aligned} \check{A}_1 &= \frac{1}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \beta' Z F^2 Z' Q^{(s)} Q^{(s)'} \beta \\ &= \frac{1}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \beta' \sum_{\alpha=1}^{r_{new}} f_\alpha^2 Z^{(\alpha)} Z^{(\alpha)'} Q^{(s)} Q^{(s)'} \beta. \end{aligned}$$

Further, a natural estimator of A_2 is

$$\begin{aligned} \check{A}_2 &= \frac{\sigma_e^2}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \text{Tr} \left(X_{new} Q^{(s)} Q^{(s)'} Q^{(s)} Q^{(s)'} X'_{new} \right) \\ &= \frac{\sigma_e^2}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \text{Tr} \left(W F Z' Q^{(s)} Q^{(s)'} Z F W' \right). \end{aligned}$$

We can easily see that

$$\text{Tr} \left(W F Z' Q^{(s)} Q^{(s)'} Z F W' \right) = \sum_{i=1}^{n_{new}} \left(\sum_{\alpha=1}^{r_{new}} f_{\alpha} Q^{(s)'} Z^{(\alpha)} W_i^{(\alpha)} \right)^2,$$

which gives

$$\check{A}_2 = \frac{\sigma_e^2}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \sum_{i=1}^{n_{new}} \left(\sum_{\alpha=1}^{r_{new}} f_{\alpha} Q^{(s)'} Z^{(\alpha)} W_i^{(\alpha)} \right)^2.$$

A natural estimator of A_3 is:

$$\begin{aligned} \check{A}_3 &= \frac{1}{n_{new}} \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)} Q^{(s)'} \beta \right)' X'_{new} X_{new} \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)} Q^{(s)'} \beta \right) \\ &= \frac{1}{n_{new}} \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} X_{new} Q^{(s)} Q^{(s)'} \beta \right)' \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} X_{new} Q^{(s)} Q^{(s)'} \beta \right). \end{aligned}$$

Using the fact that

$$X_{new} Q^{(s)} = W F Z' Q^{(s)} = \sum_{\alpha=1}^{r_{new}} f_{\alpha} Q^{(s)'} Z^{(\alpha)} W^{(\alpha)},$$

we deduce

$$\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} X_{new} Q^{(s)} Q^{(s)'} \beta = \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \sum_{\alpha=1}^{r_{new}} f_{\alpha} Q^{(s)'} Z^{(\alpha)} Q^{(s)'} \beta W^{(\alpha)}.$$

Consequently,

$$\begin{aligned} \check{A}_3 &= \frac{1}{n_{new}} \sum_{s=1}^r \sum_{\ell=1}^r \frac{d_s^2 d_{\ell}^2}{(d_s^2 + \lambda)(d_{\ell}^2 + \lambda)} \sum_{\alpha=1}^{r_{new}} f_{\alpha} Q^{(s)'} Z^{(\alpha)} Q^{(s)'} \beta W^{(\alpha)'} \sum_{\vartheta=1}^{r_{new}} f_{\vartheta} Q^{(\ell)'} Z^{(\vartheta)} Q^{(\ell)'} \beta W^{(\vartheta)} \\ &= \frac{1}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Q^{(s)'} \beta \sum_{\ell=1}^r \frac{d_{\ell}^2}{d_{\ell}^2 + \lambda} Q^{(\ell)'} \beta \sum_{\alpha=1}^{r_{new}} f_{\alpha}^2 \langle Z^{(\alpha)} Z^{(\alpha)'} Q^{(s)}, Z^{(\alpha)} Z^{(\alpha)'} Q^{(\ell)} \rangle. \end{aligned}$$

7. Extra lemmas

Lemma 1 (Bounds on $\check{\rho}_g$). *Under the same assumptions as in Theorem 3 of the main manuscript, we always have*

$$\frac{B_1}{(B_2 + B_3)^{1/2} B_4^{1/2}} \leq \check{\rho}_g \leq \rho_g^{oracle},$$

where

$$\begin{aligned}
B_1 &= \min_{1 \leq s \leq r} \frac{d_s^2}{d_s^2 + \lambda} \min_{1 \leq \alpha \leq r_{new}} f_\alpha^2 \langle ZZ' \beta, QQ' \beta \rangle, \\
B_2 &= \sigma_e^2 r r_{new} \max_{1 \leq s \leq r} \frac{d_s^2}{(d_s^2 + \lambda)^2} \max_{1 \leq \alpha \leq r_{new}} f_\alpha^2 \max_{s, \alpha} \left\| Q^{(s)'} Z^{(\alpha)} W^{(\alpha)} \right\|^2, \\
B_3 &= \max_{1 \leq s \leq r} \frac{d_s^4}{(d_s^2 + \lambda)^2} \|QQ' \beta\|^2 \max_{1 \leq \alpha \leq r_{new}} f_\alpha^2 r^2, \\
B_4 &= \|ZZ' \beta\|^2 \max_{1 \leq \alpha \leq r_{new}} f_\alpha^2.
\end{aligned}$$

Note that it is possible to replace B_2 by the quantity

$$\sigma_e^2 r_{new} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \max_{1 \leq \alpha \leq r_{new}} \left(f_\alpha^2 \langle Q^{(s)}, Z^{(\alpha)} \rangle^2 \right),$$

entailing a second lower bound for $\check{\rho}_g$.

Proof. To begin with, let us focus on the upper bound. First, we have to notice that we have the relationship

$$\check{A}_1 = \frac{1}{n_{new}} \sum_{\alpha=1}^{r_{new}} \langle f_\alpha Z^{(\alpha)} Z^{(\alpha)'} \beta, \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} f_\alpha Z^{(\alpha)} Z^{(\alpha)'} Q^{(s)} Q^{(s)'} \beta \rangle.$$

Then, applying two times the Cauchy-Schwartz inequality, we obtain

$$\begin{aligned}
\check{A}_1 &\leq \frac{1}{n_{new}} \sum_{\alpha=1}^{r_{new}} \left(\left\| f_\alpha Z^{(\alpha)} Z^{(\alpha)'} \beta \right\| \left\| \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} f_\alpha Z^{(\alpha)} Z^{(\alpha)'} Q^{(s)} Q^{(s)'} \beta \right\| \right) \\
&\leq \frac{1}{n_{new}} \left(\sum_{\alpha=1}^{r_{new}} \left\| f_\alpha Z^{(\alpha)} Z^{(\alpha)'} \beta \right\|^2 \right)^{1/2} \left(\sum_{\alpha=1}^{r_{new}} \left\| \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} f_\alpha Z^{(\alpha)} Z^{(\alpha)'} Q^{(s)} Q^{(s)'} \beta \right\|^2 \right)^{1/2} \\
&= \check{A}_4^{1/2} \frac{1}{\sqrt{n_{new}}} \left(\sum_{\alpha=1}^{r_{new}} f_\alpha^2 \left\| \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Z^{(\alpha)} Z^{(\alpha)'} Q^{(s)} Q^{(s)'} \beta \right\|^2 \right)^{1/2}.
\end{aligned}$$

We have

$$\begin{aligned}
&\sum_{\alpha=1}^{r_{new}} f_\alpha^2 \left\| \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} Z^{(\alpha)} Z^{(\alpha)'} Q^{(s)} Q^{(s)'} \beta \right\|^2 \\
&= \sum_{\alpha=1}^{r_{new}} f_\alpha^2 \left(\sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \beta' Q^{(s)} Q^{(s)'} Z^{(\alpha)} Z^{(\alpha)'} \right) \left(\sum_{\ell=1}^r \frac{d_\ell^2}{d_\ell^2 + \lambda} Z^{(\alpha)} Z^{(\alpha)'} Q^{(\ell)} Q^{(\ell)'} \beta \right) \\
&= \sum_{s=1}^r \frac{d_s^2}{d_s^2 + \lambda} \beta' Q^{(s)} \sum_{\ell=1}^r \frac{d_\ell^2}{d_\ell^2 + \lambda} \beta' Q^{(\ell)} \sum_{\alpha=1}^{r_{new}} f_\alpha^2 \langle Z^{(\alpha)} Z^{(\alpha)'} Q^{(s)}, Z^{(\alpha)} Z^{(\alpha)'} Q^{(\ell)} \rangle \\
&= n_{new} \check{A}_3.
\end{aligned}$$

Thus

$$\check{A}_1 \leq \check{A}_4^{1/2} \check{A}_3^{1/2}.$$

Since $\check{A}_2 \geq 0$, we finally obtain

$$\check{\rho}_g \leq \frac{\check{A}_1}{\check{A}_4^{1/2} \check{A}_3^{1/2}} \leq 1.$$

Let us now move on to the lower bound. We have the relationship:

$$\begin{aligned} \check{A}_2 &\leq \frac{\sigma_e^2}{n_{new}} \max_{1 \leq s \leq r} \frac{d_s^2}{(d_s^2 + \lambda)^2} \max_{1 \leq \alpha \leq r_{new}} f_\alpha^2 \left(\sum_{s=1}^r \left\| \sum_{\alpha=1}^{r_{new}} Q^{(s)'} Z^{(\alpha)} W^{(\alpha)} \right\|^2 \right) \\ &\leq \frac{r \sigma_e^2}{n_{new}} \max_{1 \leq s \leq r} \frac{d_s^2}{(d_s^2 + \lambda)^2} \max_{1 \leq \alpha \leq r_{new}} f_\alpha^2 \max_{1 \leq s \leq r} \left\| \sum_{\alpha=1}^{r_{new}} Q^{(s)'} Z^{(\alpha)} W^{(\alpha)} \right\|^2 \\ &\leq \frac{r \sigma_e^2}{n_{new}} \max_{1 \leq s \leq r} \frac{d_s^2}{(d_s^2 + \lambda)^2} \max_{1 \leq \alpha \leq r_{new}} f_\alpha^2 \max_{1 \leq s \leq r} \sum_{\alpha=1}^{r_{new}} \left\| Q^{(s)'} Z^{(\alpha)} W^{(\alpha)} \right\|^2 \\ &\leq \frac{r r_{new} \sigma_e^2}{n_{new}} \max_{1 \leq s \leq r} \frac{d_s^2}{(d_s^2 + \lambda)^2} \max_{1 \leq \alpha \leq r_{new}} f_\alpha^2 \max_{s, \alpha} \left\| Q^{(s)'} Z^{(\alpha)} W^{(\alpha)} \right\|^2. \end{aligned}$$

Coming back to the expression of \check{A}_2 , we also have:

$$\check{A}_2 \leq \frac{\sigma_e^2}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \max_{\alpha} \left(f_\alpha^2 < Q^{(s)}, Z^{(\alpha)} >^2 \right) \sum_{i=1}^{n_{new}} \left(\sum_{\omega=1}^{r_{new}} W_i^{(\omega)} \right)^2.$$

We can notice that $\sum_{i=1}^{n_{new}} \left(\sum_{\omega=1}^{r_{new}} W_i^{(\omega)} \right)^2 = \text{Tr}(WW') = \text{Tr}(W'W) = r_{new}$.

As a consequence, another bound is the following

$$\check{A}_2 \leq \frac{\sigma_e^2 r_{new}}{n_{new}} \sum_{s=1}^r \frac{d_s^2}{(d_s^2 + \lambda)^2} \max_{1 \leq \alpha \leq r_{new}} \left(f_\alpha^2 < Q^{(s)}, Z^{(\alpha)} >^2 \right).$$

On the other hand, we have

$$\begin{aligned} \check{A}_1 &\geq \min_{1 \leq s \leq r} \frac{d_s^2}{d_s^2 + \lambda} \min_{1 \leq \alpha \leq r_{new}} f_\alpha^2 \sum_{s=1}^r \left(\sum_{\alpha=1}^{r_{new}} \beta' Z^{(\alpha)} Z^{(\alpha)'} Q^{(s)} Q^{(s)'} \beta \right) \\ &= \frac{1}{n_{new}} \min_{1 \leq s \leq r} \frac{d_s^2}{d_s^2 + \lambda} \min_{1 \leq \alpha \leq r_{new}} f_\alpha^2 \beta' Z Z' Q Q' \beta \\ &= \frac{1}{n_{new}} \min_{1 \leq s \leq r} \frac{d_s^2}{d_s^2 + \lambda} \min_{1 \leq \alpha \leq r_{new}} f_\alpha^2 < Z Z' \beta, Q Q' \beta >. \end{aligned}$$

Last,

$$\begin{aligned}
\check{A}_3 &\leq \frac{1}{n_{new}} \max_{1 \leq s \leq r} \frac{d_s^4}{(d_s^2 + \lambda)^2} \max_{1 \leq s \leq r} \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2 \max_{1 \leq \alpha \leq r_{new}} f_\alpha^2 \sum_{s=1}^r \sum_{\ell=1}^r \sum_{\alpha=1}^{r_{new}} Q^{(s)'} Z^{(\alpha)} Z^{(\alpha)'} Q^{(\ell)} \\
&= \frac{1}{n_{new}} \max_{1 \leq s \leq r} \frac{d_s^4}{(d_s^2 + \lambda)^2} \max_{1 \leq s \leq r} \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2 \max_{1 \leq \alpha \leq r_{new}} f_\alpha^2 \sum_{s=1}^r \sum_{\ell=1}^r Q^{(s)'} Z Z' Q^{(\ell)} \\
&\leq \frac{1}{n_{new}} \max_{1 \leq s \leq r} \frac{d_s^4}{(d_s^2 + \lambda)^2} \max_{1 \leq s \leq r} \left\| Q^{(s)} Q^{(s)'} \beta \right\|^2 \max_{1 \leq \alpha \leq r_{new}} f_\alpha^2 \\
&\quad \times \left\{ r \max_{1 \leq s \leq r} \left\| Z Z' Q^{(s)} \right\|^2 + r(r-1) \max_{s \neq \ell} \langle Z Z' Q^{(s)}, Z Z' Q^{(\ell)} \rangle \right\}.
\end{aligned}$$

Since $Z Z'$ is an idempotent matrix and $Q^{(s)'} Q^{(s)} = 1$ for all $1 \leq s \leq r$, we have

$$\left\| Z Z' Q^{(s)} \right\|^2 \leq 1.$$

Besides, according to Cauchy-Schwartz inequality,

$$|\langle Z Z' Q^{(s)}, Z Z' Q^{(\ell)} \rangle| \leq \left\| Z Z' Q^{(s)} \right\| \left\| Z Z' Q^{(\ell)} \right\| \leq 1.$$

Finally, since $Q Q'$ is an idempotent matrix, and putting together all the above considerations, we obtain

$$\check{A}_3 \leq \frac{r^2}{n_{new}} \max_{1 \leq s \leq r} \frac{d_s^4}{(d_s^2 + \lambda)^2} \|Q Q' \beta\|^2 \max_{1 \leq \alpha \leq r_{new}} f_\alpha^2.$$

We can now easily deduce the bounds given in the statement. \square

Lemma 2. *Let us consider same assumptions as in Theorem 1 of the main manuscript. Then, the quantity $\tilde{\rho}_g$ defined in Section 5 of the main manuscript has the following expression*

$$\tilde{\rho}_g = \frac{\tilde{A}_1}{\left(\tilde{A}_2 + \tilde{A}_3\right)^{1/2} \left(\tilde{A}_4\right)^{1/2}},$$

where

$$\begin{aligned}
\tilde{A}_1 &= \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} \beta' \Sigma Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta, \quad \tilde{A}_2 = \sigma_e^2 \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} \mathbb{E} \left(\left\| Q^{(\sigma(s))} Q^{(\sigma(s))'} x_{new} \right\|^2 \right) \\
\tilde{A}_3 &= \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right)' \Sigma \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right), \quad \tilde{A}_4 = A_4.
\end{aligned}$$

Proof. After replacing the quantity $X'V^{-1}$ by $X'V^{-1}\tilde{P}\tilde{P}'$, formula (5) of Rabier et al. (2016) becomes

$$\rho_g = \frac{\beta' \text{Var}(x_{new}) X'V^{-1}\tilde{P}\tilde{P}'X\beta}{\left(\sigma_e^2 \mathbb{E} \left(\left\| x'_{new} X'V^{-1}\tilde{P}\tilde{P}' \right\|^2 \right) + \beta' X' \tilde{P} \tilde{P}' V^{-1} X \text{Var}(x_{new}) X'V^{-1}\tilde{P}\tilde{P}'X\beta \right)^{1/2} \sigma_G}.$$

As a result, let us define

$$\begin{aligned}\tilde{A}_1 &:= \beta' \text{Var}(x_{new}) X' V^{-1} \tilde{P} \tilde{P}' X \beta, \quad \tilde{A}_2 := \sigma_e^2 \mathbb{E} \left(\left\| x'_{new} X' V^{-1} \tilde{P} \tilde{P}' \right\|^2 \right), \\ \tilde{A}_3 &:= \beta' X' \tilde{P} \tilde{P}' V^{-1} X \text{Var}(x_{new}) X' V^{-1} \tilde{P} \tilde{P}' X \beta, \quad \tilde{A}_4 := A_4.\end{aligned}$$

Using the fact that $X' V^{-1} = Q \bar{D} P'$ and $\Sigma = \mathbb{E}(x_{new} x'_{new})$, we have

$$\begin{aligned}\tilde{A}_1 &= \beta' \Sigma X' V^{-1} \tilde{P} \tilde{P}' X \beta \\ &= \beta' \Sigma Q \bar{D} P' \tilde{P} \tilde{P}' X \beta.\end{aligned}$$

After some simple algebra, we obtain

$$Q \bar{D} P' \tilde{P} = \left(\frac{d_{\sigma(1)}}{d_{\sigma(1)}^2 + \lambda} Q^{(\sigma(1))}, \dots, \frac{d_{\sigma(\tilde{r})}}{d_{\sigma(\tilde{r})}^2 + \lambda} Q^{(\sigma(\tilde{r}))} \right). \quad (12)$$

Then,

$$\begin{aligned}\tilde{A}_1 &= \beta' \Sigma \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))} P^{(\sigma(s))'} \right) \left(\sum_{s=1}^{\tilde{r}} d_s P^{(s)} Q^{(s)'} \right) \beta \\ &= \beta' \Sigma \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right) \\ &= \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} \beta' \Sigma Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta.\end{aligned}$$

Let us now consider \tilde{A}_2 . We have

$$\begin{aligned}\left\| x'_{new} X' V^{-1} \tilde{P} \tilde{P}' \right\|^2 &= x'_{new} X' V^{-1} \tilde{P} \tilde{P}' \tilde{P} \tilde{P}' (X' V^{-1})' x_{new} \\ &= x'_{new} Q \bar{D} P' \tilde{P} \tilde{P}' P \bar{D} Q' x_{new}.\end{aligned}$$

According to formula (12), we obtain

$$Q \bar{D} P' \tilde{P} \tilde{P}' P \bar{D} Q' = \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} Q^{(\sigma(s))} Q^{(\sigma(s))'}$$

and

$$\begin{aligned}x'_{new} Q \bar{D} P' \tilde{P} \tilde{P}' P \bar{D} Q' x_{new} &= \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} x'_{new} Q^{(\sigma(s))} Q^{(\sigma(s))'} x_{new} \\ &= \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} \left\| Q^{(\sigma(s))} Q^{(\sigma(s))'} x_{new} \right\|^2.\end{aligned}$$

The last equality comes from the fact that $Q^{(\sigma(s))}Q^{(\sigma(s))'}$ is an idempotent matrix. To conclude, we have

$$\tilde{A}_2 = \sigma_e^2 \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} \mathbb{E} \left(\left\| Q^{(\sigma(s))} Q^{(\sigma(s))'} x_{new} \right\|^2 \right).$$

Furthermore, recall that

$$\tilde{A}_3 = \beta' X' \tilde{P} \tilde{P}' V^{-1} X \text{Var}(x_{new}) X' V^{-1} \tilde{P} \tilde{P}' X \beta.$$

Since the expression of $X' V^{-1} \tilde{P} \tilde{P}' X \beta$ is also present in \tilde{A}_1 , we easily obtain

$$\tilde{A}_3 = \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right)' \Sigma \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right).$$

□

8. Proof of Lemma 4 of the main manuscript

To begin with, let us recall the expression for \tilde{A}_1 given in Lemma 2 above:

$$\tilde{A}_1 = \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} \beta' \Sigma Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta.$$

We consider the following natural estimation \hat{A}_1 for \tilde{A}_1 :

$$\hat{A}_1 := \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} \beta' \hat{\Sigma} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta,$$

where $\hat{\Sigma} = X'X/n$ is the empirical covariance matrix.

We have

$$\begin{aligned} \hat{A}_1 &= \frac{1}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} \beta' \hat{\Sigma} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \\ &= \frac{1}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} \beta' Q D^2 Q' Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta. \end{aligned}$$

It is easy to see that

$$Q D^2 Q' Q^{(\sigma(s))} = d_{\sigma(s)}^2 Q^{(\sigma(s))}.$$

Therefore,

$$\hat{A}_1 = \frac{1}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{d_{\sigma(s)}^2 + \lambda} \beta' Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta = \frac{1}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{d_{\sigma(s)}^2 + \lambda} \left\| Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right\|^2.$$

Let us recall the expression for \tilde{A}_2 given previously:

$$\tilde{A}_2 = \sigma_e^2 \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} \mathbb{E} \left(\left\| Q^{(\sigma(s))} Q^{(\sigma(s))'} x_{new} \right\|^2 \right).$$

A natural estimation of \tilde{A}_2 is

$$\begin{aligned} \hat{A}_2 &:= \frac{\sigma_e^2}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} \sum_{i=1}^n \left\| Q^{(\sigma(s))} Q^{(\sigma(s))'} x_i \right\|^2 \\ &= \frac{\sigma_e^2}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} \text{Tr} \left(X Q^{(\sigma(s))} Q^{(\sigma(s))'} Q^{(\sigma(s))} Q^{(\sigma(s))'} X' \right) \\ &= \frac{\sigma_e^2}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} \text{Tr} \left(X Q^{(\sigma(s))} Q^{(\sigma(s))'} X' \right) \\ &= \frac{\sigma_e^2}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} \text{Tr} \left(P D Q' Q^{(\sigma(s))} Q^{(\sigma(s))'} Q D P' \right). \end{aligned}$$

Note that

$$D Q' Q^{(\sigma(s))} = d_{\sigma(s)} e_{\sigma(s)},$$

where $e_{\sigma(s)}$ denotes the $\sigma(s)$ -th vector of the canonical basis of \mathbb{R}^r . As a result,

$$\begin{aligned} \hat{A}_2 &= \frac{\sigma_e^2}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{(d_{\sigma(s)}^2 + \lambda)^2} \text{Tr} \left(P e_{\sigma(s)} e_{\sigma(s)}' P' \right) \\ &= \frac{\sigma_e^2}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{(d_{\sigma(s)}^2 + \lambda)^2} \text{Tr} \left(P' P e_{\sigma(s)} e_{\sigma(s)}' \right) \\ &= \frac{\sigma_e^2}{n} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^4}{(d_{\sigma(s)}^2 + \lambda)^2}. \end{aligned}$$

An estimation for the quantity \tilde{A}_3 is the following

$$\hat{A}_3 := \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right)' \hat{\Sigma} \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right).$$

We have the following relations

$$\begin{aligned} \hat{A}_3 &= \frac{1}{n} \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right)' X' X \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right) \\ &= \frac{1}{n} \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} X Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right)' \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} X Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right). \end{aligned}$$

Since $X = PDQ'$, we have

$$XQ^{(\sigma(s))}Q^{(\sigma(s))'}\beta = d_{\sigma(s)}Pe_{\sigma(s)}Q^{(\sigma(s))'}\beta = d_{\sigma(s)}P^{(\sigma(s))}Q^{(\sigma(s))'}\beta$$

and thus

$$\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} XQ^{(\sigma(s))}Q^{(\sigma(s))'}\beta = \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^3}{d_{\sigma(s)}^2 + \lambda} P^{(\sigma(s))}Q^{(\sigma(s))'}\beta.$$

Last, we obtain

$$\begin{aligned} \hat{A}_3 &= \frac{1}{n} \left(\sum_{\ell=1}^{\tilde{r}} \frac{d_{\sigma(\ell)}^3}{d_{\sigma(\ell)}^2 + \lambda} \beta' Q^{(\sigma(\ell))} P^{(\sigma(\ell))'} \right) \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^3}{d_{\sigma(s)}^2 + \lambda} P^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right) \\ &= \frac{1}{n} \sum_{\ell=1}^{\tilde{r}} \frac{d_{\sigma(\ell)}^3}{d_{\sigma(\ell)}^2 + \lambda} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^3}{d_{\sigma(s)}^2 + \lambda} \beta' Q^{(\sigma(\ell))} P^{(\sigma(\ell))'} P^{(\sigma(s))} Q^{(\sigma(s))'} \beta \\ &= \frac{1}{n} \sum_{\ell=1}^{\tilde{r}} \frac{d_{\sigma(\ell)}^6}{(d_{\sigma(\ell)}^2 + \lambda)^2} \beta' Q^{(\sigma(\ell))} Q^{(\sigma(\ell))'} \beta \\ &= \frac{1}{n} \sum_{\ell=1}^{\tilde{r}} \frac{d_{\sigma(\ell)}^6}{(d_{\sigma(\ell)}^2 + \lambda)^2} \left\| Q^{(\sigma(\ell))} Q^{(\sigma(\ell))'} \beta \right\|^2. \end{aligned}$$

9. Proof of Lemma 6 of the main manuscript

To simplify notations, let us put

$$\begin{aligned} u &:= \hat{A}_1, \quad \delta_1 := \hat{A}_1 - \hat{A}_1, \\ v &:= \hat{A}_2 + \hat{A}_3, \quad \delta_2 := \hat{A}_2 + \hat{A}_3 - (\hat{A}_2 + \hat{A}_3). \end{aligned}$$

With these notations the condition $\hat{\rho}_g \geq \hat{\rho}_g$ reads

$$\frac{u + \delta_1}{\sqrt{v + \delta_2}} \leq \frac{u}{\sqrt{v}},$$

which is further equivalent to

$$\delta_2 u^2 - 2u\delta_1 v - \delta_1^2 v \geq 0.$$

The discriminant in the u variable equals $\Delta = 4\delta_1^2 v(v + \delta_2)$ and is positive. The above second order inequation is thus satisfied for

$$\left| u - \frac{\delta_1}{\delta_2} v \right| \geq \frac{\delta_1}{\delta_2} \sqrt{v(v + \delta_2)}.$$

Note that the case

$$u \leq \frac{\delta_1}{\delta_2} v - \frac{\delta_1}{\delta_2} \sqrt{v(v + \delta_2)} < 0$$

is not possible, since $u \geq 0$. The only possible case is therefore

$$u \geq \frac{\delta_1}{\delta_2}v + \frac{\delta_1}{\delta_2}\sqrt{v(v + \delta_2)},$$

which further gives the desired statement when dividing by $\delta_1 \neq 0$.

10. Proof of Corollary 1 of the main manuscript

Proof of 1.

Let us denote $R := \frac{\widehat{\text{Cov}}(\tilde{Y}_{new}, Y_{new})}{\widehat{\text{Cov}}(\vec{Y}_{new}, Y_{new})}$ and $b := \frac{\text{Var}(\tilde{Y}_{new})}{\text{Var}(\vec{Y}_{new})}$. Then the condition on the estimated covariances and variances writes

$$R \geq b \left(1 + \sqrt{1 + \frac{1}{b}} \right),$$

which is equivalent to $\hat{\rho}_g \geq \hat{\rho}_g$ by Lemma 6 of the main manuscript. Moreover, if the above condition is satisfied, we deduce that

$$\frac{\hat{\rho}_g}{\hat{\rho}_g} = \frac{R}{\sqrt{b}} \geq \sqrt{b} + \sqrt{1 + b} \geq 1,$$

hence $\hat{\rho}_g \geq \hat{\rho}_g$.

Proof of 2.

When interchanging \tilde{Y}_{new} and \vec{Y}_{new} , the previously shown result implies that $\hat{\rho}_g \geq \hat{\rho}_g$ if and only if

$$\frac{1}{R} \geq \frac{1}{b} \left(1 + \sqrt{1 + b} \right),$$

and in this case we also have $\hat{\rho}_g \geq \hat{\rho}_g$.

But $\frac{b}{1 + \sqrt{1 + b}} = \sqrt{1 + b} - 1$, and hence the above condition is equivalent to $R \leq \sqrt{1 + b} - 1$, which is exactly the condition in the statement.

Proof of 3.

This result follows directly from the two previous points.

11. Proof of Lemma 7 of the main manuscript

Using Lemma 6 of the main manuscript and proceeding in the same way as in the proof of Lemma 2 of the main manuscript, we obtain

$$\begin{aligned} n\hat{A}_1 &\sim \sum_{s \in \tilde{\Omega}_1} d_s^2 \frac{n^{2\tau}}{r}, \\ n\hat{A}_2 + n\hat{A}_3 &\sim \sum_{s \in \tilde{\Omega}_1} d_s^2 \frac{n^{2\tau}}{r}, \\ n\hat{A}_4 = n\hat{A}_4 &\sim \sum_{s \in \Omega_1} d_s^2 \frac{n^{2\tau}}{r}, \end{aligned}$$

and the stated result follows.

12. Extra results

Lemma 3. *Let us consider same hypotheses as in Theorem 3 of the main manuscript. Then, a natural estimator of the quantity $\tilde{\rho}_g$ is the following:*

$$\check{\rho}_g := \frac{\check{A}_1}{\left(\check{A}_2 + \check{A}_3\right)^{1/2} \left(\check{A}_4\right)^{1/2}},$$

where

$$\begin{aligned} \check{A}_1 &= \frac{1}{n_{new}} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} \left(\sum_{\alpha=1}^{r_{new}} f_{\alpha}^2 \langle Z^{(\alpha)} Z^{(\alpha)'} \beta, Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \rangle \right), \\ \check{A}_2 &= \frac{\sigma_e^2}{n_{new}} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} \sum_{i=1}^{n_{new}} \left(\sum_{\alpha=1}^{r_{new}} f_{\alpha} Q^{(\sigma(s))'} Z^{(\alpha)} W_i^{(\alpha)} \right)^2, \\ \check{A}_3 &= \frac{1}{n_{new}} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))'} \beta \sum_{\ell=1}^{\tilde{r}} \frac{d_{\sigma(\ell)}^2}{d_{\sigma(\ell)}^2 + \lambda} Q^{(\sigma(\ell))'} \beta \left(\sum_{\alpha=1}^{r_{new}} f_{\alpha}^2 \langle Z^{(\alpha)} Z^{(\alpha)'} Q^{(\sigma(s))}, Z^{(\alpha)} Z^{(\alpha)'} Q^{(\sigma(\ell))} \rangle \right), \\ \check{A}_4 &= \check{A}_4. \end{aligned}$$

Proof. In the same way as before, we consider the estimators \check{A}_1 , \check{A}_2 and \check{A}_3 , of \hat{A}_1 , \hat{A}_2 and \hat{A}_3 , respectively:

$$\begin{aligned} \check{A}_1 &:= \frac{1}{n_{new}} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} \beta' X'_{new} X_{new} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta, \\ \check{A}_2 &:= \frac{\sigma_e^2}{n_{new}} \sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} \text{Tr} \left(X_{new} Q^{(\sigma(s))} Q^{(\sigma(s))'} Q^{(\sigma(s))} Q^{(\sigma(s))'} X'_{new} \right), \\ \check{A}_3 &= \frac{1}{n_{new}} \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right)' X'_{new} X_{new} \left(\sum_{s=1}^{\tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} Q^{(\sigma(s))} Q^{(\sigma(s))'} \beta \right). \end{aligned}$$

After some easy computations we can deduce the stated formulas. \square

Lemma 4. *Let us consider same hypotheses as in Theorem 3 of the main manuscript. Then we always have*

$$\frac{\tilde{B}_1}{\left(\tilde{B}_2 + \tilde{B}_3\right)^{1/2} \tilde{B}_4^{1/2}} \leq \check{\rho}_g \leq \rho_g^{oracle},$$

where

$$\begin{aligned} \tilde{B}_1 &= \min_{1 \leq s \leq \tilde{r}} \frac{d_{\sigma(s)}^2}{d_{\sigma(s)}^2 + \lambda} \min f_{\alpha}^2 \quad \langle ZZ' \beta, \tilde{Q} \tilde{Q}' \beta \rangle, \\ \tilde{B}_2 &= \sigma_e^2 \tilde{r} r_{new} \max_{1 \leq s \leq \tilde{r}} \frac{d_{\sigma(s)}^2}{(d_{\sigma(s)}^2 + \lambda)^2} \max f_{\alpha}^2 \max_{1 \leq s \leq \tilde{r}, \alpha} \left\| Q^{(\sigma(s))' } Z^{(\alpha)} W^{(\alpha)} \right\|^2, \\ \tilde{B}_3 &= \max_{1 \leq s \leq \tilde{r}} \frac{d_{\sigma(s)}^4}{(d_{\sigma(s)}^2 + \lambda)^2} \left\| \tilde{Q} \tilde{Q}' \beta \right\|^2 \max f_{\alpha}^2 \tilde{r}^2, \\ \tilde{B}_4 &= B_4. \end{aligned}$$

The proof relies heavily on the proof of Lemma 4 of the main manuscript, provided that we consider the expressions of $\check{A}_1, \check{A}_2, \check{A}_3$ given in Lemma 3 above.

13. Extra comments on Section 3 of the main manuscript

13.1. About the conditions

In the manuscript, we assume

$$\begin{aligned} d_1^2 &\sim n^{\psi}, \quad \text{with } 0 < \psi \leq 1, \\ d_r^2 &\sim n^{\eta}, \quad \text{with } \eta \leq \psi \leq 1 \quad \text{and } \eta \text{ and } \psi \text{ not depending on } n. \end{aligned}$$

Recall that $u_n \sim v_n$ means that $\frac{u_n}{v_n} \rightarrow 1$ when $n \rightarrow \infty$. Besides, we also consider that

$$\|QQ' \beta\|^2 \sim n^{2\tau}, \quad \text{with } \tau < \eta \text{ and } \tau \text{ not depending on } n.$$

These conditions are largely inspired from Shao and Deng (2012). However, we are mentioning the exact order of each term, since our goal in this section is to study the behavior of the quantity $\hat{\rho}_g$, which is a ratio. For instance, Condition C2 of Shao and Deng (2012) imposes $\|QQ' \beta\|^2 = O(n^{2\tau})$ and is somewhat more general than ours. On the other hand, in their Theorem 3, Shao and Deng (2012) suppose $d_1^2 = O(n)$, whereas Fan and Lv (2008) assume (in condition 4) $d_1^2 = O(n^v)$ with $v \geq 0$. Therefore, our condition on d_1^2 can be viewed as a compromise between the conditions considered in these two papers. Note that all the results in the present paper remain valid even if $\psi > 1$. Last, our condition on d_r^2 is inspired from condition C1 of Shao and Deng (2012).

13.2. About the tuning parameter λ

The setting $\lambda \rightarrow \infty$ when $p \rightarrow \infty$ is somewhat classical in genomics. The heritability h^2 of a quantitative character is only approximately known by geneticists, and it is well known that the Ridge regression can be interpreted in a Bayesian framework, assuming same variance on each regressor. As a consequence, in order to obtain an estimated value of λ , the signal (linked to h^2) is generally spread out accross all the regressors. This leads to a tuning parameter which diverges to $+\infty$ and the $\hat{\beta}_k$'s are more and more shrunked when the number of regressors increases (see for instance our section on the regularization parameter in Rabier et al. (2016)).

13.3. About the extra conditions

Let us also give a few comments regarding the conditions (C1-C2-C3-C4-C5-C6). Under (C2), the squared L^2 norm of the vector containing the largest singular values d_s for $s \in \Omega_1$ may diverge to $+\infty$ at a rate slower than λ . According to (C3), the squared L^2 norm of the vector whose components are the square of the smallest singular values may diverge to $+\infty$ at a rate slower than λ^2 . Condition (C4) assumes that the ratio $r/n^{2\tau}$ diverges faster to $+\infty$ than the tuning parameter λ . Last, (C5) and (C6) impose that the number of large singular values and the number of intermediate singular values are bounded. In other words, when $p > n$, the rank $r \leq n$ of the matrix X will diverge to $+\infty$ if and only if the number of small singular values tends to $+\infty$.

References

- 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.
- Rabier, C.E., Barre, P., Asp, T., Charmet, G. & Mangin, B. (2016). On the Accuracy of Genomic Selection. *PloS One.* **11**, (6), e0156086. doi:10.1371/journal.pone.0156086.
- Shao, J. & Deng, X. (2012). Estimation in high-dimensional linear models with deterministic design matrices. *The Annals of Statistics.* **40**, (2), 812-831.

Table 1: Comparison among different estimators of the phenotypic accuracy as a function of the number of siblings in the TRN sample (TRN and TST samples based on the same number of generations). The chromosome is of length 1M and 2 QTLs are located at 3cM and 80cM with effects +1 and -2, respectively ($n = 500$, $n_{new} = 100$, $\sigma_e^2 = 1$). Emp. Acc. refers to the empirical *phenotypic accuracy*. The standard errors for the estimates are given in brackets.

Nb Markers	Nb generations	Nb Siblings	Emp. Acc.	$\hat{\rho}_{ph}(\beta)$	$\check{\rho}_{ph}(\beta)$
100	30	0	0.6967 (0.0049)	0.6969 (0.0001)	0.6915 (0.0015)
		100	0.6941 (0.0049)	0.6890 (0.0001)	0.6772 (0.0018)
	50	0	0.6804 (0.0049)	0.6765 (0.0001)	0.6731 (0.0015)
		100	0.6868 (0.0054)	0.6572 (0.0001)	0.6806 (0.0017)
	70	0	0.6708 (0.0055)	0.6717 (0.0001)	0.6654 (0.0017)
		100	0.6708 (0.0056)	0.6717 (0.0001)	0.6654 (0.0017)
1,000	30	0	0.7015 (0.0047)	0.7155 (0.0001)	0.6889 (0.0017)
		100	0.6735 (0.0052)	0.6739 (0.0001)	0.6549 (0.0024)
	50	0	0.6683 (0.0053)	0.6597 (0.0001)	0.6576 (0.0021)
		100	0.6305 (0.0054)	0.6064 (0.0002)	0.6213 (0.0027)
	70	0	0.6713 (0.0059)	0.685 (0.0001)	0.6642 (0.0023)
		100	0.6471 (0.0056)	0.6712 (0.0002)	0.6395 (0.0020)
2,000	30	0	0.6977 (0.0056)	0.6933 (0.0001)	0.6881 (0.0016)
		100	0.6858 (0.0170)	0.7053 (0.0013)	0.6857 (0.00170)
	50	0	0.6316 (0.0060)	0.6254 (0.0001)	0.6264 (0.0025)
		100	0.6665 (0.0054)	0.6913 (0.0001)	0.6625 (0.0019)
	70	0	0.4174 (0.0078)	0.4600 (0.0004)	0.4095 (0.0041)
		100	0.6665 (0.0053)	0.6913 (0.0001)	0.6625 (0.0019)

Table 2: Comparison among different estimators of the phenotypic accuracy, in presence of a mixture of major genes and small QTLs (50 generations, $n = 500$, $n_{new} = 100$, $\sigma_e^2 = 1$). Three scenarios considered (a) 2 large QTLs with effects +0.5 and -0.6, 98 small QTLs with the same effect +0.07, (b) 2 large QTLs with effects +1 and -0.7, 98 small QTLs with the same effect +0.1, (c) 2 large QTLs with effects +2 and -2, 98 small QTLs with the same effect +0.1. The chromosome is of length 1M and the large QTLs are located at 3cM and 80cM, whereas the small QTLs are located every centimorgan (except at 3cM and 80cM). $\hat{\beta}_{LASSO}$, $\hat{\beta}_{ADLASSO}$, $\hat{\beta}_{GPLASSO}$ refer to the LASSO, Adaptive LASSO and Group LASSO estimators of β , respectively. Emp. Acc. refers to the empirical phenotypic accuracy. The standard errors for the estimates are given in brackets.

Nb markers	Method	Scenario (a)	Scenario (b)	Scenario (c)
100	Emp. Acc.	0.5479 (0.0068)	0.7012 (0.0051)	0.8074 (0.0034)
	$\hat{\rho}_{ph}(\beta)$	0.5362 (0.0001)	0.6900 (0.0001)	0.8013 (0.0001)
	$\hat{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.3792 (0.0050)	0.6096 (0.0027)	0.7614 (0.0014)
	$\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.5400 (0.0040)	0.6678 (0.0018)	0.8049 (0.0014)
	$\hat{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.3500 (0.0046)	0.5909 (0.0026)	0.7419 (0.0014)
	$\check{\rho}_{ph}(\beta)$	0.5296 (0.0030)	0.6868 (0.0026)	0.7962 (0.0019)
	$\check{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.3628 (0.0055)	0.6016 (0.0038)	0.7550 (0.0026)
	$\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.5313 (0.0045)	0.6942 (0.0029)	0.7999 (0.0022)
	$\check{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.3370 (0.0051)	0.5720 (0.0036)	0.7324 (0.0027)
1,000	Emp. Acc.	0.5867 (0.0072)	0.7374 (0.0048)	0.8307 (0.0032)
	$\hat{\rho}_{ph}(\beta)$	0.5738 (0.0002)	0.7316 (0.0001)	0.8276 (0.0001)
	$\hat{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.4187 (0.0054)	0.6575 (0.0029)	0.7935 (0.0014)
	$\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.5077 (0.0029)	0.6639 (0.0018)	0.7918 (0.0009)
	$\hat{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.4127 (0.0050)	0.6526 (0.0034)	0.7843 (0.0021)
	$\check{\rho}_{ph}(\beta)$	0.5768 (0.0032)	0.7274 (0.0027)	0.8209 (0.0018)
	$\check{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.4055 (0.0058)	0.6411 (0.0039)	0.7833 (0.0023)
	$\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.4973 (0.0042)	0.6478 (0.0033)	0.7811 (0.0021)
	$\check{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.4036 (0.0052)	0.6401 (0.0033)	0.7773 (0.0021)
2,000	Emp. Acc.	0.5446 (0.0083)	0.7063 (0.0060)	0.8038 (0.0038)
	$\hat{\rho}_{ph}(\beta)$	0.5502 (0.0001)	0.7132 (0.0001)	0.8029 (0.0001)
	$\hat{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.3710 (0.0056)	0.6297 (0.0030)	0.7633 (0.0001)
	$\hat{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.4867 (0.0026)	0.6445 (0.0015)	0.7594 (0.0001)
	$\hat{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.3578 (0.0057)	0.6197 (0.0026)	0.7488 (0.0012)
	$\check{\rho}_{ph}(\beta)$	0.5378 (0.0033)	0.6972 (0.0029)	0.7937 (0.0021)
	$\check{\rho}_{ph}(\hat{\beta}_{LASSO})$	0.3407 (0.0065)	0.5958 (0.0051)	0.7502 (0.0028)
	$\check{\rho}_{ph}(\hat{\beta}_{ADLASSO})$	0.4627 (0.0047)	0.6190 (0.0040)	0.7525 (0.0022)
	$\check{\rho}_{ph}(\hat{\beta}_{GPLASSO})$	0.3317 (0.0062)	0.5886 (0.0045)	0.7379 (0.0025)