

RESEARCH ARTICLE

WILEY Journal of CHEMOMETRICS

Model and estimators for partial least squares regression

Inge Svein Helland¹ | Solve Sæbø² | Trygve Almøy² | Raju Rimal²

¹Department of Mathematics, University of Oslo, Oslo NO-0315, Norway

²Norwegian University of Life Sciences, Ås 1430, Norway

Correspondence

Inge Svein Helland, Department of Mathematics, University of Oslo, P.O. Box 1053, Oslo NO-0316, Norway.
Email: ingeh@math.uio.no

Abstract

Partial least squares (PLS) regression has been a very popular method for prediction. The method can in a natural way be connected to a statistical model, which now has been extended and further developed in terms of an envelope model. Concentrating on the univariate case, several estimators of the regression vector in this model are defined, including the ordinary PLS estimator, the maximum likelihood envelope estimator, and a recently proposed Bayes PLS estimator. These are compared with respect to prediction error by systematic simulations. The simulations indicate that Bayes PLS performs well compared with the other methods.

KEYWORDS

Bayes PLS estimator, envelope model, partial least squares, partial least squares model, simulation

1 | INTRODUCTION

Supervised learning from multivariate data is a central problem area in applied statistics and also in chemometrics. Specifically, let our task be to predict a single variable y from a p -dimensional variable \mathbf{x} , having data on n units. From a statistical point of view, a large number of learning methods are discussed in Hastie et al.,¹ mainly under the ordinary multiple regression model. In chemometrics, partial least squares (PLS) regression is the dominating method.

Partial least squares regression has had a vigorous development in the chemometric literature since it was proposed by Wold et al.² and Martens and Næs.³ The method has been extended in several directions, and its applications have been expanded to an increasing number of fields, for instance, genomic data.⁴ Both these issues have been discussed in detail in a recent paper by Mehmood and Ahmed,⁵ where a wealth of further references may be found.

Sometimes, the issue is prediction, but very often, one also see interpretations of scoring, loading, and correlation plots; see, for instance, Martens and Martens.⁶ Such plots are not unfamiliar to statisticians in principal component connections, but they are much more used by the chemometric society, and many scientists find them informative. They are plots of the sample variants of the latent variables and parameters defined by (3), (4), and (5) below and, thus, involve consistent estimates of these quantities when $n \rightarrow \infty$ and probably also in the more general case $p/n \rightarrow 0$.

In the beginning, the PLS method was to some extent neglected or turned down by statisticians (an exception among others was Frank and Friedman⁷; see also Helland^{8,9}), but it is now included as a tool among other biased regression methods by applied statisticians. For a general discussion paper with contributions both from mathematical statisticians and chemometricians, see Sundberg.¹⁰

Indeed, there was a difference in culture between chemometricians and statisticians then, and this difference still exists to a large extent. A statement by Munck et al.¹¹ illustrates this, as seen from one side: “If chemometrics in its historical development had been limited to follow current scientific (and statistical) theories there would have been minimal progress in its wide applications today.”

Recently, the difference in culture was discussed in some detail by Martens.¹² On the one hand, Martens makes the point that the field of Chemometrics has a lot to learn from other disciplines—mathematics, statistics, and computer science.

Among other things, he says that it will not be enough to have efficient “black box” algorithms. On the other hand, he accuses statisticians in general for a predilection for “macho mathematics,” concluding in part that Chemometrics need more statistics but not more statisticians. In other parts of the paper, he talks about bridging the gap between the 2 disciplines, an effort that we whole heartedly support.

This difference in culture may in part be related to the concepts of creativity and rigor, qualities which to some extent may be called complementary. One could say that one culture puts more emphasis on creativity, the other on rigor. Of course, this is a huge simplification. First, there is a lot of creativity among statisticians, also mathematical statisticians. Secondly, one should emphasize that precise thinking also should influence practice. A case of point is the following: Chung and Keleş¹³ recently proved that the PLS regression vector is inconsistent when $p/n \rightarrow k > 0$ under a wide set of conditions. This result is probably not too well known among chemometricians; some may have a tendency to put much confidence in PLS regression when $p \sim n$ or $p > n$. It is to be emphasized that the inconsistency result in Chung and Keleş¹³ is only concerned with estimation of the regression vector. The mathematical properties of PLS as a “prediction” method when $p > n$ are largely unknown, from a statistical point of view. There is much positive empirical evidence among applied researchers on these properties, but statisticians have only started to attack this problem, since it from an analytic point of view is very difficult. In particular, see the very recent paper by Cook and Forzani,¹⁴ where asymptotic expansions allowing both n and p to be large are developed for PLS prediction with 1 component.

It is true that chemometricians have had a leading edge in the development of PLS and of certain multivariate methods, in particular, with respect to visualization etc, and they still are ahead of statisticians in this sense.

Accepting this, an important general question is what mathematical statisticians can contribute with in this development. There are relatively few papers by mathematical statisticians investigating statistical properties of the PLS regression method itself. There are however several investigations on the shrinkage properties of PLS; see Krämer¹⁵ and references there, and also Foschi¹⁶ with references. Garthwaite¹⁷ offered a simple interpretation of PLS. Stone and Brooks¹⁸ and Naik and Tsai¹⁹ discuss different generalizations of PLS; in the latter paper also, consistency of PLS is proved. In Stoica and Söderström,²⁰ an asymptotic formulae related to PLS is derived. Chun and Keleş²¹ extends consistency to the case $p/n \rightarrow 0$, introduces a sparse PLS algorithm, and compares methods by simulation. In Krämer and Sugiyama,²² the degrees of freedom of PLS regression is discussed, and this concept is used in model selection. See also references in this last paper.

In Helland and Almøy,²³ several predictors in the random x regression model were compared asymptotically as $n \rightarrow \infty$, including principal component regression (PCR) and sample PLS regression (see the next section). The conclusion was that PCR is best for very large irrelevant eigenvalues (excluded from the prediction equation), whereas PLS regression tends to be best for intermediate irrelevant eigenvalues. Because the difference is extremely small for small irrelevant eigenvalues, and because very large irrelevant eigenvalues seldom occur in practice (and if they do, they should be included in the prediction equation), it was concluded that PLS regression is the method of choice in many cases. An additional argument for PLS over PCR is that PLS involves only choosing the number of components, whereas PCR also entails deciding which of the components should be included in the prediction.

As already mentioned, Cook and Forzani¹⁴ give an asymptotic expansion of the prediction error in PLS regression, which also is informative when $p \rightarrow \infty$, but mainly limited to 1 component. Results with several components are also announced.

A vital aspect in the history of statistics is the interplay between model and estimators. Once a model is formulated, one can in principle think of several estimators in this model. A statistician will talk about a “hard” model in terms of probability distributions—at least in terms of a model equation and a statement of correlation between terms in this model. This is a concept that has had and has a great success in a number of disciplines and is at the very core of statistics as a science. Our goal in the present paper is to show that this concept can be applied—and is useful—also in connection to PLS. Specifically, our purposes are to

- stress that PLS as an algorithm can be connected to a unique statistical model (known since 1990);
- formulate 5 different ways to present this model (known in the statistical literature since 2013);
- argue that the simplest way to present the model is through the concept of relevant components—a reduction of the random x regression model;
- review briefly some statistical investigations related to PLS;
- ask if the PLS algorithm may be improved by modifying the weights;
- argue that once the model is presented, the comparison of different estimators in the model is relevant;
- present a systematic tool (`simrel`) for comparing estimators in the model with relevant components;
- present the maximum likelihood estimator in the model;

- present a Bayes estimator connected to the model;
- and compare the PLS algorithm, the maximum likelihood estimator, and the Bayes estimator in a systematic simulation study, mainly with near collinear data.

Thus, in the PLS model, one can certainly discuss other estimators than the usual PLS regression estimator, which can be seen as originating by replacing population (co)variances in the model by sample (co)variances. Two examples are the maximum likelihood estimator of Cook et al.,²⁴ see also Cook et al.^{25,26} and Cook and Zhang,²⁷ and the Bayesian estimator of Helland et al.²⁸ By simulation, both these estimators have performed well compared with PLS regression under certain conditions, but they have their disadvantages. The maximum likelihood estimator cannot be used in the case when the data matrix has rank less than p , and the Bayesian estimator requires heavy computations, in particular, when p is large.

To compare estimators, we make vital use of the recently developed simulation package `simrel`; see Sæo et al.²⁹ It is very important to have such a tool in an area where it is difficult to obtain results by purely analytical means.

We emphasize that this paper is based upon reduction of the *random* \mathbf{x} regression model. When considering latent variables from PLS, and when considering near collinearity in the observed \mathbf{x} -variables, it is natural to treat these \mathbf{x} -variables as random. It is our philosophy that this is also the best way to look upon model reduction. On the other hand, in the context of prediction, one could argue that one should condition upon the \mathbf{x} -variables and consider them as fixed. A prominent paper on PLS regression, taking fixed \mathbf{x} -variables in the basic model, is Krämer and Sugiyama,²² where further references can be found.

In recent years, there has been a rapidly growing statistical literature on the envelope model—a model generalizing the PLS model. In addition to the maximum likelihood estimation paper mentioned above, the most important papers seem to be Cook and Zhang,³⁰ where simultaneous reduction in the \mathbf{x} - and \mathbf{y} -space is proposed, and Cook and Zhang,³¹ where extensions to other regression methods than linear regression are discussed. More references can be found in these papers.

Model reduction in regression models is discussed in general from the point of view of rotations in the \mathbf{x} -space in Helland.³²

The plan of this paper is as follows: In Section 2, we formulate the model in 5 different ways, which can be shown to be equivalent. In Section 3, we define 4 different estimators in the model, including the recent Bayes PLS estimator of Helland et al.²⁸ In Section 4, we ask the question if the ordinary PLS estimator with m components can be improved by forcing the weight vector at step $m + 1$ to vanish; the answer turns out to be negative. In Section 5, we describe the simulations done for comparison of estimators with respect to prediction error, and in Section 6, we give the results of the simulations. In Section 7, we illustrate the methods on a real data set. Finally, Section 8 is a discussion section.

2 | THE MODEL: SEVERAL FORMULATIONS

Take as a point of departure the linear model

$$y = \mu_y + \beta'(\mathbf{x} - \mu_x) + \epsilon, \quad (1)$$

where β and \mathbf{x} are p -dimensional and the random predictor \mathbf{x} has mean μ_x and covariance matrix Σ_{xx} , for simplicity, assumed nonsingular here (this can be relaxed to assuming $\beta \in \text{span}(\Sigma_{xx})$ in the case where this matrix is singular; see Cook et al.,²⁴ and also C below). Independently, ϵ is distributed with mean 0 and variance σ^2 . When doing prediction from this model for near collinear data, a model reduction may be called for. Throughout this paper, a definite m -dimensional model reduction, which may be formalized in several equivalent ways, will be used. When this model holds, we say that we have an envelope model or a PLS model of dimension m or that there are m relevant components for prediction in the model.

- Given a subspace S of R^p , let \mathbf{P}_S be the projection upon S , and let \mathbf{Q}_S be the projection orthogonal to S . For simplicity, discuss the case where $\mu_x = \mathbf{0}$. Let now S be the smallest space such that (1) $\mathbf{Q}_S \mathbf{x}$ is uncorrelated with $\mathbf{P}_S \mathbf{x}$ and (2) y is uncorrelated with $\mathbf{Q}_S \mathbf{x}$ given $\mathbf{P}_S \mathbf{x}$. In this case, we may say that $\mathbf{Q}_S \mathbf{x}$ contains no linear information about y , neither directly nor through $\mathbf{P}_S \mathbf{x}$. Consider a reduction of the data to $\mathbf{P}_S \mathbf{x}$.
- Here is an algebraic characterization, which turns out to be equivalent. For a matrix \mathbf{M} , define $\mathbf{M}S$ as the space of vectors $\mathbf{M}\mathbf{z}$, as \mathbf{z} runs through S , and let S^\perp be the space perpendicular to S . Let now S be the smallest space in R^p such that (1) both $\Sigma_{xx}S \subseteq S$ and $\Sigma_{xx}S^\perp \subseteq S^\perp$ and (2) $\text{span}(\beta) \subseteq S$. In this case, we say that S is the Σ_{xx} -envelope of $\text{span}(\beta)$. It can be shown Cook et al.³³ that the envelope always exists as the smallest space with the stated properties.

C. The regression vector β can always be expanded in terms of the eigenvectors \mathbf{d}_i of Σ_{xx} :

$$\beta = \sum_{i=1}^p \gamma_i \mathbf{d}_i. \quad (2)$$

In general, when there are coinciding eigenvalues in Σ_{xx} , this expansion is not unique. However, assume that this sum can be reduced to exactly m nonzero terms: $\beta = \sum_{i=1}^m \gamma_i \mathbf{d}_i$, where the \mathbf{d}_i correspond to different eigenvalues of Σ_{xx} . We then say that there are m relevant components for prediction in the model. This reduction can be imagined to take place by 2 mechanisms: (1) Some of the γ_i 's are really zero, and (2) there are coinciding eigenvalues in Σ_{xx} . Then, one can rotate such that it is enough with 1 eigenvector for each eigenspace in the sum. In this approach, it is important that we only know that there are m nonzero terms in the sum, not which terms that are nonzero. For a closer discussion of this, see Næs and Helland³⁴ and Helland and Almøy.²³

D. Consider the population version of the well-known PLS algorithm: Take $\mathbf{e}_0 = \mathbf{x} - \mu_x$, $f_0 = y - \mu_y$, and for $a = 1, 2, \dots, m$ compute successively:

$$\mathbf{w}_a = \text{cov}(\mathbf{e}_{a-1}, f_{a-1}), \quad t_a = \mathbf{w}_a' \mathbf{e}_{a-1}, \quad (3)$$

$$\mathbf{p}_a = \text{cov}(\mathbf{e}_{a-1}, t_a) / \text{var}(t_a), \quad q_a = \text{cov}(f_{a-1}, t_a) / \text{var}(t_a), \quad (4)$$

$$\mathbf{e}_a = \mathbf{e}_{a-1} - \mathbf{p}_a t_a, \quad f_a = f_{a-1} - q_a t_a.$$

It can be proved⁹ and is important in this connection that under the reduced model C, this algorithm stops automatically after m steps when $m < p$: It stops because $\mathbf{w}_{m+1} = \text{cov}(\mathbf{e}_m, f_m) = 0$. After those m steps, we get the representations

$$\mathbf{x} = \mu_x + \mathbf{p}_1 t_1 + \dots + \mathbf{p}_m t_m + \mathbf{e}_m, \quad y = \mu_y + q_1 t_1 + \dots + q_m t_m + f_m \quad (5)$$

with the corresponding PLS population prediction

$$y_{m,PLS} = \mu_y + q_1 t_1 + \dots + q_m t_m = \mu_y + \beta_{m,PLS}' (\mathbf{x} - \mu_x).$$

Theorem 1. (Helland⁹ and Cook et al²⁴)

- (a) The 2 conditions A and B on the space S are equivalent.
- (b) The models formulated by C and D are equivalent.
- (c) When there are m relevant components for prediction, the envelope space S has dimension m , and S can be taken as $\text{span}(\mathbf{w}_1, \dots, \mathbf{w}_m) = \text{span}(\mathbf{d}_1, \dots, \mathbf{d}_m)$.
- (d) When the envelope space has dimension m , there are m relevant components for prediction.
- (e) In this case, we have $\beta_{m,PLS} = \beta$.

Proof. (a) is proved in Cook et al,²⁴ Proposition 1 and (b) in Helland.⁹ Theorem 2 Finally, (c)-(e) and the equivalence with E below are contained in Cook et al.²⁴ Proposition 5 □

In this sense, all the model formulations (A-D) are equivalent; they describe the same reduced model. In Helland⁹ and Cook et al,²⁴ a fifth equivalent formulation in terms of a Krylov sequence is also given:

E. S is also spanned by the vectors $\sigma_{xy}, \Sigma_{xx} \sigma_{xy}, \dots, \Sigma_{xx}^{m-1} \sigma_{xy}$, and m is the smallest integer such that $\beta = \Sigma_{xx}^{-1} \sigma_{xy}$ belongs to S .

The simplest way to express the model reduction implied by PLS seems to be C. In analogy with the exivalence between A and B, this can also be expressed as a reduction of the x vector. Consider again the centered case $\mu_x = \mathbf{0}$. For details, see Næs.³⁴

C'. Let \mathbf{R} be a nonrandom $p \times m$ matrix of rank m . Normalize such that $\mathbf{R}' \mathbf{R} = \mathbf{I}$. There are m relevant components $\mathbf{R}' \mathbf{x}$ for predicting y if and only if \mathbf{R} can be found such that (a) $\beta \in \text{span}(\mathbf{R})$ and (b) $\text{span}(\mathbf{R})$ is spanned by eigenvectors of Σ_{xx} .

Being a reduced model that can be motivated in so many different ways, it is definitively of interest to find a good estimator of the regression vector β under this model.

3 | ESTIMATORS IN THE PLS/ENVELOPE MODEL

Now that the PLS model is introduced, we will start to look at estimators of the parameters in this model, in particular, estimators of β , which will give prediction. Of special interest is estimators that perform well in the case of near collinear data. Some estimators are already known from the literature.

- The ordinary PLS estimator can be introduced as follows: With data (\mathbf{X}, \mathbf{y}) , take initial values $\mathbf{E}_0 = \mathbf{X} - \bar{\mathbf{x}}\mathbf{1}'$ and $\mathbf{f}_0 = \mathbf{y} - \bar{y}\mathbf{1}$. Run the population PLS algorithm for A steps with population (co)variances replaced by sample (co)variances. Ordinarily, A is found by cross-validation or by similar means. Note that from D in Section 2, the m -step PLS model is characterized by $\mathbf{w}_{m+1} = \text{cov}(\mathbf{e}_m, \mathbf{f}_m) = \mathbf{0}$. Theoretically, when $A = m$, we cannot expect the sample weights $\hat{\mathbf{w}}_{m+1}$ to be zero. However, since any continuous function of the sample covariances and variances is consistent for the same function of the population covariances and variances, since $\hat{\mathbf{w}}_{m+1}$ through the PLS algorithm is such a function and since $\mathbf{w}_{m+1} = \mathbf{0}$, we will have $\lim_{n \rightarrow \infty} \hat{\mathbf{w}}_{m+1} = \mathbf{0}$ almost surely.
- The sparse regression SPLS of Chun and Keleş²¹: This requires 2 effective tuning parameters, and it also aims at variable selection. Sparse partial least squares (SPLS) seems to be better than ordinary PLS in certain cases, also when variable selection is not an issue.
- When $\mathbf{S} = (\mathbf{X} - \bar{\mathbf{x}}\mathbf{1}')'(\mathbf{X} - \bar{\mathbf{x}}\mathbf{1}')$ has rank p , which specifically requires $n > p$, the maximum likelihood estimator of β under the multinormal envelope model was given in Cook et al.²⁴ This estimator is of course very useful, but it cannot be used for small n . Modifications of the maximum likelihood estimator, which cover also this case, were recently indicated in Cook et al.²⁵ That paper also gives a MATLAB toolbox for maximum likelihood estimation in the envelope model and in several generalizations of this model. A faster algorithm for maximum likelihood estimation is discussed in Cook and Zhang.²⁷ Even faster algorithms with modifications to small sample size $n < p$ are recently described in Cook and Zhang,³⁵ and an R-package was recently described by Cook et al.²⁶
- Under a specific rotation-invariant prior, the Bayes estimator of β under the model with m relevant components was given in Helland et al.²⁸ This estimator was shown to be close to the best equivariant estimator, but it requires heavy computation.

The estimation was performed by a Markov Chain Monte Carlo approach. Specifically, for given m , and for observed centered data \mathbf{y} and \mathbf{X} , the likelihood function is proportional to

$$f(\mathbf{y}, \mathbf{X} | \mathbf{v}, \boldsymbol{\gamma}, \mathbf{D}, \sigma^2) \propto (\sigma^2)^{-n/2} \exp \left(-\frac{1}{2\sigma^2} \left(\mathbf{y} - \mathbf{X} \sum_{i=1}^m \gamma_i \mathbf{d}_i \right)' \left(\mathbf{y} - \mathbf{X} \sum_{i=1}^m \gamma_i \mathbf{d}_i \right) \right) \times \left(\prod_{i=1}^p v_i \right)^{-n/2} \prod_{j=1}^n \exp \left(-\frac{1}{2} \mathbf{x}_j' \left(\sum_{i=1}^p \frac{1}{v_i} \mathbf{d}_i \mathbf{d}_i' \right) \mathbf{x}_j \right), \quad (6)$$

where $\mathbf{v} = [v_1, \dots, v_p]$ and $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_p]$ are the eigenvalues and the eigenvectors of the \mathbf{x} -covariance matrix Σ_{xx} and $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_m]$ are regression parameters of the PLS model.

As argued in Helland et al.,²⁸ a near optimal equivariant regressor is found as the Bayesian estimator under rotation-invariant prior for $\mathbf{d}_1, \dots, \mathbf{d}_p$ and prior $\pi(\boldsymbol{\gamma}) = \prod_i 1/\gamma_i^{1-\epsilon}$, where $1/\epsilon$ is a large uneven integer. Slightly modified scale priors are also chosen for \mathbf{v} as $\pi(\mathbf{v}) = \prod_i 1/v_i \exp(-\epsilon_v/2v_i)$ and for σ^2 as $\pi(\sigma^2) = 1/\sigma^2 \exp(-\epsilon_\sigma/2\sigma^2)$. Here, ϵ_v and ϵ_σ are some small numbers chosen to ensure properness of the posterior distribution. Estimation of model parameters may be done by means of Markov chain Monte Carlo methods. As shown in Helland et al.,²⁸ the marginal posterior distributions for σ^2 and v_i (for $i = 1, \dots, p$) are, for the given prior distributions, all inverse gamma distributions. Furthermore, the marginal posterior distributions for γ_i (for $i \in 1, \dots, m$) are approximately normally distributed. There is no closed form posterior distribution for \mathbf{D} , hence a random walk step with a Metropolis-Hastings acceptance step is necessary for the sampling from the posterior distributions of the parameters. R-code for the Bayes estimator is available at <http://www.github.com/solvsa/BayesPLS>, and further details on the Markov chain Monte Carlo implementation may be found in the supplementary documentation to Helland et al.²⁸

By simulation, both the maximum likelihood estimator c and the Bayes estimator d were shown to perform well compared to the PLS estimator a. These 2 estimators assume a multinormal distribution of the data in their derivation, but the estimators themselves are valid under more general assumptions. Both the chemometric tradition and the envelope model of Cook et al.^{24,33} demand no detailed distributional assumptions.

4 | CAN A BETTER ESTIMATOR BE FOUND BY SIMPLE MEANS?

The m step PLS model is characterized by the constraint $\mathbf{w}_{m+1} = \text{cov}(\mathbf{e}_m, f_m) = \mathbf{0}$. However, in the sample PLS algorithm, $\hat{\mathbf{w}}_{m+1}$ is a continuous random variable if the data are continuous. Hence, almost surely, $\hat{\mathbf{w}}_{m+1} \neq \mathbf{w}_{m+1} = \mathbf{0}$. This means that the estimator of the vector of PLS parameter falls outside the corresponding parameter space. On the other hand, by standard statistical theory, the maximum likelihood estimator and any Bayes estimator are always in the parameter space.

In this section, we ask the question whether we can improve the PLS algorithm in some way such that $\hat{\mathbf{w}}_{m+1} = \mathbf{0}$ for the improved algorithm. That is, we seek modified weights $\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_m$ such that $\hat{\mathbf{w}}_{m+1} = \mathbf{0}$ in the modified algorithm. Unfortunately, the answer to this question is no. This programme is only possible when \mathbf{S} is invertible, and then it by necessity leads to the least squares solution. Let $\hat{\mathbf{W}}_A = (\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_A)$ for any A .

First, we need some properties of the ordinary PLS algorithm.

Proposition 2. *At each step, the PLS weights satisfy*

$$\hat{\mathbf{w}}_{A+1} = \mathbf{s} - \mathbf{S}\hat{\mathbf{W}}_A(\hat{\mathbf{W}}_A'\mathbf{S}\hat{\mathbf{W}}_A)^{-1}\hat{\mathbf{W}}_A'\mathbf{s}, \quad (7)$$

and the A step regression vector is

$$\hat{\beta}_A = \hat{\mathbf{W}}_A(\hat{\mathbf{W}}_A'\mathbf{S}\hat{\mathbf{W}}_A)^{-1}\hat{\mathbf{W}}_A'\mathbf{s}. \quad (8)$$

Proof. These relations were proved in Helland,⁸ see equations (3.3) and (3.7) there, and were also used in Cook et al.²⁴ \square

Now, fix m . To find an algorithm such that $\hat{\mathbf{w}}_{m+1} = \mathbf{0}$, we will have to modify the weights $\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_m$.

Definition 1. For the purpose of this section, call a restricted PLS prediction any prediction method based on an estimator of β of the form (8) for $A = m$ such that

- 1.) $\hat{\mathbf{W}}_m = (\hat{\mathbf{w}}_1, \dots, \hat{\mathbf{w}}_m)$ is modified with respect to PLS in some way.
- 2.) Equation 7 holds for $A = m$ and gives $\hat{\mathbf{w}}_{m+1} = \mathbf{0}$.

Theorem 3. *An RPLS prediction method exists if and only if \mathbf{S} is invertible and $\mathbf{S}^{-1}\mathbf{s} \in \text{span}\hat{\mathbf{W}}_m$. In that case, $\hat{\beta}$ is equal to the least squares estimator $\mathbf{S}^{-1}\mathbf{s}$.*

Proof. Assume that (7) holds for $A = m$ and $\hat{\mathbf{w}}_{m+1} = \mathbf{0}$. Then, $\mathbf{s} = \mathbf{S}\hat{\mathbf{W}}_m(\hat{\mathbf{W}}_m'\mathbf{S}\hat{\mathbf{W}}_m)^{-1}\hat{\mathbf{W}}_m'\mathbf{s}$. This is possible for general \mathbf{s} only if \mathbf{S} is nonsingular, and then it is equivalent to $\mathbf{R}\sqrt{\mathbf{S}}^{-1}\mathbf{s} = \sqrt{\mathbf{S}}^{-1}\mathbf{s}$ with $\mathbf{R} = \mathbf{A}(\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}'$, where $\mathbf{A} = \sqrt{\mathbf{S}}\hat{\mathbf{W}}_m$. Since \mathbf{R} is the projector upon $\text{span}(\mathbf{A})$, this is again equivalent to $\sqrt{\mathbf{S}}^{-1}\mathbf{s} \in \text{span}(\sqrt{\mathbf{S}}\hat{\mathbf{W}}_m)$, or $\mathbf{S}^{-1}\mathbf{s} \in \text{span}(\hat{\mathbf{W}}_m)$. Then, putting $\mathbf{s} = \mathbf{S}\hat{\mathbf{W}}_m\mathbf{q}$ in (8) for some \mathbf{q} gives $\hat{\beta} = \hat{\mathbf{W}}_m\mathbf{q} = \mathbf{S}^{-1}\mathbf{s}$. \square

Thus, Theorem 3 shows clearly that it is not possible to modify the PLS weights in a nontrivial way such that the modified estimator belongs to the parameter space.

5 | DATA SIMULATIONS FOR COMPARISON OF ESTIMATORS

A comparative study of the prediction performances of the regular PLS algorithm, the maximum likelihood envelope method, the Bayes PLS, and the method of ordinary least squares (OLS) was performed on data simulated from the random regression model (1) and a real dataset measuring various properties and near infrared (NIR) spectra of diesel fuels. This and the following section will focus on simulation study in detail. In the study, we consider envelope method for predictor reduction²⁴ and use R-code discussed in Cook et al.²⁶ A detailed description of the simulation procedure can be found in Sæbø et al.²⁹ with the accompanying R-package `simrel`, but key features of the approach are presented next. The simulation set up is best explained from reexpressing model (1) in the Gaussian case as

$$\begin{bmatrix} y \\ \mathbf{x} \end{bmatrix} \sim \mathcal{N}(\mu_{yx}, \Sigma_{yx}) = \mathcal{N}\left(\begin{bmatrix} \mu_y \\ \mu_x \end{bmatrix}, \begin{bmatrix} \sigma_y^2 & \sigma_{xy}^t \\ \sigma_{xy} & \Sigma_{xx} \end{bmatrix}\right), \quad (9)$$

where σ_{xy} is a vector holding the covariances between the predictors (\mathbf{x}) and the response (y). The vector of regression coefficients β is by standard theory given as $\beta = \Sigma_{xx}^{-1} \sigma_{yx}$, which in turn can be expressed in terms of the eigenvalues v_1, \dots, v_p and the eigenvectors $\mathbf{d}_1, \dots, \mathbf{d}_p$ of Σ_{xx} :

$$\beta = \sum_{i=1}^p \frac{\mathbf{d}_i' \sigma_{yx}}{v_i} \mathbf{d}_i = \sum_{i=1}^p \gamma_i \mathbf{d}_i, \quad (10)$$

as given in Equation 2. In `simrel`, the following simplifying assumptions are made:

- ▷ It is assumed that $v_i = e^{-\eta(i-1)}$ for $i = 1, \dots, p$, implying $v_1 = 1$ (which we may assume without loss of generality) and that all subsequent eigenvalues are decaying according to the size of the parameter η . A large η gives a rapid decrease in eigenvalues, implying high level of multicollinearity in \mathbf{x} .
- ▷ It is assumed that $m \leq p$ eigenvectors are relevant for y , which means that Equation 10 (potentially) reduces to

$$\beta = \sum_{i \in \mathcal{P}} \gamma_i \mathbf{d}_i, \quad (11)$$

where m -vector \mathcal{P} is the set of indices of the relevant components (relpos) for which $\gamma_i = \mathbf{d}_i' \sigma_{xy} / v_i \neq 0$. Hence, the envelope or the relevant space has dimension m (see Theorem 1).

- ▷ Without loss of generality, it is further assumed that $\sigma_y = 1$, $\mu_y = 0$ and $\mu_x = \mathbf{0}$.

In `simrel`, the actual values of σ_{xy} were set to satisfy a prespecified value of the population coefficient of determination ρ^2 . It may be shown that under the above assumptions, $\rho^2 = \sigma_{xy}' \Sigma_{xx}^{-1} \sigma_{xy}$. This completes the specification of the parameters used in `simrel`, and in the present comparison study, a design for the simulated data sets in terms of these parameters were as defined in Table 1.

From the possible combination of the above parameters, 32 calibration sets were simulated with 5 replications of each, ie, there were 160 calibration sets (`datasets`) altogether.

6 | SYSTEMATIC COMPARISONS

A systematic comparison of the methods across the simulation designs was made on the basis of their ability to predict test samples. Since the distribution of the simulated variables is fully known, the expected mean squared error of prediction (MSEP) based on some $\hat{\beta}$ estimated from a calibration set may be found as

$$E_x [E_y(y - \hat{y})^2] = \left[\sigma^2 + E(\hat{\beta} - \beta)' \Sigma_{xx} (\hat{\beta} - \beta) \right] \frac{n+1}{n} \quad (12)$$

in the model. The expectation on the right-hand side of the above expression is estimated for each method and for each design as an average over the 5 replicated calibration sets. To study the effects of p , ρ^2 , `relpos`(\mathcal{P}), `Method`, and (η) along with their interactions, we first retrieved the minimum MSEP for each method across 1 to 10 components (assumed numbers of relevant components). In Figure 1, interaction plots for these data properties are displayed.

The effect of the third-order interaction between p , ρ^2 and `Methods`, which we see in Figure 1 (left), shows that the maximum likelihood-based estimation methods, in our case, the envelope and the OLS, perform poorly on data sets with large number of variables and low ρ^2 . Still, the performance of the envelope is better than OLS also in situations where $p = 40$ and $n = 50$, representing here $p \sim n$. The interaction plots suggest that the Bayes PLS and ordinary PLS estimation methods are better and more stable on average than the two other methods.

Similarly, the effect of third-order interaction between `relpos`, η , and `Method` in Figure 1 (right) shows that OLS method gives higher prediction error than other methods, but the effect of `relpos` is small but notable for the envelope method. Again, Bayes PLS and ordinary PLS are best.

TABLE 1 Parameters used for simulating calibration sets

Number of training samples	n	50
Number of predictor variables	p	15 and 40
Population coefficient of determination	ρ^2	0.5 and 0.9
Position of relevant components	\mathcal{P}	$\triangleright 1$, $\triangleright 1, 2$ $\triangleright 1, 3$ $\triangleright 2, 3$ and $\triangleright 1, 2, 3$
Decay factor of eigenvalues of Σ_{xx}	η	0.5 and 0.9

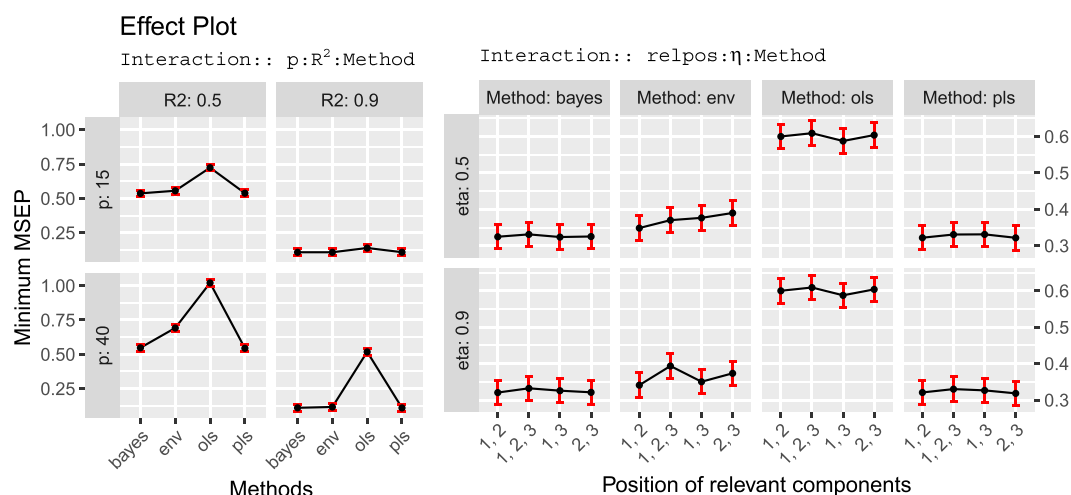


FIGURE 1 Third-order interaction effects. MSE, mean squared error of prediction; ENV, envelope; OLS, ordinary least squares; PLS, partial least squares

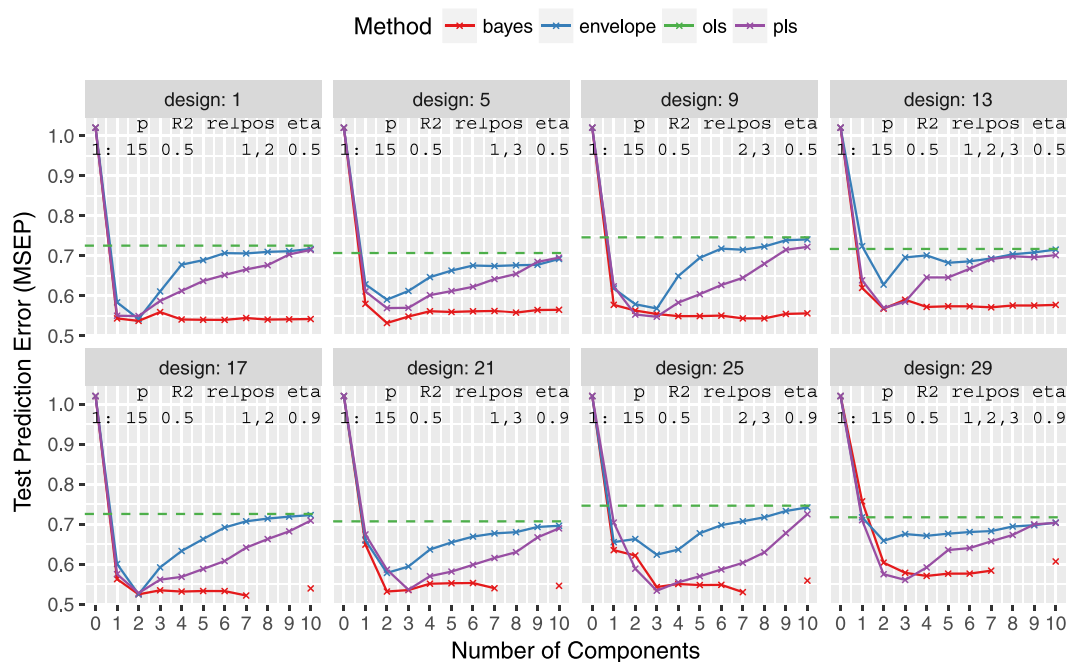


FIGURE 2 Average prediction error for designs with 15 predictor variables where coefficient of determination is 0.5. MSE, mean squared error of prediction

The prediction error plots below are organized into 4 groups: (a) $p = 15$, $\rho^2 = 0.5$; (b) $p = 15$, $\rho^2 = 0.9$; (c) $p = 40$, $\rho^2 = 0.5$; and (d) $p = 40$, $\rho^2 = 0.9$. The OLS prediction error is shown by a straight dotted line.

In group (a), with small number of variables ($p \ll n$) and noisy data ($\rho^2 = 0.5$), Figure 2 shows that all the estimation methods performed better than OLS for all designs in this group, Bayes PLS being best in nearly all cases. Some convergence problems with Bayes PLS when eigenvalues decrease rapidly can be ignored since the minimum MSE is already obtained from fewer components.

Having few variables rich with information ($\rho^2 = 0.9$), the designs in group (b) (Figure 3) leads to easy prediction with low prediction error in general for all methods. All the methods including OLS have small MSEs, but the other methods are still dominant. In most of the situations, Bayes PLS has reached minimum error with only 1 component. In this group, the performance of envelope is better than regular PLS, and the minimum error for envelope is also achieved with fewer components.

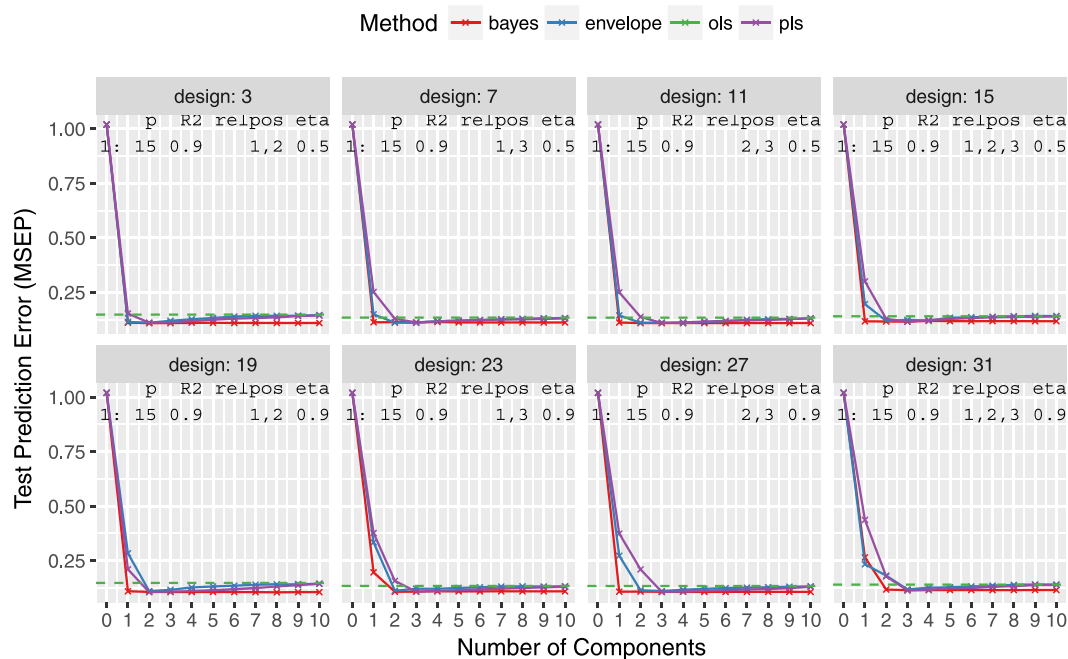


FIGURE 3 Average prediction error for designs with 15 predictor variables where coefficient of determination is 0.9. MSE, mean squared error of prediction

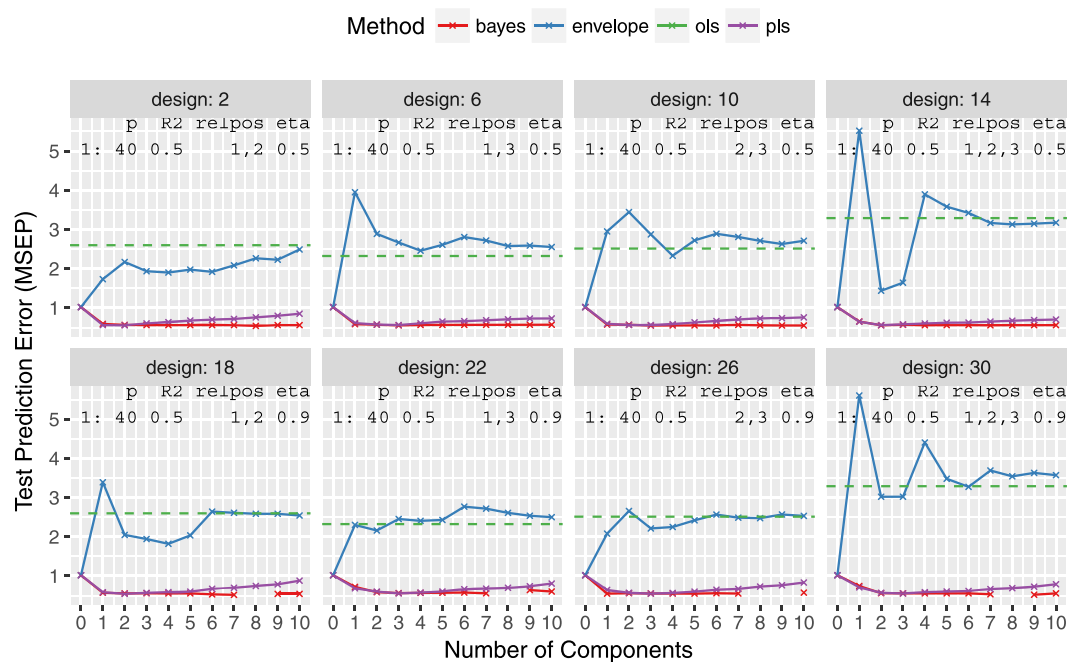


FIGURE 4 Average prediction error for designs with 40 predictor variables where coefficient of determination is 0.5. MSE, mean squared error of prediction

Low information content combined with many predictor variables characterize the designs in group (c), and prediction is in general difficult for these designs. In Figure 4, the methods based on maximum likelihood estimation performed poorly and often poorer than an average guess. Bayes PLS and regular PLS performed well, as in the previous designs.

With 40 predictors ($p \sim n$) and rich information (high ρ^2) (designs in group d), Figure 5 shows that in most of the situations (except in design 16), the envelope method has nearly attained true minimum error (0.1) and has outperformed OLS. However, its prediction error is still larger than Bayes PLS and PLS. Bayes PLS and PLS methods are highly stable

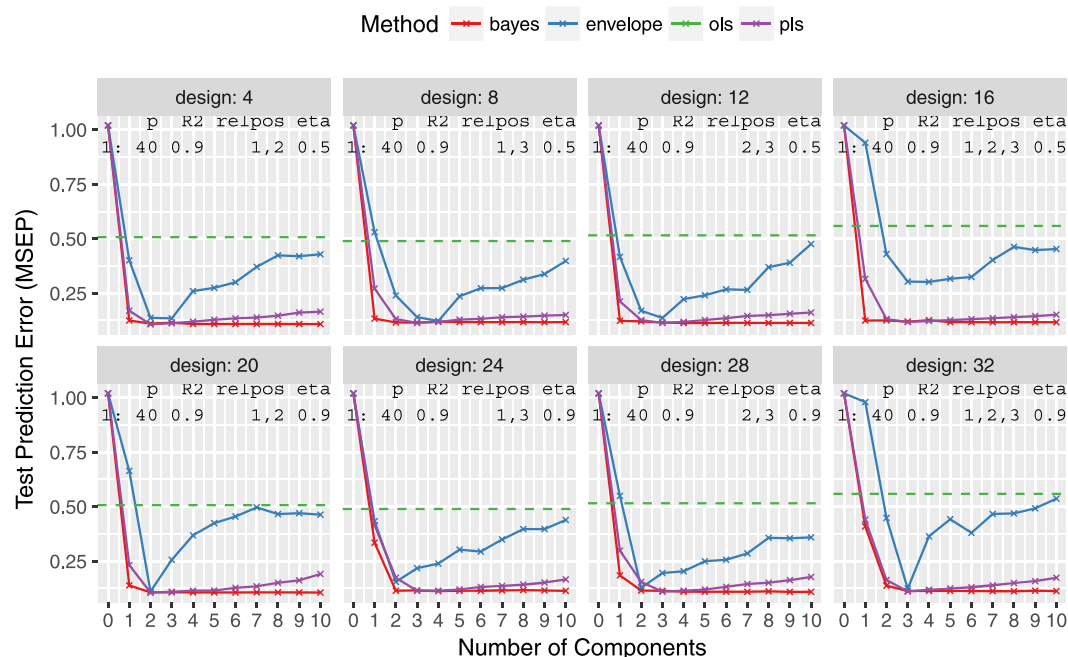


FIGURE 5 Average prediction error for designs with 40 predictor variables where coefficient of determination is 0.9. MSEP, mean squared error of prediction

and are closer to true minimum error. Further, Bayes PLS is able to obtain its minimum prediction error with only a small number of components.

In general, ordinary PLS is very stable in all situations. It is extensible (lots of variants has been developed after its introduction), easy, and less time consuming to fit than Bayes PLS and the envelope method. If the issue is to get closer prediction from squeezing information as much as possible, Bayes PLS will be a good alternative. Its performance with varying number of components is stable and better in all designs studied here. The envelope method performed better than OLS, and the performance increased for informative data ($\rho^2 = 0.9$). However, it has an increased error with additional components in many situations.

Correlation between estimated and true regression coefficients (β) along with the mean square error of estimation is presented for 4 designs in Figure 6. In case of ordinary PLS and the envelope method, the correlation for design 1 from group (a) and design 3 from group (b), both having 15 predictors, is high for small components. However, for design 2 from group (c) and design 4 from group (d), envelope methods exhibit sudden decrease in the correlation with corresponding increase in estimation error. The impressive prediction performance of Bayes PLS is also seen from the high correlation of estimated coefficients and true coefficients. In addition, the average mean square error of regression for this methods is also small compared with others for all the components.

Although having low prediction error in case of envelope estimation method, the coefficient estimates are highly unstable for different components, which we can see from its variation in correlation with true coefficients (Figure 6, top). Bayes PLS and regular PLS estimates are more stable over different replicates and for different components (Figure 6, bottom) especially when $p \sim n$. This stability agrees with the low prediction error we have discussed before.

7 | COMPARISON OF ESTIMATORS USING NIR SPECTRA OF DIESEL FUELS

Let us consider an example using a real dataset. In this example, we have used data from <http://www.eigenvector.com/data/SWRI/>, which consists of NIR spectra of diesel fuels with different properties measured such as Catane Number. Since the variables in NIR spectra are highly correlated, we have selected a subset of every 10th variable as predictors and the property Catane Number as response. After removing missing observations, the first 150 observations were used as calibration set, and the rest 231 were used as validation set.

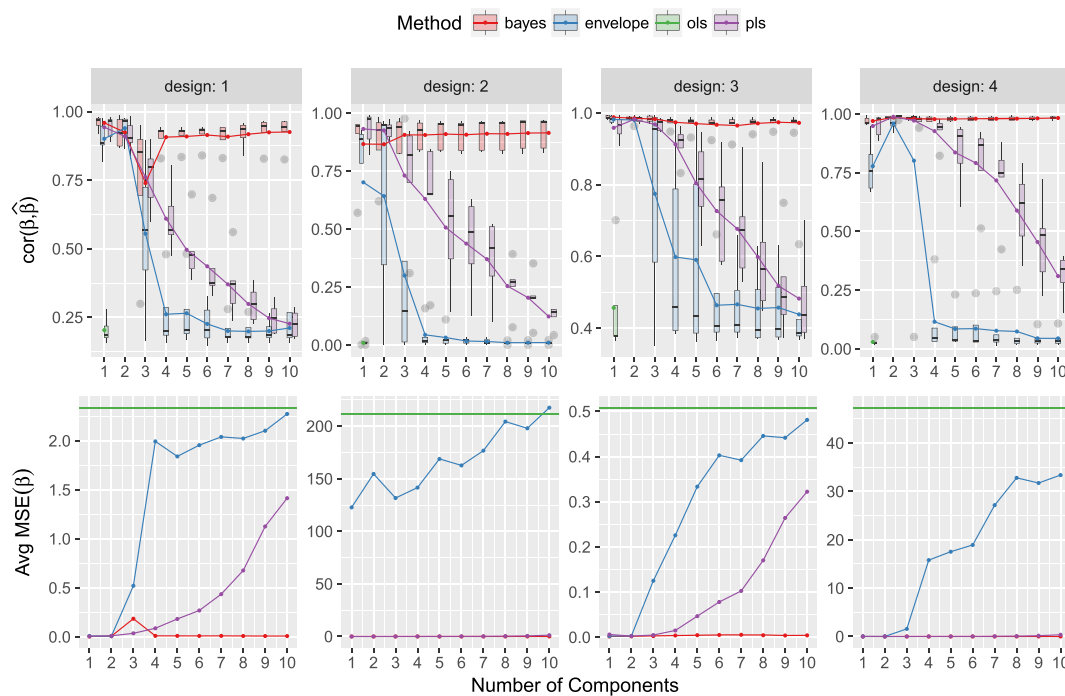


FIGURE 6 Correlation between true and estimated beta coefficient and beta estimation error. Box plots on the plots in first row show the variation in the correlation for each estimator and number of components used

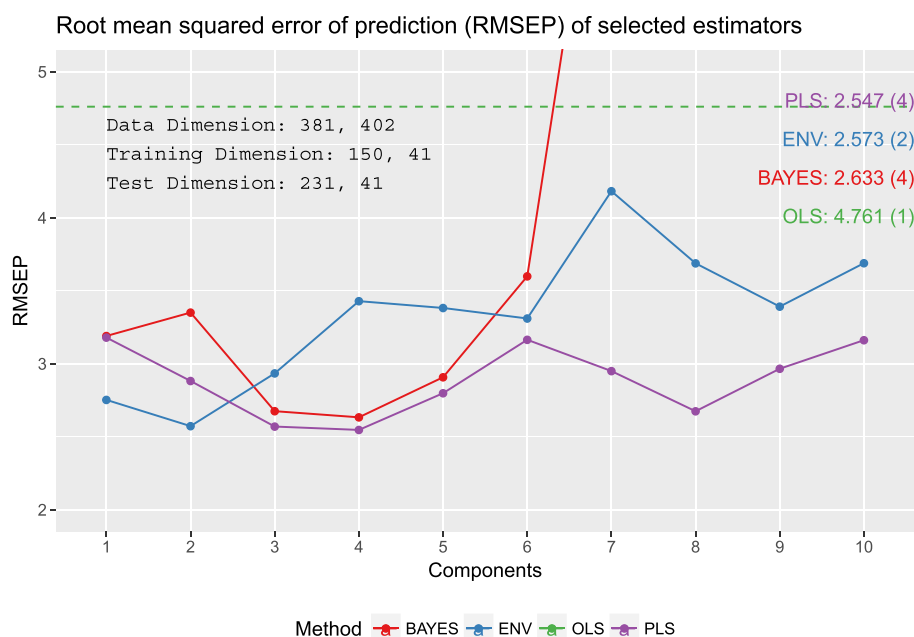


FIGURE 7 Root mean square error of prediction (RMSEP) from different estimators. Missing values were omitted in training and test datasets. ENV, envelope; OLS, ordinary least squares; PLS, partial least squares

Using the calibration set, a model with 1 to 10 components were fitted using PLS, envelope, and Bayes PLS methods. An OLS method was also fitted for reference. With each of these fitted models, the validation (test) set was used for prediction, and the root MSE was measured. Based on the prediction error, Figure 7 compares the estimators we have considered.

The results from Figure 7 show quite different results from the systematic simulation study, mainly for Bayes PLS estimation. By using 3 and 4 components, the prediction from PLS and Bayes PLS is similar and can be considered their best. Envelope model is able to attain similar prediction error just in 2 components. It is important to notice that Bayes PLS

and envelope methods here are rather sensitive to the extra number of components, which also suggest that over-fitting must be examined before using the model for predicting new observations. In the example, all the methods have significant better performance than OLS.

8 | DISCUSSION

The purpose of the present article has been to discuss the approach to PLS regression via model reduction in the random \mathbf{x} multiple regression model, and to compare estimators in this reduced model.

From simulations, the Bayes estimator under the PLS model seems to have very good properties. In virtually all of the 32 designs, the MSEP curve for Bayes PLS lies below that for ordinary PLS and also that for the maximum likelihood envelope model. A particularly desirable feature of Bayes PLS is that the MSEP curve seems to be almost flat for varying number of components. Thus, the error made by choosing a wrong number of components m by cross-validation must be expected to be small.

Envelope and Bayes PLS estimation methods, when compared with ordinary PLS methods, display better prediction performance (only when p is small for the envelope method). However, both of them have their disadvantages. The envelope method, as based on maximum likelihood, breaks down when p approaches n , while Bayes PLS has time-consuming computation, and in our simulations, it failed to converge for some cases.

However, in the results in the example using real data, the performance of Bayes PLS estimator is in contrast to its result from the simulated data. Since the predictors are highly correlated, only a few number of components are sufficient for the prediction, but when an extra number of components were used, the estimators seem to be influenced by the noise which increases with each additional component. In this respect, a more thorough study on Bayes PLS should be done for its contrast results on simulated and real dataset. A convergence issue in Bayes PLS can be suspected for the reason as seen in the example using simulated data.

For practical purposes, the ordinary PLS algorithm still seems to be a good option for prediction purposes, but from a statistical point of view, a closer study of its properties as $p \rightarrow \infty$ seems to be called for. We feel that the model approach of the present paper may give a good framework for such a study, both in terms of asymptotic expansions and in terms of further simulations. Such simulations may also include the cross-validated LASSO and other methods such as ridge regression, but note that these estimators are derived from other considerations than that of predicting the effect of relevant components.

This paper has been concentrated on the case of univariate response. We hope to discuss the multivariate case later.

ORCID

Inge Svein Helland  <http://orcid.org/0000-0002-7136-873X>

REFERENCES

1. Hastie T, Tibshirani R, Friedman J. *The Elements of Statistical Learning*, 2nd ed. Springer: New York; 2009.
2. Wold S, Martens H, Wold H. The multivariate calibration problem in chemistry solved by the PLS method. In: Ruhe A, Kågstöm B, eds. *Proc. Conf. Matrix Pencils. March 1982. Lecture Notes in Mathematics*. Heidelberg: Springer Verlag; 1983:286-293.
3. Martens H, Næs T. *Multivariate Calibration*. Chichester and New York: John Wiley & Sons; 1989.
4. Boulesteix AL, Strimmer K. Partial least squares: a versatile tool for the analysis of high-dimensional genomic data. *Briefings Bioinf.* 2007;8(1):32-44.
5. Mehmood T, Ahmed B. The diversity in the applications of partial least squares: an overview. *J Chemom.* 2016;30(1):4-17.
6. Martens H, Martens M. *Multivariate Analysis of Quality. An Introduction*. Bristol, UK: IOP Publishing; 2001.
7. Frank LE, Friedman JH. A statistical view of some chemometrics regression tools. *Technometrics.* 1993;35(2):109-135.
8. Helland IS. On the structure of partial least squares regression. *Commun Stat-Simul Comput.* 1988;17(2):581-607.
9. Helland IS. Partial least squares regression and statistical models. *Scand J Stat.* 1990;17:97-114.
10. Sundberg R. Multivariate calibration—direct and indirect regression methodology. *Scand J Stat.* 1999;26(2):161-207.
11. Munck L, Jespersen BM, Rinnan Å, et al. A physiochemical theory on the applicability of soft mathematical models—experimentally interpreted. *J Chemom.* 2010;24(7-8):481-495.
12. Martens H. Quantitative big data: where chemometrics can contribute. *J Chemom.* 2015;29:563-581.
13. Chung D, Keleş S. Sparse partial least squares classification for high dimensional data. *Stat Appl Genet Mol Biol.* 2010;9(1):17.

14. Cook RD, Forzani L. Big data and partial least-squares prediction. *Can J Stat*. 2017;46:62-78.
15. Krämer N. An overview on the shrinkage properties of partial least squares regression. *Comput Stat*. 2007;22(2):249-273.
16. Foschi P. *The geometry of PLS shrinkages*, University of Bologna; 2015.
17. Garthwaite PH. An interpretation of partial least squares. *J Am Stat Assoc*. 1994;89(425):122-127.
18. Stone M, Brooks RJ. Continuum regression: cross-validated sequentially constructed prediction embracing ordinary least squares, partial least squares and principal components regression. *J R Stat Soc Ser B (Methodological)*. 1990;52(2):237-269.
19. Naik P, Tsai CL. Partial least squares estimator for single-index models. *J R Stat Soc: Ser B (Statistical Methodology)*. 2000;62(4):763-771.
20. Stoica P, Söderström T. Partial least squares: a first-order analysis. *Scand J Stat*. 1998;25(1):17-24.
21. Chun H, Keleş S. Sparse partial least squares regression for simultaneous dimension reduction and variable selection. *J R Stat Soc: Ser B (Statistical Methodology)*. 2010;72(1):3-25.
22. Krämer N, Sugiyama M. The degrees of freedom of partial least squares regression. *J Am Stat Assoc*. 2012;106(494):697-705.
23. Helland IS, Almøy T. Comparison of prediction methods when only a few components are relevant. *J Am Stat Assoc*. 1994;89(426):583-591.
24. Cook R, Helland I, Su Z. Envelopes and partial least squares regression. *J R Stat Soc: Ser B (Statistical Methodology)*. 2013;75(5):851-877.
25. Cook D, Su Z, Yang Y, et al. envlp: A MATLAB toolbox for computing envelope estimators in multivariate analysis. *J Stat Software*. 2015;62(1):1-20.
26. Cook RD, Forzani L, Su Z. A note on fast envelope estimation. *J Multivariate Anal*. 2016;150:42-54.
27. Cook RD, Zhang X. Algorithms for envelope estimation. *J Comput Graphical Stat*. 2016;25(1):284-300.
28. Helland IS, Sæbø S, Tjelmeland H, et al. Near optimal prediction from relevant components. *Scand J Stat*. 2012;39(4):695-713.
29. Sæbø S, Almøy T, Helland IS. simrel—a versatile tool for linear model data simulation based on the concept of a relevant subspace and relevant predictors. *Chemom Intell Lab Syst*. 2015;146:128-135.
30. Cook RD, Zhang X. Simultaneous envelopes for multivariate linear regression. *Technometrics*. 2015;57(1):11-25.
31. Cook RD, Zhang X. Foundations for envelope models and methods. *J Am Stat Assoc*. 2015;110(510):599-611.
32. Helland IS. Reduction of regression models under symmetry. *Contemp Math*. 2001;287:139-154.
33. Cook RD, Li B, Chiaromonte F. Envelope models for parsimonious and efficient multivariate linear regression. *Stat Sin*. 2010;20(3):927-960.
34. Næs T, Helland IS. Relevant components in regression. *Scand J Stat*. 1993;20:239-250.
35. Cook RD, Zhang X. Fast envelope algorithms. *Stat Sin*. 2018;28(3):28.

How to cite this article: Helland IS, Sæbø S, Almøy T, Rimal R. Model and estimators for partial least squares regression. *Journal of Chemometrics*. 2018;32:e3044. <https://doi.org/10.1002/cem.3044>