



Comparison of multi-response prediction methods

Raju Rimal^{a,*}, Trygve Almøy^a, Solve Sæbø^b

^a Faculty of Chemistry and Bioinformatics, Norwegian University of Life Sciences, Ås, Norway

^b Norwegian University of Life Sciences, Ås, Norway



ARTICLE INFO

Keywords:

Model-comparison
Multi-response
Simrel

ABSTRACT

While data science is battling to extract information from the enormous explosion of data, many estimators and algorithms are being developed for better prediction. Researchers and data scientists often introduce new methods and evaluate them based on various aspects of data. However, studies on the impact of/on a model with multiple response variables are limited. This study compares some newly-developed (envelope) and well-established (PLS, PCR) prediction methods based on real data and simulated data specifically designed by varying properties such as multicollinearity, the correlation between multiple responses and position of relevant principal components of predictors. This study aims to give some insight into these methods and help the researcher to understand and use them in further studies.

1. Introduction

The prediction has been an essential component of modern data science, whether in the discipline of statistical analysis or machine learning. Modern technology has facilitated a massive explosion of data however, such data often contain irrelevant information that consequently makes prediction difficult. Researchers are devising new methods and algorithms in order to extract information to create robust predictive models. Such models mostly contain predictor variables that are directly or indirectly correlated with other predictor variables. In addition, studies often consist of many response variables correlated with each other. These interlinked relationships influence any study, whether it is predictive modelling or inference.

Modern inter-disciplinary research fields such as chemometrics, econometrics and bioinformatics handle multi-response models extensively. This paper attempts to compare some multivariate prediction methods based on their prediction performance on linear model data with specific properties. The properties include the correlation between response variables, the correlation between predictor variables, number of predictor variables and the position of relevant predictor components. These properties are discussed more in the Experimental Design section. Among others, Sæbø et al. [26] and Almøy [2] have conducted a similar comparison in the single response setting. In addition, Rimal et al. [25] have also conducted a basic comparison of some prediction methods and their interaction with the data properties of a multi-response model. The main aim of this paper is to present a comprehensive comparison of

contemporary prediction methods such as simultaneous envelope estimation (Senv) [8] and envelope estimation in predictor space (Xenv) [7] with customary prediction methods such as Principal Component Regression (PCR), Partial Least Squares Regression (PLS) using simulated dataset with controlled properties. In the case of PLS, we have used PLS1 which fits individual response separately and PLS2 which fits all the responses together. Experimental design and the methods under comparison are discussed further, followed by a brief discussion of the strategy behind the data simulation.

2. Simulation model

Consider a model where the response vector (\mathbf{y}) with m elements and predictor vector (\mathbf{x}) with p elements follow a multivariate normal distribution as follows,

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{x} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_y \\ \boldsymbol{\mu}_x \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{bmatrix} \right) \quad (1)$$

where, $\boldsymbol{\Sigma}_{xx}$ and $\boldsymbol{\Sigma}_{yy}$ are the variance-covariance matrices of \mathbf{x} and \mathbf{y} , respectively, $\boldsymbol{\Sigma}_{xy}$ is the covariance between \mathbf{x} and \mathbf{y} and $\boldsymbol{\mu}_x$ and $\boldsymbol{\mu}_y$ are mean vectors of \mathbf{x} and \mathbf{y} , respectively. A linear model based on (1) is,

$$\mathbf{y} = \boldsymbol{\beta}^t (\mathbf{x} - \boldsymbol{\mu}_x) + \boldsymbol{\varepsilon} \quad (2)$$

where, $\boldsymbol{\beta}_{m \times p}^t$ is a matrix of regression coefficients and $\boldsymbol{\varepsilon}$ is an error term

* Corresponding author.

E-mail addresses: raju.rimal@nmbu.no (R. Rimal), trygve.almoy@nmbu.no (T. Almøy), solve.sabo@nmbu.no (S. Sæbø).

such that $\epsilon \sim \mathcal{N}(0, \Sigma_{y|x})$. Here, $\beta^t = \Sigma_{yx}\Sigma_{xx}^{-1}$ and $\Sigma_{y|x} = \Sigma_{yy} - \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy}$.

In a model like (2), we assume that the variation in response y is partly explained by the predictor x . However, in many situations, only a subspace of the predictor space is relevant for the variation in the response y . This space can be referred to as the relevant space of x and the rest as irrelevant space. In a similar way, for a certain model, we can assume that a subspace in the response space exists and contains the information that the relevant space in predictor can explain (Fig. 1). Cook et al. [7] and Cook and Zhang [8] have referred to the relevant space as material space and the irrelevant space as immaterial space.

With an orthogonal transformation of y and x to latent variables w and z , respectively, by $w = Qy$ and $z = Rx$, where Q and R are orthogonal rotation matrices, an equivalent model to (1) in terms of the latent variables can be written as,

$$\begin{bmatrix} w \\ z \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_w \\ \mu_z \end{bmatrix}, \begin{bmatrix} \Sigma_{ww} & \Sigma_{wz} \\ \Sigma_{zw} & \Sigma_{zz} \end{bmatrix}\right) \quad (3)$$

where, Σ_{ww} and Σ_{zz} are the variance-covariance matrices of w and z , respectively. Σ_{zw} is the covariance between z and w . μ_w and μ_z are the mean vector of z and w respectively.

Here, the elements of w and z are the principal components of responses and predictors, which will respectively be referred to respectively as “response components” and “predictor components”. The column vectors of respective rotation matrices Q and R are the eigenvectors corresponding to these principal components. We can write a linear model based on (3) as,

$$w = \mu_w + \alpha^t(z - \mu_z) + \tau \quad (4)$$

where, α^t is a matrix of regression coefficients and τ is an error term such that $\tau \sim \mathcal{N}(0, \Sigma_{w|z})$.

Following the concept of relevant space, a subset of predictor components can be imagined to span the predictor space. These components can be regarded as relevant predictor components. Naes and Martens [22] introduced the concept of relevant components which was explored further by Helland [11], Næs and Helland [21], Helland and Almøy [13] and Helland [12]. The corresponding eigenvectors were referred to as relevant eigenvectors. A similar logic is introduced by Cook et al. [7] and later by Cook et al. [5] as an envelope which is the space spanned by the relevant eigenvectors [4, pp. 101].

In addition, various simulation studies have been performed with the model based on the concept of relevant subspace. A simulation study by Almøy [2] has used a single response simulation model based on reduced regression and has compared some contemporary multivariate estimators. In recent years Helland et al. [15], Sæbø et al. [26], Helland et al.

[14] and Rimal et al. [25] implemented similar simulation examples similar to those we are discussing in this study. This paper, however, presents an elaborate comparison of the prediction using multi-response simulated linear model data. The properties of the simulated data are varied through different levels of simulation-parameters based on an experimental design. Rimal et al. [25] provide a detailed discussion of the simulation model that we have adopted here. The following section presents the estimators being compared in more detail.

3. Prediction methods

Partial least squares regression (PLS) and Principal component regression (PCR) have been used in many disciplines such as chemometrics, econometrics, bioinformatics and machine learning, where wide predictor matrices, i.e. p (number of predictors) $>$ n (number of observation) are common. These methods are popular in multivariate analysis, especially for exploratory studies and predictions. In recent years, a concept of envelope introduced by Cook et al. [6] based on the reduction in the regression model was implemented for the development of different estimators. This study compares these prediction methods based on their prediction performance on data simulated with different controlled properties.

Principal Components Regression (PCR): Principal components are the linear combinations of predictor variables such that the transformation makes the new variables uncorrelated. In addition, the variation of the original dataset captured by the new variables is sorted in descending order. In other words, each successive component captures maximum variation left by the preceding components in predictor variables [18]. Principal components regression uses these principal components as a new predictor to explain the variation in the response.

Partial Least Squares (PLS): Two variants of PLS: PLS1 and PLS2 are used for comparison. The first one considers individual response variables separately, i.e. each response is predicted with a single response model, while the latter considers all response variables together. In PLS regression, the components are determined so as to maximize a covariance between response and predictors [10]. Among other, there are three main PLS algorithms NIPALS, SIMPLS and Kernel Algorithm all of which removes the extracted information through deflation and makes the resulting new variables orthogonal. The algorithms differ in the deflation strategy and computation of various weight vectors [1] and here we have used the kernel version of PLS. R-package pls [20] is used for both PCR and PLS methods.

Envelopes: The envelope, introduced by Cook et al. [6], was first used to define response envelope [7] as the smallest subspace in the response space and must be a reducing subspace of $\Sigma_{y|x}$ such that the span of regression coefficients lies in that space. Since a multivariate linear regression model contains relevant (material) and irrelevant (immaterial) variation in both response and predictor, the relevant part provides information, while the irrelevant part increases the estimative variation. The concept of the envelope uses the relevant part for estimation while excluding the irrelevant part consequently increasing the efficiency of the model [9].

The concept was later extended to the predictor space, where the predictor envelope was defined [5]. Further Cook and Zhang [8] used envelopes for joint reduction of the responses and predictors and argued that this produced efficiency gains that were greater than those derived by using individual envelopes for either the responses or the predictors separately. All the variants of envelope estimations are based on maximum likelihood estimation. Here we have used predictor envelope (Xenv) and simultaneous envelope (Senv) for the comparison. R-package Renvlp [19] is used for both Xenv and Senv methods.

3.1. Modification in envelope estimation

Since envelope estimators (Xenv and Senv) are based on maximum likelihood estimation (MLE), it fails to estimate in the case of wide

Relevant space within a model

A concept for reduction of regression models

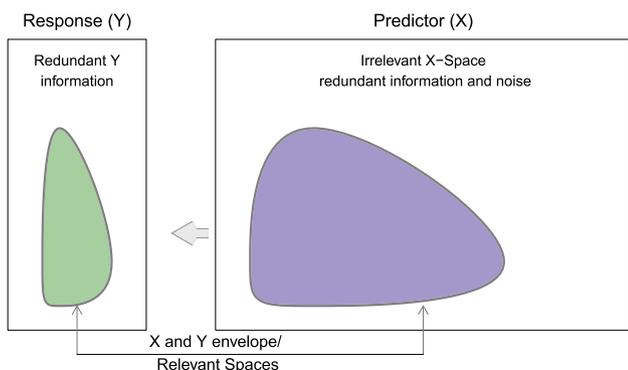


Fig. 1. Relevant space in a regression model.

matrices, i.e. $p > n$. To incorporate these methods in our comparison, we have used the principal components (\mathbf{z}) of the predictor variables (\mathbf{x}) as predictors, using the required number of components for capturing 97.5% of the variation in \mathbf{x} for the designs where $p > n$. The new set of variables \mathbf{z} were used for envelope estimation. The regression coefficients ($\hat{\alpha}$) corresponding to these new variables \mathbf{z} were transformed back to obtain coefficients for each predictor variable

$$\hat{\beta} = \mathbf{e}_k \hat{\alpha}_k$$

where \mathbf{e}_k is a matrix of eigenvectors with the first k number of components. Only simultaneous envelope allows to specify the dimension of response envelope and all the simulation is based on a single latent dimension of response, so it is fixed at two in the simulation study. In the case of Senv, when the envelope dimension for response is the same as the number of responses, it degenerates to the Xenv method and if the envelope dimension for the predictor is the same as the number of predictors, it degenerates to the standard multivariate linear regression [19].

4. Experimental design

This study compares prediction methods based on their prediction ability. Data with specific properties are simulated, some of which are easier to predict than others. These data are simulated using the R-package *simrel*, which is discussed in Sæbø et al. [26] and Rimal et al. [25]. Here we have used four different factors to vary the property of the data: a) Number of predictors (p), b) Multicollinearity in predictor variables (γ), c) Correlation in response variables (η) and d) position of predictor components relevant for the response (relpos). Using two levels of p , γ and relpos and four levels of η , 32 sets of distinct properties are designed for the simulation.

Number of predictors: To observe the performance of the methods on tall and wide predictor matrices, 20 and 250 predictor variables are simulated with the number of observations fixed at 100. Parameter p controls these properties in the *simrel* function.

Multicollinearity in predictor variables: Highly collinear predictors can be explained completely by a few components. The parameter γ in *simrel* controls the decline in the eigenvalues of the predictor variables as (5).

$$\lambda_i = e^{-\gamma(i-1)}, \gamma > 0 \text{ and } i = 1, 2, \dots, p \quad (5)$$

Here, $\lambda_i, i = 1, 2, \dots, p$ are eigenvalues of the predictor variables. We have used 0.2 and 0.9 as different levels of γ . The higher the value of γ , the higher the multicollinearity will be, and vice versa. In our simulations, the higher and lower γ values corresponded to the maximum correlation between the predictors equal to 0.990 and 0.709, respectively, in the case of $p = 20$ variables. In the case of $p = 250$ the

corresponding values for the maximum correlation were 0.998 to 0.923.

Correlation in response variables: Correlation among response variables has been explored to a lesser extent. Here we have tried to explore that part with four levels of correlation in the response variables. We have used the η parameter of *simrel* for controlling the decline in eigenvalues corresponding to the response variables as (6).

$$\kappa_j = e^{-\eta(j-1)}, \eta > 0 \text{ and } j = 1, 2, \dots, m \quad (6)$$

Here, $\kappa_j, j = 1, 2, \dots, m$ are the eigenvalues of the response variables and m is the number of response variables. We have used 0, 0.4, 0.8 and 1.2 as different levels of η . The larger the value of η , the larger will be the correlation will be between response variables and vice versa. In our simulation, the different levels of η from small to large correspond to the maximum correlation of 0, 0.442, 0.729 and 0.878 between the response variables respectively.

Position of predictor components relevant to the response: The principal components of the predictors are ordered. The first principal component captures most of the variation in the predictors. The second captures most of the remainder left by the first principal component and so on. In highly collinear predictors, the variation captured by the first few components is relatively high. However, if those components are not relevant for the response, prediction becomes difficult [13]. Here, two levels of the positions of these relevant components are used as 1, 2, 3, 4 and 5, 6, 7, 8.

Moreover, a complete factorial design from the levels of the above parameters gave us 32 designs. Each design is associated with a dataset having unique properties. Fig. 2, shows all the designs. For each design and prediction method, 50 datasets were simulated as replicates. In total, there were $5 \times 32 \times 50$, i.e. 8000 simulated datasets.

Common parameters: Each dataset was simulated with $n = 100$ number of observation and $m = 4$ response variables. Furthermore, the coefficient of determination corresponding to each response components in all the designs is set to 0.8. The informative and uninformative latent components are generated according to (3). Since Σ_{ww} and Σ_{zz} are diagonal matrices, the components are independent within w and z , but dependence between the latent spaces of x and y are secured through the non-zero elements of Σ_{wz} with positions defined by the relpos and ypos parameters. The latent components are subsequently rotated to obtain the population covariance structure of response and predictor variables. In addition, we have assumed that there is only one informative response component. Hence, the informative response component after the orthogonal rotation together with three uninformative response components generates four response variables. This spreads out the information in all simulated response variables. For further details on the simulation tool, see Ref. [25].

An example of simulation parameters for the first design is as follows:

```
simrel(
  n      = 100,          ## Training samples
  p      = 20,          ## Predictors
  m      = 4,           ## Responses
  q      = 20,          ## Relevant predictors
  relpos = list(c(1, 2, 3, 4)), ## Relevant predictor components index
  eta    = 0,           ## Decay factor of response eigenvalues
  gamma  = 0.2,         ## Decay factor of predictor eigenvalues
  R2     = 0.8,         ## Coefficient of determination
  ypos   = list(c(1, 2, 3, 4)),
  type   = "multivariate"
)
```

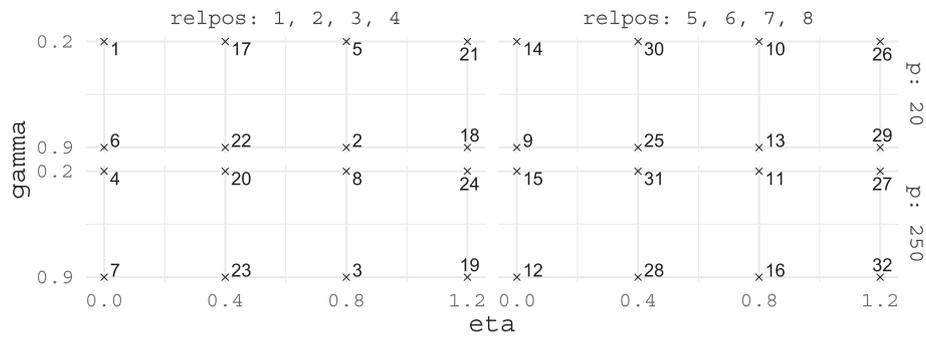


Fig. 2. Experimental Design of simulation parameters. Each point represents a unique data property.

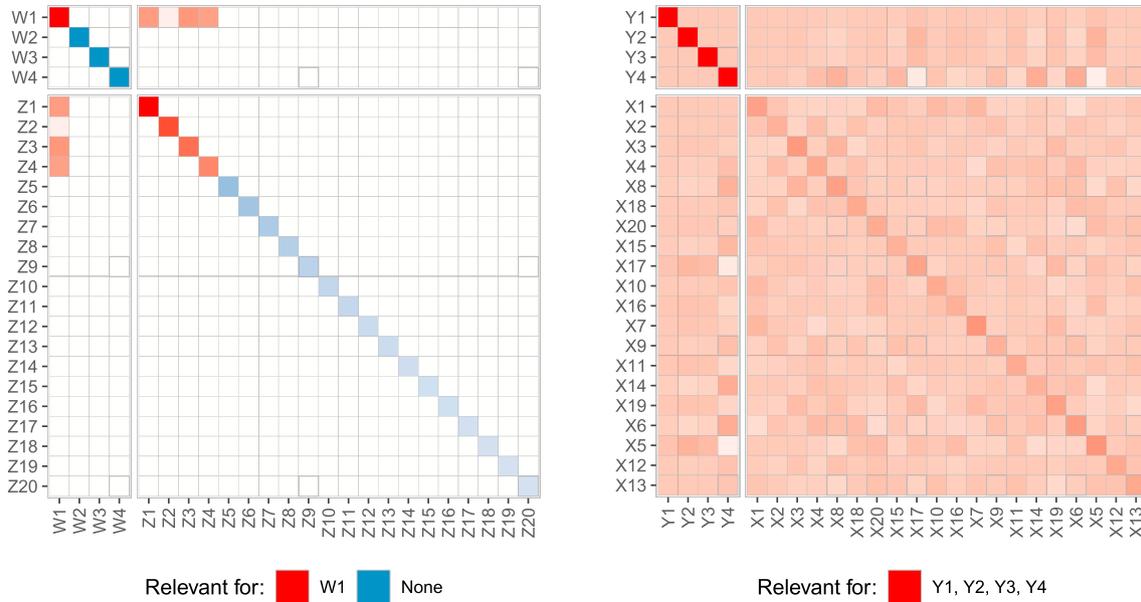


Fig. 3. (left) Covariance structure of latent components (right) Covariance structure of predictor and response.

The covariance structure of the data simulated with this design in Fig. 3 shows that the predictor components at positions 1, 2, 3 and 4 are relevant for the first response component. After the rotation with an orthogonal rotation matrix, all predictor variables are somewhat relevant for all response variables, satisfying other desired properties such as multicollinearity and coefficient of determination. For the same design, Fig. 4 (top left) shows that the predictor components 1, 2, 3 and 4 are relevant for the first response component. All other predictor components are irrelevant and all other response components are uninformative. However, due to the orthogonal rotation of the informative response component together with uninformative response components, all response variables in the population have similar covariance with the relevant predictor components (Fig. 4 (top right)). The sample covariances between the predictor components and predictor variables with response variables are shown in Fig. 4 (bottom left) and (bottom right) respectively.

A similar description can be made for all 32 designs, where each of the designs holds the properties of the data they simulate. These data are used by the prediction methods discussed in the previous section. Each prediction method is given independently simulated datasets in order to give them an equal opportunity to capture the dynamics in the data.

5. Basis of comparison

This study focuses mainly on the prediction performance of the

methods with an emphasis specifically on the interaction between the properties of the data controlled by the simulation parameters and the prediction methods. The prediction performance is measured based on the following:

- The average prediction error that a method can give using an arbitrary number of components and
- The average number of components used by the method to give the minimum prediction error

Let us define,

$$\mathcal{P} \mathcal{E}_{ijkl} = \frac{1}{\sigma_{y_j|x}^2} E \left[(\beta_j - \hat{\beta}_{ijkl})^t (\Sigma_{xx})_i (\beta_j - \hat{\beta}_{ijkl}) \right] + 1 \quad (7)$$

as a prediction error of response $j = 1, \dots, 4$ for a given design $i = 1, 2, \dots, 32$ and method $k = 1(\text{PCR}), \dots, 5(\text{Senv})$ using $l = 0, \dots, 10$ number of components. Here, $(\Sigma_{xx})_i$ is the true covariance matrix of the predictors, unique for a particular design i and $\sigma_{y_j|x}^2$ for response $j = 1, \dots, m$ is the true model error. Here prediction error is scaled by the true model error to remove the effects of influencing residual variances. Since both the expectation and the variance of $\hat{\beta}$ are unknown, the prediction error is estimated using data from 50 replications as follows,

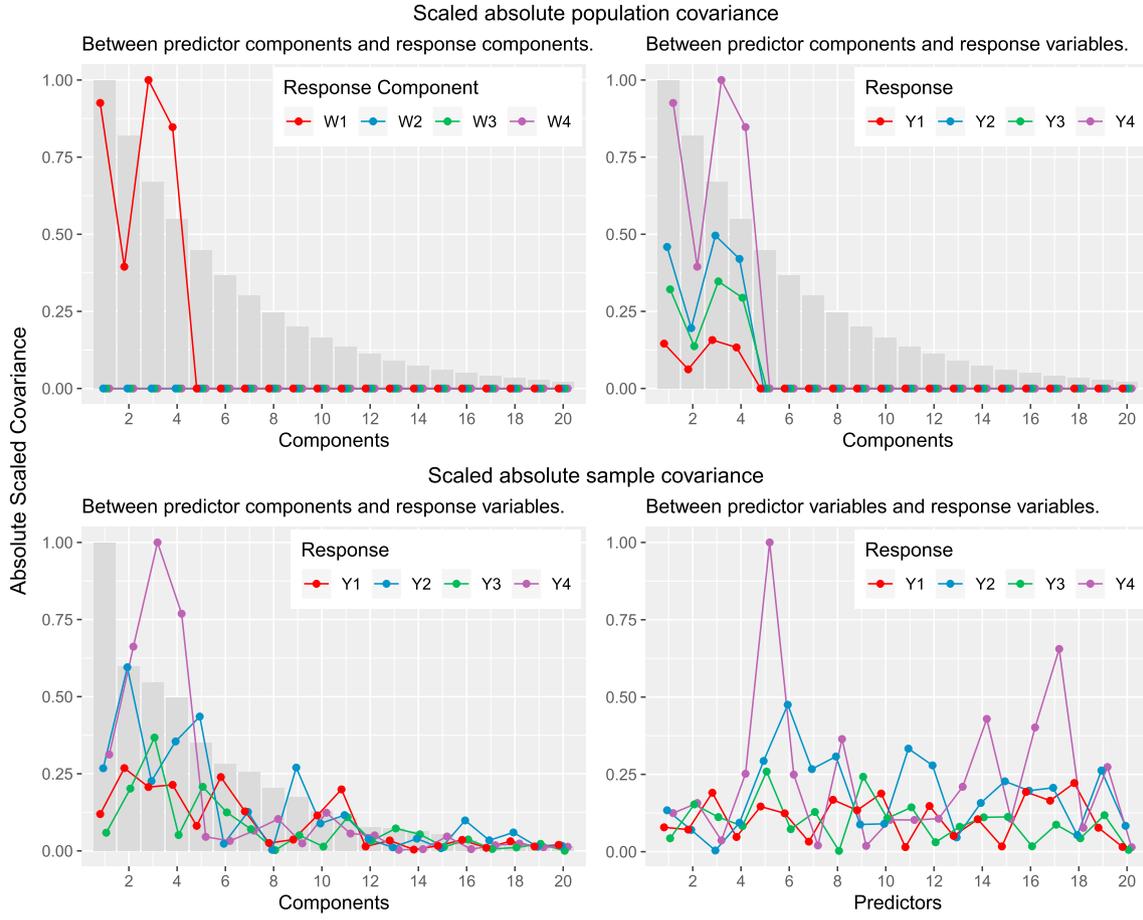


Fig. 4. Expected Scaled absolute covariance between predictor components and response components (top left). Expected Scaled absolute covariance between predictor components and response variables (top right). Sample scaled absolute covariance between predictor components and response variables (bottom left). Sample scaled absolute covariance between predictor variables and response variables (bottom right). The bar graph in the background represents eigenvalues corresponding to each component in the population (top plots) and in the sample (bottom plots). One can compare the top-right plot (true covariance of the population) with bottom-left (covariance in the simulated data) which shows a similar pattern for different components.

$$\widehat{\mathcal{P}\mathcal{E}}_{ijkl} = \frac{1}{\sigma_{y_j|x}^2} \sum_{r=0}^{50} \left[(\beta_{ij} - \widehat{\beta}_{ijklr})' (\Sigma_{xx})_i (\beta_{ij} - \widehat{\beta}_{ijklr}) \right] + 1 \quad (8)$$

where $\widehat{\mathcal{P}\mathcal{E}}_{ijkl}$ is the estimated prediction error averaged over $r = 50$ replicates.

The following section focuses on the data for the estimation of these prediction errors that are used for the two models discussed above in a) and b) of this section.

6. Data preparation

A dataset for estimating (7) is obtained from simulation which contains a) five factors corresponding to simulation parameters, b) prediction methods, c) number of components, d) replications and e) prediction error for four responses. The prediction error is computed using predictor components ranging from 0 to 10 for every 50 replicates as,

$$(\widehat{\mathcal{P}\mathcal{E}}_*)_{ijklr} = \frac{1}{\sigma_{y_j|x}^2} \left[(\beta_{ij} - \widehat{\beta}_{ijklr})' (\Sigma_{xx})_i (\beta_{ij} - \widehat{\beta}_{ijklr}) \right] + 1$$

Thus there are 32 (designs) \times 5 (methods) \times 11 (number of components) \times 50 (replications), i.e. 88000 observations corresponding to the response variables from Y1 to Y4.

Since our discussions focus on the average minimum prediction error that a method can obtain and the average number of components they use to get the minimum prediction error in each replicates, the dataset

discussed above is summarized as constructing the following two smaller datasets. Let us call them *Error Dataset* and *Component Dataset*.

Error Dataset: For each prediction method, design and response, an average prediction error is computed over all replicates for each component. Next, a component that gives the minimum of this average prediction error is selected, i.e.,

$$l_i = \operatorname{argmin}_i \left[\frac{1}{50} \sum_{r=1}^{50} (\mathcal{P}\mathcal{E}_*)_{ijklr} \right] \quad (9)$$

Using the component l_i , a dataset of $(\mathcal{P}\mathcal{E}_*)_{ijklr}$ is used as the *Error Dataset*. Let $\mathbf{u}_{(8000 \times 4)} = (u_j)$ for $j = 1, \dots, 4$ be the outcome variables measuring the prediction error corresponding to the response number j in the context of this dataset.

Component Dataset: The number of components that gives the minimum prediction error in each replication is referred to as the *Component Dataset*, i.e.,

$$l_i = \operatorname{argmin}_i [\mathcal{P}\mathcal{E}_{ijklr}] \quad (10)$$

Here l_i is the number of components that gives minimum prediction error $(\mathcal{P}\mathcal{E}_*)_{ijklr}$ for design i , response j , method k and replicate r . Let $\mathbf{v}_{(8000 \times 4)} = (v_j)$ for $j = 1, \dots, 4$ be the outcome variables measuring the number of components used for minimum prediction error corresponding to the response j in the context of this dataset.

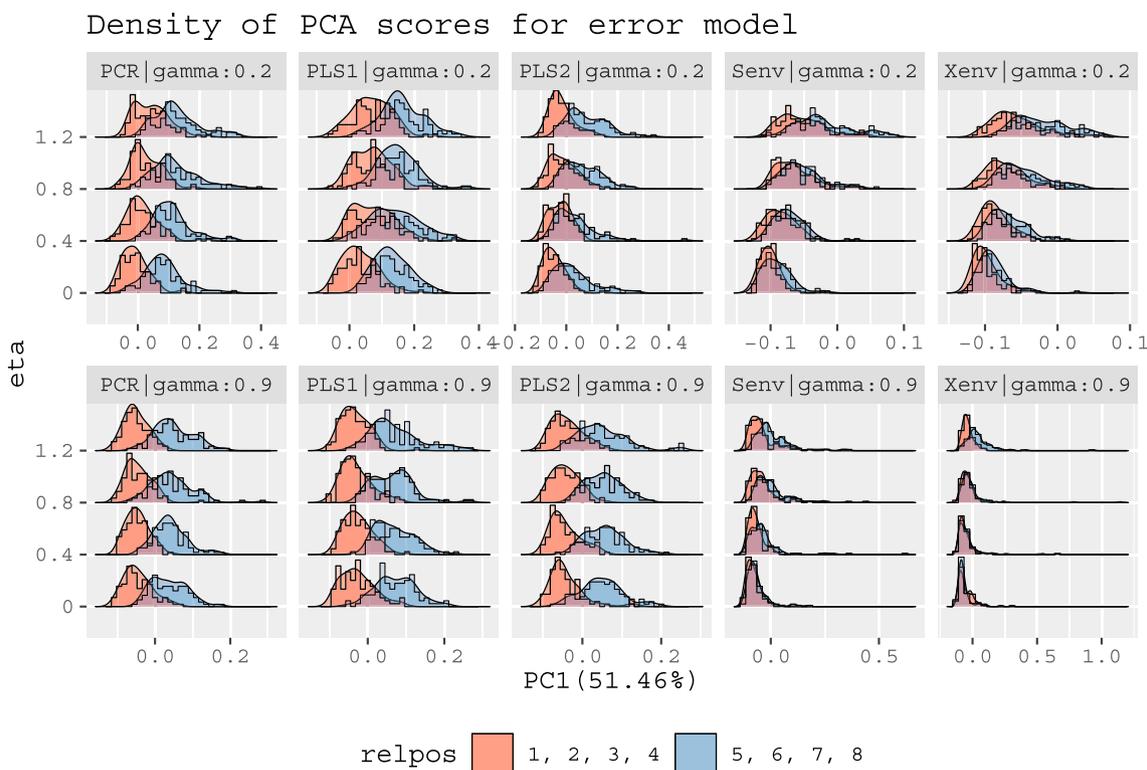


Fig. 5. Scores density corresponding to first principal component of *error dataset* (*u*) subdivided by methods, gamma and eta and grouped by relpos.

7. Exploration

This section explores the variation in the *error dataset* and the *component dataset* for which we have used Principal Component Analysis (PCA). Let t_u and t_v be the principal component score sets corresponding to PCA run on the *u* and *v* matrices respectively. The scores density in Fig. 5 corresponds to the first principal component of *u*, i.e. the first column of t_u .

Since higher prediction errors correspond to high scores, the plot shows that the PCR, PLS1 and PLS2 methods are influenced by the two levels of the position of relevant predictor components. When the relevant predictors are at positions 5, 6, 7, 8, the eigenvalues corresponding to them are relatively smaller. This also suggests that PCR, PLS1 and PLS2 depend greatly on the position of the relevant components, and the variation of these components affects their prediction performance. However, the envelope methods appeared to be less influenced by relpos in this regard.

In addition, the plot also shows that the effect of gamma, i.e., the level of multicollinearity, has a lesser effect when the relevant predictors are at positions 1, 2, 3, 4. This indicates that the methods are somewhat robust for handling collinear predictors. Nevertheless, when the relevant predictors are at positions 5, 6, 7, 8, high multicollinearity results in a small variance of these relevant components and consequently yields poor prediction. This is in accordance with the findings of Helland and Almøy [13].

Furthermore, the density curves for PCR, PLS1 and PLS2 are similar for different levels of eta, i.e., the factor controlling the correlation between responses. However, the envelope models have been shown to have distinct interactions between the positions of relevant components (relpos) and eta. Here higher levels of eta have yielded higher scores and clear separation between two levels of relpos. In the case of high multicollinearity, envelope methods have resulted in some large outliers indicating that in some cases that the methods can result in giving an unexpected prediction.

In Fig. 6, the higher scores suggest that methods have used a larger

number of components to give minimum prediction error. The plot also shows that the relevant predictor components at 5, 6, 7, 8 give larger prediction errors than those in positions 1, 2, 3, 4. The pattern is more distinct in large multicollinearity cases and PCR and PLS methods. Both the envelope methods have shown equally enhanced performance at both levels of relpos and gamma. However, for data with low multicollinearity ($\gamma = 0.2$), the envelope methods have used a lesser number of components on average than in the high multicollinearity cases to achieve minimum prediction error.

8. Statistical analysis

This section has modelled the *error data* and the *component data* as a function of the simulation parameters to better understand the connection between data properties and prediction methods using multivariate analysis of variation (MANOVA).

Let us consider a model with third order interaction of the simulation parameters (p, gamma, eta and relpos) and Methods as in (11) and (12) using datasets *u* and *v*, respectively. Let us refer to them as the *error model* and the *component model*.

Error Model:

$$\mathbf{u}_{abcdef} = \boldsymbol{\mu}_u + (p_a + \text{gamma}_b + \text{eta}_c + \text{relpos}_d + \text{Methods}_e)^3 + (\boldsymbol{\epsilon}_u)_{abcdef} \tag{11}$$

Component Model:

$$\mathbf{v}_{abcdef} = \boldsymbol{\mu}_v + (p_a + \text{gamma}_b + \text{eta}_c + \text{relpos}_d + \text{Methods}_e)^3 + (\boldsymbol{\epsilon}_v)_{abcdef} \tag{12}$$

where, \mathbf{u}_{abcdef} is a vector of prediction errors in the *error model* and \mathbf{v}_{abcdef} is a vector of the number of components used by a method to obtain minimum prediction error in the *component model*.

Although there are several test-statistics for MANOVA, all are essentially equivalent for large samples [17]. Here we will use Pillai's trace statistic which is defined as,

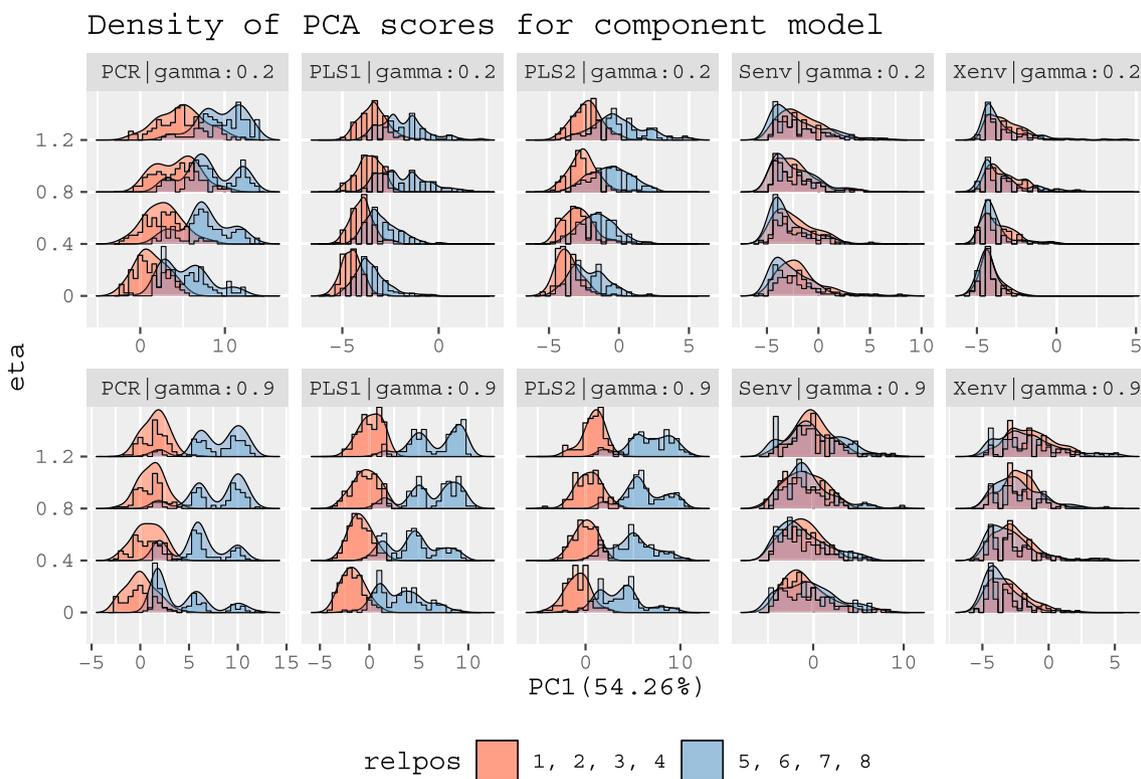


Fig. 6. Score density corresponding to the first principal component of the *component dataset* (v) subdivided by methods, gamma and eta and grouped by relpos.

$$\text{Pillai statistic} = \text{tr}[(\mathbf{E} + \mathbf{H})^{-1}\mathbf{H}] = \sum_{i=1}^m \frac{\nu_i}{1 + \nu_i} \quad (13)$$

Here the matrix \mathbf{H} holds between-sum-of-squares and sum-of-products for each of the predictors. The matrix \mathbf{E} has a within the sum of squares and sum of products for each of the predictors. ν_i represents the eigenvalues corresponding to $\mathbf{E}^{-1}\mathbf{H}$ [24].

For both the models (11) and (12), Pillai's trace statistic is used for accessing the effect of each factor and returns an F-value for the strength of their significance. Fig. 7 plots the Pillai's trace statistics as bars with corresponding F-values as text labels for both models.

Error Model: Fig. 7 (left) shows the Pillai's trace statistic for factors of the *error model*. The main effect of Method followed by relpos, eta and gamma have the largest influence on the model. A highly significant two-factor interaction of Method with gamma followed by the relpos and eta clearly shows that methods perform differently for different levels of these data properties. The significant third order interaction between Method, eta and gamma suggest that the performance of a method differs for a given level of multicollinearity and the correlation between the responses. Since only some methods consider modelling predictor and response together, the prediction is affected by the level of correlation between the responses (eta) for a given method.

Component Model: Fig. 7 (right) shows the Pillai's trace statistic for factors of the *component model*. As in the *error model*, the main effects of the Method, relpos, gamma and eta have a significantly large effect on the number of components that a method has used to obtain minimum prediction error. The two-factor interactions of Method with simulation parameters are larger in this case. This shows that the Methods and these interactions have a larger effect on the use of the number of component than the prediction error itself. In addition, a similar significant high third-order interaction as found in the *error model* is also observed in this model.

The following section will continue to explore the effects of different levels of the factors in the case of these interactions.

8.1. Effect analysis of error model

The large difference in the prediction error for the envelope models in Fig. 8 (left) is intensified when the position of the relevant predictor is at 5, 6, 7, 8. The results also show that the envelope methods are more sensitive to the levels of eta than the rest of the methods. In the case of PCR and PLS, the difference in the effect of levels of eta is small.

In Fig. 8 (right), we can see that the multicollinearity (controlled by gamma) has affected all the methods. However, envelope methods have better performance on low multicollinearity, as opposed to high multicollinearity, and PCR, PLS1 and PLS2 are robust for high multicollinearity. Despite handling high multicollinearity, these methods have higher prediction error in both cases of multicollinearity than the envelope methods.

8.2. Effect analysis of the component model

Unlike for prediction errors, Fig. 9 (left) shows that the number of components used by the methods to obtain minimum prediction error is less affected by the levels of eta. All methods appear to use on average more components when eta increases. Envelope methods are able to obtain minimum prediction error by using components ranging from 1 to 3 in both the cases of relpos. This value is much higher in the case of PCR as its prediction is based only on the principal components of the predictor matrix. The number of components used by this method ranges from 3 to 5 when relevant components are at positions 1, 2, 3, 4 and 5 to 8 when relevant components are at positions 5, 6, 7, 8.

When relevant components are at position 5, 6, 7, 8, the eigenvalues of relevant predictors become smaller and responses are relatively difficult to predict. This becomes more critical for high multicollinearity cases. Fig. 9 (right) shows that the envelope methods are less influenced by the level of relpos and are particularly better in achieving minimum prediction error using a fewer number of components than other methods.

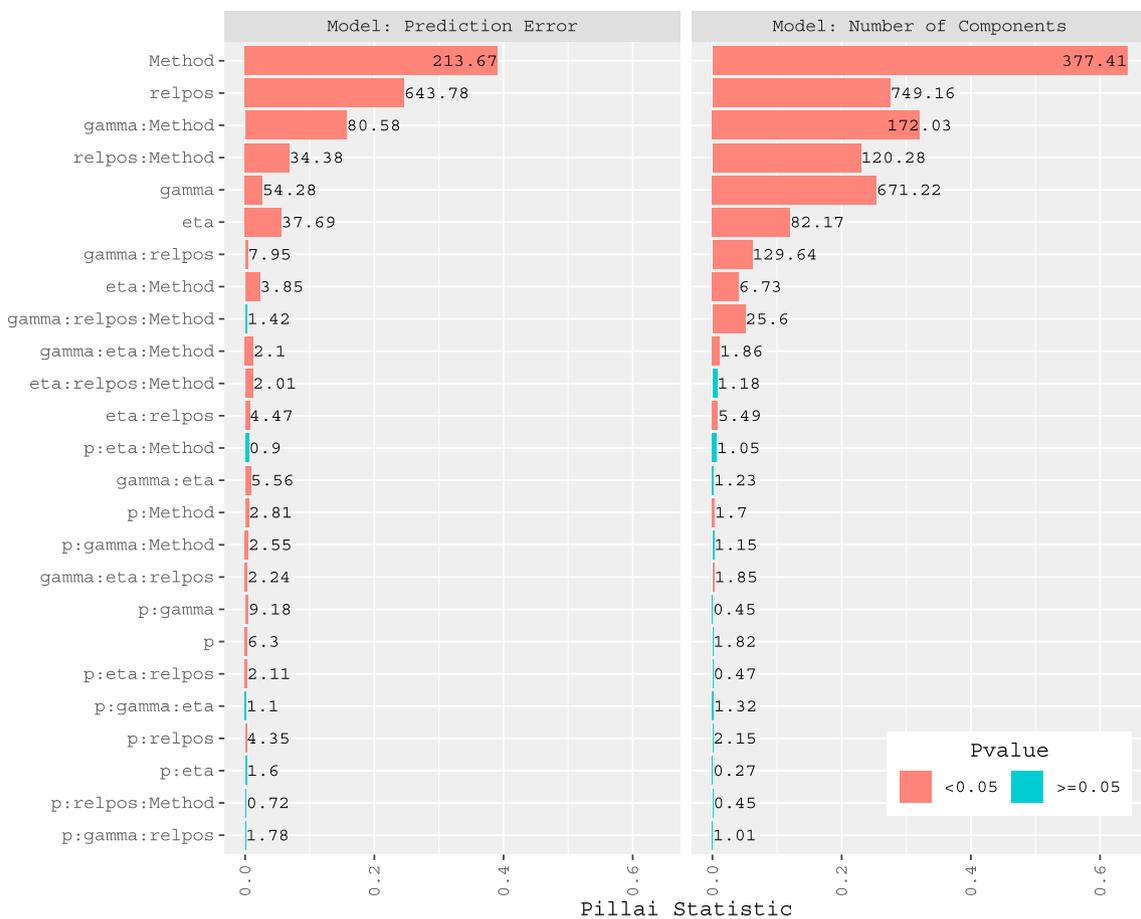


Fig. 7. Pillai Statistic and F-value for the MANOVA model. The bar represents the Pillai Statistic and the text labels are F-value for the corresponding factor.

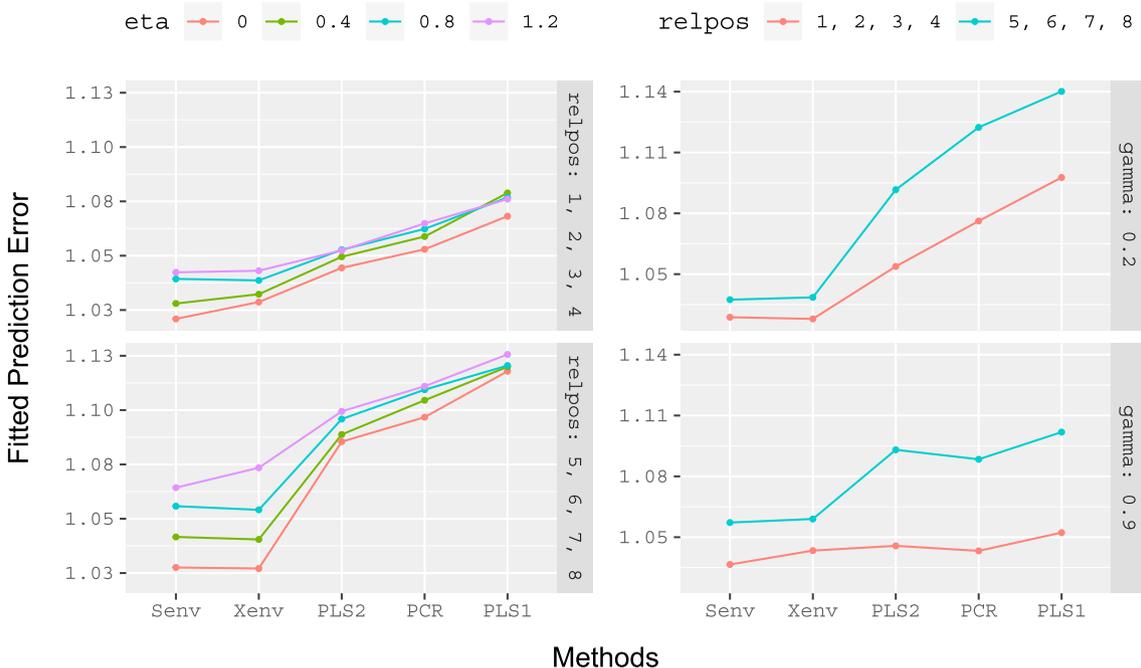


Fig. 8. Effect plot of some interactions of the multivariate linear model of prediction error.

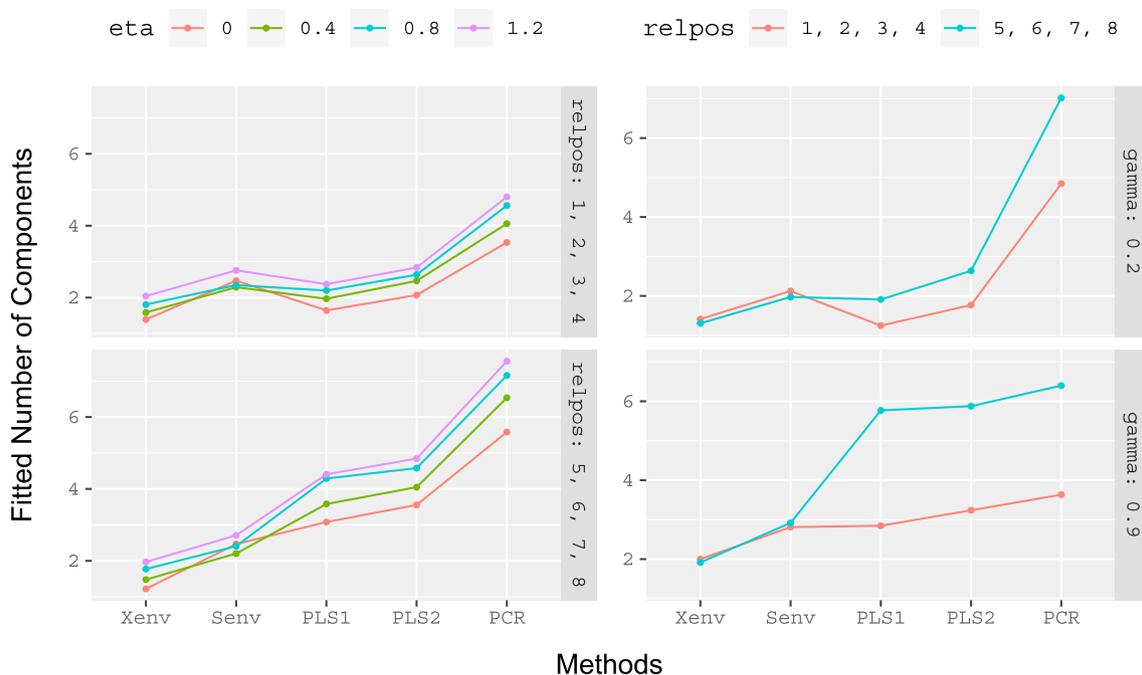


Fig. 9. Effect plot of some interactions of the multivariate linear model of the number of components to get minimum prediction error.

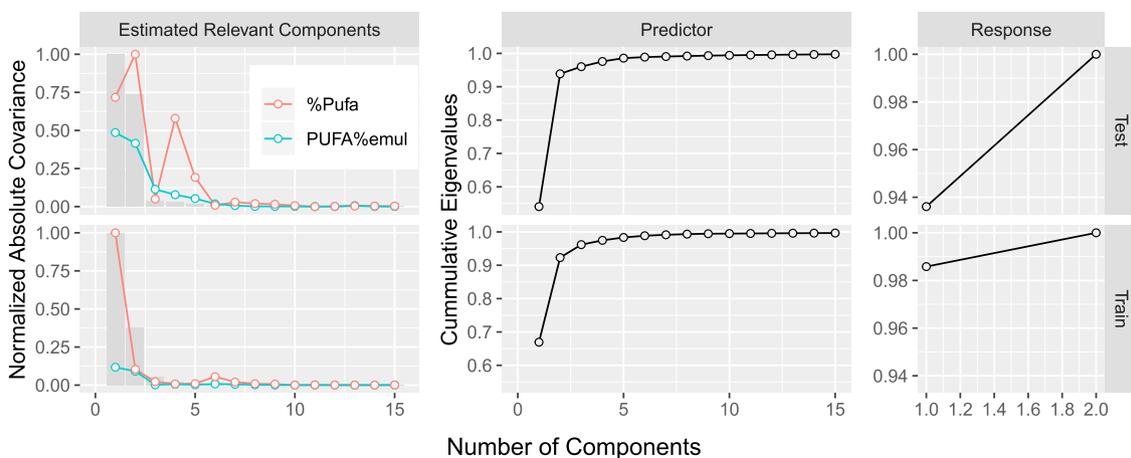


Fig. 10. (Left) Bar represents the eigenvalues corresponding to Raman Spectra. The points and line are the covariances between response and the principal components of Raman Spectra. All the values are normalized to scale from 0 to 1. (Middle) Cumulative sum of eigenvalues corresponding to predictors. (Right) The cumulative sum of eigenvalues corresponding to responses. The top and bottom row corresponds to test and training datasets respectively.

9. Examples

In addition to the analysis with the simulated data, the following two examples explore the prediction performance of the methods using real datasets. Since both examples have wide predictor matrices, principal components explaining 97.5% of the variation in them are used for envelope methods. The coefficients were transformed back after the estimation.

9.1. Raman spectra analysis of contents of polyunsaturated fatty acids (PUFA)

This dataset contains 44 training samples and 25 test samples of fatty acid information expressed as a) percentage of total sample weight and b) the percentage of total fat content. The dataset is borrowed from Næs et al. [23] where more information can be found. The samples were analysed using Raman spectroscopy from which 1096 wavelength

variables were obtained as predictors. Raman spectroscopy provides detailed chemical information from minor components in food. The aim of this example is to compare how well the prediction methods that we have considered are able to predict the contents of PUFA using these Raman spectra.

Fig. 10 (left) shows that the first few predictor components are somewhat correlated with response variables. In addition, the most variation in predictors is explained by less than five components (middle). Further, the response variables are highly correlated, suggesting that a single latent dimension explains most of the variation (right). We may therefore also believe that the relevant latent space in the response matrix is of dimension one. This resembles Design 19 (Fig. 2) from our simulation.

Using a range of components from 1 to 15, regression models were fitted using each of the methods. The fitted models were used to predict the test observation, and the root mean squared error of prediction (RMSEP) was calculated. Fig. 11 shows that PLS2 obtained a minimum

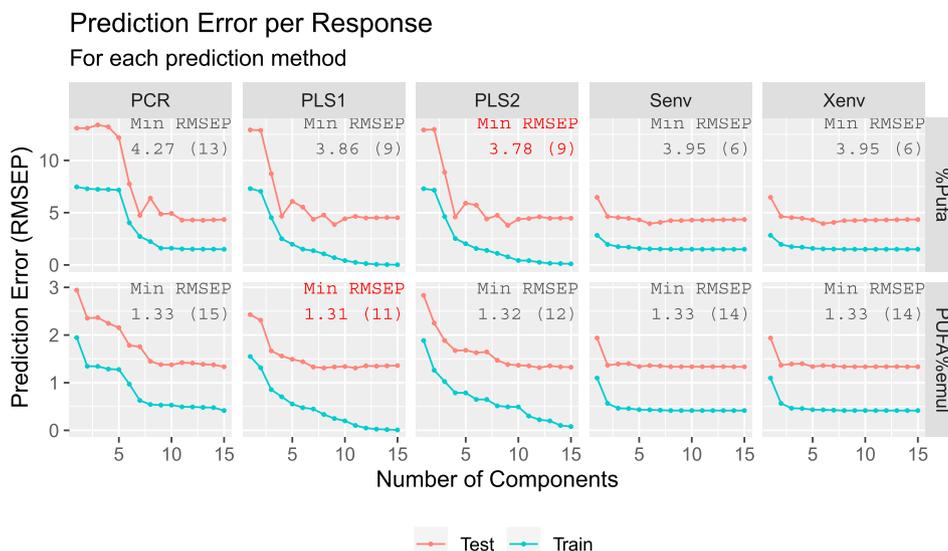


Fig. 11. Prediction Error of different prediction methods using different number of components.

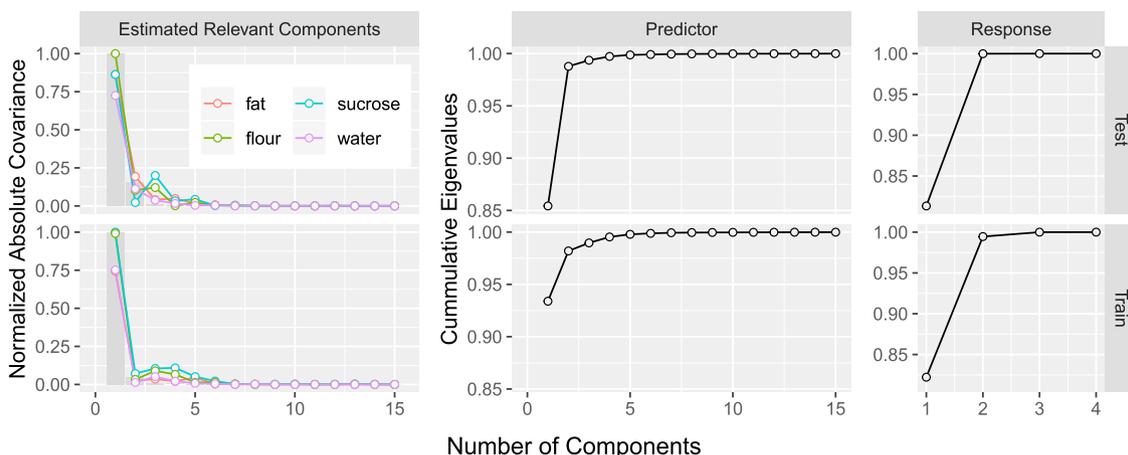


Fig. 12. (Left) Bar represents the eigenvalues corresponding to NIR Spectra. The points and line are the covariances between response and the principal components of NIR Spectra. All the values are normalized to scale from 0 to 1. (Middle) Cumulative sum of eigenvalues corresponding to predictors. (Right) The cumulative sum of eigenvalues corresponding to responses.

prediction error of 3.783 using 9 components in the case of response % Pufa, while PLS1 obtained a minimum prediction error of 1.308 using 11 components in the case of response PUFA%emul. However, the figure also shows that both envelope methods have reached to almost minimum prediction error in fewer number of components. This pattern is also visible in the simulation results (Fig. 9).

9.2. Example-2: NIR spectra of biscuit dough

The dataset consists of 700 wavelengths of NIR spectra (1100–2498 nm in steps of 2 nm) that were used as predictor variables. There are four response variables corresponding to the yield percentages of (a) fat, (b) sucrose, (c) flour and (d) water. The measurements were taken from 40 training observation of biscuit dough. A separate set of 32 samples created and measured on different occasions were used as test observations. The dataset is borrowed from Indahl [16] where further information can be obtained.

Fig. 12 (left) shows that the first predictor component has the largest variance and also has large covariance with all response variables. The second component, however, has larger variance (middle) than the succeeding components but has a small covariance with all the responses,

which indicates that the component is less relevant for any of the responses. In addition, two response components have explained most of the variation in response variables (right). This structure is also somewhat similar to Design 19, although it is uncertain whether the dimension of the relevant space in the response matrix is larger than one.

Fig. 13 (corresponding to Fig. 11) shows the root mean squared error for both test and train prediction of the biscuit dough data. Here four different methods have minimum test prediction error for the four responses. As the structure of the data is similar to that of the first example, the pattern in the prediction is also similar for all methods.

The prediction performance on the test data of the envelope methods appears to be more stable compared to the PCR and PLS methods. Furthermore, the envelope methods achieve good performance generally using fewer components, which is in accordance with Fig. 6.

10. Discussions and conclusion

Analysis using both simulated data and real data has shown that the envelope methods are more stable, less influenced by relpos and gamma and in general, performed better than PCR and PLS methods. These methods are also found to be less dependent on the number of



Fig. 13. Prediction Error of different prediction methods using different number of components.

components.

Since the facet in Figs. 5 and 6 have their own scales, despite having some large prediction errors seen at the right tail, envelope methods still have a smaller prediction error and have used a fewer number of components than the other methods.

The envelope methods may have this problem of being caught in a local optimum of the objective function. If these cases of sub-optimal convergence were identified and rerun to obtain better convergence, the envelope results may have become even better. Particularly in the case of the simultaneous envelope, since users can specify the number of dimension for the response envelope, the method can leverage the relevant space of response while PCR, PLS and Xenv are constrained to play only on predictor space.

Furthermore, we have fixed the coefficient of determination (R^2) as a constant throughout all the designs. Initial simulations (not shown) indicated that low R^2 affects all methods in a similar manner and that the MANOVA is highly dominated by R^2 . Keeping the value of R^2 fixed has allowed us to analyze other factors properly.

Two clear comments can be made about the effect of correlation of response on the prediction methods. The highly correlated response has shown the highest prediction error in general and the effect is most distinct in envelope methods. Since the envelope methods identify the relevant space as the span of relevant eigenvectors, the methods are able to obtain the minimum average prediction error by using a lesser number of components for all levels of eta.

To our knowledge, the effect of correlation in the response on PCR and PLS methods has been explored only to a limited extent. In this regards, it is interesting to see that these methods have applied a large number of components and returned a larger prediction error than envelope methods in the case of highly correlated responses. To fully understand the effect of eta, it is necessary to study the estimation performance of these methods with different numbers of components.

In addition, since using principal components or actual variables as

predictors in envelope methods has shown similar results, we have used principal components that have explained 97.5% of the variation, as mentioned previously, in the cases of envelope methods for the designs where $p > n$. Using 97.5% is slightly arbitrary here, but for the chosen simulation designs this proportion captured a fair amount of variations in predictor variables and also reduce the dimension significantly while enabling us to use envelope methods in all settings. The analyst should choose this number to balance the explained amount of variation to the number of components which is practical for model fitting using the envelope model. The methodology used to adapt envelopes to settings in which $p > n$ is, in fact, the same as that used by PLS: reduce by principal components, run the method, and then back transform to the original scale. The minor relative impact of p shown in Fig. 7 suggests that this adaptation method is useful.

The results from this study will help researchers to understand these methods for their performance in various linear model data and encourage them to use newly developed methods such as the envelopes. Since this study has focused entirely on prediction performance, further analysis of the estimative properties of these methods is required. A study of estimation error and the performance of methods on the non-optimal number of components can give a deeper understanding of these methods.

A shiny application [3] is available at <http://therimalaya.shinyapps.io/Comparison> where all the results related to this study can be visualized. In addition, a GitHub repository at <https://github.com/therimalaya/03-prediction-comparison> can be used to reproduce this study.

Acknowledgment

We are grateful to Inge Helland for his inputs on this paper throughout the period. His guidance on the envelope models and his review of the paper helped us greatly. Our gratitude also goes to thank Kristian Lillan, Ulf Indahl, Tormod Næs, Ingrid Måge and the team for

providing the data for analysis. We are also thankful to the reviewers for their comments which helped us to improve this paper.

References

- [1] A. Alin, Comparison of pls algorithms when number of objects is much larger than number of variables, *Stat. Pap.* 50 (4) (2009) 711–720. <https://doi.org/10.1007/s00362-009-0251-7>.
- [2] T. Almøy, A simulation study on comparison of prediction methods when only a few components are relevant, *Comput. Stat. Data Anal.* 21 (1) (jan 1996) 87–107.
- [3] W. Chang, J. Cheng, J. Allaire, Y. Xie, J. McPherson, Shiny: Web Application Framework for R. R Package Version 1.2.0, 2018. <https://CRAN.R-project.org/package=shiny>.
- [4] R.D. Cook, *An Introduction to Envelopes : Dimension Reduction for Efficient Estimation in Multivariate Statistics*, first ed., John Wiley & Sons, Hoboken, NJ, 2018, 2018.
- [5] R.D. Cook, I.S. Helland, Z. Su, Envelopes and partial least squares regression, *J. R. Stat. Ser. B Stat. Methodol.* 75 (5) (2013) 851–877.
- [6] R.D. Cook, B. Li, F. Chiaromonte, Dimension reduction in regression without matrix inversion, *Biometrika* 94 (3) (aug 2007) 569–584.
- [7] R.D. Cook, B. Li, F. Chiaromonte, Envelope models for parsimonious and efficient multivariate linear regression, *Stat. Sin.* 20 (3) (2010) 927–1010.
- [8] R.D. Cook, X. Zhang, Simultaneous envelopes for multivariate linear regression, *Technometrics* 57 (1) (2015) 11–25.
- [9] R.D. Cook, X. Zhang, Algorithms for envelope estimation, *J. Comput. Graph. Stat.* 25 (1) (2016) 284–300.
- [10] S. de Jong, SIMPLS: an alternative approach to partial least squares regression, *Chemometr. Intell. Lab. Syst.* 18 (3) (mar 1993) 251–263.
- [11] I.S. Helland, Partial least squares regression and statistical models, *Scand. J. Stat.* 17 (2) (1990) 97–114.
- [12] I.S. Helland, Model reduction for prediction in regression models, *Scand. J. Stat.* 27 (1) (mar 2000) 1–20.
- [13] I.S. Helland, T. Almøy, Comparison of prediction methods when only a few components are relevant, *J. Am. Stat. Assoc.* 89 (426) (1994) 583–591.
- [14] I.S. Helland, S. Sæbø, T. Almøy, R. Rimal, S. Sæbø, T. Almøy, R. Rimal, Model and estimators for partial least squares regression, *J. Chemom.* 32 (9) (sep 2018), e3044.
- [15] I.S. Helland, S. Sæbø, H.K. Tjelmeland, Near optimal prediction from relevant components, *Scand. J. Stat.* 39 (4) (mar 2012) 695–713.
- [16] U. Indahl, A twist to partial least squares regression, *J. Chemom.* 19 (1) (2005) 32–44.
- [17] R. Johnson, D. Wichern, *Applied Multivariate Statistical Analysis (Classic Version)*. Pearson Modern Classics for Advanced Statistics Series. Pearson Education Canada, 2018. <https://books.google.no/books?id=QBqlswEACAAJ>.
- [18] I.T. Jolliffe, *Principal Component Analysis*, second ed., 2002.
- [19] M. Lee, Z. Su, Renvlp: Computing Envelope Estimators. R Package Version 2.5, 2018. <https://CRAN.R-project.org/package=Renvlp>.
- [20] B.-H. Mevik, R. Wehrens, K.H. Liland, Pls: Partial Least Squares and Principal Component Regression. R Package Version 2.7-0, 2018. <https://CRAN.R-project.org/package=pls>.
- [21] T. Næs, I.S. Helland, Relevant components in regression, *Scand. J. Stat.* 20 (3) (1993) 239–250.
- [22] T. Næs, H. Martens, Comparison of prediction methods for multicollinear data, *Commun. Stat. Simulat. Comput.* 14 (3) (jan 1985) 545–576.
- [23] T. Næs, O. Tomic, N.K. Afseth, V. Segtnan, I. Måge, Multi-block regression based on combinations of orthogonalisation, pls-regression and canonical correlation analysis, *Chemometr. Intell. Lab. Syst.* 124 (2013) 32–42.
- [24] A.C. Rencher, *Methods of Multivariate Analysis*, vol. 492, John Wiley & Sons, 2003.
- [25] R. Rimal, T. Almøy, S. Sæbø, A tool for simulating multi-response linear model data, *Chemometr. Intell. Lab. Syst.* 176 (may 2018) 1–10.
- [26] S. Sæbø, T. Almøy, I.S. Helland, Simrel - a versatile tool for linear model data simulation based on the concept of a relevant subspace and relevant predictors, *Chemometr. Intell. Lab. Syst.* 146 (2015) 128–135.