

# Comparison of Multivariate Prediction Methods

Raju Rimal<sup>a,\*</sup>, Trygve Almøy<sup>a</sup>, Solve Sæbø<sup>b</sup>

<sup>a</sup>Faculty of Chemistry and Bioinformatics, Norwegian University of Life Sciences, Ås, Norway

<sup>b</sup>Prorector, Norwegian University of Life Sciences, Ås, Norway

---

## 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 model with multiple response model is limited. This study compares some newly-developed (envelope) and well-established (PLS, PCR) prediction methods based on simulated data specifically designed by varying properties such as multicollinearity, correlation between multiple responses and amount of information content in predictor variables. This study aims to give some insight on these methods and help researcher to understand and use them for further study.

*Keywords:* model-comparison, multi-response, simrel

---

## 1. Introduction

Prediction has been an essential components of modern data science, weather it is statistical analysis or machine learning. Modern technology has facilitated a massive explosion of data, however, such data often contain irrelevant information consequently making prediction difficult. Researchers are devising new methods and algorithms in order to extract information to create robust predictive models. Mostly such models contain

---

\*Corresponding Author

Email addresses: `raju.rimal@nmbu.no` (Raju Rimal), `trygve.almoy@nmbu.no` (Trygve Almøy), `solve.sabo@nmbu.no` (Solve Sæbø)

predictor variables that are directly or indirectly correlated with other predictor variables. In addition studies often constitute of many response variables correlated with each other. These interlinked relationships influence any study, whether it is predictive modeling or inference.

Modern inter-disciplinary research fields such as chemometrics, econometrics and bioinformatics are handling 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 includes correlation between response variables, 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. \(2015\)](#) and [Almøy \(1996\)](#) have made a similar comparison in the single response setting. In addition, [Rimal et al. \(2018\)](#) has also made a basic comparison on 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) ([Cook and Zhang, 2015](#)) and envelope estimation in predictor space (Xenv) ([Cook et al., 2010](#)) with customary prediction methods such as Principal Component Regression (PCR), Partial Least Squares Regression (PLS ([PLS1 and PLS2](#))) using simulated dataset with controlled properties. An 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,  $\boldsymbol{\Sigma}_{yy}$  and  $\boldsymbol{\Sigma}_{xx}$  are the variance-covariance matrices of  $\mathbf{y}$  and  $\mathbf{x}$ , respectively,  $\boldsymbol{\Sigma}_{xy}$  is

the covariance between  $\mathbf{x}$  and  $\mathbf{y}$  and  $\boldsymbol{\mu}_y$  and  $\boldsymbol{\mu}_x$  are mean vectors of  $\mathbf{x}$  and  $\mathbf{y}$ , respectively. A linear model based on (1) is,

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

where,  $\boldsymbol{\beta}^t$  is a matrix of regression coefficients and  $\boldsymbol{\epsilon}$  is an error term such that  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \boldsymbol{\Sigma}_{y|x})$ . Here,  $\boldsymbol{\beta}^t = \boldsymbol{\Sigma}_{yx}\boldsymbol{\Sigma}_{xx}^{-1}$  and  $\boldsymbol{\Sigma}_{y|x} = \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{yx}\boldsymbol{\Sigma}_{xx}^{-1}\boldsymbol{\Sigma}_{xy}$

In a model like (2), we assume that the variation in response  $\mathbf{y}$  is partly explained by the predictor  $\mathbf{x}$ . However, in many situations, only a subspace of the predictor space is relevant for the variation in the response  $\mathbf{y}$ . This space can be referred to as the relevant space of  $\mathbf{x}$  and the rest as irrelevant space. In the similar way, for a certain model, we can assume that a subspace in the response space exists which contains the information that the relevant space in predictor can explain (Figure-1). Cook et al. (2010) and Cook and Zhang (2015) have referred to the relevant space as material space, and the irrelevant space as immaterial space.

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

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{z} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \boldsymbol{\mu}_w \\ \boldsymbol{\mu}_z \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{ww} & \boldsymbol{\Sigma}_{wz} \\ \boldsymbol{\Sigma}_{zw} & \boldsymbol{\Sigma}_{zz} \end{bmatrix} \right) \quad (3)$$

where,  $\boldsymbol{\Sigma}_{ww}$  and  $\boldsymbol{\Sigma}_{zz}$  are the variance-covariance matrices of  $\mathbf{w}$  and  $\mathbf{z}$ , respectively.  $\boldsymbol{\Sigma}_{zw}$  is the covariance between  $\mathbf{z}$  and  $\mathbf{w}$ .  $\boldsymbol{\mu}_w$  and  $\boldsymbol{\mu}_z$  are mean vector of  $\mathbf{z}$  and  $\mathbf{w}$  respectively.

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

$$\mathbf{w} = \boldsymbol{\mu}_w + \boldsymbol{\alpha}^t(\mathbf{z} - \boldsymbol{\mu}_z) + \boldsymbol{\tau} \quad (4)$$

# Relevant space within a model

A concept for reduction of regression models

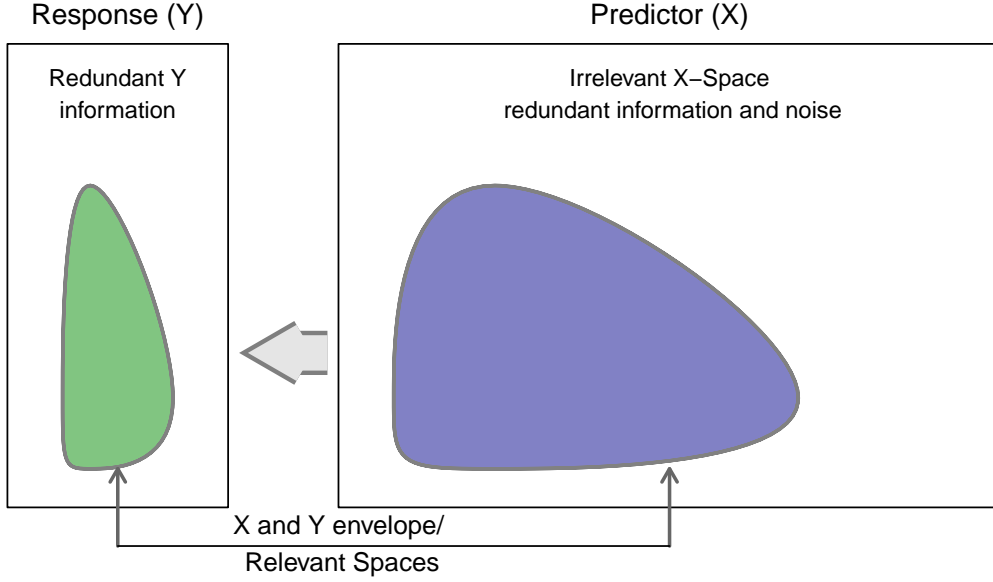


Figure 1: Relevant space in a regression model

where,  $\alpha^t$  is a matrix of regression coefficients and  $\tau$  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 (1985) introduced the concept of relevant components which was explored further by Helland (1990), Næs and Helland (1993), Helland and Almøy (1994) and Helland (2000). The corresponding eigenvectors were referred to as relevant eigenvectors. A similar logic is introduced by Cook et al. (2010) and later by Cook et al. (2013) as an envelope which is the space spanned by the relevant eigenvectors (Cook, 2018, 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 (1996) has used a single response simulation model based on reduced regression and has compared some contemporary multivariate estimators. In the recent years Helland et al. (2012), Sæbø

et al. (2015), Helland et al. (2018) and Rimal et al. (2018) implemented similar simulation examples as 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 parameter based on an experimental design. Rimal et al. (2018) has a detail discussion about the simulation model that we have opted here. The following section presents the estimators under comparison in more detail.

### 3. Prediction Methods

Partial least squares regression (PLS) and Principal component regression (PCR) has 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) is common. These methods are popular in multivariate analysis, especially for exploratory studies and prediction. In recent years, a concept of envelope introduced by Cook et al. (2007) based on reduction in regression model has been implemented for the development of different estimators. In this study, we will follow estimation 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 are sorted in descending order. In other words, each successive components captures maximum variation left by the preceding components in predictor variables (Jolliffe, 2002). Principal components regression uses these principal components as a new predictors to explain the variation in the response.

**Partial Least Squares (PLS):** Two variants of PLS: PLS1 and PLS2 will be 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 such as to maximize a covariance between response and predictors (de Jong, 1993).

**Envelopes:** The envelope, introduced by Cook et al. (2007), was first used to define response envelope (Cook et al., 2010) as a smallest subspace in the response space 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 irrelevant part increases the estimative variation. The concept of envelope uses the relevant part for estimation while excluding the irrelevant part consequently increasing the efficiency of the model (Cook and Zhang, 2016).

The concept was later extended to the predictor space, where the predictor envelope was defined (Cook et al., 2013). Further Cook and Zhang (2015) uses envelopes for joint reduction of the responses and predictors and argued to produce efficiency gains greater than using individual envelopes either of the response and predictors. All the variants of envelope estimations are based on maximum likelihood estimation. Here in this study we will also use predictor envelope (Xenv) and simultaneous envelope (Senv) for the comparison.

### 3.1. Modification in envelope estimation

Since envelope estimators (Xenv and Senv) are based on maximum likelihood estimation (MLE), it fails to estimate in case of wide matrices, i.e.  $p > n$ . In order 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  $\mathbf{x}$ . The new set of variables,  $\mathbf{z}$ , were used for envelope estimation. The regression coefficients ( $\hat{\boldsymbol{\alpha}}$ ) corresponding to these new variables  $\mathbf{z}$  were transformed back to obtain coefficients for each predictor variable as,

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

where,  $\mathbf{e}_k$  is a matrix of eigenvectors with first  $k$  number of components.

#### 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. (2015) and Rimal et al. (2018). Here we will use four different factors to vary the property of the data: a) Number of predictors ( $p$ ), b) Multicollinearity in predictor variables ( $\gamma$ ), c) Correlation in response variables ( $\eta$ ) and d) position of predictor components relevant for the response ( $relpos$ ). Using two levels of  $p$ ,  $\gamma$  and  $relpos$  and four levels of  $\eta$ , 32 set of distinct properties are designed for the simulation.

**Number of predictors:** In order to observe the performance of the methods on tall and wide predictor matrices, 20 and 250 predictor variables are simulated. Parameter  $p$  controls this properties in the `simrel` function.

**Multicollinearity in predictor variables:** Highly collinear predictors can be explained completely by few components. The parameter  $\gamma$  ( $\gamma$ ) in `simrel` controls 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. Here we have used 0.2 and 0.9 as different levels of  $\gamma$ . The higher the value of  $\gamma$ , the higher will be the multicollinearity and vice versa.

**Correlation in response variables:** Correlation among response | | multicollinearity in response | | variables is a less explored area. Here we have tried to explore that part with 4 levels of correlation in the response variables. We have used the  $\eta$  ( $\eta$ ) parameter of `simrel` for controlling the decline in eigenvalues corresponding to the response variables as (6).

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

Here,  $\kappa_i, i = 1, 2, \dots, m$  are the eigenvalues of the response variables and  $m$  is the number of response variables. Here we have used 0, 0.4, 0.8 and 1.2 as different levels of  $\eta$ . The larger the value of  $\eta$ , the larger will be the correlation between response variables and vice versa.

**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 the most in the rest that is left by the first principal components 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 (Helland and Almøy, 1994). Here, two levels of the positions of these relevant components are used: 1, 2, 3, 4 and 5, 6, 7, 8.

Further, a complete factorial design from the levels of the above given parameters gave us 32 designs. Each design is associated with a dataset having unique properties. Figure~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 dataset simulated.

**Common parameters:** Each dataset was simulated with  $n = 100$  number of observation and  $m = 4$  response variables. Further, the coefficient of determination corresponding to each response components in all the designs is set to and 0.8. In addition, we have assumed that there is only one informative response component. Hence, the informative response component is rotated orthogonally together with three uninformative response components to generate four response variables. This spread out the information in all simulated response variables. For further details on the simulation tool see (Rimal et al., 2018).



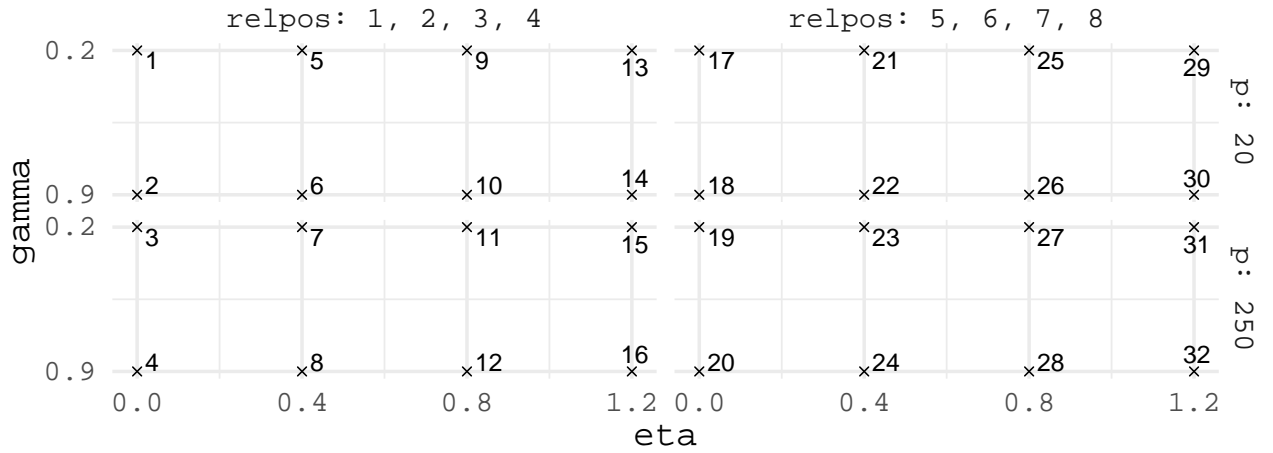


Figure 2: Experimental Design of simulation parameters. Each point represents an unique data property.

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"
)
```

Figure 3 shows the covariance structure of the data simulated with this design. The figure shows that the predictor components at position 1, 2, 3 and 4 are relevant for the first response component. After the rotation with orthogonal rotation matrix, all predictors are somewhat relevant for all response variables, fulfilling other desired properties like multicollinearity and coefficient of determination. For same design, Figure 4(top left) shows

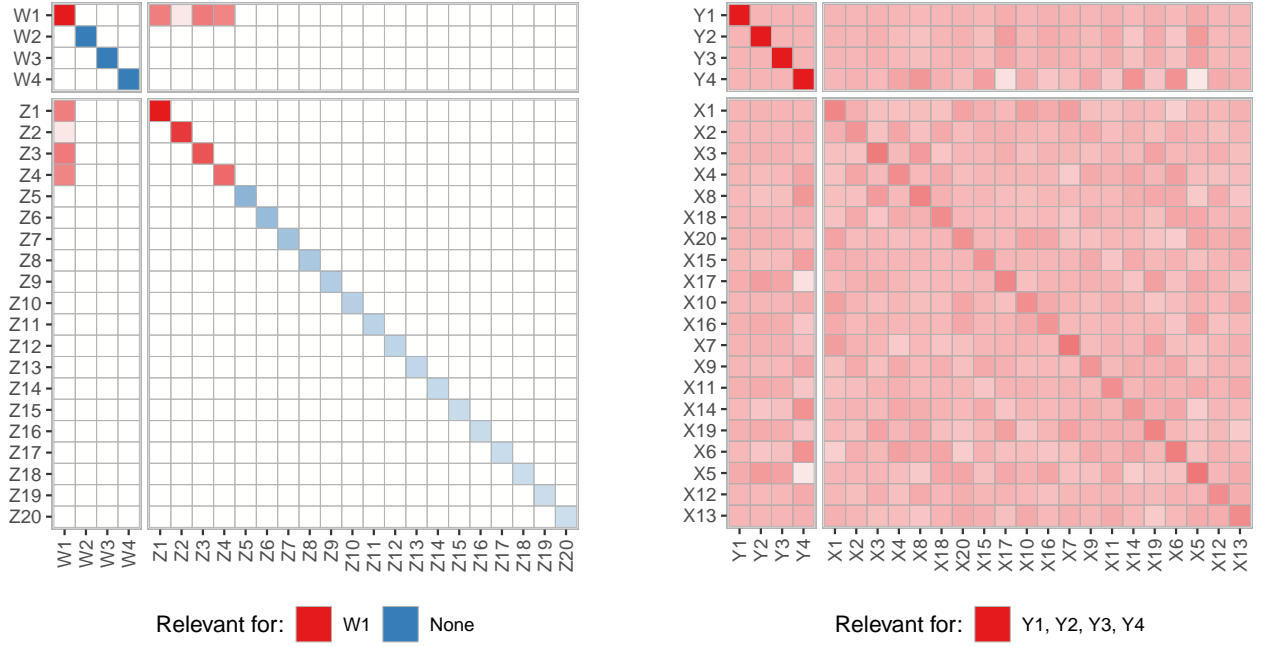


Figure 3: (left) Covariance structure of latent components. (right) Covariance structure of predictor and response

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 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 (Figure 4(top right)). The sample covariances between the predictors components and predictor variables with response variables are in Figure 4 (bottom left) and (bottom right) respectively.

A similar discussion can be made on all 32 designs where each of the design holds the properties of the data they simulate. These data are used by the prediction methods discussed in previous section. Each prediction methods are given independent dataset simulated in order to give them equal opportunity to understand the dynamics in the data.

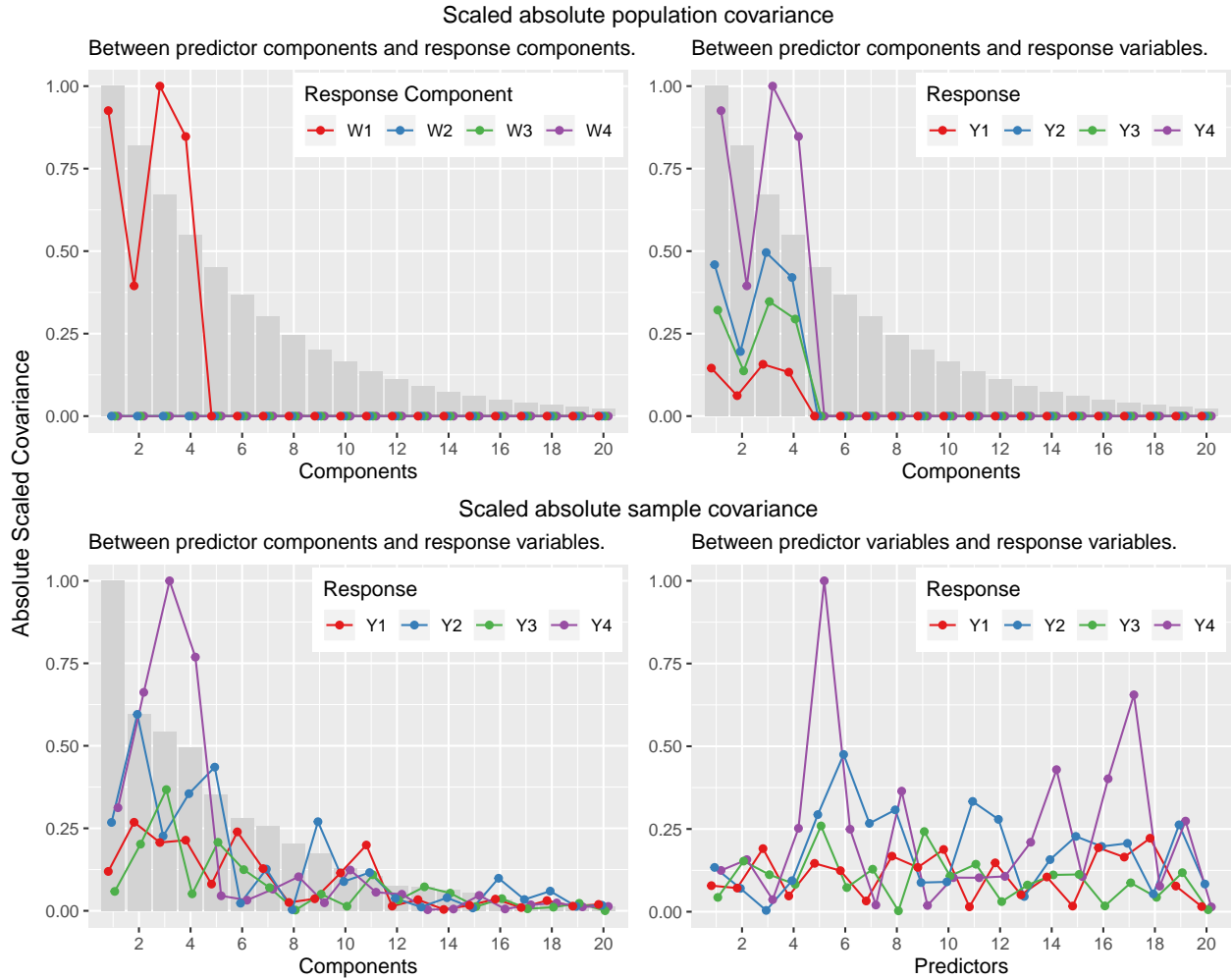


Figure 4: Expected Scaled absolute covariance between predictor components and response components (top left). Expected Scaled absolute covariance between predictor components and response variables (top right). Sample scaled absolute covariance between predictor components and response variables (bottom left). Sample scaled absolute covariance between predictor variables and response variables (bottom right). The bar in the background are eigenvalues corresponding to each components in population (top plots) and in 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.

## 5. Basis of comparison

This study focuses mainly on the prediction performance of the methods and 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 on the following basis:

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

Let us define,

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

as a prediction error of response  $j = 1, \dots, 4$  for a given design  $i = 1, 2, \dots, 32$  and method  $k = 1(PCR), \dots, 5(Senv)$  using  $l = 0, \dots, 10$  number of components. Here,  $(\Sigma_{xx})_i$  is the true covariance matrix of predictor unique for a particular design  $i$  and  $\sigma_{y_{ij}|x}$  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 from the prediction error. Since both the expectation and the variance of  $\hat{\beta}$  are unknown, the prediction error are estimated using data from 50 replications as follows,

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

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

Following section focuses on the data for the estimation of these prediction error 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 five factors corresponding to simulation parameters, prediction methods, number of components, replications and prediction error for four responses. The prediction error is computed using 0 to 10 predictor components for each 50 replicates as,

$$\left(\widehat{\mathcal{PE}}_{\circ}\right)_{ijklr} = \frac{1}{\sigma_{y_{ij}|x}^2} \left[ \left(\beta_{ij} - \hat{\beta}_{ijklr}\right)^t (\Sigma_{xx})_i \left(\beta_{ij} - \hat{\beta}_{ijklr}\right) \right] + 1$$

.

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

Since we will focus our discussion 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 to construct 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 components. Next, a component that gives the minimum of this average prediction error is selected, i.e.,

$$l_{\circ} = \underset{l}{\operatorname{argmin}} \left[ \frac{1}{50} \sum_{i=1}^{50} (\mathcal{PE}_{\circ})_{ijklr} \right] \quad (9)$$

Using the component  $l_{\circ}$ , a dataset of  $(\mathcal{PE}_{\circ})_{ijkl_{\circ}r}$  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 the response  $j$  in the context of this dataset.

**Component Dataset:** Components that gives the minimum prediction error in each replication is used as *Component Dataset*, i.e.,

$$l_{\circ} = \underset{l}{\operatorname{argmin}} [\mathcal{PE}_{ijklr}] \quad (10)$$

Here  $l_o$  is the number of component that gives minimum prediction error  $(\mathcal{PE}_o)_{ijklr}$  for design  $i$ , response  $j$ , method  $k$  and replicate  $r$ .  $l_o$  is used as the *component Dataset*. 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 the response  $j$  in the context of this dataset.

## 7. Exploration

This section focuses on exploring the variation in *error dataset* and *component dataset* for which we will use Principal Component Analysis (PCA). Let  $\mathbf{t}_u$  and  $\mathbf{t}_v$  be the principal component scores corresponding to the  $\mathbf{u}$  and  $\mathbf{v}$  respectively. Figure-5 plots the scores density corresponding to first principal component of  $\mathbf{u}$ , i.e. first column of  $\mathbf{t}_u$ .

Since higher prediction error results in high scores the plot shows that the PCR, PLS1 and PLS2 methods are influenced by two levels of position of relevant predictor components. When the relevant predictors are at position 5, 6, 7, 8, the eigenvalues corresponding to them becomes smaller. This also suggest that PCR, PLS1 and PLS2 depends highly on the relevant components and the variation of these components affect their performance. However, the envelope methods have less influence of `relpos` in this regard.

In addition, the plot also shows that the effect of `gamma`, i.e., the level of multicollinearity, has smaller effect when the relevant predictors are at position 1, 2, 3, 4. This indicates that the methods are somewhat robust to handle collinear predictors. Although, when the relevant predictors are at position 5, 6, 7, 8 high multicollinearity results in small variance of these relevant components and consequently gives poor prediction.

Further, the density curve for PCR, PLS1 and PLS2 are similar for different levels of `eta`, i.e., the factor controlling the correlation between responses. However, this is not true for envelope models. The envelope methods have shown to have distinct interaction between position of relevant components and `eta`. Here higher levels of `eta` is giving larger scores and clear separation between two level of `relpos`. This behavior is expected in the simultaneous envelope as the method has claimed to model relevant (material) response space.

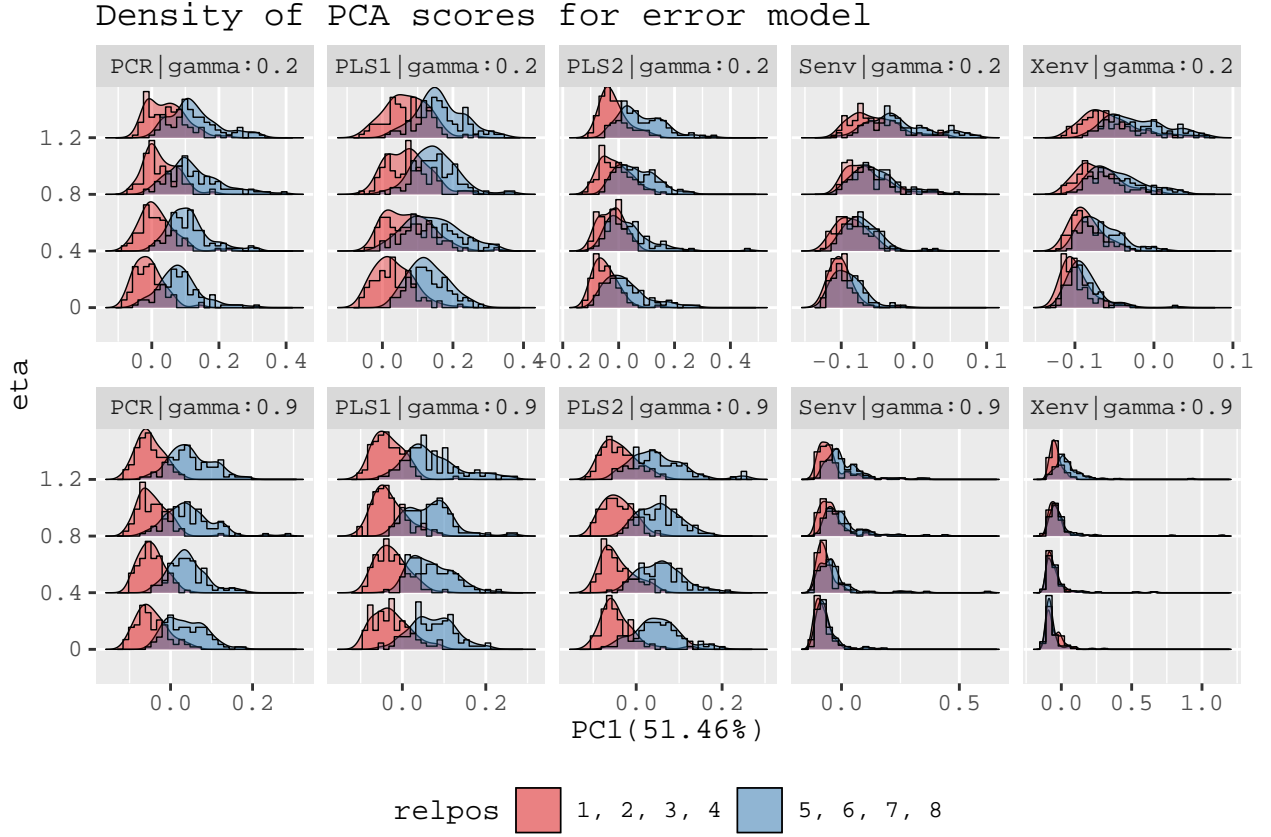


Figure 5: Scores density corresponding to first principal component of *error* dataset (**u**) subdivided by methods, gamma and eta and grouped by relpos.

However in the case of high multicollinearity, envelope methods have resulted in some large outliers. This suggests that in the case of multicollinearity, the methods can give unexpected prediction.

In the Figure 6, the higher scores suggest the Methods have used more number of components to give minimum prediction error. The plot also shows that the relevant predictor components at 5, 6, 7, 8 gives larger prediction error than those which are at the position 1, 2, 3, 4. The pattern is more distinct in large multicollinearity case and PCR and PLS methods. Both the envelope methods have shown equally better performance at both levels of relpos and gamma. However in low multicollinear predictors envelope methods have used fewer number of components in average than high multicollinear predictors to give minimum prediction error.

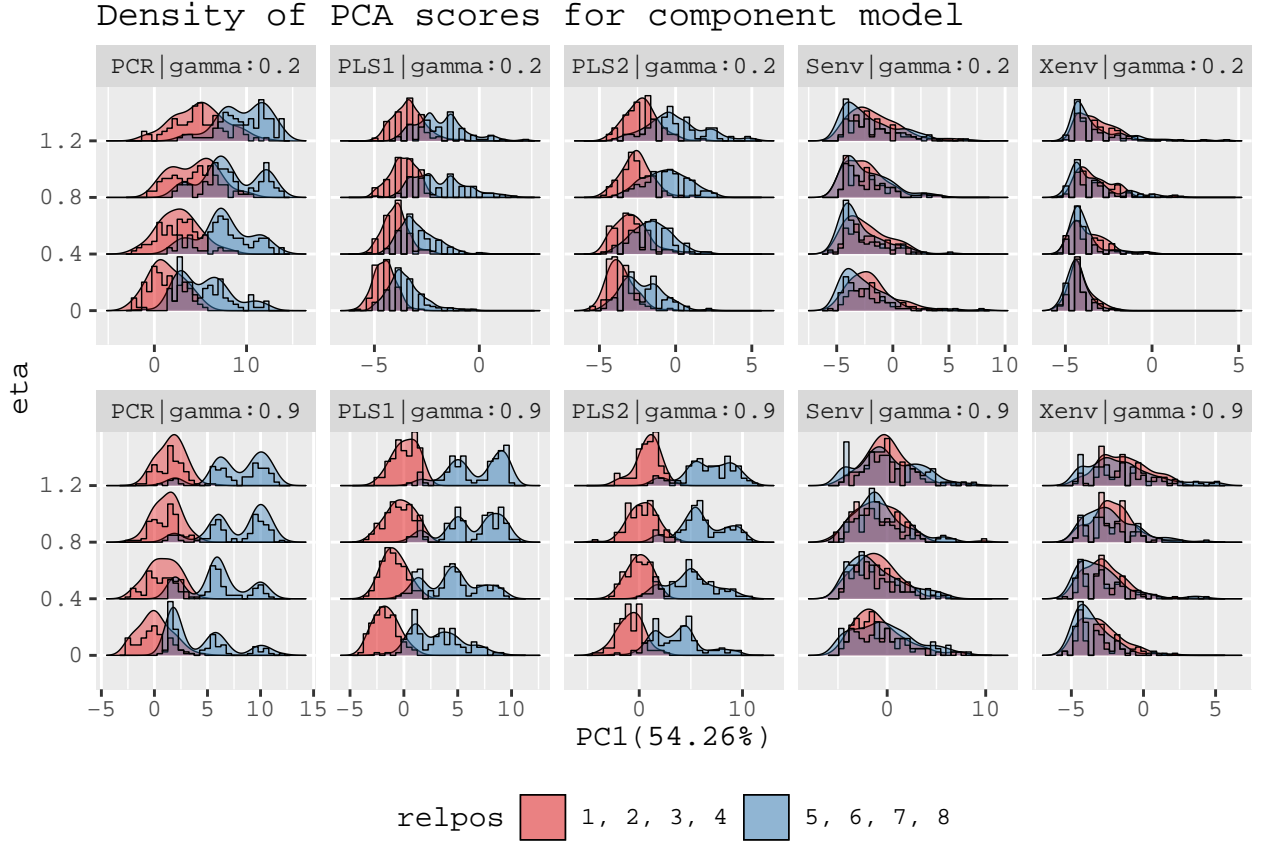


Figure 6: Score density corresponding to first principal component of *component dataset* (**v**) subdivided by methods, gamma and eta and grouped by relpos.

## 8. Statistical Analysis

Let us consider a model with third order interaction of simulation parameters (*p*, *gamma*, *eta* and *relpos*) and Methods as in (11) and (12) using *error dataset u* and **v** respectively. Let us call them *error model* and *\_component model\*\** respectively.

**Error Model:**

$$\mathbf{u}_{abcdef} = \mu_u + (p_a + \text{gamma}_b + \text{eta}_c + \text{relpos}_d + \text{Methods}_e)^3 + (\varepsilon_u)_{abcdef} \quad (11)$$

**Component Model:**

$$\mathbf{v}_{abcdef} = \mu_v + (p_a + \text{gamma}_b + \text{eta}_c + \text{relpos}_d + \text{Methods}_e)^3 + (\varepsilon_v)_{abcdef} \quad (12)$$



where,  $\mathbf{u}_{abcdef}$  is a vector of prediction errors in *error model* and  $\mathbf{v}_{abcdef}$  is a vector of number of components used by a method to obtain minimum prediction error in *component model*. Although there are several test-statistic for MANOVA, for large sample all are essentially equivalent (Johnson and Wichern, 2018). Here we will use Pillai's trace statistic which is defined as,

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

Here the matrix  $\mathbf{H}$  has between sum of squares and sum of products for each of the predictors. The matrix  $\mathbf{E}$  has a within sum of squares and sum of products for each of the predictors.  $v_i$  represents the eigenvalues corresponding to  $\mathbf{E}^{-1}\mathbf{H}$  (Rencher, 2003). For both the models (11) and (12), Pillai's trace statistic is used for accessing the effect of each factors and F-value for the strength of their significance. Figure 7 plots the Pillai's trace statistic as bar with corresponding F-values as text label for both models.

**Error Model:** Figure 7 (left) shows the Pillai's trace statistic for factors of the *error model*.

The main effects of the Method has largest influence on the model followed by relpos, eta and gamma. A highly significant interaction of Method with eta, relpos and gamma clearly shows that methods perform differently for different levels of these data properties. Further, the significant third order interaction between Method, eta and gamma suggests that a method perform differently 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 response (eta) for a given method.

**Component Model:** Figure 7 (right) shows the Pillai's trace statistic for factors of the *component model*. As in *error model*, the main effects of the Method, relpos, gamma and eta have significantly large effect on number of components that a method has used to get minimum prediction error. The interaction of Method with simulation parameters are significantly larger in this case. This shows that the Methods and these interactions have larger effect on the use of number of component than the

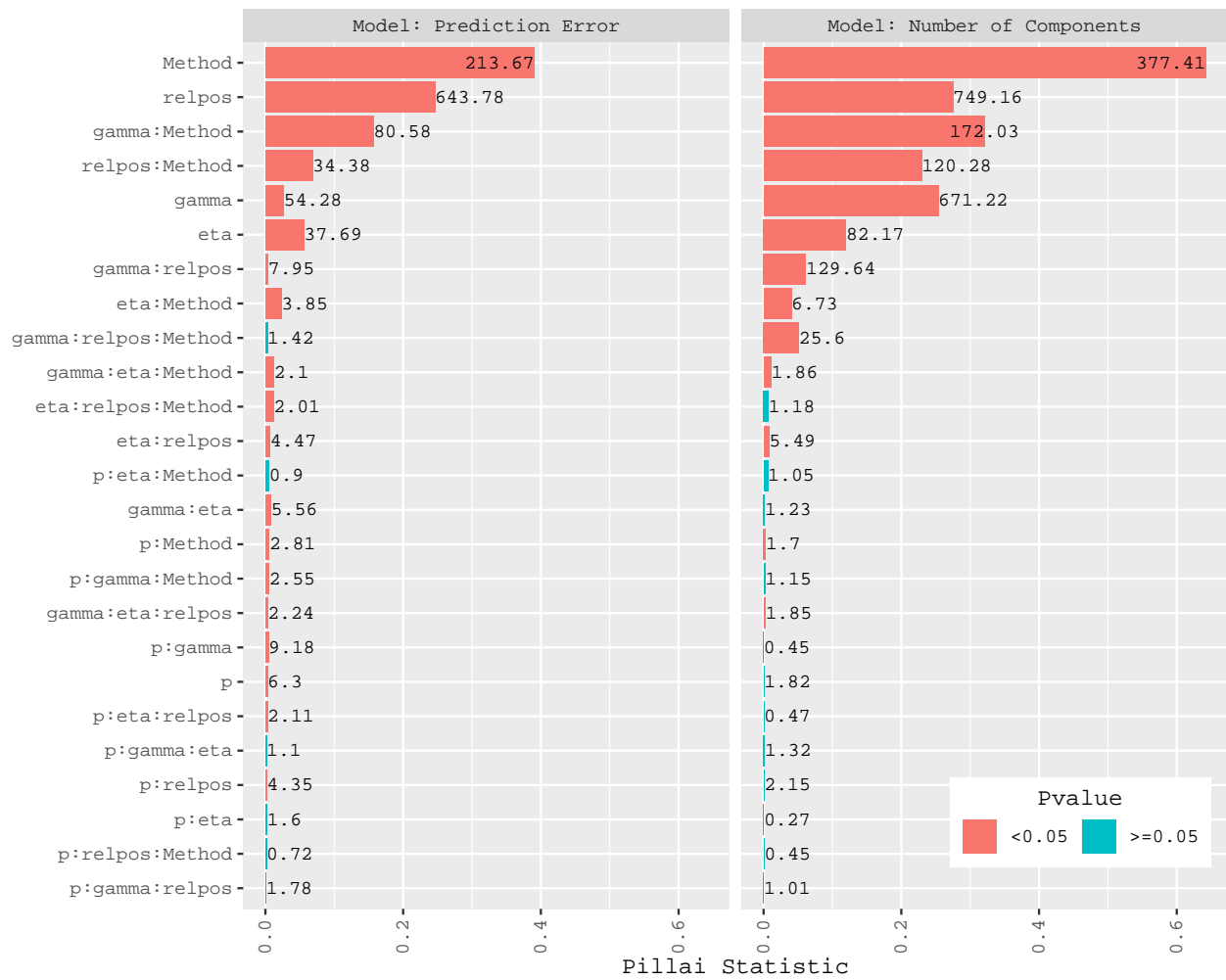


Figure 7: Pillai Statistic and F-value for the MANOVA model. The bar represents the Pillai Statistic and the text labels are F-value for corresponding factor.

prediction error itself. In addition, a similar significant high third order interaction as in *error model* is also observed in this model.

Following section will continue demystifying the effect of different levels of the factors in the case of these interactions.

### 8.1. Effect Analysis of Error Model

Figure 8 (left) shows large difference in prediction error in envelope models is observed. The difference is intensified when position of relevant predictor at at 5, 6, 7, 8. The results also shows that the envelope methods respond to the levels of eta more 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 the Figure 8 (right), we can see that the multicollinearity has affected all the methods however, PCR, PLS1 and PLS2 are more robust for the condition than envelopes. Rather these methods have shown better performance when high multicollinearity is present in the data. Envelope methods on the other hand are better in handling the model when relevant position are at 5, 6, 7, 8 in both high and low multicollinearity cases.

### 8.2. Effect Analysis of Component Model

Unlike prediction error, Figure 9 (left) shows that the number of components used by the methods to give minimum prediction error is less affected by the levels of eta. However, simultaneous envelope has shown distinct reaction when there is no correlation between the response variables. 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 predictors. The number of components used by this method ranges from 3 to 5 when relevant predictors are at 1, 2, 3, 4 and 5 to 8 when relevant predictors are at 5, 6, 7, 8.

We can also see that at relpos 1, 2, 3, 4 for PLS1 and PLS2 have used fewer components than simultaneous envelope. However, in the case when relevant components are at position 5, 6, 7, 8, simultaneous envelope manage to get smaller prediction error using

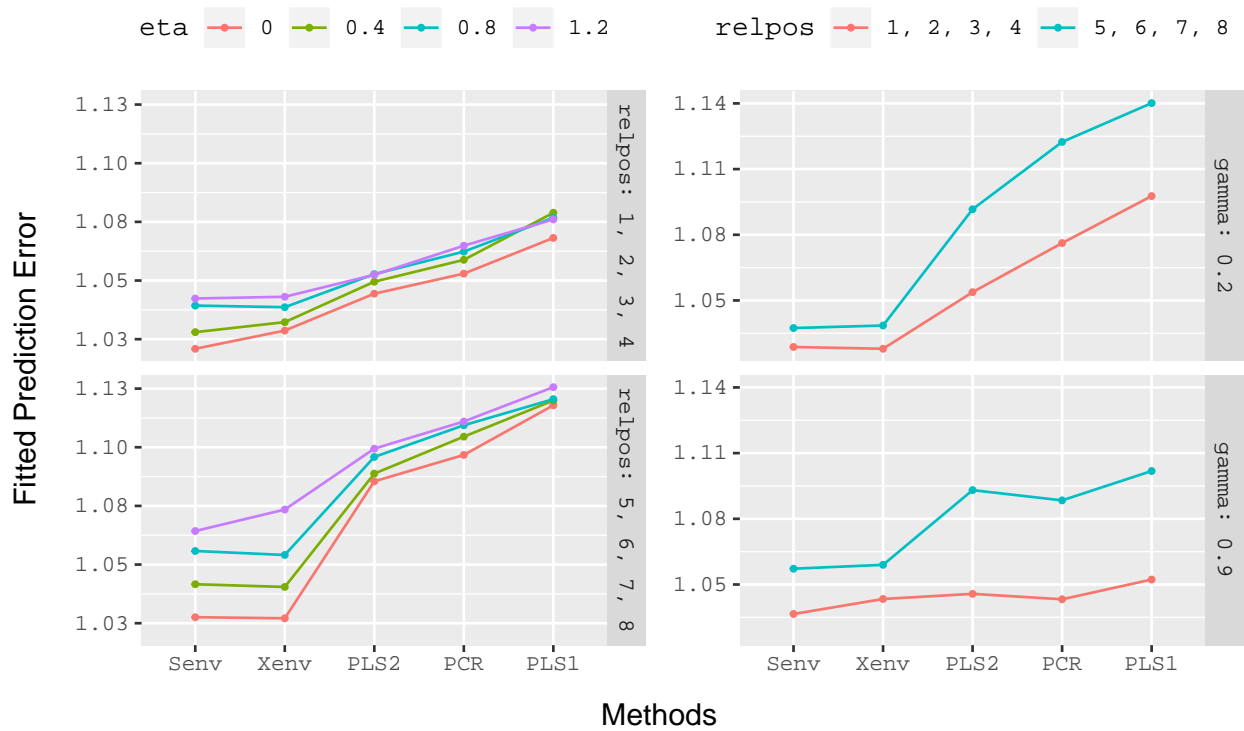


Figure 8: Effect plot of some interactions of the multivariate linear model of prediction error

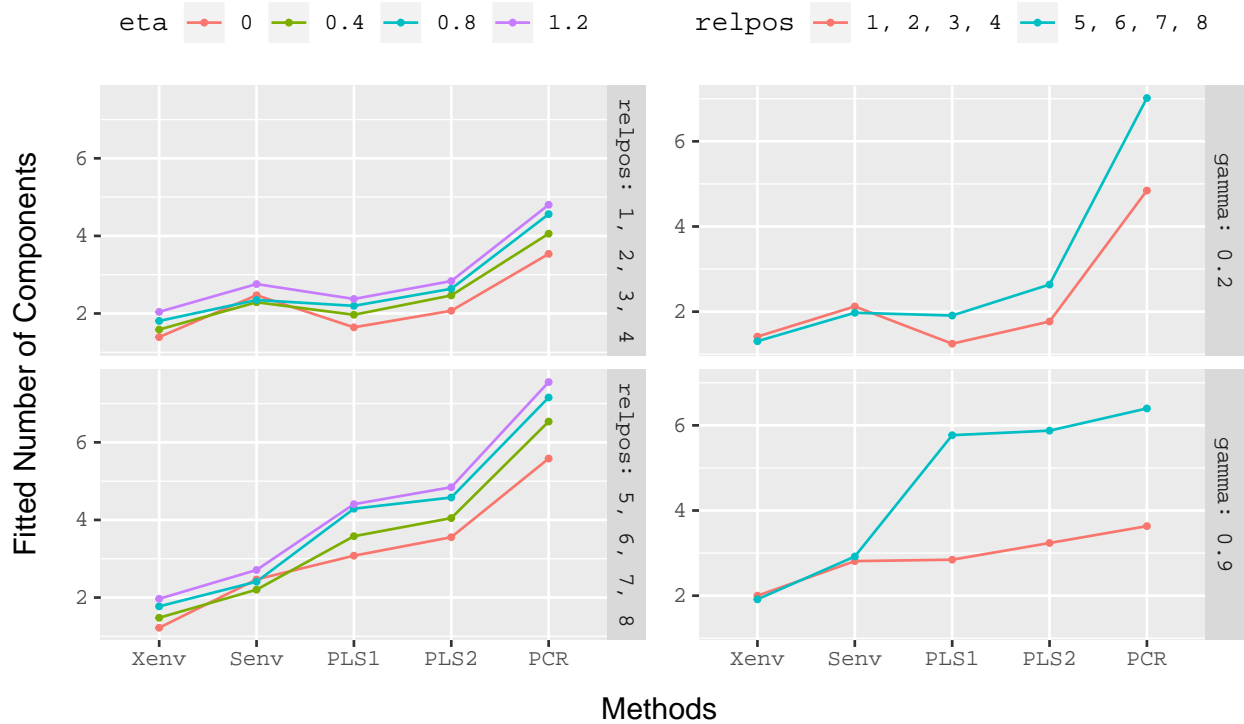


Figure 9: Effect plot of some interactions of the multivariate linear model of number of components to get minimum prediction error

fewer number of components than that of PLS models. This is the case when eigenvalues of relevant predictors are small and are difficult to predict.

In the interaction effect of gamma, relpos and Method (Figure 9 (right)), PCR, PLS1 and PLS2 methods have used fairly large number of components for the model with high multicollinearity and relevant predictors at position 5, 6, 7, 8. The number of components used by envelope methods in both the cases of relpos is similar although slightly more for the models with high multicollinearity than models with low multicollinearity.

In the case of PCR, the prediction relies heavily on the relevant components, the number of components used by the method is as expected in both the levels of relpos. This reinforce the results from the density plot (Figure 6) in previous section.

## 9. Discussions and Conclusion

On one hand, the envelope methods have shows better prediction performance with fewer number of components. However, for data with high multicollinearity, the methods have shown some cases of unstable prediction. Here for the prediction error, we have used the components that gives the miniumum average prediciton error for all replicates. This hints that for some replicates, the methos has used non-optimal number of components and consequently resulting in high prediction error.

On the other hand, PLS1, PLS2 and PCR methods have smaller prediction error in the model with high multicollinearity which suggest their robustness in the particular case. However, they are poor in modelling information at relevant predictors with small variance. These methods have used larger number of components, in general, than envelopes. Although resulting high prediction error than envelope in most situation, these methods are more stable especially in the cases of high multicollinearity.

Further, we have fixed the coefficient of determination ( $R^2$ ) constant throughtout all the design. Since low  $R^2$  affected all the methods in similar manner and the MANOVA is highly dominated by  $R^2$ . Keeping the value of  $R^2$  fixed have allowed us to analyze other factors properly.

In addition, since using principal components or actual variables as predictors in envelope methods have 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. As the envelope methods are based on MLE and this can be an alternative way of using the methods in data with wide predictors. The results from this study will help researcher to understand these methods for different nature of data. We encourage researcher to use newly developed methods such as envelope based on the nature of data they are working on.

Since, this study have focused entirely on prediction performance, further analysis of their estimative properties of these methods is required. A study of estimation error and the behaviour of methods on non-optimal number of components can give deeper understanding of these methods.

A shiny application (Chang et al., 2018) 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.

- Almøy, T., jan 1996. A simulation study on comparison of prediction methods when only a few components are relevant. *Computational Statistics & Data Analysis* 21 (1), 87–107.
- Chang, W., Cheng, J., Allaire, J., Xie, Y., McPherson, J., 2018. shiny: Web Application Framework for R. R package version 1.2.0.  
URL <https://CRAN.R-project.org/package=shiny>
- Cook, R. D., 2018. An introduction to envelopes : dimension reduction for efficient estimation in multivariate statistics, 1st Edition. Hoboken, NJ : John Wiley & Sons, 2018.
- Cook, R. D., Helland, I. S., Su, Z., 2013. Envelopes and partial least squares regression. *Journal of the Royal Statistical Society. Series B: Statistical Methodology* 75 (5), 851–877.
- Cook, R. D., Li, B., Chiaromonte, F., aug 2007. Dimension reduction in regression without matrix inversion. *Biometrika* 94 (3), 569–584.
- Cook, R. D., Li, B., Chiaromonte, F., 2010. Envelope Models for Parsimonious and Efficient Multivariate Linear Regression. *Statistica Sinica* 20 (3), 927–1010.
- Cook, R. D., Zhang, X., 2015. Simultaneous envelopes for multivariate linear regression. *Technometrics* 57 (1), 11–25.

- Cook, R. D., Zhang, X., 2016. Algorithms for Envelope Estimation. *Journal of Computational and Graphical Statistics* 25 (1), 284–300.
- de Jong, S., mar 1993. SIMPLS: An alternative approach to partial least squares regression. *Chemometrics and Intelligent Laboratory Systems* 18 (3), 251–263.
- Helland, I. S., 1990. Partial least squares regression and statistical models. *Scandinavian Journal of Statistics* 17 (2), 97–114.
- Helland, I. S., mar 2000. Model Reduction for Prediction in Regression Models. *Scandinavian Journal of Statistics* 27 (1), 1–20.
- Helland, I. S., Almøy, T., 1994. Comparison of prediction methods when only a few components are relevant. *Journal of the American Statistical Association* 89 (426), 583–591.
- Helland, I. S., Sæbø, S., Almøy, T., Rimal, R., Sæbø, S., Almøy, T., Rimal, R., sep 2018. Model and estimators for partial least squares regression. *Journal of Chemometrics* 32 (9), e3044.
- Helland, I. S., Sæbø, S., Tjelmeland, H. K., mar 2012. Near Optimal Prediction from Relevant Components. *Scandinavian Journal of Statistics* 39 (4), 695–713.
- Johnson, R., Wichern, D., 2018. *Applied Multivariate Statistical Analysis (Classic Version)*. Pearson Modern Classics for Advanced Statistics Series. Pearson Education Canada.  
URL <https://books.google.no/books?id=QBqlswEACAAJ>
- Jolliffe, I. T., 2002. *Principal Component Analysis*, Second Edition.
- Næs, T., Helland, I. S., 1993. Relevant components in regression. *Scandinavian Journal of Statistics* 20 (3), 239–250.
- Naes, T., Martens, H., jan 1985. Comparison of prediction methods for multicollinear data. *Communications in Statistics - Simulation and Computation* 14 (3), 545–576.
- Rencher, A. C., 2003. *Methods of multivariate analysis*. Vol. 492. John Wiley & Sons.
- Rimal, R., Almøy, T., Sæbø, S., may 2018. A tool for simulating multi-response linear model data. *Chemometrics and Intelligent Laboratory Systems* 176, 1–10.
- Sæbø, S., Almøy, T., Helland, I. S., 2015. Simrel - A versatile tool for linear model data simulation based on the concept of a relevant subspace and relevant predictors. *Chemometrics and Intelligent Laboratory Systems* 146, 128–135.