

# Comparison of Multivariate Estimation 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

Prediction performance often does not reflect the estimation behaviour of a method. High error in estimation not necessarily results in high prediction error but can lead to an unreliable prediction when test data are in a different direction than the training data. In addition, the effect of a variable becomes unstable and can not be interpreted in such situations. Many research fields are more interested in these estimates than performing prediction. This study compares some newly-developed (envelope) and well-established (PLS, PCR) prediction methods using simulated data with specifically designed properties such as multicollinearity, the correlation between multiple responses and position of principal components of predictor that are relevant for the response. This study aims to give some insight into these methods and help the researcher to understand and use them for further study. *Write some specifics from the results to show what we have found.*

*Keywords:* model-comparison, multi-response, simrel, estimation, estimation error, meta modeling

---

## 1. Introduction

Estimation of parameters in a regression model is an integral part of many research study. Research fields such as social science, econometrics, psychology and medical study are

---

\*Corresponding Author

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

more interested in measuring the impact of certain indicator or variable rather than performing prediction. Such studies have a large influence on people's perception and also help in policy making and decisions.

Technology has facilitated researcher to collect a large amount of data however often times, such data either contains irrelevant information or are highly collinear. Researchers are devising new estimators to extract information and identify their inter-relationship. Some estimators are robust towards fixing the multicollinearity problem while some are targeted to model only the relevant information contained in the response variable.

This study extends the (Rimal et al., 2019) and compares some well-established estimators such as Principal Components Analysis (PCA), Partial Least Squares (PLS) together with two new methods based on envelope estimation: Envelope estimation in predictor space (Xenv) (Cook et al., 2010) and simultaneous estimation of envelope (Senv) (Cook and Zhang, 2015). The estimation process of these methods is discussed in [Methods] section. The comparison tests the estimation performance of these methods using multi-response simulated data from a linear model with controlled properties. The properties include the number of predictors, level of multicollinearity, the correlation between different response variables and the position of relevant predictor components. These properties are explained in Experimental Design section together with the strategy behind the simulation and data model.

## 2. Simulation Model

- Reduction of the regression model
- Include the figure from previous paper
- How the covariance and coefficients are related
- From the construction of the covariance matrix of latent variables to the simulated data
- How and what simulation parameters are related to properties of data

### 3. Estimation Methods

A regression model is written as,

$$\underset{(1 \times m)}{\mathbf{y}} = \underset{(1 \times p)(p \times m)}{\mathbf{x}\boldsymbol{\beta}} + \underset{(1 \times m)}{\boldsymbol{\varepsilon}} \quad (1)$$

where  $\mathbf{y}$  is a vector of  $m$  responses measured about their means,  $\mathbf{x}$  is a vector of  $p$  predictors measured about their means,  $\boldsymbol{\beta}$  is a matrix of regression coefficients and  $\boldsymbol{\varepsilon}$  is a vector of independent error terms with constant variance  $\Sigma_{y|x}$ . In ordinary least squares, coefficient  $\boldsymbol{\beta}$  is estimated as,

$$\underset{(p \times m)}{\hat{\boldsymbol{\beta}}} = \left( \underset{(p \times n)(n \times p)}{\mathbf{x}^t \mathbf{x}} \right)^{-1} \underset{(p \times n)(n \times m)}{\mathbf{x}^t \mathbf{y}} \quad (2)$$

#### 3.1. Principal Components Regression

Principal Components are new set of variables from the transformation of original dataset such that they are uncorrelated with each other and the variation in the original data are ordered from first to last of these new variables. Let us define a transformation of  $\mathbf{x}$  as,

$$\underset{(1 \times p)}{\mathbf{x}} = \underset{(1 \times k)(k \times p)}{\mathbf{z} \mathbf{R}^t} \quad (3)$$

where  $\mathbf{R}$  is the eigenvectors corresponding to the covariance of  $\mathbf{x}$  and  $\mathbf{z}$  are the principal components. A regression model can be defined in terms of  $\mathbf{z}$  as,

$$\underset{(1 \times m)}{\mathbf{y}} = \underset{(1 \times k)(k \times m)}{\mathbf{z} \boldsymbol{\alpha}} + \underset{(1 \times m)}{\boldsymbol{\varepsilon}} \quad (4)$$

Since the variation is ordered in  $\mathbf{z}$ , only  $k \leq p$  columns of  $\mathbf{z}$  are used so that  $p - k$  uninformative components are not used for modeling. The regression coefficient of (4) can be estimated as,

$$\underset{(k \times m)}{\hat{\boldsymbol{\alpha}}} = \left( \underset{(k \times n)(n \times k)}{\mathbf{z}^t \mathbf{z}} \right)^{-1} \underset{(k \times n)(n \times m)}{\mathbf{z}^t \mathbf{y}} \quad (5)$$

Using (3) in (2), we get,

$$\hat{\beta}_{(p \times m)} = \begin{bmatrix} \mathbf{R}\mathbf{z}^t & \mathbf{z}\mathbf{R}^t \end{bmatrix}_{(p \times k)(k \times n)(n \times k)(k \times p)}^{-1} \begin{bmatrix} \mathbf{R}\mathbf{z}^t & y \end{bmatrix}_{(p \times k)(k \times n)(n \times m)}$$

### 3.2. Partial Least Squares Regression

- How beta coefficients are constructed
- How it is dependent on the variance-covariance matrices
- In what way PLS1 and PLS2 differ

### 3.3. Envelope Estimations

- How beta coefficients are constructed
- How it is dependent on the variance-covariance matrices
- In what way Xenv and Senv differ

## 4. Experimental Design

An R (R Core Team, 2018) package `simrel` (Rimal et al., 2018; Sæbø et al., 2015) is used to simulate data. For the simulation the number of observation is kept fixed at  $n = 100$  and following four simulation parameters are used to obtain the data with wide range of properties.

**Number of predictors:** In order to cover both tall ( $n > p$ ) and wide ( $p > n$ ) cases,  $p = 20$  and  $p = 250$  number of predictors are simulated.

**Multicollinearity in predictor variables:** A parameter gamma ( $\gamma$ ) in simulation controls the exponential decline of eigenvalues ( $\lambda_i, i = 1, \dots, p$ ) corresponding to predictor variables as,

$$\lambda_i = e^{-\gamma(i-1)}, \gamma > 0 \text{ and } i = 1, 2, \dots, p \quad (6)$$

Two levels 0.2 and 0.8 of gamma are used for simulation so that level 0.2 simulates the data with low multicollinearity and 0.8 simulates the data with high multicollinearity.

**Position of relevant components:** Initial principal components of a non-singular covariance matrix are larger than the later one. If the principal components corresponding to predictors with larger variation is not relevant for a response, this will just increase noise in the model. Here we will use two different levels of position index of predictor components: a) 1, 2, 3, 4 and b) 5, 6, 7, 8. Predictor components irrelevant for a response makes prediction difficult (Helland and Almøy, 1994). When combined with multicollinearity, this factor can create both easy and difficult model for both estimation and prediction.

**Correlation in response variables:** Many estimators also uses the structure of response for their estimation. Here the correlation between the responses are varied through a simulation parameter eta ( $\eta$ ). The parameter controls the exponential decline of eigenvalues  $\kappa_j, j = 1, \dots, m$  (number of responses) corresponding to response variables as,

$$\eta_i = e^{-\kappa(j-1)}, \kappa > 0 \text{ and } j = 1, 2, \dots, m \quad (7)$$

Four levels 0, 0.4, 0.8 and 1.2 of eta are used in the so that level 0 simulates the data with uncorrelated response variables while 1.2 simulates the highly correlated response variables.

Here we have assumed that there is only one informative response component. In the final dataset, all predictors together span the same space as the relevant predictor components and all response together span the same space as the one informative response component. In addition, coefficient of determination is fixed at 0.8 for all dataset.

A complete factorial design is adopted using different levels of factors discussed above to create 32 design (Figure 1) each of which gives dataset with unique properties. From each of these design and each estimation method, 50 different datasets are simulated so that each of them have same true population structure. In total,  $5 \times 32 \times 50$  i.e., 8000 datasets are simulated.

The simulation properties are directly reflected in the simulated data. For example, in



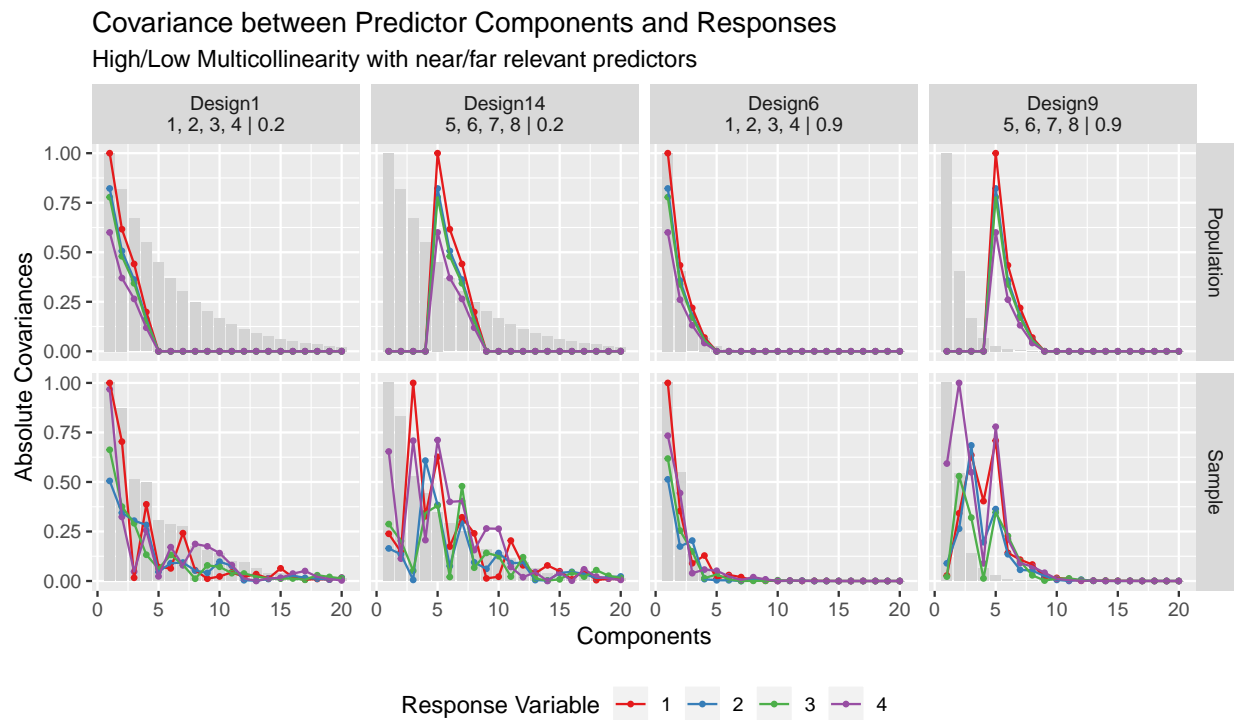


Figure 2: Covariance between predictor components and response variables in population (top) and in the simulated data (bottom) for four different designs. The Bar in the background represents the variance of corresponding components.

for response  $j = 1, \dots, 4$  in a given design  $i = 1, 2, \dots, 32$  and method  $k = 1(PCR), \dots, 5(Senv)$  using  $l = 0, \dots, 10$  number of components. Since both the expectation and the variance of  $\hat{\beta}$  are unknown, the prediction error are estimated using data from 50 replications as follows,

$$\widehat{\mathcal{E}\mathcal{E}_{ijkl}} = \sum_{r=0}^{50} \left[ \left( \beta_{ij} - \hat{\beta}_{ijklr} \right)^t \left( \beta_{ij} - \hat{\beta}_{ijklr} \right) \right] \quad (9)$$

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

## 6. Exploration

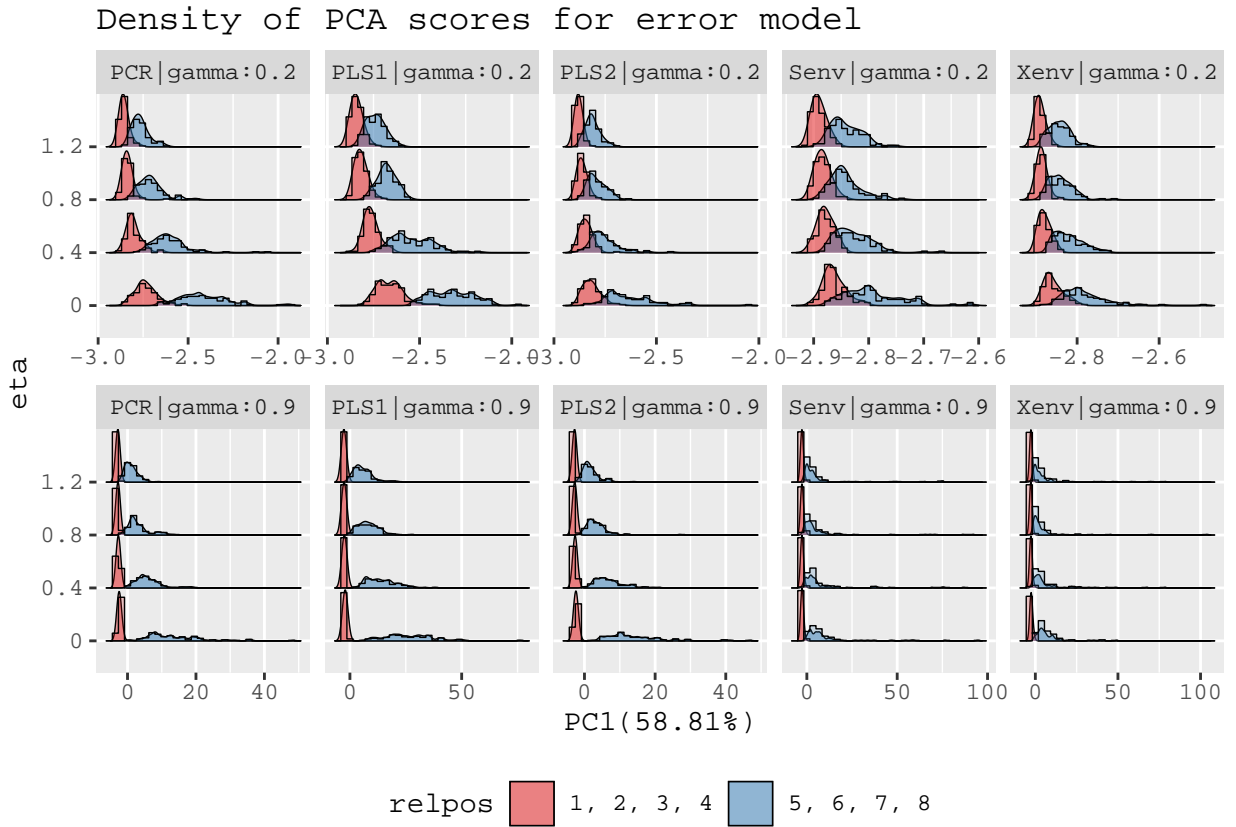


Figure 3: Scores density corresponding to first principal component of *error dataset* (**u**) subdivided by methods,  $\gamma$  and  $\eta$  and grouped by relpos.

- A similar exploration as previous paper but can be different in case of new idea



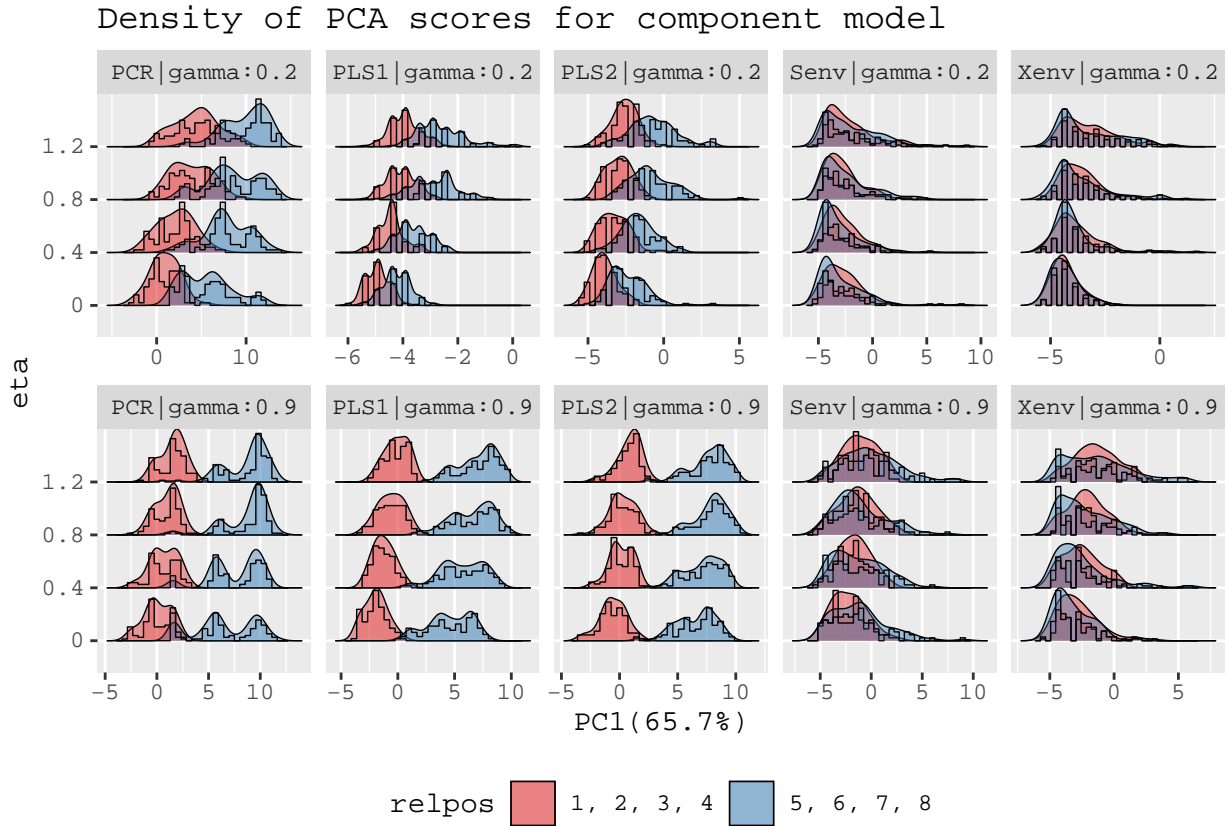


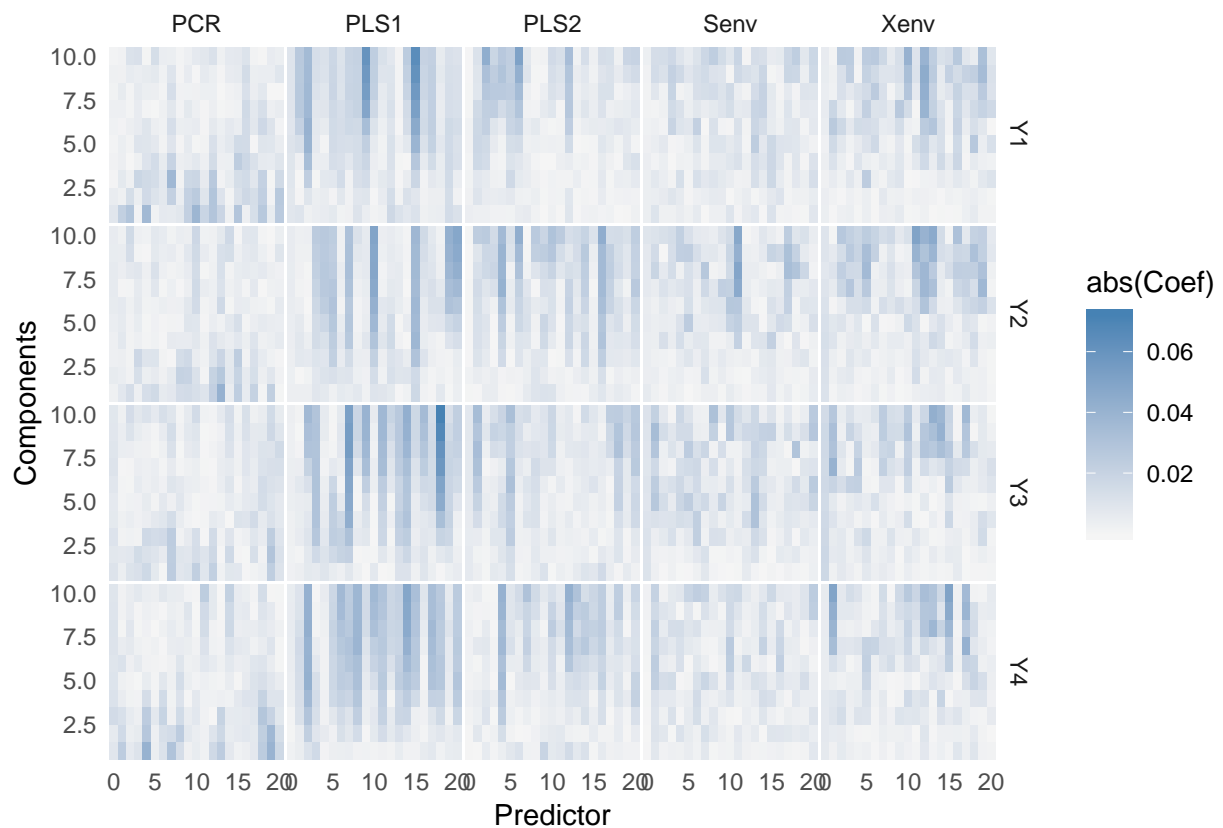
Figure 4: Score density corresponding to first principal component of *component dataset* ( $\mathbf{v}$ ) subdivided by methods, gamma and eta and grouped by relpos.

### 6.1. Dataset for Analysis

- Preparation of dataset for MANOVA analysis (i.e. minimum estimation error using arbitrary number of components)
- A component dataset is also created for testing the use of components by each of these methods

### 6.2. Regression Coefficients

- In case of some idea on comparing regression coefficients through some statistical way, this can be included here
- Otherwise can also be done just by using plots



### 6.3. Prediction and Estimation Error

- Explore both estimation error and number of components and try to bind them with the prediction error for the similar case

## 7. Analysis

- A MANOVA model is fitted using the dataset prepared in previous section

### 7.1. Error Analysis

- Effect analysis of estimation error model
- Tie up these results with prediction error in previous paper

### 7.2. Component Analysis

- Effect analysis of number of component model
- Tie up these results in previous paper

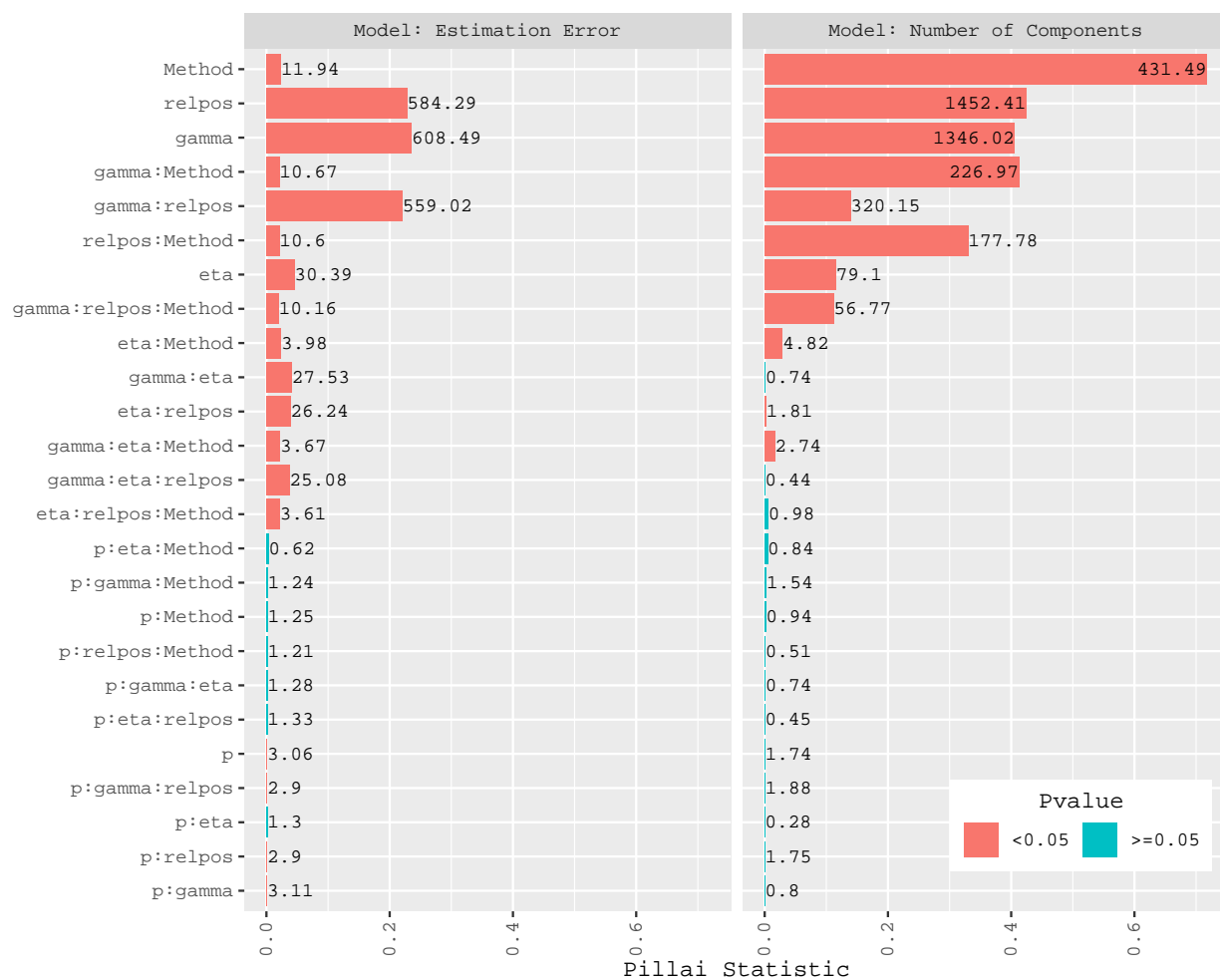


Figure 5: Pillai Statistic and F-value for the MANOVA model. The bar represents the Pillai Statistic and the text labels are F-value for corresponding factor.

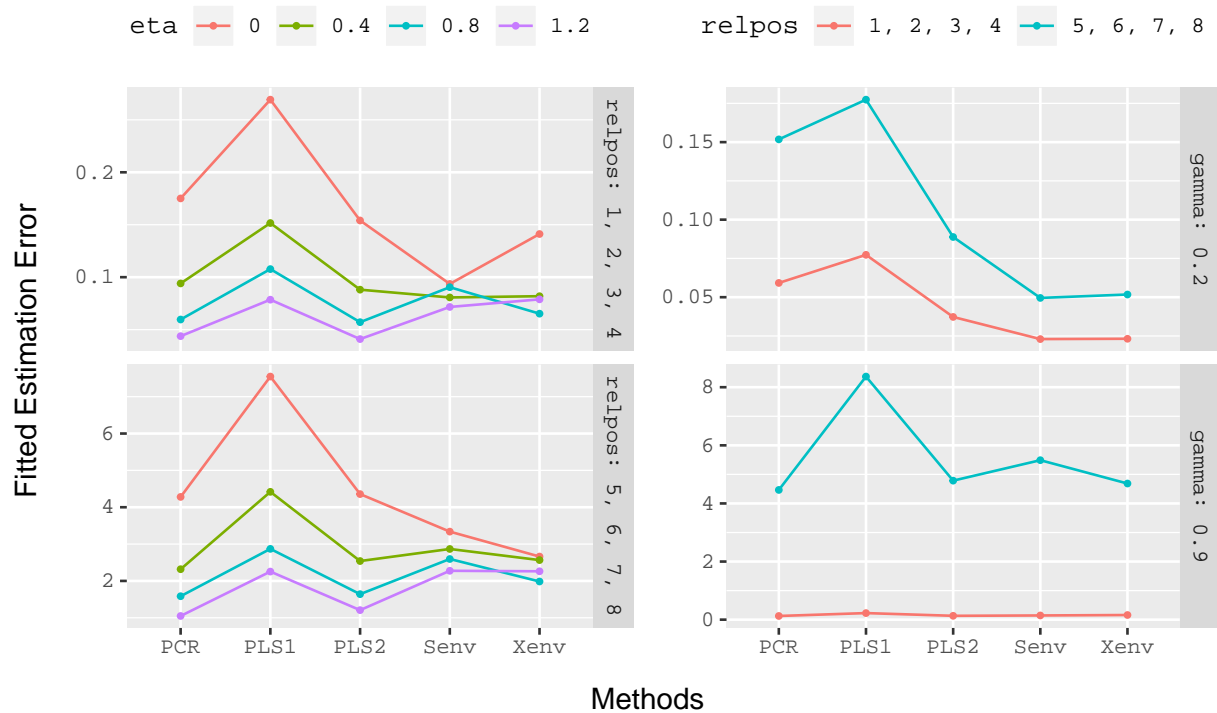


Figure 6: Effect plot of some interactions of the multivariate linear model of estimation error

## 8. Discussion and Conclusion

- A similar discussion but based more on why the methods worked in the way we have seen in the results in previous sections
- Some concluding remarks and limitations (or a gate for further exploration)

## References

- Cook, R.D., Li, B., Chiaromonte, F., 2010. Envelope Models for Parsimonious and Efficient Multivariate Linear Regression. *Statistica Sinica* 20, 927–1010.
- Cook, R.D., Zhang, X., 2015. Simultaneous envelopes for multivariate linear regression. *Technometrics* 57, 11–25. doi:[doi:10.1080/00401706.2013.872700](https://doi.org/10.1080/00401706.2013.872700).
- Helland, I.S., Almøy, T., 1994. Comparison of prediction methods when only a few components are relevant. *Journal of the American Statistical Association* 89, 583–591. doi:[doi:10.1080/01621459.1994.10476783](https://doi.org/10.1080/01621459.1994.10476783).
- R Core Team, 2018. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing. Vienna, Austria. URL: <https://www.R-project.org/>.

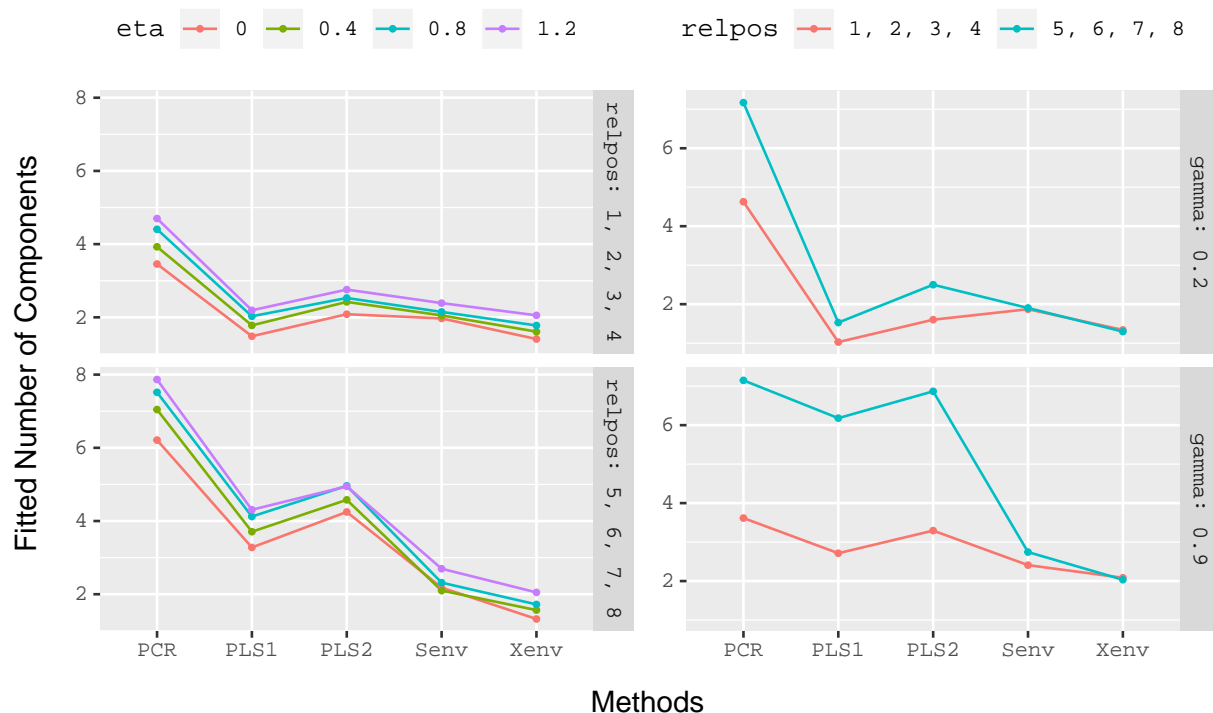


Figure 7: Effect plot of some interactions of the multivariate linear model of number of components to get minimum prediction error

Rimal, R., Almøy, T., Sæbø, 2019. Comparison of multivariate prediction methods.

Rimal, R., Almøy, T., Sæbø, S., 2018. A tool for simulating multi-response linear model data. *Chemometrics and Intelligent Laboratory Systems* 176, 1–10. doi:[doi:10.1016/j.chemolab.2018.02.009](https://doi.org/10.1016/j.chemolab.2018.02.009).

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. doi:[doi:10.1016/j.chemolab.2015.05.012](https://doi.org/10.1016/j.chemolab.2015.05.012).