

A tool for simulating multi-response linear model data

Raju Rimal, Trygve Almøy, Solve Sæbø^{1,*}

Abstract

Data science is generating enormous amounts of data, and new and advanced analytical methods are constantly being developed to cope with the challenge of extracting information from such “big-data”. Researchers often use simulated data to assess and document the properties of these new methods, and in this paper we present `simrel-m`, which is a versatile and transparent tool for simulating linear model data with an extensive range of adjustable properties. The method is based on the concept of relevant components Helland and Almøy [11], which is equivalent to the envelope model by Dennis Cook. It is a multi-response extension of R-package `simrel` which is available in R-package repository CRAN, and as `simrel` the new approach is essentially based on random rotations of latent relevant components to obtain a predictor matrix \mathbf{X} , but in addition we introduce random rotations of latent components spanning a response space in order to obtain a multivariate response matrix \mathbf{Y} . The properties of the linear relation between \mathbf{X} and \mathbf{Y} are defined by a small set of input parameters which allow versatile and adjustable simulations. Sub-space rotations also allow for generating data suitable for testing variable selection methods in multi-response settings. The method is implemented as an update to the R-package `simrel`.

Keywords: `simrel` package in r, data simulation, linear model, `simrel-m`,

Introduction

Technological advancement has opened a door for complex and sophisticated scientific experiments that was not possible before. Due to this change, enormous amounts of raw data are generated which contains massive information but difficult to excavate. Finding information and performing scientific research on these raw data has now become another problem. In order to tackle this situation new methods are

^{*}Dep. of Chemistry and Food Science, NMBU, Ås (nmbu.no)

being developed. However, before implementing any method, it is essential to test its performance and explore its properties. Often, researchers use simulated data for the purpose which itself is a time-consuming process. The main focus of this paper is to present a simulation method, along with an r-package called `simrel-m`, that is versatile in nature and yet simple to use.

The simulation method we are presenting here is based on the principle of relevant space for prediction [11] **which assumes that there exists a subspace in the complete space of response variables that is spanned by a subset of eigenvectors of predictor variables**. Our r-package based on this principle lets the user specify various population properties such as which latent components are relevant for a latent subspace of the responses \mathbf{y} and the collinearity structure of \mathbf{x} . This enables the possibility to construct data for evaluating estimation methods and methods developed for variable selection.

Among several publications on simulation, Ripley [15] and Gamerman and Lopes [6] has exhaustively discussed the topic. In addition, many publications have implemented simulated data in order to investigate new estimation methods and prediction strategies [see: 4, 5, 13]. However, most of the simulations in these studies were developed to address their specific problem. A systematic tool for simulating linear model data with single response, which could serve as a general tool for all such comparisons, was presented in Sæbø et al. [17] and as the r-package `simrel`. This paper extends `simrel` in order to simulate linear model data with multivariate response

with an r-package `simrel-m`.

Statistical Model

In this section we describe the model and the model parameterization which is assumed throughout this paper. We assume:

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

where, \mathbf{y} is a response vector with m response variables y_1, y_2, \dots, y_m with mean vector $\boldsymbol{\mu}_y$, and \mathbf{x} is vector of p predictor variables with mean vector $\boldsymbol{\mu}_x$. Further,

$\boldsymbol{\Sigma}_{yy}$ is the variance-covariance matrix of \mathbf{y}

$\boldsymbol{\Sigma}_{xx}$ is the variance-covariance matrix of variables \mathbf{x}

$\boldsymbol{\Sigma}_{xy}$ is the matrix of covariance between \mathbf{x} and \mathbf{y}

Standard theory in multivariate statistics may be used to show that \mathbf{y} conditioned on \mathbf{x} corresponds to the linear model,

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

where, $\boldsymbol{\beta}^t$ is a matrix of regression coefficient and $\boldsymbol{\varepsilon}$ is an error term such that $\boldsymbol{\varepsilon} \sim$

$N(0, \Sigma_{y|x})$. The properties of the linear model (2) can be expressed in terms of covariance matrices in (1).

Regression Coefficients The vector of regression coefficients is given by

$$\beta = \Sigma_{xx}^{-1} \Sigma_{xy}$$

Coefficient of Determination The diagonal elements of the coefficient of determination matrix ρ_y^2 gives the amount of variation in each response variable that is explained by \mathbf{x} .

$$\rho_y^2 = \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \Sigma_{yy}^{-1}$$

Conditional variance The conditional variance covariance matrix of \mathbf{y} given \mathbf{x} is,

$$\Sigma_{y|x} = \Sigma_{yy} - \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}.$$

The diagonal elements of this matrix equals the theoretical least square error of prediction for each of the response variables.

Let us define a transformation of \mathbf{x} and \mathbf{y} as, $\mathbf{z} = \mathbf{R}\mathbf{x}$ and $\mathbf{w} = \mathbf{Q}\mathbf{y}$. Here, $\mathbf{R}_{p \times p}$ and $\mathbf{Q}_{m \times m}$ are rotation matrices that rotate \mathbf{x} and \mathbf{y} to yield \mathbf{z} and \mathbf{w} , respectively. The model (1) can be re-expressed in terms of these transformed variables as:

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{z} \end{bmatrix} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (3)$$

$$\begin{aligned} &= N \left(\begin{bmatrix} \boldsymbol{\mu}_w \\ \boldsymbol{\mu}_z \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{ww} & \boldsymbol{\Sigma}_{wz} \\ \boldsymbol{\Sigma}_{zw} & \boldsymbol{\Sigma}_{zz} \end{bmatrix} \right) \\ &= N \left(\begin{bmatrix} \mathbf{Q}\boldsymbol{\mu}_y \\ \mathbf{R}\boldsymbol{\mu}_x \end{bmatrix}, \begin{bmatrix} \mathbf{Q}\boldsymbol{\Sigma}_{yy}\mathbf{Q}^t & \mathbf{Q}\boldsymbol{\Sigma}_{yx}\mathbf{R}^t \\ \mathbf{R}\boldsymbol{\Sigma}_{xy}\mathbf{Q}^t & \mathbf{R}\boldsymbol{\Sigma}_{xx}\mathbf{R}^t \end{bmatrix} \right) \end{aligned} \quad (4)$$

In addition, a linear model relating \mathbf{w} and \mathbf{z} can be written as,

$$\mathbf{w} = \boldsymbol{\mu}_w + \boldsymbol{\alpha}^t (\mathbf{z} - \boldsymbol{\mu}_z) + \boldsymbol{\tau} \quad (5)$$

where $\boldsymbol{\alpha}$ is the regression coefficient vector for the transformed model and $\boldsymbol{\tau} \sim N(\mathbf{0}, \boldsymbol{\Sigma}_{w|z})$. Further, if both \mathbf{Q} and \mathbf{R} are orthonormal matrices such that $\mathbf{Q}^t \mathbf{Q} = \mathbf{I}_q$ and $\mathbf{R}^t \mathbf{R} = \mathbf{I}_p$, the inverse transformation can be defined as,

$$\begin{aligned} \boldsymbol{\Sigma}_{yy} &= \mathbf{Q}^t \boldsymbol{\Sigma}_{ww} \mathbf{Q} & \boldsymbol{\Sigma}_{yx} &= \mathbf{Q}^t \boldsymbol{\Sigma}_{wz} \mathbf{R} \\ \boldsymbol{\Sigma}_{xy} &= \mathbf{R}^t \boldsymbol{\Sigma}_{zw} \mathbf{Q} & \boldsymbol{\Sigma}_{xx} &= \mathbf{R}^t \boldsymbol{\Sigma}_{zz} \mathbf{R} \end{aligned} \quad (6)$$

From this, we can find a direct connection between different population properties of (2) and (5).

Regression Coefficients

$$\begin{aligned}\alpha &= \Sigma_{wz} \Sigma_{zz}^{-1} = Q \Sigma_{YZ} \mathbf{R}^t [R \Sigma_{xx} \mathbf{R}^t]^{-1} \\ &= Q \left[\Sigma_{yx} \Sigma_{xx}^{-1} \right] \mathbf{R}^t = Q \beta \mathbf{R}^t\end{aligned}$$

Conditional Variance Further, the conditional variance covariance matrix of \mathbf{w} given \mathbf{z} is,

$$\begin{aligned}\Sigma_{w|z} &= Q \Sigma_{yy} \mathbf{Q}^t - Q \Sigma_{yx} \mathbf{R}^t [R \Sigma_{xx} \mathbf{R}^t]^{-1} R \Sigma_{xy} \mathbf{Q}^t \\ &= Q \Sigma_{yy} \mathbf{Q}^t - Q \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \mathbf{Q}^t \\ &= Q \left[\Sigma_{yy} - \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \right] \mathbf{Q}^t \\ &= Q \Sigma_{y|x} \mathbf{Q}^t\end{aligned}$$

Coefficient of Determination The coefficient of determination matrix (5) is,

$$\begin{aligned}\rho_w^2 &= \Sigma_{wz} \Sigma_{zz}^{-1} \Sigma_{zw} \Sigma_{ww}^{-1} \\ &= \mathbf{Q}^t \Sigma_{yx} \mathbf{R}^t (R \Sigma_{xx} \mathbf{R}^t)^{-1} R \Sigma_{xy} \mathbf{Q}^t (Q \Sigma_{yy}^{-1} \mathbf{Q}^t) \\ &= \mathbf{Q}^t \left[\Sigma_{yx} \Sigma_{xx} \Sigma_{xy} \Sigma_{yy}^{-1} \right] \mathbf{Q} \\ &= Q \rho_Y^2 \mathbf{Q}^t\end{aligned}$$

From the eigenvalue decomposition principle, if $\Sigma_{xx} = \mathbf{R} \Lambda \mathbf{R}^t$ and $\Sigma_{yy} = \mathbf{Q} \Omega \mathbf{Q}^t$ then \mathbf{z} and \mathbf{w} can be interpreted as principal components of \mathbf{x} and \mathbf{y} respectively. In this paper, these principal components will be termed as *predictor components* and *response components* respectively. Here, Σ and Ω are diagonal matrices of eigenvalues of Σ_{xx} and Σ_{yy} , respectively.

Relevant Components

Consider a single response linear model with p predictors.

$$y = \mu_y + \boldsymbol{\beta}^t (\mathbf{x} - \mu_x) + \epsilon$$

where, $\epsilon \sim N(0, \sigma^2)$ and \mathbf{x} are random predictors. Following the principal of relevant space and irrelevant space which are discussed extensively in Helland and Almøy [11], Helland [10], Helland et al. [13], Cook et al. [5], Sæbø et al. [17] and Helland et al. [12], we can assume that there exists a subspace of the full predictor space which is relevant for y . An orthogonal space to this space does not contain any information about y and is considered as irrelevant. Here, the y –relevant subspace of \mathbf{x} is spanned by a subset of the principal components defined by the eigenvectors of the covariance matrix of \mathbf{x} , i.e. $\boldsymbol{\Sigma}_{xx}$.

This concept can be extended to m responses so that the subspace of \mathbf{x} is relevant for a subspace of \mathbf{y} . This corresponds to the concept of simultaneous envelopes [4] where relevant (material) and irrelevant (immaterial) space were discussed for both response and predictor variables.

Model Parameterization

In order to construct a fully specified and unrestricted covariance matrix of \mathbf{z} and \mathbf{w} for the model in equation (4), we need to identify $1/2(p + m)(p + m + 1)$ unknown

parameters. However, Σ_{zz} is a diagonal matrix of eigenvalues of Σ_{xx} , hence the number of parameters may be reduced by $1/2p(p-1)$ parameters. In addition, the covariances of \mathbf{z} that are not relevant for \mathbf{w} will be zero. This further reduces the number of unknown parameters. For the purpose of simulation, we implement some assumptions to re-parameterize and simplify the model. This enables us to construct a wide range of model properties from only few key parameters.

Parameterization of Σ_{zz} If we let the rotation matrix \mathbf{R} correspond to the eigenvectors of Σ_{xx} , then \mathbf{z} becomes the set of principal components of \mathbf{x} . In that case Σ_{zz} is a diagonal matrix with eigenvalues $\lambda_1, \dots, \lambda_p$. Further, we adopt the same parametric representation as Sæbø et al. [17] for these eigenvalues:

$$\lambda_j = e^{-\gamma(j-1)}, \gamma > 0 \text{ and } j = 1, 2, \dots, p$$

Here, as γ increases, the decline of eigenvalues becomes steeper, hence the parameter γ controls the level of multicollinearity in \mathbf{x} .

Parameterization of Σ_{ww} Here we assume that \mathbf{w}' 's are independent and mutinormal distributed with variances 1, hence $\Sigma_{ww} = \mathbf{I}_m$.

Parameterization of Σ_{zw} After parameterization of Σ_{zz} and Σ_{ww} , we are left with $m \times p$ number of unknowns corresponding to Σ_{zw} . Some of the elements of Σ_{zw} may be equal to zero, which implies that the given \mathbf{z} is irrelevant for the given variable \mathbf{w} . The non-zero elements define which of the \mathbf{z} are relevant for \mathbf{w} . We

typically refer to the indices of these z variables as the positions of relevant components. In order to re-parameterize this covariance matrix, it is necessary to discuss the position of relevant components in detail.

Position of relevant components

Let k_1 components be relevant for w_1 , k_2 components be relevant for w_2 and so on. Let the positions of these components be given by the index sets $\mathcal{P}_1, \mathcal{P}_2, \dots, \mathcal{P}_m$ respectively. Further, the covariance between w_j and z_i is non-zero only if z_i is relevant for w_j . If σ_{ij} is the covariance between w_j and z_i then $\sigma_{ij} \neq 0$ if $i \in \mathcal{P}_j$ where $i = 1, \dots, p$ and $j = 1, \dots, m$ and $\sigma_{ij} = 0$ otherwise.

In addition, the true regression coefficients for w_j (ref:(5)) is given by:

$$\alpha_j = \Lambda^{-1} \sigma_{ij} = \sum_{i \in \mathcal{P}_j} \frac{\sigma_{ij}}{\lambda_i}, \quad j = 1, 2, \dots, m$$

The positions of the relevant components have heavy impact on prediction. Helland and Almøy [11] have shown that if the relevant components have large eigenvalues (variances), which here implies small index values in \mathcal{P}_j , prediction of \mathbf{y} from \mathbf{x} is relatively easy and if the eigenvalues (variances) of relevant components are small, the prediction becomes difficult, given that the coefficient of determination and other model parameters are held constant. For example, if the first and second components, z_1 and z_2 , are relevant for w_1 and fifth and sixth components, z_5 and z_6 , are relevant

for w_2 , it is relatively easier to predict w_1 than w_2 , other properties being similar. This is so, because the first and second principal components have larger variances than the fifth and sixth components.

Although the covariance matrix may depend on few relevant components, we can not choose these covariances freely since we also need to satisfy following two conditions:

- The covariance matrices Σ_{zz} , Σ_{ww} and Σ must be positive definite
- The covariance σ_{ij} must satisfy user defined coefficient of determination

We have the relation,

$$\rho_w^2 = \Sigma_{zw}^t \Sigma_{zz}^{-1} \Sigma_{zw} \Sigma_{ww}^{-1}$$

Applying our above given assumptions that, $\Sigma_{ww} = \mathbf{I}_m$ and $\Sigma_{zz} = \Lambda$, we obtain,

$$\begin{aligned} \rho_w^2 &= \Sigma_{zw}^t \Lambda^{-1} \Sigma_{zw} \mathbf{I}_m \\ &= \begin{bmatrix} \sum_{i=1}^p \sigma_{i1}^2 / \lambda_i & \dots & \sum_{i=1}^p \sigma_{i1} \sigma_{im} / \lambda_i \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^p \sigma_{i1} \sigma_{im} / \lambda_i & \dots & \sum_{i=1}^p \sigma_{im}^2 / \lambda_i \end{bmatrix} \end{aligned}$$

Furthermore, we assume that there are no overlapping relevant components for any two w , i.e, $\mathcal{P}_j \cap \mathcal{P}_{j*} = \emptyset$ or $\sigma_{ij} \sigma_{ij*} = 0$ for $j \neq j*$. The additional unknown parameters in the diagonal of ρ_w^2 should agree with user specified coefficients of determination for \mathbf{w} . i.e, ρ_{wj}^2 is,

$$\rho_{wj}^2 = \sum_{i=1}^p \frac{\sigma_{ij}^2}{\lambda_i}$$

Here, only the relevant components have non-zero covariances with w_j , so,

$$\rho_{wj}^2 = \sum_{i \in \mathcal{P}_j} \frac{\sigma_{ij}^2}{\lambda_i}$$

For some user defined ρ_{jw}^2 the σ_{ij}^2 is determined as follows,

1. Sample k_j values from a uniform distribution $\mathcal{U}(-1, 1)$ distribution. Let them be denoted $\mathcal{S}_{\mathcal{P}_1}, \dots, \mathcal{S}_{\mathcal{P}_{k_j}}$.

2. Define,

$$\sigma_{ij} = \text{Sign}(\mathcal{S}_i) \sqrt{\frac{\rho_{wj}^2 |\mathcal{S}_i|}{\sum_{k \in \mathcal{P}_j} |\mathcal{S}_k|}} \lambda_i$$

for $i \in \mathcal{P}_j$ and $j = 1, \dots, m$

Data Simulation

From the above given parameterizations and the user defined choices of model parameters, a fully defined and known covariance matrix $\mathbf{\Sigma}$ of (\mathbf{w}, \mathbf{z}) is given. For the simulation of a single observation of (\mathbf{w}, \mathbf{z}) let us define $\mathbf{g} = \mathbf{\Sigma}^{-1/2} \mathbf{u}$ such that $\text{cov}(\mathbf{g}) = \mathbf{\Sigma}$. Here $\mathbf{\Sigma}^{-1/2}$ is obtained from Choleskey decomposition of $\mathbf{\Sigma}$ and \mathbf{u} is simulated from standard normal distribution and has covariance $\text{cov}(\mathbf{u}) = \mathbf{I}$.

Similarly, in order to simulate n observations, we define $\mathbf{G}_{n \times (m+p)} = \mathbf{U}\mathbf{\Sigma}^{-1/2}$. Here the first m columns of \mathbf{G} will serve as \mathbf{W} and remaining p columns will serve as \mathbf{Z} . Further, each row of \mathbf{G} will be a vector sampled independently from the joint normal distribution of (\mathbf{w}, \mathbf{z}) . Finally, these simulated matrices \mathbf{W} and \mathbf{Z} are orthogonally rotated in order to obtain \mathbf{Y} and \mathbf{X} respectively. The following section discuss about these rotation matrices in detail.

Rotation of predictor space

Initially, let us consider an example where a regression model with $p = 10$ predictors (\mathbf{x}) and $m = 4$ responses (\mathbf{y}). Let's assume that only three response components (w_1, w_2 and w_3) are needed to describe all four response variables. Further, let the index sets $\mathcal{P}_1 = \{1, 2\}, \mathcal{P}_2 = \{3, 4\}$ and $\mathcal{P}_3 = \{5, 6\}$ define the positions of the predictor components of \mathbf{x} that are relevant for w_1, w_2 and w_3 respectively. Let $\mathcal{S}_1, \mathcal{S}_2$ and \mathcal{S}_3 be the orthogonal spaces spanned by each set of predictor components. These spaces together span $\mathcal{S}_k = \mathcal{S}_1 \oplus \mathcal{S}_2 \oplus \mathcal{S}_3$ which is the minimum relevant space and equivalent to the \mathbf{x} -envelope as discussed by Cook et al. [5].

Moreover, let $q_1 = 3, q_2 = 3$ and $q_3 = 2$ be the number of predictor variables we want to be relevant for w_1, w_2 and w_3 respectively. Then $q_1 = 3$ predictors may be obtained by rotating the predictor components in \mathcal{P}_1 along with one more irrelevant component. Similarly, $q_2 = 3$ predictors, relevant for w_2 , can be obtained by rotating predictor components in \mathcal{P}_2 along with one more irrelevant component and finally,

$q_3 = 2$ predictors, relevant for w_3 , can be obtained by rotating the components in \mathcal{P}_3 without any additional irrelevant component. Let the space spanned by the q_1, q_2 and q_3 number of predictors be $\mathcal{S}_{q_1}, \mathcal{S}_{q_2}$ and \mathcal{S}_{q_3} . Together they span a space $\mathcal{S}_q = \mathcal{S}_{q_1} \oplus \mathcal{S}_{q_2} \oplus \mathcal{S}_{q_3}$. This space is bigger than \mathcal{S}_k since in the process two irrelevant components were included in the rotations. Here, \mathcal{S}_k is orthogonal to \mathcal{S}_{p-k} and \mathcal{S}_q is orthogonal to \mathcal{S}_{p-q} . Generally speaking, here we are splitting the complete variable space \mathcal{S}_p into two orthogonal spaces – \mathcal{S}_k relevant for \mathbf{w} and \mathcal{S}_{p-k} irrelevant for \mathbf{w} .

In the previous section, we discussed about the construction of a covariance matrix for the latent structure. Figure~1 (left) shows a similar structure resembling the example here. The three colors represent the relevance with the three latent response components (w_1, w_2 and w_3). Here we can see that z_1 and z_2 (first and second predictor components of \mathbf{x}) have non-zero covariance with w_1 (first latent component of response \mathbf{y}). In the similar manner other non-zero covariances are self-explanatory.

In order to simulate predictor variables (\mathbf{x}), we construct matrix \mathbf{R} which then is used for orthogonal rotation of the predictor components \mathbf{z} . This defines a new basis for the same space as is spanned by the predictor components. In principle, there are many possible options for defining a rotation matrix. Among them, the eigenvector matrix of $\mathbf{\Sigma}_{xx}$ can be a candidate. However, in this reverse engineering approach both rotation matrices \mathbf{R} and \mathbf{Q} along with the covariance matrices $\mathbf{\Sigma}_{xx}$ are unknown. So, we are free to choose any \mathbf{R} that satisfied the properties of a real valued rotation matrix, i.e $\mathbf{R}^{-1} = \mathbf{R}^t$ so that \mathbf{R} is orthonormal and its determinant becomes

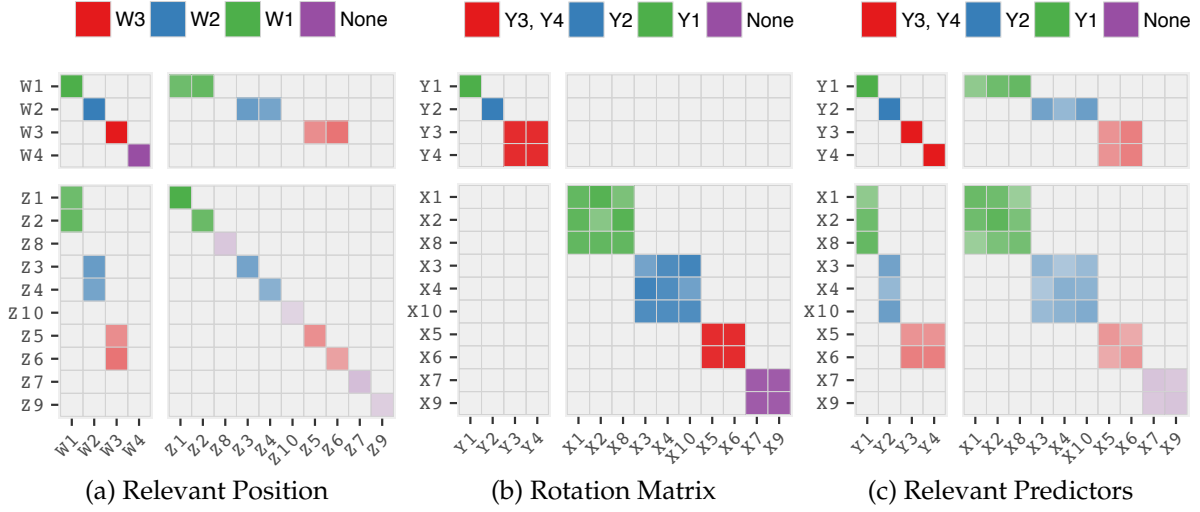


Figure 1: Simulation of predictor and response variables after orthogonal transformation of predictor and response components by rotation matrices Q and R shown as the upper left and the lower right block matrices in (b).

± 1 . Here the rotation matrix \mathbf{R} should be block diagonal as in figure~1 (middle) in order to rotate spaces $\mathcal{S}_1, \mathcal{S}_2 \dots$ separately. Figure~2 (left) shows the simulated predictor components \mathbf{z} that we are following in our example where we can see that the components z_1 and z_2 (relevant for w_1) is getting rotated together with an irrelevant component z_8 . The resultant predictors (Figure 2, right) x_1, x_2 and x_8 will hence also be relevant for w_1 . In the figure, we can see that components z_7, z_8, z_9 and z_{10} are not relevant for any responses before rotation, however, the x_8, x_{10} predictors become relevant after rotation keeping x_7 and x_9 still irrelevant.

Among several methods [1, 9] for generating random orthogonal matrix, in this paper we are using orthogonal matrix \mathcal{Q} obtained from QR-decomposition of a matrix filled with standard normal variates. The rotation here can be a) restricted and b) unrestricted. The latter rotates all components \mathbf{z} together and makes all predictor

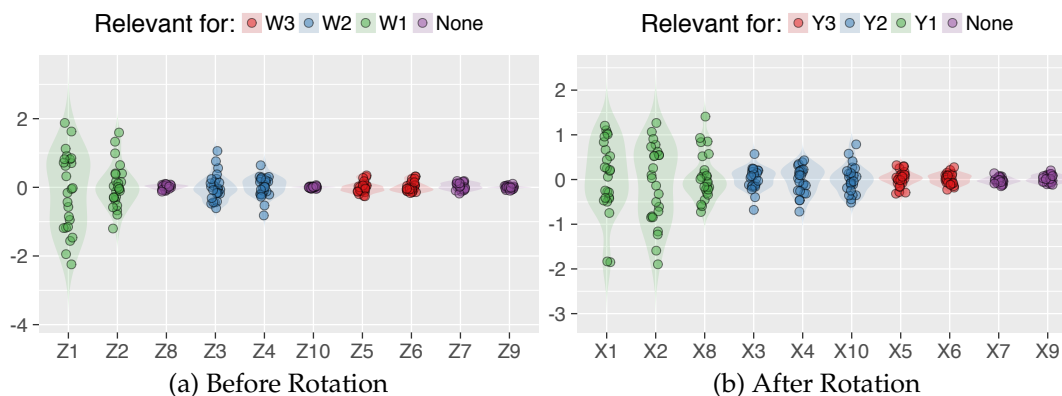


Figure 2: Simulated Data before

variables somewhat relevant for all response components. However, the former performs a block-wise rotation so that it rotates certain selected predictor components together. This gives control for specifying certain predictors as relevant for selected responses, which was discussed in our example above. This also allows us to simulate irrelevant predictors such as x_7 and x_9 which can be detected during variable selection procedures.

Rotation of response space

The previous example has four response variables with only three informative components w_1, w_2 and w_3 . During the rotation procedure, the response space is also rotated along with the predictor space. Figure 1 shows that the informative response component w_3 is rotated together with the uninformative response component w_4 so that the predictors which were relevant for w_3 will be relevant for response variables y_3 and y_4 . Similarly, response components w_1 and w_2 are rotated separately so that predictors relevant for w_1 and w_2 will only be relevant for y_1 and y_2 respectively, which we can

see in Figure~2. In the r-package *simrel-m*, the combining of the response components is specified by a parameter *ypos*.

Implementation

This section demonstrates an application of *simrel-m* in order to compare different estimation methods on the basis of prediction error. For the comparison, we have considered four well established estimation methods.

- a) Ordinary Least Squares (OLS),
- b) Principal Component Regression (PCR),
- c) Partial Least Squares predicting individual response variable separately (PLS1)
and
- d) Partial Least Squares predicting all response variables together (PLS2).

We have also considered four relatively new estimation methods in multi-response regression:

- a) Canonically Powered Partial Least Squares regression (CPPLS) [14],
- b) Canonical Partial Least Squares regression (CPLS) [14],
- c) Envelope estimation in predictor space (*xenv*) [2],
- d) Envelope estimation in response space (*yenv*) [3] and
- e) Simultaneous estimation of *x*- and *y*-envelope (*senv*) [4]

From the possible combinations of two levels of coefficient of determination (R^2) and two levels of gamma (the factor that controls the multicollinearity in predictor variables), four simulation designs (design 1-4) were prepared. Replicating each design 20 times, 80 datasets with five response variables ($m = 5$) and 16 predictor variables ($p = 16$) were simulated using the method discussed in this paper. It was also assumed that three principle components of response variables (w_1, w_2 and w_3) completely describe the variation present in five response variables ($y_1 \dots y_5$). The four designs are presented in the Table~1. All datasets contained 100 sampled observations and out of 16 predictor variables, three disjoint set of five predictor variables are relevant for response components w_1, w_2 and w_3 . Further, predictor components z_1 and z_6 were relevant for response component w_1 , predictor components z_2 and z_5 were relevant for response component w_2 and predictor component z_3 was relevant for response component w_3 . In addition, following the discussion about rotation of response space, w_1 was rotated together with w_4 and w_2 was rotated together with w_5 .

Table 1: Parameter setting of simulated data for model comparison

	Design1	Design2	Design3	Design4
Decay of eigenvalues (γ)	0.2	0.8	0.2	0.8
Coef. of Determination $(\rho_{w_j}^2)$	0.8, 0.8, 0.4	0.8, 0.8, 0.4	0.4, 0.4, 0.4	0.4, 0.4, 0.4

For each method, an estimate of test prediction error was computed as,

$$\alpha_{m \times m} = (\hat{\beta} - \beta)^t \Sigma_{xx} (\hat{\beta} - \beta) + \Sigma_{Y|X}$$

where, $\hat{\beta}$ is an estimate of true regression coefficient β and Σ_{xx} is the true covariance structure of the predictor variables obtained from `simrel-m`. Also, $\Sigma_{Y|X}$ is the true minimum error of the model. Here $\hat{\beta}$ varies accross different estimation methods while the remaining terms are same for each dataset design. Further, an overall prediction error of all responses is measured by the trace of α .

The minimum prediction error (measured as discussed above) for nine estimation methods averaged over 20 replications of four designs are shown in Table 2. The table also gives the number of predictor components a method has used in order to obtain the minimum of average prediciton error.

Table 2: Minimum average prediction error (number of components corresponding to minimum prediction error, minimum prediction error) (For Yenv, the number of response components is given)

Model	Design: 1	Design: 2	Design: 3	Design: 4
<i>CPLS</i>	(3, 3.24)	(4, 3.22)	(3, 4.09)	(3, 4.05)
<i>CPPLS</i>	(3, 3.21)	(3, 3.17)	(3, 4.11)	(3, 4.04)
<i>OLS</i>	(1, 3.6)	(1, 3.58)	(1, 4.57)	(1, 4.5)
<i>PCR</i>	(7, 3.28)	(6, 3.19)	(6, 4.08)	(6, 4.04)
<i>PLS</i>	(5, 3.29)	(6, 3.19)	(3, 4.11)	(6, 4.06)

Model	Design: 1	Design: 2	Design: 3	Design: 4
<i>PLS1</i>	(2, 3.32)	(5, 3.2)	(1, 4.16)	(5, 4.07)
<i>Senv</i>	(4, 3.17)	(5, 3.14)	(3, 4.35)	(5, 4.28)
<i>Xenv</i>	(5, 3.23)	(6, 3.2)	(5, 4.1)	(6, 4.11)
<i>Yenv</i>	(3, 3.24)	(3, 3.23)	(3, 4.29)	(3, 4.24)

Table 2 shows that the simultaneous envelope has prediction error of 3.17 and 3.14 in design 1 (with 4 components) and design 2 (with 5 components), respectively, which is smaller than other methods. However the method was not able to show the same performance in design 3 and design 4. The PCR model has the smallest prediction error (4.08) from 6 components in design 3 and Canonically Powered PLS has minimum prediction error (4.04) from 3 components in design 4. In design 3, we can also see that the Canonical PLS method has second best performance with only three components. The number of components vary accross different replicated dataset but the component corresponding to minimum prediction error is discussed here. A detailed picture of prediction error for each estimation method obtained for each additional component is shown in Figure 3. Although designs 2 and 4 have higher levels of multicollinearity, the performance of the estimation methods is indifferent to its effect. Since all the methods, except OLS, are based on shrinking of estimates, they are less influenced by the multicollinearity problem.

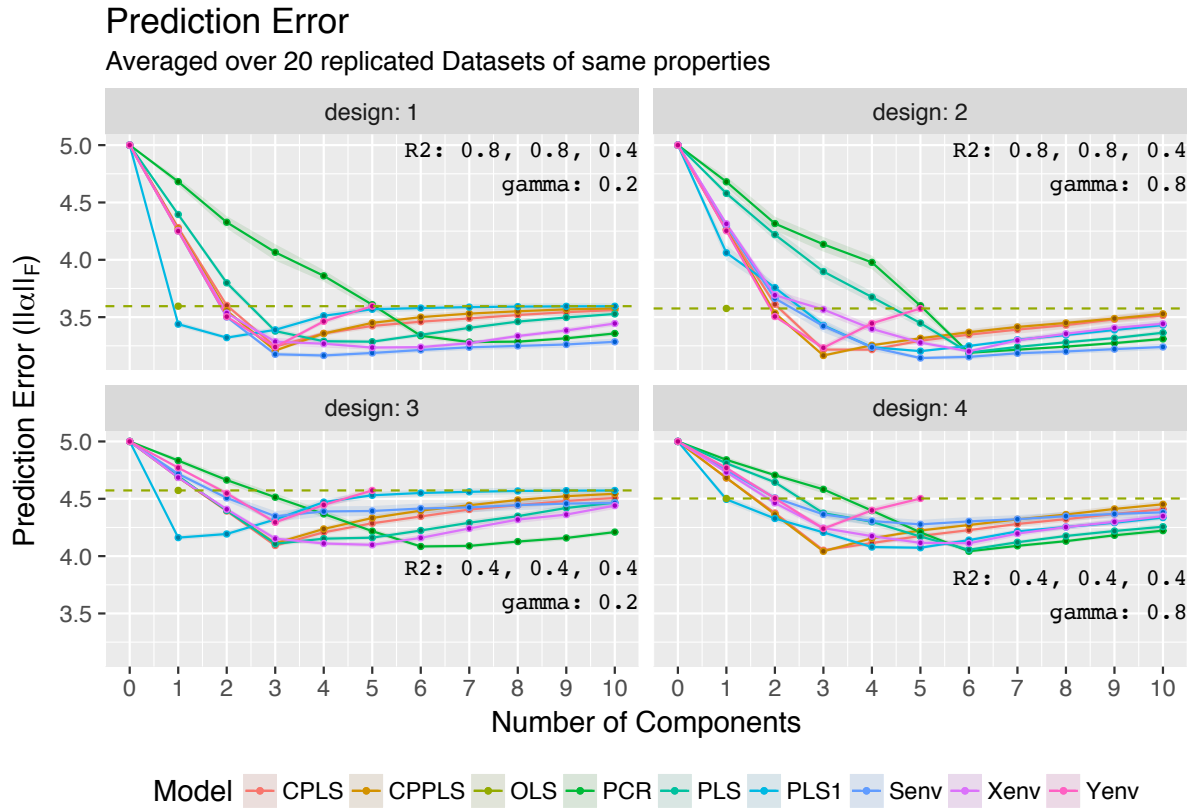


Figure 3: Minimum of Average Prediction Error

The analysis presented in Figure 3 has addressed some questions such as how methods works when there exist a true reduced dimension in response space, but also arised other questions like why they perform differently. For example, what is the reason for the decreasing relative performance of the simultaneous envelope method as the ρ^2 values are reduced? Does this depend on the dimensions and spheric shape of the y envelopes? Since, the example is merely intended as a demonstration of how `simrel-m` can be used in scientific study, a more elaborative study would be necessary to answer such questions, but for this purpose `simrel-m` would be a powerful tool.

Web Interface

In order to give an alternative interface for `simrel-m`, we have created a shiny app which allows users to input the simulation parameters through different input fields. Figure 4 shows a screenshot of the application. The application contains three main sections through which the user can interact with this simulation approach. A random seed can be selected using section Figure 4 (a) so that a particular set of data can be re-simulated if needed. Figure 4 (b) has all the input panels where the user dependent parameters for simulation can be entered. Here the user also has the option to simulate univariate (uses `simrel` package in CRAN), bivariate (not yet available in CRAN) and multivariate simulation (`simrel-m`). In addition, a simulated R-object is comprised by the simulated data can be download as Rdata format (section (e) in Figure 4). The object holds the simulated data along with other properties such as coefficient of determination for each response, true regression coefficients and rotation matrices.

All `simrel-m` parameters can be entered using a simple user interface where a vectors are separated with comma(,) and lists are separated with semicolon(;). For instance, the relevant position discussed in the implementation section of this paper can be entered as 1, 6; 2, 5; 3, 4 which is equivalent to R syntax `list(c(1, 6), c(2, 5), c(3, 4))`.

The application not only allows users to simulate data, but also gives some insight into simulated data properties. The example used in the screenshot has simulated

200 training and 50 test samples with 15 predictor and 4 response variables. There are two latent variables (response components) that completely span the informative response space. Five predictor variables are relevant for the first response component and explains 80 percent of the variation in x . In addition, the first and second predictor components span the same space as spanned by these five relevant predictors. Similarly, another set of four predictor variables are relevant for second response component and explains 70 percent of the variation. Here, third, fourth and sixth predictor components span the same space as spanned by these four relevant predictors. Further, the first response component is rotated together with a normally distributed random vector to obtain first and third response variable and second response components are rotated together with another normally distributed random vector to obtain second and fourth response variable.

Section (c) in Figure 4 contains three plots – a) true regression coefficients b) relevant components and c) estimated relevant components. In the first plot we can see that predictor variables (1, 2, 8, 9 and 13) are relevant for the first and third response variable by their non-zero coefficients, whereas predictor variables (3, 4, 6 and 15) are relevant for the second and fourth response variable. The second plot shows the covariances between the response components and the predictor components along with the corresponding eigenvalues in the background (bar plot). In the plot the absolute value of the covariances after scaling with the largest covariance are shown. As in our parameter setting, the plot shows that first and second predictor components

have non-zero covariance with first response component and third, fourth and sixth predictor components have non-zero covariance with second response component. The third plot is the estimated covariance between predictor components and the response variables, for the simulated data. Since the first and third response components are rotated together, in the plot, the covariance between predictor components and first and third response variable is following similar pattern. This also suggests that the predictor components which were relevant for first response components gets relevant for first and third response variables after rotation.

Along with these main sections, section (d) in the figure contains additional analysis performed with the simulated data such as its estimation with different methods. This section is intended for educational purposes to show how changing the data properties influences the performances of different estimation and prediction methods.

Many scientific studies [4, 13, 16] are using simulated data in order to compare their findings with others or assess its properties. In many of these situations, a user-friendly and versatile simulation tool like `simrel-m` can play an important role. Gangsei et al. [7]; Gangsei et al. [8]; and Sæbø et al. [17] are some examples where the univariate and bivariate form of `simrel` have been used for such purposes.

References

- [1] Anderson TW, Olkin I, Underhill LG. Generation of random orthogonal matrices. *SIAM Journal on Scientific and Statistical Computing* 1987;8(4):625–629.
- [2] Cook RD, Li B, Chiaromonte F. Envelope models for parsimonious and efficient multivariate linear regression. *Statistica Sinica* 2010;p. 927–960.
- [3] Cook RD, Zhang X. Foundations for envelope models and methods. *Journal of the American Statistical Association* 2015;110(510):599–611.
- [4] Cook RD, Zhang X. Simultaneous envelopes for multivariate linear regression. *Technometrics* 2015;57(1):11–25.
- [5] Cook R, Helland I, Su Z. Envelopes and partial least squares regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 2013;75(5):851–877.
- [6] Gamerman D, Lopes HF. Markov chain Monte Carlo: stochastic simulation for Bayesian inference. CRC Press; 2006.
- [7] Gangsei LE, Almøy T, Sæbø S. Theoretical evaluation of prediction error in linear regression with a bivariate response variable containing missing data. *Communications in Statistics-Theory and Methods* 2016;(just-accepted).
- [8] Gangsei L, Almøy T, Sæbø S. Linear regression with bivariate response variable containing missing data. An empirical Bayes strategy to increase prediction

precision Submitted manuscript to Communications in Statistics–Simulation and Computation 2016;.

- [9] Heiberger RM. Algorithm AS 127: Generation of random orthogonal matrices. Journal of the Royal Statistical Society Series C (Applied Statistics) 1978;27(2):199–206.
- [10] Helland IS. Model Reduction for Prediction in Regression Models. Scandinavian Journal of Statistics 2000 Mar;27(1):1–20. <http://dx.doi.org/10.1111/1467-9469.00174>.
- [11] Helland IS, Almøy T. Comparison of prediction methods when only a few components are relevant. Journal of the American Statistical Association 1994;89(426):583–591.
- [12] Helland IS, Sæbø S, Almøy T, Rimal R. Model and estimators for partial least squares; 2017.
- [13] Helland IS, Sæbø S, Tjelmeland H, et al. Near optimal prediction from relevant components. Scandinavian Journal of Statistics 2012;39(4):695–713.
- [14] Indahl UG, Liland KH, Næs T. Canonical partial least squares—a unified PLS approach to classification and regression problems. Journal of Chemometrics 2009;23(9):495–504.
- [15] Ripley BD. Stochastic simulation, vol. 316. John Wiley & Sons; 2009.

- [16] Sæbø S, Almøy T, Flatberg A, Aastveit AH, Martens H. LPLS-regression: a method for prediction and classification under the influence of background information on predictor variables. *Chemometrics and Intelligent Laboratory Systems* 2008;91(2):121–132.
- [17] Sæbø S, Almøy T, Helland IS. simrel – A versatile tool for linear model data simulation based on the concept of a relevant subspace and relevant predictors. *Chemometrics and Intelligent Laboratory Systems* 2015;.

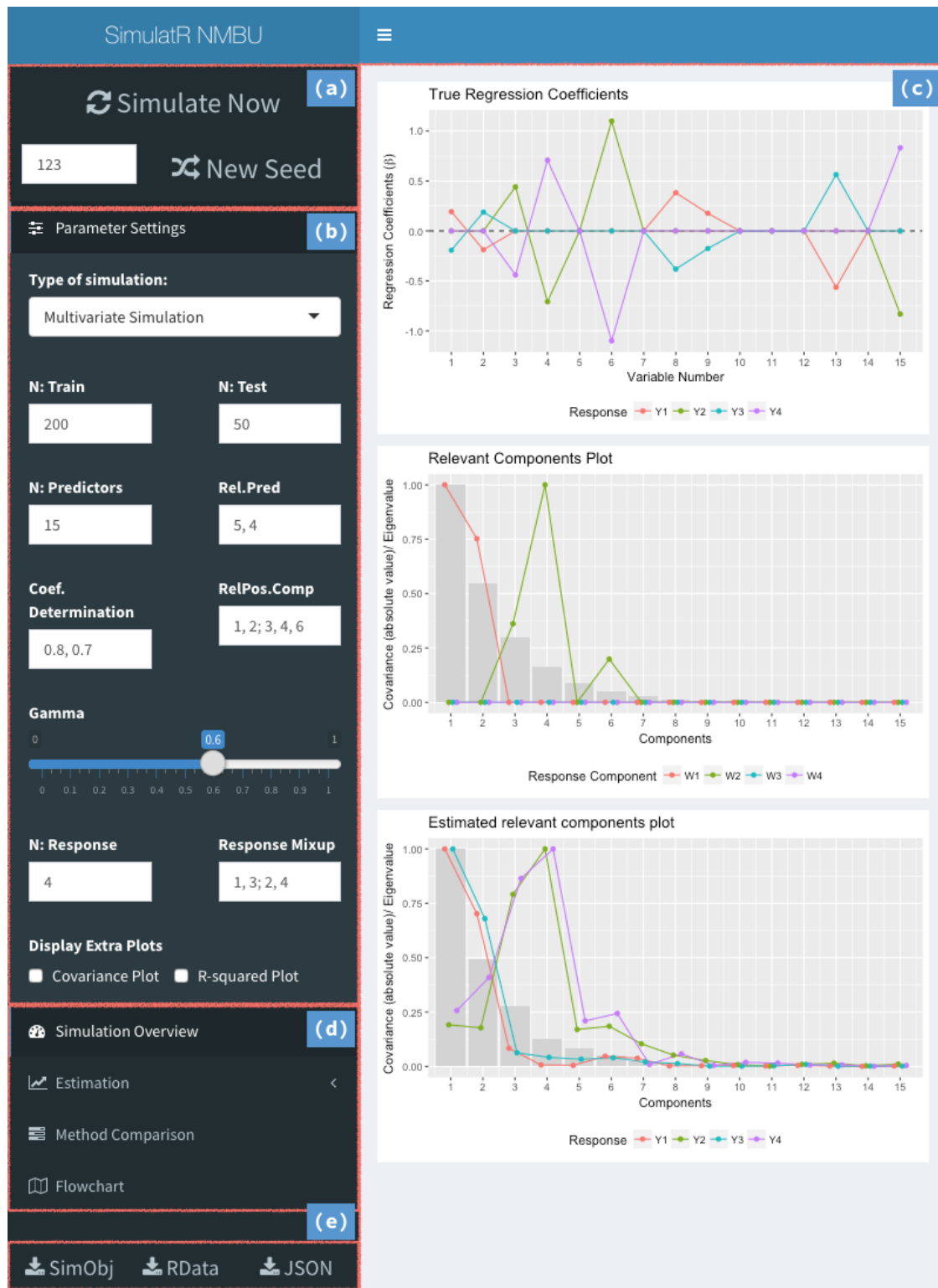


Figure 4: Application interface of 'simulatR'. (a) Seed and simulation button (b) Parameter control panel (c) Properties of simulated data (d) Additional analysis (e) Download option of simulated data