

**Comparison of Some Dimensionality  
Reduction Methods for Prediction**

**G.M. Merola and B. Abraham**  
*University of Waterloo*

**RR-99-02**  
January 1999

# COMPARISON OF SOME DIMENSIONALITY REDUCTION METHODS FOR PREDICTION

G.M. Merola and B. Abraham  
University of Waterloo

## Abstract

In this paper we briefly discuss the common dimensionality reduction methods such as principal components analysis, canonical correlation, partial least squares and reduced rank regression. We also consider some recent techniques such as maximum overall redundancy (MOR), weighted MOR and iteratively weighted reduced rank regression (IWRRR). Performance of these methods are compared by a simulation study.

**Keywords:** Canonical correlation regression, partial least squares, principal component regression, maximum overall redundancy, reduced rank regression, simulation.

## 1. Introduction

The growing use of computers in industry and the availability of inexpensive computer storage devices have created the need for dealing with large data sets in many industrial processes. Many processes such as chemical reactors are equipped with sensors connected to computers that can provide hundreds of measurements taken on many process variables ( $\mathbf{x}$ ) every few seconds and on output characteristics ( $\mathbf{y}$ ) sometimes less frequently. The variables are often highly correlated as well. The availability of the  $\mathbf{x}$  measurement can be used to monitor the process itself and as a diagnostic tool for causes of out-of-control values of the  $\mathbf{y}$ -variables. The structure of such data calls for an approach which looks for a lower dimensional subspace in which the process can be monitored and from which  $\mathbf{y}$  can be predicted.

In section 2, we outline the usual dimension reduction methods (DRMs). Section 3 gives a simulation study to compare the various methods. Section 4 gives some concluding remarks.

## 2. Dimensionality Reduction Methods (DRMs)

Let  $\mathbf{X}$  be an  $(n \times p)$  matrix containing  $n$  independent measurements on  $p$  process (explanatory) variables and  $\mathbf{Y}$  be an  $(n \times q)$  matrix of  $n$  independent observations on  $q$  response (quality) variables. We assume that columns of these matrices are mean centred. It is also common practice to scale the columns to unit length. It is assumed that the data follow the multivariate linear model

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E} \quad (2.1)$$

where  $\mathbf{B}$  is the  $(p \times q)$  matrix of regression coefficients and  $\mathbf{E}$  is a matrix of errors.

### 2.1 Principal Components Analysis (PCA)

The DRMs that we consider determine a set of  $d$  ( $\leq p$ ) orthogonal latent variables,  $\mathbf{T}$ , which form a subspace of  $L(\mathbf{X})$ , the space spanned by the columns of  $\mathbf{X}$ . Then the matrix  $\mathbf{X}$  can be represented as

$$\mathbf{X} = \mathbf{T}\mathbf{P} + \mathbf{F} \quad (2.2)$$

with the constraint that  $\mathbf{T}'\mathbf{F} = 0$ . If we take

$$\hat{\mathbf{P}} = (\mathbf{T}'\mathbf{T})^{-1}\mathbf{T}'\mathbf{X} \quad (2.3)$$

then this requirement is satisfied. Now take  $\mathbf{T} = \mathbf{X}\mathbf{A}$  where  $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_d)$  is a  $(p \times d)$  matrix of coefficients and consider

$$\|\mathbf{X} - \mathbf{T}\hat{\mathbf{P}}\|^2 = \|\mathbf{X} - \mathbf{X}\mathbf{A}(\mathbf{A}'\mathbf{X}'\mathbf{X}\mathbf{A})^{-1}\mathbf{A}'\mathbf{X}'\mathbf{X}\|^2. \quad (2.4)$$

Subject to the condition that the  $d$  latent variables are orthogonal ie.  $\mathbf{A}'\mathbf{X}'\mathbf{X}\mathbf{A} = \mathbf{I}$ , the expression (2.4) is minimized w.r.t. the matrix  $\mathbf{A}$  by the eigen-vectors of  $\mathbf{X}'\mathbf{X}$  corresponding

to the  $d$  largest eigen-values (see Rao (1964), Hotelling (1935)). These latent variables are referred to as the principal components (PC) of  $\mathbf{X}$ .

Principal Component Regression (PCR) utilizes an appropriate number of these PC's to predict  $\mathbf{y}$ . It should be noted that the PC's are not obtained using an optimality criterion for the prediction of  $\mathbf{y}$  but for the reconstruction of  $\mathbf{X}$  (prediction of  $\mathbf{x}$ ).

## 2.2 Reduced Rank Regression (RRR)

The linear model in the latent variables  $\mathbf{T}$  can be written as

$$\mathbf{Y} = \mathbf{T}\mathbf{Q} + \mathbf{E}^* \quad (2.5)$$

where  $\mathbf{Q}$  is a  $(d \times q)$  matrix of regression coefficients. Taking  $\mathbf{T} = \mathbf{X}\mathbf{A}$  as before

$$\mathbf{Y} = \mathbf{X}\mathbf{A}\mathbf{Q} + \mathbf{E}^* = \mathbf{X}\mathbf{M} + \mathbf{E}^*. \quad (2.6)$$

The linear relationship between the explanatory variables  $\mathbf{X}$  and the responses  $\mathbf{Y}$  is expressed by the  $(p \times d)$  matrix  $\mathbf{M} = \mathbf{A}\mathbf{Q}$  of rank  $d$ . Thus, the use of DRMs in prediction can be regarded as a regression with rank deficient matrix of coefficients. Model (2.5) is known as the reduced rank regression (RRR) model.

The residual sum of squares (RSS) for this model can be written as

$$\| \mathbf{Y} - \mathbf{X}\mathbf{A}(\mathbf{A}'\mathbf{X}'\mathbf{X}\mathbf{A})^{-1}\mathbf{A}'\mathbf{X}'\mathbf{Y} \|^2. \quad (2.7)$$

This is to be minimized with respect to  $\mathbf{A}$  subject to the condition  $\mathbf{A}'\mathbf{X}'\mathbf{X}\mathbf{A} = \mathbf{I}$ . It is known (see Rao (1964), Izenman (1975)) that the optimal latent vectors for this RRR model are given by the eigen-vectors associated with the  $d$  largest eigen-values obtained from

$$\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}\mathbf{Y}'\mathbf{T} = \mathbf{T}\mathbf{\Phi} \quad (2.8)$$

where  $\mathbf{\Phi}$  is the diagonal matrix of eigen-values. Suppose now that  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\mathbf{B}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$  is the ordinary least square estimate of  $\mathbf{Y}$ , then  $\hat{\mathbf{Y}}'\mathbf{T} = \mathbf{Y}'\mathbf{T}$  and hence (2.8) becomes

$$\hat{\mathbf{Y}}\hat{\mathbf{Y}}'\mathbf{T} = \mathbf{T}\mathbf{\Phi}. \quad (2.9)$$

This means that the latent variables obtained are the principal components of the orthogonal projection of  $\mathbf{Y}$  on to  $\mathbf{X}$ . It should be noted that these latent vectors are obtained from an optimality criteria (minimization of RSS) for the prediction of  $\mathbf{y}$ .

### 2.3 Canonical Correlation Analysis (CCA)

In this procedure pairs of latent variables  $(\mathbf{r}_j, \mathbf{t}_j) = (\mathbf{Y}\mathbf{d}_j, \mathbf{X}\mathbf{a}_j)$  are obtained by maximizing the objective function

$$(\mathbf{a}'_j \mathbf{X}' \mathbf{Y} \mathbf{d}_j)^2 / [\mathbf{a}'_j \mathbf{X}' \mathbf{X} \mathbf{a}_j \mathbf{d}'_j \mathbf{Y}' \mathbf{Y} \mathbf{d}_j] \quad (2.10)$$

with respect to  $\mathbf{d}_j$  and  $\mathbf{a}_j$  such that  $\mathbf{a}'_j \mathbf{a}_j = 1 = \mathbf{d}'_j \mathbf{d}_j$  and  $\mathbf{a}'_j \mathbf{X}' \mathbf{X} \mathbf{a}_i = 0 \quad i < j$  (see Hotelling (1936), Magnus and Neudecker (1988)). For predicting  $\mathbf{y}$  the first  $d$  latent variables  $\mathbf{T} = \mathbf{X}\mathbf{A}$  are used as regressors. From the objective function it is clear that subsets of latent variables chosen this way are not optimal for the prediction of the  $\mathbf{y}$  variables. Hence, we will not consider this in our comparisons.

### 2.4 Partial Least Squares (PLS)

Wold (1982) proposed an algorithm to compute pairs of latent variables by maximizing their covariance. Let us denote  $\mathbf{T}_{(k-1)}$  the latent variables in the  $\mathbf{X}$ -space after  $(k-1)$  pairs of latent variables have been determined. Then

$$\mathbf{X}^{(k)} = \mathbf{X} - \mathbf{T}_{(k-1)} (\mathbf{T}'_{(k-1)} \mathbf{T}_{(k-1)})^{-1} \mathbf{T}'_{(k-1)} \mathbf{X}$$

represent the residuals from the  $\mathbf{X}$  matrix at this stage and is referred to as the deflated  $\mathbf{X}$ -matrix. Then the coefficients of the  $k$ -th pair of latent variables are obtained by maximizing

$$\mathbf{d}'_k \mathbf{Y}' \mathbf{X}^{(k)} \mathbf{a}_k \quad \text{such that} \quad \mathbf{a}'_k \mathbf{a}_k = 1 = \mathbf{d}'_k \mathbf{d}_k. \quad (2.11)$$

The  $k$ -th pair of latent variables are given by

$$\mathbf{r}_k = \mathbf{Y}\mathbf{d}_k \quad \text{and} \quad \mathbf{t}_k = \mathbf{X}^{(k)} \mathbf{a}_k. \quad (2.12)$$

Several versions of this algorithm have been proposed, see for example Gelaldi and Kowalski (1986), Hoskuldsson (1988), Helland (1988), Nomikos and MacGregor (1993), de Jong

(1993) and Schmidli (1995). Merola (1998) gave an improved algorithm which is given in the Appendix.

## 2.5 Maximum Overall Redundancy (MOR)

Earlier in this section we saw that

- (i) the first  $d$  principal components of  $\mathbf{X}$  are the  $d$  latent variables in  $L(\mathbf{X})$  giving the best reconstruction of  $\mathbf{X}$ .
- (ii) the first  $d$  RRR-latent variables (PC of  $\hat{\mathbf{Y}}$ ) are optimal for predicting  $\mathbf{y}$ .

In the context of multivariate process control it is important to predict  $\mathbf{y}$  as well as to reconstruct the  $\mathbf{X}$  matrix from the latent subspace. However there is a trade off between these two objectives and PLS is a compromise between these without showing any particular optimality. Now we like to obtain a set of latent variables which meet both of these objectives. Merola (1998) has shown that, if the two objectives are equally important, then the latent variables are given by the eigen solution to

$$\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{Y}\mathbf{Y}' + \mathbf{X}\mathbf{X}')\mathbf{T} = \mathbf{T}\mathbf{\Lambda} \quad (2.13)$$

where  $\mathbf{\Lambda}$  is a diagonal matrix of eigen-values. It should be noted that, if  $\mathbf{X}'\mathbf{X}$  does not have an inverse we can replace that with a generalised inverse. Now equation (2.13) can be written as

$$(\hat{\mathbf{Y}}\hat{\mathbf{Y}}' + \mathbf{X}\mathbf{X}')\mathbf{T} = \mathbf{T}\mathbf{\Lambda} \quad (2.14)$$

where  $\hat{\mathbf{Y}}$  is the projection of  $\mathbf{Y}$  onto  $L(\mathbf{X})$ . Thus the resulting latent variables are the eigen-vectors corresponding to the  $d$  largest eigen-values of  $\hat{\mathbf{Y}}\hat{\mathbf{Y}}' + \mathbf{X}\mathbf{X}'$  which is the sum of the matrices generating the latent variables in RRR and PCA respectively. We refer to this procedure as the Maximum Overall Redundancy (MOR).

### Weighted MOR

We can generalize the MOR procedure to obtain a set of latent variables which are the eigen-vectors of

$$(1 - \alpha)\hat{\mathbf{Y}}\hat{\mathbf{Y}}' + \alpha\mathbf{X}\mathbf{X}' \quad (2.15)$$

which is a convex combination of the matrices in RRR and PCA. This will be referred to as the Weighted MOR (WMOR) method. It is easy to see that  $\alpha = 0, .5, 1$  correspond to RRR, MOR and PCR solutions respectively. As  $\alpha$  becomes smaller the prediction of  $\mathbf{y}$  gets more weight and vice versa. Now we consider some special weighting schemes. Let  $\ell(\mathbf{y}) = \text{tr}[(\mathbf{Y}'\mathbf{Y})^2]$ ,  $\ell_1(\mathbf{y}) = \text{tr}(\mathbf{Y}'\mathbf{Y})$ ,  $\ell(\hat{\mathbf{y}}) = \text{tr}[(\hat{\mathbf{Y}}'\hat{\mathbf{Y}})^2]$ ,  $\ell_1(\hat{\mathbf{y}}) = \text{tr}(\hat{\mathbf{Y}}'\hat{\mathbf{Y}})$ ,  $\ell(\mathbf{x}) = \text{tr}[(\mathbf{X}'\mathbf{X})^2]$ ,  $\ell_1(\mathbf{x}) = \text{tr}(\mathbf{X}'\mathbf{X})$ .

Based on these we define

$$\begin{aligned} \text{(i)} \quad \alpha_1 &= \sqrt{\ell(\mathbf{y})}/(\sqrt{\ell(\mathbf{y})} + \sqrt{\ell(\mathbf{x})}), \quad \text{(ii)} \quad \alpha_2 = \ell_1(\mathbf{y})/(\ell_1(\mathbf{y}) + \ell_1(\mathbf{x})), \\ \text{(iii)} \quad \alpha_3 &= \sqrt{\ell(\hat{\mathbf{y}})}/(\sqrt{\ell(\mathbf{x})} + \sqrt{\ell(\hat{\mathbf{y}})}), \quad \text{(iv)} \quad \alpha_4 = \ell_1(\hat{\mathbf{y}})/(\ell_1(\hat{\mathbf{y}}) + \ell_1(\mathbf{x})) \end{aligned} \quad (2.16)$$

As can be seen, the weights are based on the norms of  $\mathbf{Y}$ ,  $\hat{\mathbf{Y}}$  and  $\mathbf{X}$ . The procedure corresponding to  $\alpha_i$  will be referred to as  $WMOR_i$  ( $i = 1, 2, 3, 4$ ) and these will be compared with the other DRMs in the next section.

### Iterative Weighting

One unique feature of PLS is the deflation of the  $\mathbf{X}$  matrix at each iteration. In the other DRMs including MOR and WMOR it is possible to obtain the solutions simultaneously because the constraints can be reduced to the form  $\mathbf{T}'\mathbf{T} = \mathbf{I}$ . The idea of deflating the  $\mathbf{X}$  space after each latent component is obtained can be exploited to assign weights iteratively to the RRR solution matrix. These weights would represent the relative ‘‘importance’’ of each  $\mathbf{x}$  variable. We consider a matrix of diagonal weights  $\mathbf{W}_{[k]}$  in which each weight  $w_{i[k]}$  expresses the proportion of  $\mathbf{x}_i$  that still remains to be explained. Thus we take

$$w_{i[k+1]} = \frac{\mathbf{x}'_i \mathbf{x}_i - \hat{\mathbf{x}}'_{i[k]} \hat{\mathbf{x}}_{i[k]}}{\mathbf{x}'_i \mathbf{x}_i}$$

where  $\hat{\mathbf{x}}_{i[k]}$  is the rank  $k$  reconstruction of  $\mathbf{x}$  obtained with the first  $k$  latent variables. For  $k = 1$  we let  $\hat{\mathbf{x}}_{i[0]} = \mathbf{0}$ ,  $\forall i = 1, \dots, p$ , hence  $\mathbf{W}_{[1]} = \mathbf{I}_p$ . The weights  $w_{i[k]}$  converge to zero when  $\hat{\mathbf{x}}_{i[k]} = \mathbf{x}_i$ , which happens for  $k \leq p$ . When this happens the variable  $\mathbf{x}_i$  is deleted from the objective function. To obtain orthogonal solutions we take the solutions to be

$$\mathbf{t}_k = \mathbf{F}_{[k]} \mathbf{a}_k$$

where  $\mathbf{F}_{[k]} = \mathbf{X} - \hat{\mathbf{X}}(\mathbf{T}_{(k)})$  is the  $\mathbf{X}$  matrix deflated of the previous components. Hence, the  $k$ -th latent component is the projection of  $\mathbf{X}\mathbf{a}_k$  onto the space orthogonal to the previous components  $\mathbf{t}_1, \dots, \mathbf{t}_{k-1}$ . In other words we take a Gram-Schmidt orthogonalization of the matrix  $\mathbf{X}\mathbf{A}$ . We will refer to this method as Iteratively Weighted Reduced Rank Regression (IWRRR) and an algorithm to implement this procedure is given in the Appendix. This is somewhat heuristic and is hard to justify on rigorous optimization arguments; however the success of PLS and PCR together with the non-popularity of RRR in some applications indicate that the rigorous minimization of the sample Residual Sum of Squares may not lead to better predictive techniques.

### 3. Simulation study

We compare the performance of different DRMs in prediction through a simulation study. Each simulation corresponds to generating  $(n + s)$  independent observations of the  $\mathbf{x}$  and  $\mathbf{y}$  variables. The first  $n$  of these observations constitute the training sample with which the parameters of the models are determined and  $s$  the test sample used for prediction. For a given structure of the data,  $N$  pseudo-random samples are generated following the prescription. Different DRMs are then performed on each sample and the distributions of the results over the  $N$  repetitions are used for comparison. For comparing the performance over the training sample, we consider the measure of goodness-of-fit Average Residual Sum of Squares (ARSS) defined by

$$ARSSy(k, m) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^q [\mathbf{y}_{ij} - \hat{\mathbf{y}}_{ij}(m\mathbf{T}_{(k)})]^2 \quad (3.1)$$

where  $k = 1, \dots, p$  is the number of latent components used,  $p$  the number of  $\mathbf{x}$  variables,  $q$  the number of responses and  $m$  the method. For the fit of the explanatory variables the ARSS takes the form

$$ARSSx(k, m) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p [\mathbf{x}_{ij} - \hat{\mathbf{x}}_{ij}(m\mathbf{T}_{(k)})]^2$$



A measure of joint goodness-of-fit is given by

$$ARSS_t(k, m) = \frac{ARSSx(k, m)}{p} + \frac{ARSSy(k, m)}{q}$$

In some cases predictive DRMs suffer from the *Robin Hood* (RH) *effect*, that is the effect by which certain responses that are well predicted by ordinary least squares (OLS) are made substantially worse to achieve modest improvement in those that are poorly predicted. We take the ratio of the RSS for individual responses fitted with each method and the corresponding RSS of the OLS fits, and consider the average of these ratios over the  $q$  responses. This is defined by

$$Ia(k, m) = \frac{1}{q} \sum_{j=1}^q \frac{\sum_{i=1}^n (y_{ij} - \hat{y}_{ij}(k, m))^2}{\sum_{i=1}^n (y_{ij} - \hat{y}_{ij}(OLS))^2} = \frac{1}{q} \sum_{j=1}^q \frac{RSS(\mathbf{y}_j, k, m)}{RSS(\mathbf{y}_j, OLS)}$$

This index measures the extent of the RH effect on each method, the higher the index is the worse the method is affected.

As a measure of predictive efficiency we consider the average Prediction Error Sum of Squares (PRESS) over the test sample. These are defined in a way analogous to the ARSS indices. The average value of these quantities over the  $N$  simulations is then used for comparing the different methods.

CCR is not included for it is known to have a poor predictive performance.

The random variables are all generated as pseudo-random Normal variables using the linear congruential generator built in the Splus 3.4 package. In Tables and Figures we will denote the WMORi methods as WMRi and IWRRR as IWRR for ease of representation.

### The Model

We consider a reduced rank model in which both sets of variables consist of linear combinations of common latent variables with added independent noises. The model is given below

$$\mathbf{X} = \mathbf{TP} + \mathbf{F}, \quad \mathbf{Y} = \mathbf{TQ} + \mathbf{E} \quad (3.2)$$

where  $T$  is the matrix of latent variables,  $P$  and  $Q$  the matrices of loadings and  $E$  and  $F$  the matrices of errors. In this study we take  $q = 3$  responses and  $p = 6$  explanatory variables generated from 2 latent variables. The signal to noise ratio (SNR) is chosen to be 3 for both sets of data.  $E$  and  $F$  are random noises with diagonal covariance matrices so that the SNR is 3 for every variable. The simulation consists of  $n = 50$  observations for the training sample and  $s = 10$  for the test sample replicated  $N = 500$  times. The squared canonical correlation coefficients for this covariance structure are:

Table 3.1: Squared Canonical Correlation Coefficients

$\rho^2$	0.95070	0.86199	0.00006
----------	---------	---------	---------

As expected, there are two common directions of high correlation and an almost orthogonal one.

Tables 3.2-3.4 give the average values of the ARSS indices for different numbers of latent components. With respect to the average ARSSy, RRR dominates all other methods and PCR is always worse than all others. PLS values, with the exception of the first latent component, is always slightly higher than the values of the WMORs. MOR has slightly higher ARSSy than PLS but for two components, that is for the right number of components. For ARSSst all methods except RRR and PCR show nearly the same values. MOR and the WMORs seem to have a slight edge over PLS.

The distributions of the ARSS over the simulated samples for 2 latent variables are given in Figure 3.1.

Table 3.2: Average ARSSy in the training sample

ARSSy	PLS	MOR	WMR1	WMR2	WMR3	WMR4	RRR	IWRR	PCR
1 comp	0.579	0.585	0.535	0.541	0.520	0.529	0.495	0.495	0.665
2 comps	0.375	0.371	0.366	0.367	0.365	0.366	0.360	0.375	0.380
3 comps	0.359	0.363	0.358	0.359	0.356	0.357	0.352	0.360	0.374
4 comps	0.353	0.359	0.355	0.355	0.354	0.355	0.495	0.354	0.367
5 comps	0.352	0.354	0.353	0.353	0.353	0.353	0.360	0.352	0.359
6 comps	0.352	0.352	0.352	0.352	0.352	0.352	0.352	0.352	0.352

Table 3.3: Average ARSSx in the training sample

ARSSx	PLS	MOR	WMR1	WMR2	WMR3	WMR4	RRR	IWRR	PCR
1 comp	1.169	1.159	1.242	1.228	1.288	1.258	1.499	1.499	1.124
2 comps	0.477	0.480	0.487	0.486	0.492	0.489	0.560	0.477	0.476
3 comps	0.354	0.338	0.346	0.345	0.351	0.348	0.419	0.354	0.334
4 comps	0.233	0.209	0.215	0.214	0.218	0.216	1.499	0.234	0.206
5 comps	0.116	0.096	0.099	0.098	0.100	0.099	0.560	0.116	0.095
6 comps	0.000	0.000	0.000	0.000	0.000	0.000	0.419	0.000	0.000

Table 3.4: Average ARSS<sub>t</sub> in the training sample

ARSS <sub>t</sub>	PLS	MOR	WMR1	WMR2	WMR3	WMR4	RRR	IWRR	PCR
1 comp	0.388	0.388	0.385	0.385	0.388	0.386	0.415	0.415	0.409
2 comps	0.205	0.204	0.203	0.203	0.204	0.203	0.213	0.205	0.206
3 comps	0.179	0.177	0.177	0.177	0.177	0.177	0.187	0.179	0.180
4 comps	0.157	0.154	0.154	0.154	0.154	0.154	0.415	0.157	0.157
5 comps	0.137	0.134	0.134	0.134	0.134	0.134	0.213	0.137	0.135
6 comps	0.117	0.117	0.117	0.117	0.117	0.117	0.187	0.117	0.117

Examining these distributions we conclude that all methods give almost the same results. RRR has a slight edge over the other methods for ARRS<sub>y</sub>. However its ARRS<sub>x</sub> is higher than the others leading to a higher overall ARSS<sub>t</sub>. As expected, the WMOR methods have the lowest values of ARSS<sub>t</sub>. MOR is performing slightly better than PLS with respect to ARSS<sub>y</sub> and ARSS<sub>t</sub>.

It has been observed that PLS often suffers from the Robin Hood effect (Breiman and Friedman (1997)). In this study, fitting with two latent variables gives almost the same  $Ia$  values for every method (Figure 3.2), but when “over-fitting” with three latent variables we notice that the  $Ia$  for PCR, MOR and PLS are larger than the WMOR methods. This implies that the addition of the third latent variable in WMOR decreases proportionally all RSS of the responses while in PCR, PLS and MOR there are some responses that are not well fitted with respect to the OLS “best” fits.

Tables 3.5-3.7 give the average PRESS values. All methods but RRR practically give the same PRESS<sub>y</sub> for the 2 latent component predictions. It is interesting to see how the

method that minimizes ARSS<sub>y</sub> performs worse with respect to PRESS<sub>y</sub>, even in a situation in which the data are extremely well behaved and follow a latent model.

Table 3.5: Average PRESS for the  $y$  variables

PRESS <sub>y</sub>	PLS	MOR	WMR1	WMR2	WMR3	WMR4	RRR	IWRR	PCR
1 comp	1.680	1.709	1.590	1.601	1.560	1.576	1.539	1.539	1.896
2 comps	1.150	1.157	1.169	1.167	1.176	1.171	1.225	1.150	1.147
3 comps	1.203	1.176	1.190	1.188	1.196	1.191	1.238	1.200	1.158
4 comps	1.225	1.189	1.204	1.201	1.211	1.207		1.224	1.170
5 comps	1.234	1.213	1.223	1.222	1.226	1.224		1.235	1.198
6 comps	1.238	1.238	1.238	1.238	1.238	1.238		1.238	1.238

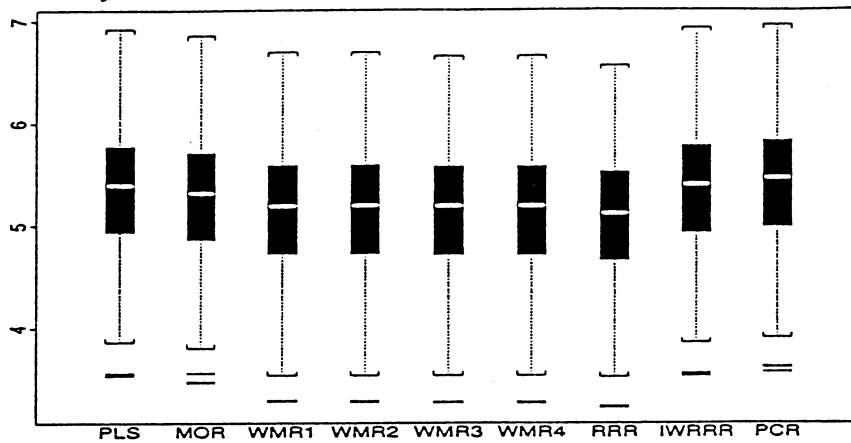
Table 3.6: Average PRESS for the  $x$  variables

PRESS <sub>x</sub>	PLS	MOR	WMR1	WMR2	WMR3	WMR4	RRR	IWRR	PCR
1 comp	2.547	2.530	2.717	2.681	2.820	2.749	3.294	3.294	2.439
2 comps	1.060	1.067	1.085	1.081	1.097	1.088	1.274	1.060	1.059
3 comps	0.821	0.903	0.889	0.891	0.881	0.887	0.963	0.826	0.905
4 comps	0.561	0.696	0.673	0.676	0.661	0.669		0.558	0.702
5 comps	0.286	0.410	0.400	0.401	0.396	0.399		0.287	0.413
6 comps	0.000	0.000	0.000	0.000	0.000	0.000		0.000	0.000

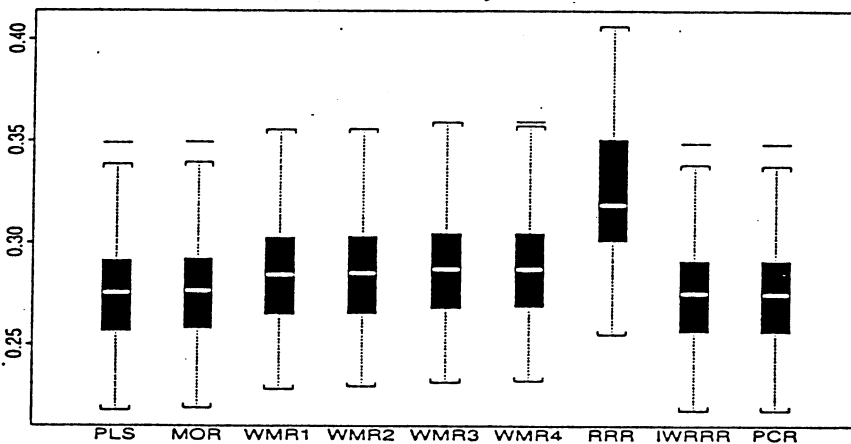
Table 3.7: Average Total PRESS

PRESS <sub>t</sub>	PLS	MOR	WMR1	WMR2	WMR3	WMR4	RRR	IWRR	PCR
1 comp	0.985	0.991	0.983	0.981	0.990	0.983	1.062	1.062	1.038
2 comps	0.560	0.564	0.571	0.569	0.575	0.572	0.621	0.560	0.559
3 comps	0.538	0.542	0.545	0.544	0.545	0.545	0.573	0.538	0.537
4 comps	0.502	0.512	0.514	0.513	0.514	0.514		0.501	0.507
5 comps	0.459	0.473	0.474	0.474	0.475	0.475		0.459	0.468
6 comps	0.413	0.413	0.413	0.413	0.413	0.413		0.413	0.413

ARSS<sub>y</sub>: 2 latent components, SNR<sub>y</sub>=SNR<sub>x</sub>=3 noises uncor.



ARSS<sub>x</sub>: 2 components, SNR<sub>y</sub>=SNR<sub>x</sub>=3 noises uncor.



ARSS<sub>t</sub>: 2 components, SNR<sub>y</sub>=SNR<sub>x</sub>=3 noises uncor.

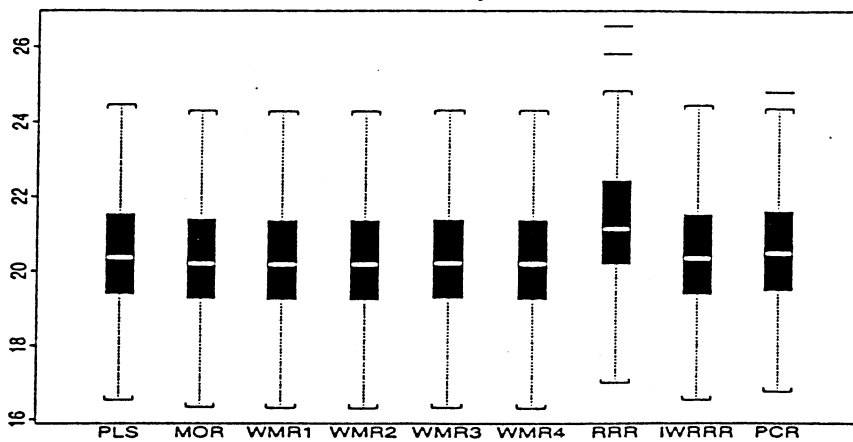
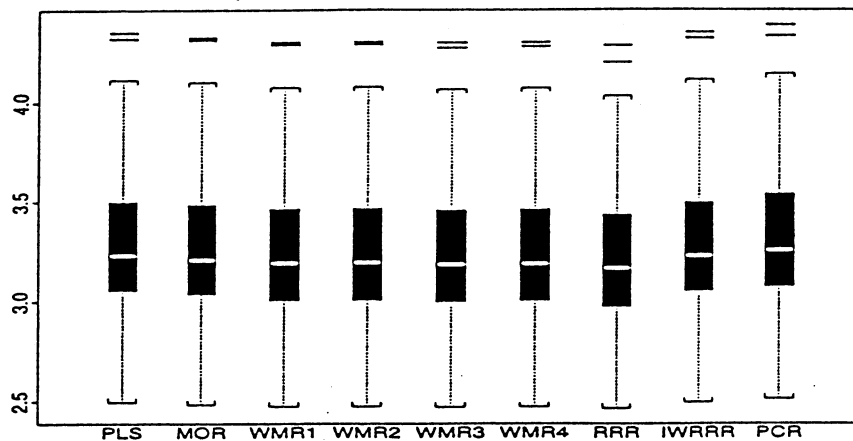


Figure 3.1 Distribution of  $ARSS_y$ ,  $ARSS_x$  and  $ARSS_t$

Ia: 2 components, SNRy=SNRx=3, noise uncor.



Ia: 3 components, SNRy=SNRx=3, noises uncorr.

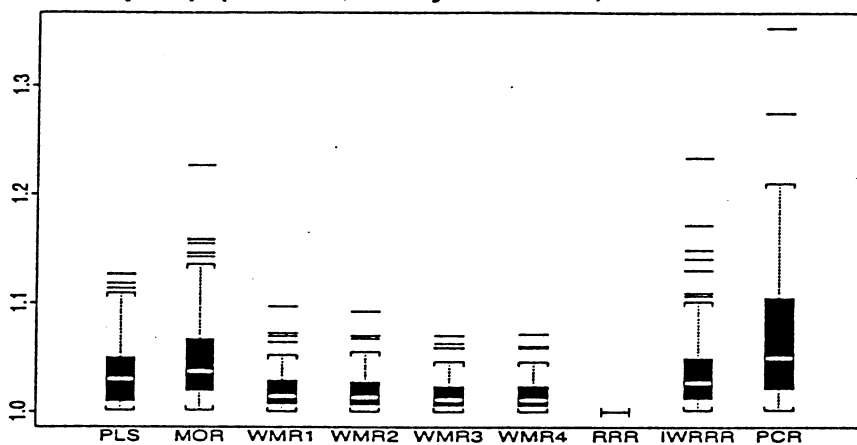


Figure 3.2 *Ia* Indices for the *y* variables

Note that for all methods the PRESS<sub>y</sub> increases when the model is over-fitted. Such an increase is more marked for PLS. The methods that give the best “predictions” of  $\mathbf{x}$  variables are PCR and PLS. The WMOR methods give good reconstructions of the explanatory variables while the values of RRR stand out for being higher than the others. With respect to PRESS<sub>t</sub> all methods but RRR give very close results.

## 4. Concluding Remarks

In this paper we briefly reviewed some common DRMs such as PCR, RRR, CCR and PLS. We also discussed the relatively new techniques MOR, WMOR and IWRRR. We note here that there have been attempts to look for a common framework for these procedures (see Burnham et al. (1995)). Merola (1998) and Merola and Abraham (1998) have discussed a common objective function from which the objective functions of the various DRMs can be obtained by changing some parameters.

From the simulation study we observe that:

- (i) In the training sample (ARSS) RRR is the best for predicting  $\mathbf{y}$  and PCR is the worst; for predicting  $\mathbf{x}$  PCR is the best and RRR is the worst. MOR, WMORs and PLS are reasonably good for both.
- (ii) In the post sample (PRESS) MOR, WMORs and PLS are doing well for prediction of  $\mathbf{x}$  as well as  $\mathbf{x}$  and  $\mathbf{y}$  jointly. RRR has a slight edge only in predicting  $\mathbf{y}$ .
- (iii) It is interesting to note that the behaviour of ARSS (training sample) and PRESS (post sample) are somewhat different. Also it should be noted that although PLS is not an “optimal” method it is doing very well in post sample predictions.

## Acknowledgements

B. Abraham was partially supported by a grant from the Natural Sciences and Engineering Research Council of Canada.

## References

- Breiman, L. and Friedman, J.H. (1997). Predicting multivariate responses in multivariate regression. *J. Royal Stat. Soc. B*, 59(1), 3-54.
- Burnham, A.J., Viveros, R., and MacGregor, J.F. (1995). Frameworks for latent variable multivariate regression. *J. of Chemometrics*, 20.
- de Jong, S. (1993). Simpls: an alternative approach to partial least squares regression. *Chemom. and Intell. Lab. Systems*, 18:251-263.
- Gelaldi, P. and Kowalski, B.R. (1986). Partial least-squares regression: A tutorial. *Analytica Chimica Acta*, 185:1-17.
- Helland, I.S. (1988). On the structure of partial least squares. *Comm. Stat.-sim*, 17(2):581-607.
- Hoskuldsson, P. (1988). Pls regression methods. *J. of Chemometrics*, 2:211-228.
- Hotelling, H. (1935). The most predictable criterion. *J. Educ. Psychol.*, pages 139-142.
- Hotelling, H. (1936). Relation between two sets of variates. *Biometrika*, 28:321-377.
- Izenman, A.J. (1975). Reduced-rank regression for the multivariate bilinear model. *J. of Multivariate Analysis*, 5:248-264.
- Merola, G.M. (1998). Dimensionality reduction methods in multivariate prediction. Unpublished Ph.D. Thesis. Department of Statistics and Actuarial Science, University of Waterloo, Waterloo, Canada.
- Merola, G.M. and Abraham, B. (1998). An objective function approach for dimensionality reduction methods in prediction. Research Report, IIQP, University of Waterloo, Waterloo, Canada.



- Magnus, J.R. and Neudecker, H. (1988). *Matrix Differential Calculus with Applications in Statistics and Econometrics*. Wiley.
- Nomikos, P. and MacGregor, J.F. (1993). Monitoring of batch processes using multi-way principal component analysis. *A.I.Ch.E. Jour.*
- Rao, C.R. (1964). The use and interpretation of principal component analysis in applied research. *Sankhya A*, 26:329-358.
- Schmidli, H. (1995). *Reduced Rank Regression*. Contributions to Statistics. Physica-Verlag.
- Wold, H. (1982). Soft modelling, the basic design and some extensions. In Joresorg, K. and Wold, H., editors, *Systems Under Indirect Observation*, volume II, pages 589-591. Wiley and Sons.

## APPENDIX

We assume that the data matrices  $X$  and  $Y$  are column mean centred and the columns are scaled to unit variance.

Let  $svd(X) = U\Gamma V'$  where  $U$  and  $V$  are orthonormal matrices and  $\Gamma$  is diagonal with positive, non-decreasing, diagonal entries.

Also let  $svd(Y) = W\Gamma Z'$  where  $W$ ,  $Z$ , and  $\Gamma$  are defined similar to  $U$ ,  $V$  and  $\Gamma$  respectively.

### Partial Least Squares

- 0) Initialize  $X_0 = X$ ,  $Y_0 = Y$
- 1)  $svd(X'_{k-1}Y) = J\Delta\Gamma'$
- 2)  $c_k = j_1$  (first column of  $J$ )
- 3)  $t_k = X_{k-1}c_k$
- 4)  $c_k \leftarrow [c_k / \sqrt{t'_k t_k}] \sqrt{n}$
- 5)  $t_k \leftarrow [t_k / \sqrt{t'_k t_k}] \sqrt{n}$
- 6)  $X_k = (I_n - \frac{t_k t'_k}{n}) X_{k-1}$
- 7) if  $\text{sum}(\text{diag}(X'_k X_k)) < \epsilon$  go to 8  
else go to 1
- 8)  $N = XC$   
 $\text{qr}(N) = QR$   
 $A = CR^{-1}$

After the coefficient  $A$  is computed,  $\hat{M}_k$  is computed as

$$\hat{M}_k = a_k t'_k Y + \hat{M}_{k-1}, \quad k = 1, \dots, d, \quad \hat{M}_0 = 0$$

### Iteratively Weighted Reduced Rank Regression (IWRRR)

- 0)  $W_1 = I_p$   $F_1 = X$  Initialization
- 1) Compute svd ( $W_i(F_i^T F_i)^{-1} F_i^T Y Y^T F_i$ )
- 2)  $a_i = A_{(1)}$  Computation of coefficients and scores
- 3)  $t_i = F_i a_i / \| F_i a_i \|$
- 5)  $H_i = t_i t_i'$  Projection matrix
- 6)  $\hat{X}_i = H_i X$
- 8)  $F_{i+1} \leftarrow F_i - \hat{X}_i$  Estimates and deflation
- 9)  $W_{[i+1]} = \text{diag}\{F_{i+1}' F_{i+1}\} [\text{diag}\{X' X\}]^{-1}$  Computation of the weights and
- 10) if  $\| W_{[i+1]} \| > \epsilon$  go to 2; else exit stopping rule