# Efficient Cross-validatory Computations and Influence Measures for Principal Component and Partial Least Squares Decompositions with Applications in Chemometrics

ProQuest Number: 10017686

All rights reserved

INFORMATION TO ALL USERS
The quality of this reproduction is dependent upon the quality of the copy submitted.

In the unlikely event that the author did not send a complete manuscript
and there are missing pages, these will be noted. Also, if material had to be removed,
a note will indicate the deletion.



ProQuest 10017686

# Acknowledgements

# Abstract

The application of the efficient rank-one perturbation algebra of Bunch, Nielsen and Sorensen [BNS78] to the leave-one-out cross-validation of principal component regression is described. Similarly, we consider the restriction of the efficient leave-one-out cross-validatory algebra for continuum regression proposed by Stone and Brooks [SB90][SB92] to partial least squares regression. Implementations of both cross-validatory procedures are presented in the numerical analysis package GAUSS [Apt92], together with procedures for the computation of principal component and partial least squares regression equations. We describe the application to the leave-one-out cross-validatory assessment of principal component and partial least squares prediction equations, using near infrared spectroscopic data. The methodologies are compared with the existing procedures for efficiency and numerical accuracy. We derive influence measures from the cross-validatory computations and consider an application in food analysis on the use of near infrared spectroscopy for the calibration of total oil content of white mustard seeds. Finally, a published paper on the SIMCA [Wol76] method of classification which is based on the principal component decomposition and usually involves a cross-validatory assesment is bound in with this thesis, following the bibliography.

# Contents

# List of Tables

# List of Figures

# Notation Glossary

| | |
|---|---|
| $\widetilde{\mathbf{Y}}$ | column vector of observed responses |
| $\tilde{y}_i$ | observed response for the $i^{\text{th}}$ observation |
| $\widetilde{\mathbf{X}}$ | matrix of observed predictor variables |
| $\tilde{\mathbf{x}}_i$ | row vector of observed predictors for the $i^{\text{th}}$ observation |
| $\mathbf{1}$ | column vector of 1's |
| $\mathbf{Y}$ | column vector of mean-centred responses |
| $y_i$ | mean-centred response for the $i^{\text{th}}$ observation |
| $\mathbf{X}$ | matrix of mean-centred predictor variables |
| $\mathbf{x}_i$ | row vector of mean-centred predictors for the $i^{\text{th}}$ observation |
| $r_x$ | rank of $\mathbf{X}$ |
| $\mathbf{S}$ | sample covariance matrix |
| $\mathbf{s}$ | sample covariance vector |
| $\mathbf{C}$ | cross-products matrix |
| $\mathbf{c}$ | cross-products vector |
| $\mathbf{Q}$ | matrix of component coefficients |
| $\mathbf{q}_j$ | column vector that contains the $j^{\text{th}}$ component coefficients |
| $\mathbf{\Lambda}$ | diagonal matrix of principal component variances |

| | |
|---|---|
| $\lambda_j$ | the variance of the $j^{\text{th}}$ principal component |
| $\mathbf{E}$ | diagonal matrix of principal component eigenvalues |
| $e_j$ | the eigenvalue of the $j^{\text{th}}$ principal component |
| $\mathbf{U}$ | matrix of component scores |
| $\mathbf{u}_j$ | column vector that contains the $j^{\text{th}}$ component scores |
| $\mathbf{W}$ | matrix of standardized principal components |
| $\mathbf{I}$ | identity matrix |
| $\mathbf{g}_j$ | the $j^{\text{th}}$ standard basis vector |
| $n$ | number of observations |
| $p$ | number of predictor variables |
| $\gamma$ | number of components selected for the construction of a prediction equation |
| $\mathbf{z}$ | normalized principal component scores for a cross-validated observation |
| $(i)$ | a subscript used to indicate quantities that have been downdated for the omission of the $i^{\text{th}}$ observation |
| $\mathbf{P}(l, m)$ | permutation matrix |
| $\mathbf{J}(l, m, \theta)$ | rotation matrix |
| $\nu$ | the cross-validatory constant $n/(n-1)$ |
| $\rho$ | the total eigenvalue downdate |
| $\theta$ | an angle |
| $\mu_j$ | the $j^{\text{th}}$ normalized eigenvalue downdate |
| $\psi$ | a rational function |
| $\phi$ | a rational function |

| | |
|---|---|
| $\tau$ | a constant of order unity |
| $\eta$ | the machine epsilon |
| $\| \ \|$ | the Euclidean vector norm |
| $\| \ \|_2$ | the matrix 2-norm |
| $\| \ \|$ | the absolute value |

# Chapter 1

# Introduction

Much interest in multivariate statistics has centred recently on the use and development of statistical methodology in the field of chemometrics. This is particularly so for the calibration of linear regression equations for prediction, using high dimensional predictor data. Typical examples are in the application of near infrared spectroscopy to the compositional analysis of foodstuffs [OFMD84]. Many measurements in food chemistry are difficult or expensive to obtain with classical methods from analytical chemistry. Spectroscopic analysis on the other hand is often quick and less expensive to perform and the spectral data may contain information relevant to the measure of interest. Near infrared spectroscopy in particular is a relatively inexpensive and simple method of analysis which has often been applied with success in food analysis. The objective of these analyses is to replace existing procedures for the determination of a measure of interest with a prediction rule obtained from a properly calibrated regression equation, using the spectral data. It is at this point that statistical methodology is applied to construct and optimize prediction equations for use in practical applications. Many of these applications have been restricted to linear regression of the measurements to be calibrated on the matrix of predictor data.

The spectra generated from a spectroscopic analysis are usually stored as measurements on a fixed grid of equidistant wavelengths. The size of the predictor matrix is thus determined by the size of the calibration set and the number of wavelengths at which the spectra are measured. The latter is potentially unlimited, while sample sizes are usually restricted. Thus, in these applications, the number of measurements available for prediction is normally much larger than the number of samples available. More important than the dimensionality problem, these calibration problems are usually complicated by the fact that it is typically the small-variation perturbations which are correlated with the measure of interest. In most applications this precludes ordinary or minimum-length least squares for the construction of a linear prediction equation, as small sources of variation will inflate the variability of the prediction rule. Thus, if we want to maintain the convenience of linear prediction equations, a problem of bias-variance trade-off with respect to least squares calibration needs to be addressed. While there are several possibilities to address this problem, least squares regression on principal or partial least squares components has often been employed, using cross-validatory assessment, as described by Martens and Naes [MN89] as well as Brown [Bro93], for example.

These methodologies have depended heavily on advances in technology and instrumentation. This is particularly true for improvements in computational facilities which have made efficient computations for large data matrices possible. Thus, techniques like principal component or partial least squares regression have become routine calculations in the calibration of many prediction equations involving high dimensional data. Furthermore, we can employ the pioneering work of such researchers as Householder [Hou58], Wilkinson [Wil65] and Golub and Van Loan [GL89] for the efficient computation of principal component regression equations. Similarly, for partial least squares regression, efficient computations have recently been proposed by Stone and Brooks [SB90] for both the partial least squares decomposition and cross-validatory computations, embedded in a more general algebra for continuum regression.

To date, however, no such efficient computations have been made available for the leave-one-out cross-validation of principal component regression equations. Naive implementations proceed from a random or systematic division of the data into subsamples, referred to as validation samples. For each subsample, the remainder of the data constitutes a construction sample. The cross-products matrices are then recomputed for each construction sample followed by a complete principal component analysis of the new cross-products matrix. With such implementations, the cost of a full leave-one-out cross-validation may be extremely high. An algorithm proposed by Bunch, Nielsen and Sorensen [BNS78], based on work by Golub [Gol73] and further investigated by DeGroat [DeG90], provides an elegant solution to the problem of leave-one-out calculations. By applying this algorithm to the cross-validation of principal components, one can calculate a full leave-one-out cross-validation from a single principal component analysis on the whole data. Rather than recomputing the principal components from scratch after leaving out an observation, we consider the change in the principal components that results from a leave-one-out perturbation. This computation is based on the downdating of the eigenvalues that is associated with the omission of an observation. The method is conceptually elegant, as we only calculate the decomposition that we are actually interested in. We can store this decomposition for use in further analysis.

As for principal component regression, we can implement an efficient leave-one-out cross-validatory algebra for partial least squares. This is obtained through the restriction of the efficient leave-one-out cross-validatory

computations for continuum least squares [SB90] to the partial least squares case. Both these algorithms are implemented in the numerical analysis package GAUSS [Apt92, Version 3.0]. Finally, we consider the derivation of influence measures in the leave-one-out cross-validatory algebra and consider their use in the calibration of linear prediction equations with the methods considered.

Throughout, we restrict ourselves to the decomposition of the covariance matrix, for both principal component and partial least squares regression and their cross-validatory assessment. The algebra described does not apply to the correlation matrix. This has also been observed by Sundberg [Sun93] and Stone and Brooks [SB90][SB92] in the cross-validation of continuum regression. In such situations where a scaled analysis may be preferable, one could scale the data prior to analysis in a manner which is not dependent on the data. One could then apply the cross-validatory algebras which will be presented here for principal and partial least squares components.

## 1.1 Preliminaries

We formalise some of the previous discussion. Assume that a sample of $n$ observations $(\tilde{y}_i, \tilde{x}_{i1}, \ldots, \tilde{x}_{ip})$, $i = 1, \ldots, n$, has been generated from a spectroscopic analysis, where

$$\widetilde{\mathbf{Y}} = (\tilde{y}_1, \ldots, \tilde{y}_n)^T$$

is a column vector of measurements on a variable of interest and

$$\widetilde{\mathbf{X}} = \begin{pmatrix} \tilde{\mathbf{x}}_1 \\ \vdots \\ \tilde{\mathbf{x}}_n \end{pmatrix}$$

is a corresponding matrix of predictor data, with $\tilde{\mathbf{x}}_i = (\tilde{x}_{i1}, \ldots, \tilde{x}_{ip})$, $i = 1, \ldots, n$. It is often more easy to work with mean-centred data and hence we introduce the notations

$$\mathbf{Y} = (\tilde{y}_1 - \overline{\widetilde{\mathbf{Y}}}, \ldots, \tilde{y}_n - \overline{\widetilde{\mathbf{Y}}})^T = (y_1, \ldots, y_n)^T$$

where $\overline{\overline{\mathbf{Y}}} = \sum_{i=1}^{n} \tilde{y}_i / n$ and

$$\mathbf{X} = \begin{pmatrix} \tilde{\mathbf{x}}_1 - \overline{\overline{\mathbf{X}}} \\ \vdots \\ \tilde{\mathbf{x}}_n - \overline{\overline{\mathbf{X}}} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{pmatrix}$$

with $\overline{\overline{\mathbf{X}}} = \sum_{i=1}^{n} \tilde{\mathbf{x}}_i / n$. The sample covariance matrix is defined as

$$\mathbf{S} = \frac{\mathbf{X}^T \mathbf{X}}{n-1}$$

and likewise, we may define a sample covariance vector

$$\mathbf{s} = \frac{\mathbf{X}^T \mathbf{Y}}{n-1}.$$

We will be concerned with the application of full leave-one-out cross-validation [Sto74] to the construction of a linear prediction equation of the form

$$\overline{\overline{\mathbf{Y}}}\mathbf{1} + \mathbf{X}\boldsymbol{\alpha},$$

for the univariate set of measurements $\widetilde{\mathbf{Y}}$, from the observed explanatory data $\widetilde{\mathbf{X}}$, using principal component [Jol86, page 129][MN89, page 97] and partial least squares [MN89, page 116] regression, where $\boldsymbol{\alpha}$ is a column vector of regression coefficients and $\mathbf{1}$ is a column vector of 1's. These biased regression methods adhere to the Stone-Brooks [SB90] formulation of continuum regression, in that they sequentially derive a set of regressors $\mathbf{u}_j = \mathbf{X}\mathbf{q}_j$, $j = 1, \ldots, r_x$, where $r_x = \min(n-1, p)$ is the rank of the predictor matrix $\mathbf{X}$, such that the regressors have positive variances and are uncorrelated, with each vector of coefficients $\mathbf{q}_j$, $j = 1, \ldots, r_x$, a column vector of unit length. The coefficient vectors $\mathbf{q}_j$, $j = 1, \ldots, r_x$, define a new system of coordinates in the observed multivariate predictor space. We will refer to these directions as component coefficients. Likewise, we may refer to the vector of coordinate values $\mathbf{u}_j$, $j = 1, \ldots, r_x$, of the observations in this new coordinate system as scores or components.

Each method constructs, for each $0 < \gamma < r_x$, a biased prediction formula

$$\widehat{\overline{\mathbf{Y}}}_{(\gamma)} = \overline{\overline{\mathbf{Y}}}\mathbf{1} + b_1 \mathbf{u}_1 + \cdots + b_\gamma \mathbf{u}_\gamma$$

from an ordinary least squares regression of $\widetilde{Y}$ on $u_1, \ldots, u_\gamma$. Hence, the only difference between both methodologies is in the optimality properties of the derived components. Therefore, the above framework of biased regression is complete apart from two points. First, and obviously, there is the optimality criterion to be employed in the construction of the components. Secondly, an appropriate stopping rule must be used to choose the number of regressors for optimal prediction.

Stone and Brooks considered a full leave-one-out cross-validatory choice of the number of regressors in their original paper on continuum regression [SB90, page 245].

An application of full leave-one-out cross-validation to the above framework for biased regression involves, for each datum $(\tilde{y}_i, \tilde{x}_i)$, $i = 1, \ldots, n$, the computation of the leave-one-out cross-validated prediction equation

$$\widehat{\widetilde{Y}}_{(i,\gamma)} = \overline{\overline{Y}}_{(i)}\mathbf{1} + b_{1(i)}\mathbf{u}_{1(i)} + \cdots + b_{\gamma(i)}\mathbf{u}_{\gamma(i)},$$

that is obtained when the datum is left out of the calculations. This equation is then applied to the left-out datum, to obtain the cross-validated predicted value

$$\widehat{\tilde{y}}_{i(i,\gamma)} = \overline{\overline{Y}}_{(i)} + b_{1(i)}u_{i1(i)} + \cdots + b_{\gamma(i)}u_{i\gamma(i)}$$

for $\tilde{y}_i$, where

$$u_{ij(i)} = (\tilde{x}_i - \overline{\overline{X}}_{(i)})\mathbf{q}_{j(i)},$$

with $\mathbf{q}_{j(i)}$ the $j^{\text{th}}$ vector of component coefficients when the $i^{\text{th}}$ observation has been removed from the data. A full leave-one-out cross-validatory assessment may then be calculated for each $0 < \gamma < r_x$ from the differences between the observed values $\tilde{y}_i$, $i = 1, \ldots, n$, and the cross-validated predicted values $\widehat{\tilde{y}}_{i(i,\gamma)}$, $i = 1, \ldots, n$. A suitable criterion is the prediction error sum of squares

$$\text{PRESS}_\gamma = \sum_{i=1}^{n}(\tilde{y}_i - \widehat{\tilde{y}}_{i(i,\gamma)})^2,$$

or equivalently, the mean prediction error sum of squares

$$\text{MSEP}_\gamma = \frac{\sum_{i=1}^{n}(\tilde{y}_i - \widehat{\tilde{y}}_{i(i,\gamma)})^2}{n},$$

which may also be referred to as the mean squared error of prediction. This implementation of cross-validation is referred to as full leave-one-out cross-validation, as each observation is removed from the data on its own, retaining

all other observations for the calculation of the prediction equation for the left-out datum.

Cross-validation makes efficient use of the data as it does not require a separate validation set. Full leave-one-out cross-validation is particularly efficient in its use of the data, as each cross-validated prediction equation is calibrated from the maximum number of observations and hence, at least from a purely heuristic point of view, leave-one-out cross-validation uses the data to the full. Unfortunately, the computational cost of naive implementations of full leave-one-out cross-validation can be prohibitively high. Hence, the above described efficiency does not extend to the computational aspects of leave-one-out cross-validation.

Stone and Brooks have considered the efficient cross-validation of continuum regression and, by implication, of partial least squares. Unfortunately, the Stone-Brooks computations for continuum regression do not extend to the principal component limit. Therefore, we will consider the efficient numerical computation of full leave-one-out cross-validatory calculations for principal component regression. We may then consider the application of the Stone-Brooks cross-validatory algebra to partial least squares regression and draw comparisons.

It is clear that the essential step in the cross-validatory computations is the calculation of the cross-validated components $\mathbf{q}_j$, $j = 1, \ldots, r_x$, irrespective of whether we use principal component or partial least squares decomposition. We will refer to these as the downdated component coefficients and to the process of computing the cross-validated components as downdating. Thus, the problem of constructing efficient methods for the leave-one-out cross-validation of principal component and partial least squares decompositions reduces to the downdating of the components under leave-one-out perturbations of the data. The leave-one-out cross-validatory algebra will be simplified if we introduce the cross-products matrix

$$\mathbf{C} = \mathbf{X}^T\mathbf{X},$$

such that $\mathbf{S} = \mathbf{C}/(n - 1)$ and the cross-products vector

$$\mathbf{c} = \mathbf{X}^T\mathbf{Y},$$

such that $\mathbf{s} = \mathbf{c}/(n - 1)$.

These methods may be criticised on a number of points. First, although they deal effectively with the dimensionality problem they ignore the natural ordering of wavelengths completely. Potentially more damaging criticism may be found in the approach to the bias-variance trade-off. For principal component regression, the above implementation has been justly criticised (e.g. Jolliffe [Jol82]) as the high-variance components may not be the best predictors for the response considered. Thus, although rejecting small-variance components is optimal in terms of rejecting those components that increase the variability of the prediction rule most, the prediction equation obtained may not be optimal when the ordering of components with respect to correlation with the response measure differs from that obtained on the basis of predictor variance alone. Some reordering of components may be needed so as to allow for relative differences in correlation between principal components. Many authors have suggested a ranking according to correlation with the measure of interest [Jol86].

Ordering of components is intrinsic in the partial least squares approach but the covariance criterion itself may be sub-optimal. This problem is addressed by continuum regression as it seeks to find the optimal optimisation criterion through a cross-validatory assessment. Thus, continuum regression attempts to achieve a similar result to that which would be achieved through selection of principal components. The above version of principal component regression emerges as one limit of continuum regression and hence, this implementation is of interest from a purely computational point of view. Furthermore, it is a straightforward exercise to change an implementation of the subsequent cross-validatory algebra to allow for a different ordering of the principal component scores in regression, using a permutation of the principal component indices. For these reasons, and to simplify notations, we restrict ourselves to the previously described implementation .

## 1.2 Mustard data

We will use an example from near infrared spectroscopy to illustrate the efficient cross-validatory computations. The example is concerned with 50 samples of white mustard seed and was provided by the Campden Food and Drink Research Association [EJH91]. In this application, interest resides in the total oil content of the samples. After determination of these measures,

Figure 1.1: The near infrared spectra of fifty samples of white mustard seed.

the data were stored as the response variable $\widetilde{\mathbf{Y}}$ to be predicted. A near infrared spectroscopic analysis was also carried out for the same samples, measuring the absorbance $\log(1/R)$ at 700 wavelengths, ranging from 1100 nanometres to 2500 nanometres at intervals of 2 nanometres. Figure 1.1 shows a plot of the measured near infrared spectra for these samples of white mustard seed. These data were stored as the predictor data $\widetilde{\mathbf{X}}$ and are a fairly typical example of the use of near infrared spectroscopy in the analysis of foods. While the distinct spectra are of a smooth and predominantly parallel nature, their interpretation is only of secondary interest. True interest lies in the calibration of the measure to be predicted. Interpretation of the predictor data often comes only at a subsequent stage, after calibration, possibly in the form of the interpretation of the principal or partial least squares components.

Principal component and partial least squares regression are applied to these data to construct linear predictors for the oil content from the observed spectra, as explained in chapters 2 and 3 respectively. Leave-one-out cross-validation is used to choose the number of components so as to optimize the mean squared error of prediction.

# Chapter 2

# Principal Component Regression

## 2.1 Principal Component Decomposition and Cross-validation

### 2.1.1 Principal Component Decomposition

Principal components are derived sequentially, such that each component maximizes the sample variance of the corresponding principal component scores, subject to the side conditions that the principal component scores are uncorrelated between distinct components and have positive sample variance, while the principal component coefficient vectors are of unit length [Hot33][Jol86, page 1]. Stone and Brooks [SB90, page 241] remark that principal component regression is peculiar as the components are constructed without reference to the observed response. This has also been noted by Gnanadesikan and Kettenring [GK72, page 110] who refer to principal component decomposition as a purely internal data analysis technique.

**Principal Component Algebra**

The algebra of principal component decomposition is well known and has been described in many texts on multivariate statistics [And58][Krz88][Jol86]. As this algebra parallels that of continuum regression, and therefore also that of partial least squares, we will include it here for completeness.

The first component coefficient vector $\mathbf{q}_1$ must maximize the sample variance $\mathbf{q}^T \mathbf{S} \mathbf{q}$ of the linear combination $\mathbf{X}\mathbf{q}$, among all $\mathbf{q}$ with $\mathbf{q}^T \mathbf{q} = 1$. Hence, using the method of Lagrange multipliers, we must maximize

$$L_1(\mathbf{q}) = \mathbf{q}^T \mathbf{S} \mathbf{q} - \lambda(\mathbf{q}^T \mathbf{q} - 1),$$

and the first component coefficient vector is the solution of the equation

$$\frac{d}{d\mathbf{q}} L_1(\mathbf{q}) = \mathbf{S}\mathbf{q} - \lambda \mathbf{q} = 0.$$

Therefore, $\mathbf{q}$ satisfies the equation

$$\mathbf{S}\mathbf{q} = \lambda \mathbf{q},$$

and hence, $\mathbf{q}_1$ is an eigenvector of the square symmetric matrix $\mathbf{S}$. Furthermore, the variance of the principal component scores $\mathbf{u}_1 = \mathbf{X}\mathbf{q}_1$ is $\mathbf{q}_1^T \mathbf{S} \mathbf{q}_1$,

which equals the eigenvalue of $\mathbf{q}_1$, and hence, as $\mathbf{u}_1$ must have maximum variance, $\mathbf{q}_1$ must be the first eigenvector of $\mathbf{S}$.

Let us postulate that the first $k$ principal component coefficient vectors $\mathbf{q}_j$, $j = 1, \ldots, k$, are the first $k$ eigenvectors of $\mathbf{S}$, such that each $j^{\text{th}}$ component coefficient vector equals the $j^{\text{th}}$ eigenvector. The variance of the principal component scores $\mathbf{u}_j = \mathbf{X}\mathbf{q}_j$ is equal to the $j^{th}$ eigenvalue of $\mathbf{S}$.

Consider the computation of $\mathbf{q}_{k+1}$, with $1 < k < r_x$. We must maximize the criterion $\mathbf{q}^T\mathbf{S}\mathbf{q}$, subject to the side conditions $\mathbf{q}^T\mathbf{q} = 1$ and $\mathbf{q}_j^T\mathbf{S}\mathbf{q} = 0$, $j = 1, \ldots, k$. Hence, we must maximize the Lagrangian equation

$$L_{k+1}(\mathbf{q}) = \mathbf{q}^T\mathbf{S}\mathbf{q} - \lambda(\mathbf{q}^T\mathbf{q} - 1) - \alpha_1\mathbf{q}_1^T\mathbf{S}\mathbf{q} - \cdots - \alpha_k\mathbf{q}_k^T\mathbf{S}\mathbf{q}.$$

As the first $k$ components are eigenvectors, the side conditions $\mathbf{q}_j^T\mathbf{S}\mathbf{q} = 0$, $j = 1, \ldots, k$, reduce to $\mathbf{q}_j^T\mathbf{q} = 0$, $j = 1, \ldots, k$, and the Lagrangian equation becomes

$$L_{k+1}(\mathbf{q}) = \mathbf{q}^T\mathbf{S}\mathbf{q} - \lambda(\mathbf{q}^T\mathbf{q} - 1) - \alpha_1\mathbf{q}_1^T\mathbf{q} - \cdots - \alpha_k\mathbf{q}_k^T\mathbf{q}.$$

Thus, $\mathbf{q}_{k+1}$ is the solution of equation

$$\frac{d}{d\mathbf{q}}L_{k+1}(\mathbf{q}) = \mathbf{S}\mathbf{q} - \lambda\mathbf{q} - \alpha_1\mathbf{q}_1 - \cdots - \alpha_k\mathbf{q}_k = 0.$$

Using $\mathbf{q}_j^T\mathbf{q} = 0$, $j = 1, \ldots, k$, we find that $\mathbf{q}_j^T\mathbf{S}\mathbf{q} = \alpha_j$, $j = 1, \ldots, k$. However, $\mathbf{q}_j^T\mathbf{S}\mathbf{q} = \mathbf{q}^T\mathbf{S}\mathbf{q}_j = 0$, $j = 1, \ldots, k$, and hence $\alpha_1 = \cdots = \alpha_k = 0$. Therefore, $\mathbf{q}_{k+1}$ satisfies the equation

$$\mathbf{S}\mathbf{q} = \lambda\mathbf{q}.$$

Hence, as for the first $k$ components, we find that $\mathbf{q}_{k+1}$ is an eigenvector of $\mathbf{S}$. Furthermore, with the same argument as used for $\mathbf{q}_1$, we find that $\mathbf{q}_{k+1}$ is the $(k+1)^{th}$ eigenvector, with variance equal to the $(k+1)^{th}$ eigenvalue of $\mathbf{S}$.

This proves by induction that the principal component coefficient vectors are the eigenvectors of the square symmetric matrix $\mathbf{S}$, such that

$$\mathbf{S} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T,$$

where $\mathbf{Q} = (\mathbf{q}_1, \ldots, \mathbf{q}_{r_x}) = [(q_{ij})]$ is the matrix of principal component coefficients with

$$\mathbf{Q}^T\mathbf{Q} = \mathbf{I},$$

and $\boldsymbol{\Lambda}$ is a diagonal matrix with the principal component variances $\lambda_1 > \cdots > \lambda_{r_x}$ on the diagonal [Krz88, page 63] [Jol86, page 3] [Seb84, page 176] [Rao73, page 592].

## Principal Component Computations

We have shown that principal component decomposition is the orthogonal decomposition of the sample covariance matrix. Hence, any method for the computation of the orthogonal decomposition of a real square symmetric matrix may be used for the computation of the principal component decomposition.

A standard method for the computation of such a decomposition is the symmetric QR algorithm, as described by Golub and Van Loan [GL89, chapter 8], Wilkinson [Wil65, chapter 8] and Wilkinson and Reinsch [WR71]. The principal component scores may be obtained from

$$\mathbf{U} = \mathbf{XQ},$$

where $\mathbf{U} = (\mathbf{u}_1, \ldots, \mathbf{u}_{r_x}) = [(u_{ij})]$ with

$$\mathbf{U}^T\mathbf{U} = \mathbf{E},$$

and $\mathbf{E}$ is the diagonal matrix with the eigenvalues $e_1 > \cdots > e_{r_x}$ of the cross-products matrix $\mathbf{C}$ on the diagonal. We have

$$\boldsymbol{\Lambda} = \frac{\mathbf{E}}{n-1},$$

or equivalently, $(\lambda_1, \ldots, \lambda_{r_x}) = (e_1, \ldots, e_{r_x})/(n-1)$ and

$$\mathbf{X} = \mathbf{UQ}^T.$$

The last equation points to an alternative method for the computation of the principal component decomposition. We may compute the singular value decomposition

$$\mathbf{X} = \mathbf{WE}^{\frac{1}{2}}\mathbf{Q}^T,$$

where $\mathbf{W} = \mathbf{XQE}^{-\frac{1}{2}}$, such that $\mathbf{W}^T\mathbf{W} = \mathbf{I}$. We have $\mathbf{U} = \mathbf{WE}^{\frac{1}{2}}$. The Golub and Reinsch method [GL89, pages 239 and 430] [GK65] [GR70] [WR71] is

an efficient algorithm for the computation of the singular value decompo-
sition that avoids the errors involved in forming the cross-products matrix.
Another attractive feature of the singular value decomposition is that the
principal component coefficients, the principal component eigenvalues and
the principal component scores may be computed simultaneously. A com-
prehensive discussion of the efficient and numerically stable computation of
the singular value decomposition may be found in the book by Golub and
Van Loan [GL89].

From a practical point of view, we can always obtain the principal com-
ponent decomposition in the regular case $(n - 1 > p)$ from the eigendecom-
position of the cross-products matrix

$$\mathbf{C} = \mathbf{QEQ}^T.$$

In the singular case $(n - 1 < p)$, it is more efficient to compute the prin-
cipal component decomposition from the eigendecomposition of the matrix
$\mathbf{XX}^T$:

$$\mathbf{XX}^T = \mathbf{WEW}^T.$$

The principal components are then obtained from

$$\mathbf{Q} = \mathbf{X}^T\mathbf{UE}^{-1} = \mathbf{X}^T\mathbf{WE}^{-\frac{1}{2}}.$$

## 2.1.2   Principal Components and Cross-validation

The earliest description of principal component decomposition dates back
to 1901 and is due to Pearson who considered a purely geometric definition
of principal components [Pea01]. This approach is essentially one of find-
ing good lower dimensional approximations to high dimensional data, using
simple linear combinations and the Euclidean distance as a goodness-of-fit
criterion. The statistical definition of principal component decomposition
as a transformation of a set of random variates is due to Hotelling [Hot33].
Thus, the above definition and derivation of principal components as a set
of uncorrelated variates may be referred to Hotelling and has since become a
standard way of introducing principal component analysis, as in Anderson's
book on multivariate statistical analysis [And58] or Jolliffe's reference book
on principal component analysis [Jol86], for example. Many different char-
acterisations and properties of principal components have since been found.

Rao's article on the 'Use and Interpretation of Principal Component Analysis in Applied Research' [Rao64] is a good summary and Jolliffe's book provides a more recent treatment with an extensive account of properties and applications.

The numerical computation of principal components has long remained a problem, limiting practical applications of principal components in multivariate analysis. The solution of this problem has depended essentially on the identification of principal component decomposition as an eigendecomposition problem. Seber, for example, has defined principal component analysis directly as the eigendecomposition of the sample covariance matrix [Seb84, page 176], ignoring the more conventional introduction of principal components in the statistical literature.

The pioneering work of Householder and Wilkinson as well as the advent of modern computer technology provided the breakthrough necessary for the efficient computation of principal component decompositions. Indeed, the development of numerical methods for principal component decomposition has gone hand in hand with developments in the computer industry. The fundamental reference text on the numerical treatment of eigenproblems is 'The Algebraic Eigenvalue Problem' by Wilkinson in 1965 [Wil65]. Many of the fundamental algorithms that stem from this work have been published in ALGOL versions in the 'Handbook of Automatic Computation, Volume 2, Linear Algebra', edited by Wilkinson and Reinsch [WR71]. The numerical computation of eigendecompositions and related topics continues to attract much attention. A recent account of the fundamentals as well as many references can be found in the second edition of 'Matrix Computations', by Golub and Van Loan [GL89].

Most of the methods described in these books have been made available through the numerical software library EISPACK [SBI+70][GBDM72]. An automatic distribution system accessible through electronic mail has been implemented with the NETLIB facility, which operates free of charge. A description of this electronic distribution service may be found in the book from Golub and Van Loan, as well as in an article from Dongarra and Grosse [DS87]. A complementary library of routines that deals with numerical methods for linear algebra is available in the LINPACK [DBMS78] package. An updated library LAPACK has now been implemented which is intended to supersede both the EISPACK and LINPACK set of routines and is available through the same channels.

LAPACK and EISPACK, as well as LINPACK, are written in FORTRAN. User friendly packages have been written which devise an easy and intuitive matrix language to make these procedures available to a wider public. GAUSS [Apt92] and MATLAB [Mat92] are the main implementations.

Cross-validation is a much more recent development in statistics. The first comprehensive treatment is due to Stone in his extensive paper on 'Cross-validatory Choice and Assessment of Statistical Predictions' [Sto74]. Many of the ideas leading to the unified theory explained in this article had been around for much longer though. Stone refers to Mosteller and Tukey [MT68] for the first accurate description of full leave-one-out cross-validation.

Given the development of principal component computations, it is not surprising that the application of full leave-one-out cross-validation to principal component analysis has been cumbersome. This is primarily due to computational problems, as the cost of naive implementations makes such calculations unpractical. A further reason is concerned with the application of Stone's paradigm to the cross-validation of principal components.

The first application of cross-validation to principal component decomposition has been in chemometrics rather than in statistics. This has been described by Svante Wold in 1976 for the SIMCA method of biased discriminant analysis [Wol76]. Wold's approach adheres to Stone's and Mosteller's paradigm in that it is based on the complete omission of observations from the data. Wold splits the data in a limited number of subgroups to reduce the workload, rather than computing a full leave-one-out analysis. Furthermore, his computations employ the NIPALS algorithm, which is essentially an implementation of the power method, rather than any of the efficient methods available at that time. This approach was subsequently generalized in an article from the same author, to the cross-validation of principal component and factor analysis models, in 1978 [Wol78]. The methodology is identical to that used for the SIMCA method, except that principal component decomposition is generalized to a framework of 'data modelling structures' which may also be applied to factor analysis. Again, computations are based on the NIPALS algorithm and the deletion of subgroups rather than individual observations.

The first and often quoted contribution in statistics is in an article from Eastment and Krzanowski, in 1982 [EK82]. This article considers the cross-validatory choice of the number of components in a principal component representation of the data and largely ignores the previous work due to Wold.

Eastment and Krzanowski take the view that for applications of principal component decomposition to obtain low dimensional representations of data, an application of Stone's paradigm must be based on the omission of individual cells from the data matrix, rather than on the complete omission of a datum. A computational scheme for this alternative cross-validatory procedure is proposed and Monte Carlo simulations are considered to verify the performance of the algorithm. Two subsequent articles by Krzanowski consider further simulations as well as an application to the Alate Adelges data [Krz83][Krz87].

This approach has largely set the tone for the statistical treatment of principal component cross-validation in statistics. Jolliffe [Jol86, page 99] follows this framework completely in a discussion on the cross-validation of principal components, as does Seber [Seb84, page 187]. An exception may be found in the paper by Stone and Brooks on continuum regression [SB90].

Most of these applications of cross-validation to principal component analysis have been concerned with the choice of the number of components necessary for an adequate representation of data. In chemometrics, the methods developed by Wold have subsequently been used in the calibration of principal component prediction equations, with similar computational approaches. Hence, these methods are based on the power algorithm and consider subsets of the data, rather than a full leave-one-out approach. A description of these methods may be found in the book from Martens and Naes on 'Multivariate Calibration' [MN89, page 254]. Essentially, the same approach to the cross-validation of principal components has been considered by Stone and Brooks in their treatment of continuum regression [SB90], but with a leave-one-out implementation. In the limit to ordinary least squares, their implementation of cross-validation will coincide with the classical implementation of cross-validation for least squares. Unfortunately, this algebra breaks down in the principal component limit.

## 2.2 Leave-One-Out Principal Component Modifications

### 2.2.1 Leave-One-Out Perturbations

The downdated principal component matrix $\mathbf{Q}_{(i)}$ for the deletion of the $i^{\text{th}}$ datum is obtained from the orthogonal eigendecomposition

$$\mathbf{C}_{(i)} = \mathbf{Q}_{(i)}\mathbf{E}_{(i)}\mathbf{Q}_{(i)}^T$$

of the downdated cross-products matrix $\mathbf{C}_{(i)} = \mathbf{X}_{(i)}^T\mathbf{X}_{(i)}$, where

$$\mathbf{X}_{(i)} = \begin{pmatrix} \tilde{\mathbf{x}}_1 - \overline{\overline{\mathbf{X}}}_{(i)} \\ \vdots \\ \tilde{\mathbf{x}}_{i-1} - \overline{\overline{\mathbf{X}}}_{(i)} \\ \tilde{\mathbf{x}}_{i+1} - \overline{\overline{\mathbf{X}}}_{(i)} \\ \vdots \\ \tilde{\mathbf{x}}_n - \overline{\overline{\mathbf{X}}}_{(i)} \end{pmatrix},$$

with $\overline{\overline{\mathbf{X}}}_{(i)}$ the mean of the leave-one-out downdated predictor data. We have

$$\mathbf{C}_{(i)} = \mathbf{C} - \nu\mathbf{x}_i^T\mathbf{x}_i,$$

with $\nu = n/(n-1)$ [Seb84, page 15] and hence, employing the principal component decomposition $\mathbf{C} = \mathbf{Q}\mathbf{E}\mathbf{Q}^T$, the leave-one-out modification problem reduces to the eigendecomposition of the matrix

$$\mathbf{E} - \rho\mathbf{z}^T\mathbf{z},$$

where $\mathbf{z} = \mathbf{x}_i\mathbf{Q}/\|\mathbf{x}_i\mathbf{Q}\|$ and $\rho = \nu\|\mathbf{x}_i\mathbf{Q}\|^2 = \nu\|\mathbf{x}_i\|^2$.

From the principal component decomposition

$$\mathbf{E} - \rho\mathbf{z}^T\mathbf{z} = \mathbf{V}\mathbf{E}_{(i)}\mathbf{V}^T,$$

we obtain the downdated matrix $\mathbf{E}_{(i)}$ with the eigenvalues $e_{1(i)}, \ldots, e_{r_x(i)}$ on the diagonal. The downdated principal components may then be derived from

$$\mathbf{Q}_{(i)} = \mathbf{Q}\mathbf{V}.$$

This determines the framework for leave-one-out principal component downdating. The downdating principal components $\mathbf{V}$ will be computed in two stages. With this terminology we distinguish the components $\mathbf{V}$ which are applied to downdate $\mathbf{Q}$ from the downdated matrix $\mathbf{Q}_{(i)}$ itself. The downdated principal component eigenvalues are computed first. The principal components $\mathbf{V}$ are then obtained as a function of the downdated eigenvalues $e_{1(i)}, \ldots, e_{r_x(i)}$, the eigenvalues $e_1, \ldots, e_{r_x}$ and the principal components $\mathbf{Q}$.

## Deflation

It is important that we study those situations for which the downdating problem is trivial, so as to identify a well-behaved problem. This will also reduce the workload in any implementation, as those components that do not downdate may be removed from the computations immediately.

We may always assume that $\rho$ is strictly positive. Indeed, $\rho = 0$ is equivalent to $\mathbf{z} = \mathbf{0}$, which corresponds to the trivial downdating problem in which none of the principal component eigenvalues and principal components are downdated. The geometric interpretation of this phenomenon is that the cross-validated observation is situated at the mean of the data. Hence this observation does not contribute to the sum of squared projections of the data on the principal component directions and, as a consequence, it can not influence any of the principal component directions.

The second situation is that in which $z_l = 0$ for some $1 \leq l \leq r_x$. In this case $e_l = e_{l(i)}$ and $\mathbf{q}_l = \mathbf{q}_{l(i)}$, as

$$(\mathbf{E} - \rho \mathbf{z}^T \mathbf{z})\mathbf{g}_l = \mathbf{E}\mathbf{g}_l = e_l \mathbf{g}_l,$$

where $\mathbf{g}_l = (0, \ldots, 0, 1, 0, \ldots, 0)$ is the $l^{\text{th}}$ standard basis vector, and hence, $\mathbf{v}_l = \mathbf{g}_l$. This follows as the ordering of the downdated components is unaffected, as may be seen from the present deflation algebra and theorem 2.2, which will be discussed subsequently. If there is a $z_m \neq 0$ for some

$l < m \leq r_x$, then we can find a permutation $\mathbf{P}(l, m)$ of $r_x$ elements,

$$
\mathbf{P}(l, m) = \begin{bmatrix}
1 & \cdots & 0 & & \cdots & & 0 & \cdots & 0 \\
\vdots & \ddots & \vdots & & & & \vdots & & \vdots \\
 & & 1 & & & & & & \\
0 & \cdots & 0 & & \cdots & & 1 & \cdots & 0 \\
 & & & 1 & & & & & \\
\vdots & & \vdots & & \ddots & & \vdots & & \vdots \\
 & & & & & 1 & & & \\
0 & \cdots & 1 & & \cdots & & 0 & \cdots & 0 \\
\vdots & & \vdots & & & & \vdots & 1 & \vdots \\
 & & & & & & & & \ddots \\
0 & \cdots & 0 & & \cdots & & 0 & \cdots & 1
\end{bmatrix}
\begin{matrix}
\\ \\ \\ l \\ \\ \\ \\ m \\ \\ \\
\end{matrix}
$$

$$\qquad\qquad\qquad l \qquad\qquad\qquad m$$

that interchanges the positions of $z_l$ and $z_m$ in $\mathbf{z}$, while leaving all other elements of $\mathbf{z}$ unchanged. This amounts to a permutation of the principal component directions, and hence, the same permutation must also be applied to the matrix of principal component coefficients $\mathbf{Q}$ as well as to the diagonal matrix of eigenvalues $\mathbf{E}$. We have that

$$\mathbf{P}(l, m)^T \mathbf{E} \mathbf{P}(l, m)$$

is a diagonal matrix. As a consequence, all zero elements of $\mathbf{z}$ may be permuted to the 'bottom' of the eigensystem.

The last case is that for which $e_l = e_m$ for some $1 \leq l < m \leq r_x$. In this case we can find a Jacobi or Givens rotation [GL89, pages 201 and 463] of $r_x$

elements,

$$
\mathbf{J}(l,m,\theta) = \begin{bmatrix} 1 & \cdots & 0 & & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & & \vdots & & \vdots \\ & & 1 & & & & & \\ 0 & \cdots & \cos(\theta) & & \cdots & \sin(\theta) & \cdots & 0 \\ & & & 1 & & & & \\ \vdots & & \vdots & & \ddots & \vdots & & \vdots \\ & & & & 1 & & & \\ 0 & \cdots & -\sin(\theta) & & \cdots & \cos(\theta) & \cdots & 0 \\ & & & & & & 1 & \\ \vdots & & \vdots & & & \vdots & & \ddots & \vdots \\ 0 & \cdots & 0 & & \cdots & 0 & \cdots & 1 \end{bmatrix} \begin{matrix} \\ \\ \\ l \\ \\ \\ \\ m \\ \\ \\ \end{matrix}
$$

$$
\qquad\qquad\qquad l \qquad\qquad\qquad m
$$

such that the $m^{\text{th}}$ component of

$$
\mathbf{z}\mathbf{J}(l,m,\theta)
$$

is zero. As for the previous case, this rotation must be applied to the diagonal matrix of eigenvalues $\mathbf{E}$ and to the matrix of eigenvectors $\mathbf{Q}$. As for the previous case, we find that

$$
\mathbf{J}(l,m,\theta)^T\mathbf{E}\mathbf{J}(l,m,\theta)
$$

is diagonal. We may now apply the previous analysis and permute the rotated $m^{\text{th}}$ dimension to the 'bottom' of the eigensystem.

From a repeated application of the above steps we may find an orthogonal transformation matrix $\mathbf{T}$ such that $\mathbf{T}^T\mathbf{E}\mathbf{T} = \text{diag}(d_1,\ldots,d_{r_x})$, with $d_1 > \ldots > d_k$ and $\boldsymbol{\zeta} = \mathbf{z}\mathbf{T}$ with $\zeta_{k+1} = \cdots = \zeta_{r_x} = 0$. We have shown that the downdating of the last $r_x - k$ components of this transformed eigensystem is trivial. The downdating of the first $k$ components is a subproblem of the same type as before, but with distinct eigenvalues and no zero principal component scores for the left-out observation. Hence, we may restrict ourselves to the decomposition of such eigensystems.

## The eigensystem of $E - \rho z^T z$

We will assume that all eigenvalues are distinct, such that $e_1 > \cdots > e_{r_x}$ and that $z$ has no zero components.

Bunch, Nielsen and Sorensen [BNS78] considered the following theorem for the computation of the matrix of downdating principal components $V = (v_1, \cdots, v_{r_x})$.

**Theorem 2.1** *If* $A$ *is a square invertible real* $n$ *by* $n$ *matrix,* $u$ *and* $v$ *are real* $n$ *by* $1$ *vectors and* $\rho$ *is a real number different from zero, then the statements*

$$(A + \rho u v^T)x = b,$$

$$\begin{bmatrix} A & u \\ v^T & -\frac{1}{\rho} \end{bmatrix} \begin{bmatrix} x \\ \theta \end{bmatrix} = \begin{bmatrix} b \\ 0 \end{bmatrix}$$

*and*

$$x = A^{-1}b - \theta A^{-1}u, \quad \xi\theta = v^T A^{-1}b, \quad with \ \xi = \frac{1}{\rho} + v^T A^{-1}u,$$

*are equivalent.*

This theorem originates from the original paper by Bunch and Rose [BR74] on the 'partitioning, tearing and modification of sparse linear systems'. An application of this theorem shows that $V$ may be computed directly without any further iterative computations, such as inverse iteration [Wil65][WR71, page 418][PW79][GL89, page 383], once the downdated eigenvalues have been obtained.

Consider the computation of component $v_k$. Deflation ensures that $D_k = E - e_{k(i)}I$ is invertible, or equivalently, $E$ and $E - \rho z^T z$ do not have any eigenvalues in common. Indeed, if we were to assume that $e_{k(i)}$ was an eigenvalue of $E$, for instance, then $e_{k(i)} = e_l$, for some $1 \le l \le r_x$ and

$$0 = g_l[(E - e_{k(i)}I)v_k - \rho(zv_k)z^T] = -\rho(zv_k)z_l.$$

As $\rho$ and $z_l$ are non-zero, we have $zv_k = 0$ and therefore $Ev_k = e_{k(i)}v_k$. However, as $E$ has distinct eigenvalues, this implies that $v_k$ is in the vector space spanned by $g_l$ and hence, $z_l = zv_k = 0$, which is in contradiction to the assumptions.

We may now apply theorem 2.1 to the equation

$$(D_k - \rho z^T z)v_k = 0,$$

with $\mathbf{A} = \mathbf{D}_k$, $\mathbf{u} = \mathbf{v} = \mathbf{z}^{\mathbf{T}}$ and $\mathbf{b} = \mathbf{0}$. We find that the equations of theorem 2.1 are consistent if and only if $\xi = 0$ and $\theta$ arbitrary. Hence, we find that

$$\mathbf{v}_k \propto (\mathbf{E} - e_{k(i)}\mathbf{I})^{-1}\mathbf{z}^T,$$

and therefore

$$\mathbf{v}_k = \frac{(\mathbf{E} - e_{k(i)}\mathbf{I})^{-1}\mathbf{z}^T}{\|(\mathbf{E} - e_{k(i)}\mathbf{I})^{-1}\mathbf{z}^T\|},$$

as $\mathbf{v}_k^T \mathbf{v}_k = 1$.

A simpler argument that avoids the use of the above theorem has been considered by Golub and Van Loan [GL89, page 462]. Each downdating eigenpair $(e_{k(i)}, \mathbf{v}_k)$ satisfies the equation

$$(\mathbf{E} - \rho\mathbf{z}^T\mathbf{z})\mathbf{v}_k = e_{k(i)}\mathbf{v}_k,$$

and hence,

$$(\mathbf{E} - e_{k(i)}\mathbf{I})\mathbf{v}_k - \rho(\mathbf{z}\mathbf{v}_k)\mathbf{z}^T = 0.$$

As we work with deflated eigensystems, $\mathbf{E} - e_{k(i)}\mathbf{I}$ is invertible and $\mathbf{z}\mathbf{v}_k \neq 0$. Hence, $\mathbf{v}_k$ is in the span of

$$(\mathbf{E} - e_{k(i)}\mathbf{I})^{-1}\mathbf{z}^T$$

and we find

$$\mathbf{v}_k = \frac{(\mathbf{E} - e_{k(i)}\mathbf{I})^{-1}\mathbf{z}^T}{\|(\mathbf{E} - e_{k(i)}\mathbf{I})^{-1}\mathbf{z}^T\|}$$

as before.

The downdated principal component coefficients may now be computed from

$$\mathbf{q}_{k(i)} = \mathbf{Q}\mathbf{v}_k,$$

and hence, the downdating principal component coefficients $\mathbf{v}_j$, $j = 1, \ldots, r_x$ are obtained once the downdated eigenvalues $e_{1(i)}, \cdots e_{r_x(i)}$ have been computed. The downdating of the principal component eigenvalues is therefore central to the leave-one-out cross-validation of principal components.

The following theorem is fundamental to the computation of the downdated principal component eigenvalues.

**Theorem 2.2** *The eigenvalues of* $\mathbf{E} - \rho \mathbf{z}^T \mathbf{z}$, *with* $\mathbf{E} = \mathrm{diag}(e_1, \cdots, e_{r_x})$, $e_1 > \cdots > e_{r_x}$, $\rho \neq 0$ *and* $\|\mathbf{z}\| = 1$, *such that no components of* $\mathbf{z}$ *are zero, are*

$$e_{j(i)} = e_j - \rho \mu_j, j = 1, \ldots, r_x,$$

*with* $\sum_{j=1}^{r} \mu_j = 1$ *and* $0 < \mu_j < 1$, $j = 1, \ldots, r_x$. *The eigenvalues of* $\mathbf{E} - \rho \mathbf{z}^T \mathbf{z}$ *strictly separate those of* $\mathbf{E}$, *such that*

$$e_1 > e_{1(i)} > e_2$$
$$\vdots$$
$$e_j > e_{j(i)} > e_{j+1}$$
$$\vdots$$
$$e_{r_x-1} > e_{r_x-1(i)} > e_{r_x}$$
$$e_{r_x} > e_{r_x(i)} > e_{r_x} - \rho.$$

Whereas the direct computation of the downdating principal component coefficients is a fairly recent discovery, the properties of leave-one-out eigenvalue downdating have been known for much longer. One of the first comprehensive descriptions may be found in Wilkinson's book on 'The Algebraic Eigenvalue Problem', which considers the eigenvalues of a sum of two symmetric matrices in the context of perturbation theory [Wil65, page 95]. For a general sum of two symmetric matrices only a minimax characterisation is derived. For the perturbation of a symmetric matrix by another of rank one, the downdated eigenvalues may be expressed as the solutions of a characteristic polynomial. This polynomial has a simple structure which can be exploited to obtain the properties of the perturbed eigenvalues. The above theorem is a summary of these results for the downdating case.

The same findings have also been considered by Thompson [Tho76] as well as by Golub [Gol73] and Bunch, Nielsen and Sorensen [BNS78]. A good summary of this theory may be found in the book from Golub and Van Loan [GL89, page 462]

The quantities $\rho \mu_j$, $j = 1, \ldots, r_x$, in the above theorem may be referred to as the eigenvalue downdates. The downdated variances $\lambda_{j(i)}$ are obtained from

$$\lambda_{j(i)} = \frac{n-1}{n-2} \lambda_j - \frac{\rho \mu_j}{n-2}.$$

## 2.2.2 Efficient Rank-One Principal Component Variance Perturbations

As before, we may assume that all eigenvalues are distinct, such that $e_1 > \cdots > e_{r_x}$ and that $\mathbf{z}$ has no zero components.

### The secular equation

The downdated eigenvalues are those values $\lambda$ for which the equation

$$(\mathbf{E} - \rho \mathbf{z}^T \mathbf{z})\mathbf{v} = \lambda \mathbf{v}$$

has a non-trivial solution. Hence, the downdated eigenvalues are the solutions of the equation

$$f(\lambda) = \det(\mathbf{E} - \rho \mathbf{z}^T \mathbf{z} - \lambda \mathbf{I}) = 0.$$

This equation is ill-conditioned due to catastrophic cancellation [Gol73][GL89, page 62]. An equivalent function that is suitable for computation can be found by rearranging the equation. We have

$$f(\lambda) = \det(\mathbf{E} - \lambda \mathbf{I}) \det(\mathbf{I} - \rho(\mathbf{E} - \lambda \mathbf{I})^{-1} \mathbf{z}^T \mathbf{z})$$

$$= \prod_{j=1}^{r_x}(e_j - \lambda)\left(1 - \rho \sum_{j=1}^{r_x} \frac{z_j^2}{e_j - \lambda}\right),$$

as $\mathbf{z}^T \mathbf{z}$ is a matrix of rank one [Seb84, page 518]. Furthermore, because of deflation, the downdated eigenvalues are strictly smaller than the eigenvalues $e_1, \ldots, e_{r_x}$ and hence, the downdated eigenvalues are the solutions of the equation

$$\omega(\lambda) = 1 - \rho \sum_{j=1}^{r_x} \frac{z_j^2}{e_j - \lambda} = 0.$$

The function $\omega(\lambda)$ has been referred to as the secular function by Golub [Gol73, page 327]. The same function was described by Wilkinson [Wil65, page 96] in the related context of bordered diagonal matrices and perturbed symmetric matrices.

When downdating the $i^{\text{th}}$ eigenvalue, the secular equation may be rewritten as a function of the proportion change $\mu$ from the total eigenvalue downdate $\rho$ for the principal component eigenvalue $e_i$:

$$\omega_i(\mu) = 1 + \sum_{j=1}^{r_x} \frac{z_j^2}{\delta_j - \mu} = 0,$$

where $\delta_j = (e_i - e_j)/\rho$, $j = 1, \ldots, r_x$ and $\lambda = e_i - \rho\mu$. Hence, the modification of the eigenvalues is obtained by solving the secular equation for the proportion change.

Bunch, Nielsen and Sorensen proposed a variant of the Newton method to solve the secular equation, in the following manner. We may rewrite the secular equation as

$$\omega_i(\mu) = 1 + \psi_i(\mu) + \phi_i(\mu),$$

where $\psi_i(\mu) = \sum_{j=1}^{i} \frac{z_j^2}{\delta_j - \mu}$ and $\phi_i(\mu) = \sum_{j=i+1}^{r_x} \frac{z_j^2}{\delta_j - \mu}$. Hence, we must solve the equation

$$-\psi_i(\mu) = 1 + \phi_i(\mu)$$

for $\mu$, with $0 < \mu < \min(1 - \sum_{j=1}^{i-1} \mu_j, \delta)$, where $\delta = \delta_{i+1}$. The function $-\psi_i(\mu)$ is a decreasing function on this interval, while $1 + \phi_i(\mu)$ is increasing. Both functions are positive and convex. They intersect at the solution.

This geometry is exploited to obtain an efficient method for the approximation of the roots. The dominant terms of the functions $\psi_i$ and $\phi_i$ in the interval $[0, \delta_{i+1}]$ are $z_i^2/(\delta_i - \mu)$ and $z_{i+1}^2/(\delta_{i+1} - \mu)$ respectively. These are simple rational functions and hence, an efficient procedure of the Newton-Raphson type may be obtained from rational function approximations to the true functions. The computations proceed from an initial approximation $t_0$ to the true value. Successive corrections to the true value are then calculated with these rational function approximations at each current approximation $t_\kappa$. The above equation is solved with the rational function approximations and a new approximation $t_{\kappa+1}$ is obtained. Two different choices for the rational functions have been considered in the literature.

## The quadratic method

The quadratic method approximates the true functions $\psi_i$ and $\phi_i$ with the rational functions

$$\widehat{\psi}_i(t) = \frac{l}{m - t} \quad \text{for} \quad \psi_i,$$

and

$$\widehat{\phi}_i(t) = r + \frac{s}{\delta - t} \quad \text{for} \quad \phi_i,$$

at each current approximation $t_\kappa$. This choice is due to Bunch, Nielsen and Sorensen in their original implementation of the methodology [BNS78]. From

the usual first order fitting conditions, we find that

$$l = \frac{\psi_i^2(t_\kappa)}{\psi_i'(t_\kappa)},$$

$$m = t_\kappa + \frac{\psi_i(t_\kappa)}{\psi_i'(t_\kappa)},$$

$$r = \phi_i(t_\kappa) - \Delta\phi_i'(t_\kappa),$$

and

$$s = \Delta^2 \phi_i'(t_\kappa),$$

where $\Delta = \delta - t_\kappa$. A new approximation $t_{\kappa+1}$ is obtained by solving the quadratic equation

$$\frac{-l}{m-t} = 1 + r + \frac{s}{\delta - t}.$$

Bunch, Nielsen and Sorensen [BNS78] have shown that the correct iteration formula may be written as

$$t_{\kappa+1} = t_\kappa + \frac{2b}{a + \sqrt{a^2 - 4b}},$$

where

$$a = (\Delta(1 + \phi_i(t_\kappa)) + \frac{\psi_i^2(t_\kappa)}{\psi_i'(t_\kappa)})/c + \frac{\psi_i(t_\kappa)}{\psi_i'(t_\kappa)},$$

$$b = \frac{\Delta\omega_i(t_\kappa)\psi_i(t_\kappa)}{\psi_i'(t_\kappa)c},$$

with

$$c = 1 + \phi_i(t_\kappa) - \Delta\phi_i'(t_\kappa).$$

The details and reasons for arranging the equations in this way have been explained in the original paper from Bunch, Nielsen and Sorensen. They are concerned only with the numerical stability of the resulting algorithm. Bunch, Nielsen and Sorensen provide extensive proof that this optimisation scheme has a guaranteed convergence to the solution if the starting value $t_0$ is chosen in the interval $]0, \mu_i[$. Such a value may be obtained from a solution to the equation

$$\frac{z_i^2}{t} = \frac{z_{i+1}^2}{\delta - t} + 1 + \sum_{j \neq i, i+1} \frac{z_j^2}{\delta_j - \delta}.$$

A numerically stabilized formula for the correct solution is

$$t_0 = \frac{2\delta z_i^2}{a + \sqrt{a^2 - 4\delta z_i^2(1 + \sum\limits_{j \neq i, i+1} \frac{z_j^2}{\delta_j - \delta})}},$$

where $a = z_i^2 + z_{i+1}^2 + \delta(1 + \sum\limits_{j \neq i, i+1} \frac{z_j^2}{\delta_j - \delta})$. The initial value for the last eigenvalue is special and may be obtained from a solution to the equation

$$\frac{z_{r_x}^2}{t} = 1 + \sum_{j=1}^{r_x} \frac{z_j^2}{\delta_j - 1}.$$

The solution is

$$t_0 = \frac{z_{r_x}^2}{1 + \sum_{j=1}^{r_x} z_j^2 / (\delta_j - 1)}.$$

The iteration for the last eigenvalue simplifies to

$$t_{\kappa+1} = t_\kappa + \left(\frac{1 + \psi_i(t_\kappa)}{\psi_i'(t_\kappa)}\right) \psi_i(t_\kappa)$$

as $\phi \equiv 0$.

**The linear method**

The linear method considers the approximations

$$\widehat{\psi}_i(t) = \frac{l}{m - t} \quad \text{for} \quad \psi_i,$$

and

$$1 + \widehat{\phi}_i(t) = \frac{r}{s - t} \quad \text{for} \quad 1 + \phi_i,$$

at each current approximation $t_\kappa$. This choice is due to DeGroat [DeG90]. A first order approximation to the true functions is obtained when

$$l = \frac{\psi_i^2(t_\kappa)}{\psi_i'(t_\kappa)},$$

$$m = t_\kappa + \frac{\psi_i(t_\kappa)}{\psi_i'(t_\kappa)},$$

$$r = \frac{(1 + \phi_i(t_\kappa))^2}{\phi_i'(t_\kappa)},$$

and

$$s = t_\kappa + \frac{1 + \phi_i(t_\kappa)}{\phi'_i(t_\kappa)}.$$

A new approximation $t_{\kappa+1}$ may then be found from the solution of the linear equation

$$\frac{-l}{m-t} = \frac{r}{s-t}.$$

The solution is

$$t_{\kappa+1} = t_\kappa + \frac{\omega_i(t_\kappa)(1 + \phi_i(t_\kappa))\psi_i(t_\kappa)}{(1 + \phi_i(t_\kappa))^2 \psi'_i(t_\kappa) + \psi_i^2(t_\kappa)\phi'_i(t_\kappa)}.$$

DeGroat shows that this updating scheme will converge when an appropriate initial value is chosen. The correct initial value is the solution of the equation

$$1 + \phi_i(0) = -\psi_i(\delta) + \frac{z_i^2}{t} - \frac{z_i^2}{\delta}, \quad \text{if} \quad 1 + \phi_i(0) \geq -\psi_i(\delta),$$

and

$$1 + \phi_i(0) = -\psi_i(\delta) + \frac{z_{i+1}^2}{t - \delta} + \frac{z_{i+1}^2}{\delta}, \quad \text{if} \quad 1 + \phi_i(0) < -\psi_i(\delta).$$

Hence, the initial value may be found from

$$t_0 = \frac{\delta z_i^2}{z_i^2 + \delta a}, \quad \text{if} \quad a \geq 0,$$

and

$$t_0 = \frac{\delta}{1 - z_{i+1}^2/\delta a}, \quad \text{if} \quad a < 0,$$

where $a = 1 + \phi_i(0) + \psi_i(\delta)$. As before, the rationale behind this choice is in the geometry of the approximating rational functions, as has been described by DeGroat. The last eigenvalue is special. The initial value is the solution of the equation

$$1 = -\psi_{r_x}(1) + \frac{z_{r_x}^2}{t} - z_{r_x}^2,$$

The solution is

$$t_0 = \frac{z_{r_x}^2}{1 + \sum_{j=1}^{r_x} z_j^2/(\delta_j - 1) + z_{r_x}^2}.$$

Iterations for the last eigenvalue are as for the quadratic method.

## Convergence

Both the quadratic and the linear method have a quadratic rate of convergence, such that,

$$| t_{\kappa+1} - \mu_i | < \tau | t_\kappa - \mu_i |^2,$$

with $\tau$ a positive constant of order unity and independent of the iteration. The proof is in the literature [BNS78][DeG90].

## 2.2.3  Implementation

### Deflation

To implement the deflation analysis, we need to specify criteria to decide when an element $z_l$ of $z$ is close to zero and when two adjacent eigenvalues $e_l$ and $e_{l+1}$ are close enough to be considered equal.

To solve this problem, Dongarra and Sorensen [DG87] ask the question when an eigenpair of $\mathbf{E}$ may be considered a good approximation of the corresponding eigenpair of $\mathbf{E} - \rho z^T z$. Hence, we consider

$$\|(\mathbf{E} - \rho z^T z)\mathbf{g}_l - e_l \mathbf{g}_l\| = |\rho z_l| \|z\| = \rho |z_l|.$$

This leads to a criterion of the type

$$|z_l| < \eta(e_1 + \rho)/\rho,$$

where $\eta$ is the machine precision. In contrast, Golub and Van Loan [GL89, page 463] suggest the criterion

$$|z_l| < \eta(e_1 + \rho).$$

We will employ this criterion and similarly, for the comparison of the eigenvalues, we will use

$$e_l - e_{l+1} < \eta(e_1 + \rho).$$

### Cancellation

It is paramount that numeric cancellation [GL89, page 62] is minimized in the computation of the eigenvalue downdates as well as in the eigenvector

downdating. To obtain stable numerical computations, we may exploit the normalization of the eigenvalue downdating problem. Let us write

$$\mathbf{D}_k = -\rho \mathbf{F}_k,$$

where $\mathbf{F}_k = \text{diag}(\Delta_1^k, \cdots, \Delta_{r_x}^k)$ with $\Delta_j^k = \delta_j - \mu_k$. The vector $\mathbf{F}_k^{-1}\mathbf{z}^T$ has the same direction as $\mathbf{D}_k^{-1}\mathbf{z}^T$ and hence, the downdated eigenvectors may also be obtained from

$$\mathbf{v}_k = \frac{\mathbf{F}_k^{-1}\mathbf{z}^T}{\|\mathbf{F}^{-1}\mathbf{z}^T\|}$$

and

$$\mathbf{q}_{k(i)} = \mathbf{Q}\mathbf{v}_k.$$

This will minimize cancellation in the iterations, as the successive eigenvalue updates may be applied directly to each current approximation of $\mathbf{F}_k$, and hence we do not have to compute the downdated eigenvalues explicitly. Indeed, at each iteration, we only compute the update $\Delta(t_\kappa)$, such that $t_{\kappa+1} = t_\kappa + \Delta(t_\kappa)$ and hence, these updates may be applied to the current approximation $\mathbf{F}_{k\kappa}$ of $\mathbf{F}_k$, such that

$$\mathbf{F}_{k\kappa+1} = \mathbf{F}_{k\kappa} - \Delta(t_\kappa)$$

with $\mathbf{F}_{k0} = t_0 - \Delta_j$. If necessary, the downdated eigenvalues may be computed separately, from the updating sequence

$$g_{\kappa+1} = g_\kappa - \Delta(t_\kappa)$$

with $g_0 = e_k/\rho - t_0$. If $t_\kappa$ is a good approximation of $\mu_k$, then $\rho g_\kappa$ will approximate $e_{k(i)}$ with the same relative accuracy.

### Convergence Criteria

We must decide on an appropriate stopping rule for the approximation of each eigenvalue downdate. A good stopping rule is to accept the current approximation $t_\kappa$ of the normalized downdated eigenvalue when

$$\frac{\|\mathbf{C}_{(i)}\mathbf{q}_{k(i)\kappa} - e_{k(i)\kappa}\mathbf{q}_{k(i)\kappa}\|}{\|\mathbf{C}_{(i)}\|_2} < \eta.$$

We have

$$\|\mathbf{C}_{(i)}\mathbf{q}_{k(i)_\kappa} - e_{k(i)_\kappa}\mathbf{q}_{k(i)_\kappa}\| = \|\mathbf{Q}^T\mathbf{C}_{(i)}\mathbf{Q}\mathbf{Q}^T\mathbf{q}_{k(i)_\kappa} - e_{k(i)_\kappa}\mathbf{Q}^T\mathbf{q}_{k(i)_\kappa}\|$$
$$= \|(\mathbf{E} - \rho\mathbf{z}^T\mathbf{z})\mathbf{v}_{k_\kappa} - e_{k(i)_\kappa}\mathbf{v}_{k_\kappa}\|.$$

Hence, we must ensure that

$$\|(\mathbf{E} - \rho\mathbf{z}^T\mathbf{z})\mathbf{v}_{k_\kappa} - e_{k(i)_\kappa}\mathbf{v}_{k_\kappa}\| < \eta\|\mathbf{E} - \rho\mathbf{z}^T\mathbf{z}\|_2.$$

Utilising the previous updating scheme, we find

$$\|(\mathbf{E} - \rho\mathbf{z}^T\mathbf{z})\mathbf{v}_{k_\kappa} - e_{k(i)_\kappa}\mathbf{v}_{k_\kappa}\| = \|\mathbf{E}\mathbf{v}_{k_\kappa} - \rho(\mathbf{z}\mathbf{v}_{k_\kappa})\mathbf{z}^T - (e_k - \rho t_\kappa)\mathbf{v}_{k_\kappa}\|$$
$$= \frac{\rho\|\mathbf{F}_{k_\kappa}(\mathbf{F}_{k_\kappa}^{-1}\mathbf{z}^T) + (\mathbf{z}^T\mathbf{F}_{k_\kappa}^{-1}\mathbf{z})\mathbf{z}^T\|}{\|\mathbf{F}_{k_\kappa}^{-1}\mathbf{z}^T\|}$$
$$= \frac{\rho\|\mathbf{z}\|}{\|\mathbf{F}_{k_\kappa}^{-1}\mathbf{z}^T\|} \mid 1 + \mathbf{z}^T\mathbf{F}_{k_\kappa}^{-1}\mathbf{z} \mid,$$

and hence, as $1 + \mathbf{z}^T\mathbf{F}_{k_\kappa}^{-1}\mathbf{z} = 1 + \psi(t_\kappa) + \phi(t_\kappa)$, $\|\mathbf{z}\| = 1$ and $\|\mathbf{F}_{k_\kappa}^{-1}\mathbf{z}^T\|^2 = \psi'(t_\kappa) + \phi'(t_\kappa)$, we must have

$$\mid 1 + \psi(t_\kappa) + \phi(t_\kappa) \mid < \eta\frac{\|(\mathbf{E} - \rho\mathbf{z}^T\mathbf{z})\|_2}{\rho}\sqrt{\psi'(t_\kappa) + \phi'(t_\kappa)}.$$

This convergence criterion may be employed for both the quadratic and the linear methods, as the quantities $\omega(t_\kappa)$, $\phi'(t_\kappa)$, and $\psi'(t_\kappa)$ are calculated anyway in the iteration. Both the quadratic and linear methods have a quadratic rate of convergence, and hence, the criterion

$$\mid t_{\kappa+1} - t_\kappa \mid < \varepsilon \mid t_\kappa \mid$$

may be employed as a second test on convergence, where $\varepsilon$ is the required relative accuracy.

Hence, a good stopping rule is to accept $t_\kappa$ if

$$\mid 1 + \psi(t_\kappa) + \phi(t_\kappa) \mid < \tau\eta\frac{v}{\rho}\sqrt{\psi'(t_\kappa) + \phi'(t_\kappa)},$$

and

$$\mid t_{\kappa+1} - t_\kappa \mid < \varepsilon \mid t_\kappa \mid,$$

with $v = \max(1, e_1 - \rho)$, as $\|\mathbf{E}\|_2 = e_1$.

## 2.3 Principal Component Regression and Influence Measures

### 2.3.1 Principal Component Regression

The application of full leave-one-out principal component cross-validation to the computation of the PRESS statistic in principal component regression proceeds along the following lines.

We must compute the cross-validated predicted values $\widehat{\widetilde{y}}_{i(i,\gamma)}$, for $i = 1, \ldots, n$ and $0 \le \gamma \le r_x$. We have

$$\widehat{\widetilde{y}}_{i(i,\gamma)} = \overline{\overline{\mathbf{Y}}}_{(i)} + \sum_{k=1}^{\gamma} (\widehat{\widetilde{y}}_{i(i,k)} - \widehat{\widetilde{y}}_{i(i,k-1)}) \quad \text{for} \quad 0 < \gamma \le r_x,$$

and

$$\widehat{\widetilde{y}}_{i(i,\gamma)} = \overline{\overline{\mathbf{Y}}}_{(i)} \quad \text{for} \quad \gamma = 0.$$

The principal components are uncorrelated and hence,

$$\widehat{\widetilde{y}}_{i(i,k)} - \widehat{\widetilde{y}}_{i(i,k-1)} = \frac{\sum_{j \ne i} u_{jk(i)} y_{j(i)}}{\sum_{j \ne i} u_{jk(i)}^2} u_{ik(i)}, \tag{2.1}$$

where $y_{j(i)} = \widetilde{y}_j - \overline{\overline{\mathbf{Y}}}_{(i)}$. We have

$$\widetilde{\mathbf{x}}_j - \overline{\overline{\mathbf{X}}}_{(i)} = \mathbf{x}_j + \sum_{l=1}^{p} \widetilde{\mathbf{x}}_l/n - \sum_{l \ne i} \widetilde{\mathbf{x}}_l/(n-1)$$

$$= \mathbf{x}_j + ((n-1) \sum_{l=1}^{p} \widetilde{\mathbf{x}}_l - n \sum_{l \ne i} \widetilde{\mathbf{x}}_l)/(n(n-1))$$

$$= \mathbf{x}_j + \widetilde{\mathbf{x}}_i/(n-1) - \sum_{l=1}^{p} \widetilde{\mathbf{x}}_l/(n(n-1))$$

$$= \mathbf{x}_j + \widetilde{\mathbf{x}}_i/(n-1) - \overline{\overline{\mathbf{X}}}/(n-1)$$

$$= \mathbf{x}_j + \mathbf{x}_i/(n-1),$$

for any $j = 1, \ldots, n$, and hence,

$$u_{ik(i)} = \nu \mathbf{x}_i \mathbf{q}_{k(i)}$$

for $j = i$ and

$$u_{jk(i)} = \mathbf{x}_j \mathbf{q}_{k(i)} + \mathbf{x}_i \mathbf{q}_{k(i)} / (n - 1)$$

if $j \neq i$. This allows us to simplify the computations further, as

$$\mathbf{X}\mathbf{q}_{k(i)} = \mathbf{X}\mathbf{Q}\mathbf{v}_k = \mathbf{U}\mathbf{v}_k,$$

and hence, we do not need to downdate the principal component coefficients explicitly in principal component regression. Instead, we apply the downdating principal component coefficients $\mathbf{v}_k$ directly to the principal component scores $\mathbf{U}$. Also, we have

$$e_{k(i)} = \sum_{j \neq i} u_{jk(i)}^2,$$

and hence, we do not need to compute the denominators (equation 2.1) explicitly if we maintain the downdated eigenvalues in the downdating algorithm.

## 2.3.2   Influence Measures

The concept of influence is relatively new in statistics and is intrinsically linked to the notion of full leave-one-out cross-validation. Indeed, in the discussion of Stone's paper on cross-validation [Sto74], Atkinson remarks that it 'should be possible to extract some information on the presence of outliers' in leave-one-out cross-validatory calculations.

In principal component decomposition, the analysis of influence poses serious computational difficulties. The following quote from Jolliffe [Jol86, page 187] more or less sums up the state of affairs to date.

> The intuitive definition of the influence of an observation on a statistic such as the $k^{\text{th}}$ eigenvalue $[e_k]$, or eigenvector $[\mathbf{q}_k]$, of a sample covariance matrix is simply the change in $[e_k]$ (or $[\mathbf{q}_k]$), perhaps normalized in some way, when the observation is deleted from the sample. The problem with influence defined in this way is that there may be no closed form for the influence function, and it needs to be computed afresh for each different sample.

It is of interest that Jolliffe considers the leave-one-out perturbations of the principal component eigenvalues. From a heuristic point of view, these

measures seem appropriate, as the principal components are defined so as to sequentially maximize sums of squared projections, subject to orthogonality constraints. This approach is vindicated by the present cross-validatory algebra which identifies the principal component eigenvalue perturbations as fundamental in the leave-one-out cross-validatory algebra. Moreover, these influence measures have been identified as downdates with intuitively appealing properties. Hence, the eigenvalue downdates emerge as natural influence measures in a principal component analysis.

Due to the computational costs of naive implementations of leave-one-out principal component cross-validation, the exact computation of principal component influence measures is cumbersome. For this reason, statisticians have concentrated on the derivation of adequate approximations for principal component influence measures instead. Such derivations have been based on a theoretical influence function, which postulates a form for the influence of an observation, as explained by Radhakrishnan and Kshirsagar [RK81], for example. With this approach, approximate influence measures have been obtained for the principal component eigenvalues, by Radhakrishnan and Kshirsagar as well as Critchley [Cri85], for the decomposition of the covariance matrix. For scaled analyses, based on the decomposition of the correlation matrix, Jolliffe [Jol86, page 187] points to work carried out by Calder [Cal84].

The present approach shows that, for the case of the covariance matrix, principal component influence measures may be derived with an efficient procedure and hence we need not worry about the lack of closed form solutions. Furthermore, we may exploit the eigenvalue downdating properties to refine and extend the principal component influence analysis suggested by Jolliffe.

From theorem 2.2, we know that the total principal component eigenvalue downdate

$$\rho = \sum_{j=1}^{r_x} (e_j - e_{j(i)})$$

is a sum of influence measures, and hence it defines an influence measure in its own right. We have

$$\rho = \nu \|\mathbf{x}_i \mathbf{Q}\|^2 = \nu \|\mathbf{x}_i\|^2$$

and therefore, these measures are easily obtained, without any need to resort to a Newton-Raphson-type approximation method. In essence, the total

eigenvalue downdate $\rho$ is determined by the Euclidean distance $\|\mathbf{x}_i\|$ of the downdated observation from the sample mean, $\nu$ being a cross-validatory constant. Hence, the length of an observation is a principal component influence measure. This explains the geometric intuition that observations that are far removed from the sample mean of the data have a large influence on the principal component decomposition, irrespective of their geometric orientation. We may always consider this influence measure, whether the principal component decomposition is used in the context of principal component regression or not.

Theorem 2.2 shows that the total principal component downdate $\rho$ must be complemented with the proportion principal component eigenvalue downdates $\mu_1, \ldots, \mu_{r_x}$ to describe the principal component eigenvalue downdating completely. This provides us with the obvious normalisation of the principal component eigenvalue downdates. Essentially, we may distinguish two sets of principal component eigenvalue influence measures. The first is the total downdate $\rho$. The second set of statistics are the proportion change statistics

$$\mu_j = \frac{e_j - e_{j(i)}}{\rho}, j = 1, \ldots, r_x.$$

We may wonder what statistics would be appropriate to study the influence of observations on the principal component directions. The choice of the influence measures is less clear cut here. The cross-validatory algebra shows that the eigenvalue downdating is crucial in the cross-validatory computations. Nevertheless, the downdated eigenvalues are not sufficient on their own to compute the downdated principal component coefficients. These are obtained from linear combinations of the principal component coefficients $\mathbf{Q}$. Hence, an influence analysis of the principal component directions requires a choice of the features of the eigenvector downdating that we wish to study. The most obvious measure seems to be the angle between each principal component direction $\mathbf{q}_k$ and the downdated direction $\mathbf{q}_{k(i)}$, as has already been considered by Pack *et al.* [PJM88] in the decomposition of the correlation matrix. We have

$$\cos(\theta_{k(i)}) = \mathbf{q}_k^T \mathbf{q}_{k(i)},$$

with $\theta_{k(i)}$ the angle between $\mathbf{q}_k$ and $\mathbf{q}_{k(i)}$, as $\|\mathbf{q}_k\| = \|\mathbf{q}_{k(i)}\| = 1$. The angle between $\mathbf{q}_k$ and $\mathbf{q}_{k(i)}$ is identical to the angle between $\mathbf{g}_k$ and $\mathbf{v}_k$, as $\mathbf{Q}$ defines

an orthonormal transformation and hence

$$\cos(\theta_{k(i)}) = \frac{z_k/(e_k - e_{k(i)})}{\sqrt{\sum_{j=1}^{r_x} z_j^2/(e_j - e_{k(i)})^2}}.$$

The norm $\|\mathbf{q}_k - \mathbf{q}_{k(i)}\|$ is an equivalent measure, as

$$\|\mathbf{q}_k - \mathbf{q}_{k(i)}\|^2 = 2 - 2\cos(\theta_{k(i)}),$$

and hence, it is a function of $\cos(\theta_{k(i)})$ only.

## 2.3.3 Computational Cost and Performance

We compare the computational cost of the efficient algorithms for principal component cross-validation with the costs involved in naive implementations of leave-one-out cross-validation. With respect to the theoretical considerations on the computational cost of the algebra for the efficient methods (discussed subsequent to the results on simulations and real data), we ignore the initial computations on the complete data in this comparison, as these are only done once.

**Simulations**

The linear and quadratic methods of efficient leave-one-out cross-validation for principal component regression have been implemented in the numerical analysis package GAUSS [Apt92] (appendix, page 118). The performance of these implementations has been compared with an implementation of the naive approach, in the same package, using simulated data.

The implementation of the naive method was based on the native GAUSS procedure EIGHV (GAUSS Command Reference [Apt92, volume 2, page 1180]) for the orthogonal decomposition of a square symmetric matrix. The regression coefficients of the downdated principal component scores were calculated with the formula

$$b_{j(i)} = \frac{\mathbf{u}_{j(i)}^T \mathbf{Y}_{(i)}}{\mathbf{u}_{j(i)}^T \mathbf{u}_{j(i)}}, \quad j = 1, \ldots, \gamma.$$

This formula reduces the computations to a minimum for the naive approach, as we exploit the fact that the principal component scores are uncorrelated.

Hence, each subsequent value of the PRESS statistic can be computed from a correction to the previous value, as is done in the efficient method.

Data were generated with the GAUSS pseudo-random number generator RNDU (GAUSS Command Reference [Apt92, volume 2, page 1452]), which simulates numbers from the uniform distribution. The distribution itself is irrelevant to the present purposes, as we only require these simulations to assess the computation times of the implementations. We are not concerned with the statistical properties of the methods involved. These simulations were run on an IBM 70/386 PC equipped with 4 Mb RAM and running at 25 mhz. Timing was done with the GAUSS procedure HSEC (GAUSS Command Reference [Apt92, volume 2, page 1242]). In these simulations, the number of observations was chosen to be small and kept fixed, for both the regular and the singular case, so that the computation times will primarily reflect the costs due to the number of predictor variables.

For the regular case $(n-1 \geq p)$, matrices of 100 observations by 5, 10, 15, 20, 25, 30, 35, 40, 45 and 50 predictor variables were generated. For the responses, a vector of 100 uniformly distributed numbers was generated for each simulated matrix of predictor variables. The three implementations of full leave-one-out cross-validation were then applied to each of these simulated regression problems, each time cross-validating all principal components and calculating the PRESS statistics for all factors. Figure 2.1 shows the computation times in seconds for each method versus the number of predictor variables. Results are plotted on the log scale for the number of predictor variables, so as to enhance the effect of the different methods when a small number of predictor variables is generated. The naive method starts to perform worse than the quadratic and linear methods when the number of predictor variables exceeds 20. For 40 predictor variables, a halving of the computational time is achieved. For 50 variables, the computation times are 1079 seconds for the naive approach, 517 seconds for the quadratic method and 411 seconds for the linear method. We did not consider regression problems with more predictor variables, as the memory requirements for the naive method become prohibitive.

The same procedure was applied to the singular case $(n-1 < p)$, but with matrices of 12 observations by 15, 20, 25, 50, 75, 100, 500, 1000 and 2500 predictor variables respectively. For each simulated problem, we calculated the PRESS statistics for all 10 factors with each of our three implementations of full leave-one-out cross-validation (we lose two dimensions as the

Figure 2.1: Computation times in seconds for full leave-one-out cross-validation of principal component regression in the regular case with 100 observations.

Figure 2.2: Computation times in seconds for full leave-one-out cross-validation of principal component regression in the singular case with 12 observations.

data are mean-centred and because of leave-one-out cross-validation). Figure 2.2 shows the results. The results are again plotted on the log scale for the number of predictor variables. The cost of both the linear and quadratic method appears to be largely unaffected by the number of predictor variables considered. There is a minor increase in the computational time for both the quadratic and linear methods beyond 1000 predictor variables. This is probably due to the cost of computing the initial principal component decomposition. The naive method suffers from a linear increase in cost as the number of predictor variables increases. With our implementation, the naive method broke down, when more than 2500 predictor variables were generated, due to insufficient memory. For the quadratic and linear procedures, it was possible to consider problems with up to 10000 predictor variables, at no extra computational cost. As expected [DeG90], the linear method was invariably quicker than the quadratic method. Both methods are significantly faster than the naive procedure when a large number of variables is considered.

**Real Data**

A limited number of real data sets from the statistical literature on principal component regression was also analyzed. The conclusions were no different from those above. A full leave-one-out cross-validatory analysis of the ethanol data (11 observations and 101 predictor variables, as described by Stone and Brooks [SB90]) was accomplished in 4.1 seconds for the naive method, 4 seconds for the quadratic method and 2.9 seconds for the linear method. For the wheat data (24 observations and 6 variables, published by Fearn [Fea83] and reanalysed by Stone and Brooks), calculations were done in 1.37 seconds for the naive method, 4.5 for the quadratic method and 3.4 for the linear method. As a last example, we calculated a full leave-one-out cross-validation of the NIR data that was considered by Stone and Brooks (page 252, the data originate from Osborne *et al.* [OFMD84]). We considered all 40 observations (including outlier) with measurements at 601 wavelengths. The computational times were 1005 seconds for the naive approach, 115 seconds for the quadratic method and 87 seconds for the linear method.

To understand the findings from these simulations, we consider a theoretical study of the computational costs involved in these procedures. To

simplify the analysis, we will only concern ourselves with orders of magnitude of costs. Furthermore, following Bunch and Nielsen in their analysis of the cost of updating the singular value decomposition [BN78], we only consider the number of multiplications as a measure of efficiency. As before, the regular and the singular case are considered separately.

## The regular case

Naive implementations of leave-one-out principal component cross-validation are based on the recomputation of the cross-products matrices after the deletion of an observation. There is an order of $np^2$ multiplications involved in the computation of each cross-products matrix. Hence, the total number of multiplications due to the computation of the cross-products matrices is of order $n^2p^2$.

For each cross-products matrix, the principal component decomposition must then be computed from scratch. The number of multiplications involved in the computation of the principal component coefficients and the principal component scores depends on the method that is chosen for these computations. Golub and Van Loan [GL89, page 423], Wilkinson [Wil65] and Wilkinson and Reinsch [WR71, part 2] are the standard references on the computation of the symmetric eigenvalue problem. For the symmetric QR algorithm [GL89, page 423], the computations will be of order $p^3$. Calculating the principal component scores from the principal component coefficients will involve $np^2$ multiplications. Thus the total cost of the cross-validatory procedure will be of order $n^2p^2 + np^3$.

Several authors in the statistical literature refer to the singular value decomposition as a suitable method for the joint computation of the principal component coefficients and the principal component scores [Jol86, page 235] [Seb84, page 506]. An implementation of the naive approach with the Golub and Reinsch method [GL89, page 430] for the singular value decomposition will involve an order of $np^2 + p^3$ multiplications for each datum [GL89, page 239]. Performing this for all data requires an order of $n^2p^2 + np^3$ multiplications. Hence, such implementations will involve roughly the same computational costs as compared to those based on the QR algorithm.

With the efficient procedures, the number of multiplications involved in downdating the principal component coefficients is of order $p^3$, for each datum. Hence, the overall cost of the computations for the principal component

coefficients is of order $np^3$. Therefore, when the number of observations is large compared to the number of predictor variables, this approach should be significantly faster than the naive implementations of principal component coefficient downdating. This is because we avoid a computation of order $n^2p^2$. Thus we achieve an order of magnitude improvement in the downdating of the principal component coefficients.

Unfortunately, for applications to principal component regression, this gain is largely lost in the computation of the PRESS statistics. Indeed, the leave-one-out cross-validation of principal component regression involves the computation of the principal component scores, which requires an order of $n^2p^2$ multiplications. This does not mean that we no longer achieve the improvement in the calculation of the principal component coefficients, but rather that this improvement is drowned in a subsequent calculation which is of the same order of magnitude as for the naive procedures. Hence, we can only halve the number of multiplications in the regular case.

| Method<br>Object | Naive | Efficient |
|---|---|---|
| $\mathbf{S}_{(i)}$ | $np^2$ | not applicable |
| $\mathbf{Q}_{(i)}$ | $p^3$ | $p^3$ |
| $\mathbf{U}_{(i)}$ | $np^2$ | $np^2$ |
| Total | $np^2 + p^3$ | $np^2 + p^3$ |

Table 2.1: The contributions of a single observation to the order of the number of multiplications needed for the downdating of a principal component decomposition in the regular case with the leave-one-out approach.

## The singular case

In the singular case, the naive approach is based on the computation of the matrices $\mathbf{X}_{(i)}\mathbf{X}_{(i)}^T$. There is an order of $pn^2$ multiplications involved in the computation of each such matrix. Thus, the total number of multiplications necessary for the computation of these matrices is of the order $pn^3$. With the QR algorithm, the number of multiplications necessary for the computation

| Method Object | Naive | Efficient |
|---|---|---|
| $\mathbf{S}_{(i)}$ | $np^2$ | not applicable |
| $\mathbf{Q}_{(i)}$ | $p^3$ | not applicable |
| $\mathbf{U}_{(i)}$ | $np^2$ | $np^2$ |
| Subtotal | $np^2 + p^3$ | $np^2$ |
| PRESS | $np$ | $np$ |
| Total | $np^2 + np + p^3$ | $np^2 + np$ |

Table 2.2: The contributions of a single observation to the order of the number of multiplications needed for the computation of the prediction error sum of squares for principal component regression in the regular case, using the leave-one-out approach.

of the principal component decomposition of one such matrix is of the order $n^3$. The principal component scores can be obtained from this decomposition, which requires an order of $n^2$ multiplications. We must also compute the principal component coefficients, which involves an order of $pn^2 + n^2$ multiplications. Thus, the total number of multiplications necessary to obtain the cross-validated principal component analysis is of the order $pn^3 + n^4$.

The efficient rank-one modification procedure involves an order of $n^4$ multiplications for the computation of the downdated principal component scores for all observations and an order $pn^3 - pn^2$ [BN78] multiplications for the calculation of the principal component coefficients. In principal component regression, we only need to consider the downdating of the scores. As can be seen, the computational cost of calculating these quantities does not depend on $p$. As a consequence, the computational cost of the efficient approach is constant when the number of observations is fixed, irrespective of the number of predictor variables that are considered.

In applications to principal component regression, the calculation of the PRESS statistic will involve an order of $n^3$ multiplications and hence, this gain is maintained in the full leave-one-out cross-validation of principal component regression. Thus we achieve an order of magnitude reduction for the cross-validation of principal component regression in the singular case and

the computational cost due to the number of variables is eliminated completely.

| Method<br>Object | Naive | Efficient |
|---|---|---|
| $\mathbf{S}_{(i)}$ | $pn^2$ | not applicable |
| $\mathbf{Q}_{(i)}$ | $pn^2 + n^2$ | $pn^2 + n^2$ |
| $\mathbf{U}_{(i)}$ | $n^3$ | $n^3$ |
| Total | $pn^2 + n^3$ | $pn^2 + n^3$ |

Table 2.3: The contributions of a single observation to the order of the number of multiplications needed for the downdating of a principal component decomposition in the singular case with the leave-one-out approach.

| Method<br>Object | Naive | Efficient |
|---|---|---|
| $\mathbf{S}_{(i)}$ | $pn^2$ | not applicable |
| $\mathbf{Q}_{(i)}$ | $pn^2 + n^2$ | not applicable |
| $\mathbf{U}_{(i)}$ | $n^3$ | $n^3$ |
| Subtotal | $pn^2 + n^3$ | $n^3$ |
| PRESS | $n^2$ | $n^2$ |
| Total | $pn^2 + n^3$ | $n^3$ |

Table 2.4: The contributions of a single observation to the order of the number of multiplications needed for the computation of the prediction error sum of squares for principal component regression in the singular case, using the leave-one-out approach.

Finally, tables 2.1 up to 2.4 summarize these results. The tables are concerned with costs induced by each single observation in the full leave-one-out cross-validation. Hence, the appropriate quantities must be multiplied by $n$ for comparison with the relevant quantities in the text.

## 2.3.4 Accuracy

Detailed analyses on the accuracy of the efficient methods for the leave-one-out downdating of the principal component decomposition are in the numerical literature. In fact, this is precisely what these publications are concerned with: to provide efficient methods for rank-one modifications of the eigendecomposition of square symmetric matrices such that a minimum and prespecified level of accuracy is achieved, irrespective of the condition of the problem.

Most articles consider the performance of the rank-one modification algorithm within the 'divide and conquer' method. This is a relatively new method to compute the eigendecomposition of large square symmetric matrices when the classical approaches based on the QR decomposition break down. Thus, the 'divide and conquer' has been one of the first applications of the rank-one modification methods for principal component decompositions. Cuppen's article [Cup81] on the 'divide and conquer' method for the symmetric tridiagonal eigenproblem deserves special mention here for an extensive mathematical analysis of the perturbation properties of the rank-one modification method, as well as for a report on experiments with simulated data. This compares the eigendecompositions obtained from a numerically stable algorithm using the QR method with those obtained from the 'divide and conquer' approach, based on the rank-one algorithm. Likewise, Dongarra and Sorensen [DS87] report theoretical work on the rank-one modification algorithm and simulations for the 'divide and conquer' method. The latter compare the performance between a number of different computing platforms. Finally, Bunch, Nielsen and Sorensen [BNS78] discuss theoretical considerations and simulations based on the uniform distribution for the rank-one modification algorithm only.

Cuppen concludes that the 'divide and conquer' method, and by consequence the rank-one modification algorithm, is numerically stable and efficient. The conclusions from other publications are similar.

Not withstanding this published material, we have considered the accuracy of the GAUSS implementations for the leave-one-out cross-validation of principal component regression in the simulations discussed in the previous section. For both the regular and the singular case, the absolute difference between the PRESS statistics calculated by the efficient and naive procedures was never larger than $10e^{-12}$ and almost always smaller than $10e^{-15}$.

| $\gamma$ | A | B | C |
|---|---|---|---|
| 1 | 1.3322676e-15 | 1.3322676e-15 | 0 |
| 2 | 0 | 8.8817842e-16 | 8.8817842e-16 |
| 3 | 3.9968029e-15 | 3.9968029e-15 | 0 |
| 4 | 2.5535130e-15 | 2.7755576e-15 | 2.2204460e-16 |
| 5 | 7.7715612e-16 | 5.5511151e-16 | 2.2204460e-16 |
| 6 | 9.1038288e-15 | 8.9928065e-15 | 1.1102230e-16 |
| 7 | 1.7708057e-14 | 1.7819080e-14 | 1.1102230e-16 |
| 8 | 2.6534330e-14 | 2.6645353e-14 | 1.1102230e-16 |
| 9 | 1.1601831e-14 | 1.1379786e-14 | 2.2204460e-16 |
| 10 | 3.0753178e-14 | 3.0808689e-14 | 5.5511151e-17 |
| 11 | 3.1419312e-14 | 3.1530334e-14 | 1.1102230e-16 |
| 12 | 3.7192471e-14 | 3.7136960e-14 | 5.5511151e-17 |
| 13 | 1.7469359e-13 | 1.7480462e-13 | 1.1102230e-16 |
| 14 | 2.4363844e-13 | 2.4374947e-13 | 1.1102230e-16 |
| 15 | 2.1105340e-13 | 2.1116442e-13 | 1.1102230e-16 |
| 16 | 1.1390888e-13 | 1.1385337e-13 | 5.5511151e-17 |
| 17 | 5.0903726e-13 | 5.0887072e-13 | 1.6653345e-16 |
| 18 | 9.6717079e-13 | 9.6711528e-13 | 5.5511151e-17 |
| 19 | 3.8802295e-13 | 3.8791192e-13 | 1.1102230e-16 |
| 20 | 5.7504002e-13 | 5.7498450e-13 | 5.5511151e-17 |
| 21 | 2.8941294e-12 | 2.8941294e-12 | 0 |
| 22 | 1.1469714e-12 | 1.1469714e-12 | 0 |
| 23 | 3.5503267e-12 | 3.5501602e-12 | 1.6653345e-16 |
| 24 | 8.8636321e-12 | 8.8637431e-12 | 1.1102230e-16 |
| 25 | 5.8018590e-12 | 5.8017480e-12 | 1.1102230e-16 |
| 26 | 7.9243834e-12 | 7.9243279e-12 | 5.5511151e-17 |
| 27 | 1.0680068e-11 | 1.0679957e-11 | 1.1102230e-16 |
| 28 | 2.4901303e-11 | 2.4901303e-11 | 0 |
| 29 | 1.6747770e-11 | 1.6747881e-11 | 1.1102230e-16 |
| 30 | 2.6088021e-12 | 2.6086910e-12 | 1.1102230e-16 |
| 31 | 2.2481184e-11 | 2.2481017e-11 | 1.6653345e-16 |
| 32 | 8.5478291e-12 | 8.5477181e-12 | 1.1102230e-16 |
| 33 | 9.7003516e-12 | 9.7003516e-12 | 0 |
| 34 | 7.7162721e-12 | 7.7161610e-12 | 1.1102230e-16 |
| 35 | 2.4789393e-11 | 2.4789504e-11 | 1.1102230e-16 |
| 36 | 5.8201555e-11 | 5.8201666e-11 | 1.1102230e-16 |
| 37 | 7.3549278e-11 | 7.3549056e-11 | 2.2204460e-16 |
| 38 | 5.4517169e-11 | 5.4517169e-11 | 0 |
| 39 | 9.7468478e-11 | 9.7468256e-11 | 2.2204460e-16 |
| 40 | 1.5084756e-10 | 1.5084745e-10 | 1.1102230e-16 |
| 41 | 1.9282498e-10 | 1.9282476e-10 | 2.2204460e-16 |
| 42 | 1.9318691e-10 | 1.9318669e-10 | 2.2204460e-16 |
| 43 | 1.5413038e-10 | 1.5413004e-10 | 3.3306691e-16 |
| 44 | 1.2476176e-10 | 1.2476153e-10 | 2.2204460e-16 |
| 45 | 1.4194101e-10 | 1.4194068e-10 | 3.3306691e-16 |
| 46 | 1.3419166e-10 | 1.3419144e-10 | 2.2204460e-16 |
| 47 | 1.3983348e-10 | 1.3983337e-10 | 1.1102230e-16 |
| 48 | 1.5002810e-10 | 1.5002810e-10 | 0 |

Table 2.5: Absolute differences between the mean prediction error sum of squares for principal component regression of the Mustard data calculated with (A) the linear method and the naive approach, (B) the quadratic method and the naive approach and (C) the linear and the quadratic method.

Finally, we have compared the performance of the efficient implementations with that from a naive implementation of principal component regression, using the principal component regression algorithm in the appendix (page 117) on the Mustard data described in chapter 1. The results in table 2.5 show the absolute differences in the computed mean squared error of prediction statistics. They are slightly worse than those obtained from the simulations. This is probably due to the condition of the dispersion matrix, which is likely to be worse than that obtained for random simulations. There is a gradual decrease in accuracy as further components of increasingly smaller variance are added to the prediction equation, although the results are of sufficient accuracy throughout. Interestingly, both the linear and quadratic methods coincide in the results obtained, as these are indistinguishable up to machine precision.

## 2.4 Efficient Leave-One-Out Principal Component Regression Cross-validation Applied to Near Infrared Spectroscopy

### 2.4.1 Cross-validation

**GAUSS Implementation**

The appendix contains complete transcripts of GAUSS implementations of efficient leave-one-out cross-validation for principal component regression, for both the quadratic (page 118) and the linear method (page 129). A complete leave-one-out cross-validatory analysis can be computed with this code in a normal GAUSS session, in the following manner.

A GAUSS matrix x must be created that contains the predictor data. Each row must contain the data from a single observation and each column must represent one predictor variable. Similarly, a GAUSS column vector y must be created that contains the observed responses for these observations, in the same order. The statements

$$\{\mathtt{msep, dd, ddiff, ycvpred, l, q}\} = \mathtt{pcrcvq(y, x, w)};$$

will then generate the leave-one-out analysis with the quadratic method, and

similarly,

$$\{\text{msep}, \text{dd}, \text{ddiff}, \text{ycvpred}, \text{l}, \text{q}\} = \text{pcrcvl}(\text{y}, \text{x}, \text{w});$$

generates the same results with the linear method. w is a positive integer that specifies the number of components that must be derived in the cross-validation.

The output from these procedures is as follows. The matrix msep will contain the PRESS and MSEP statistics in the second and third columns respectively. The first column is the number of components in the prediction equation. Thus, the top row will contain the results when there are no principal components in the prediction equation and only the mean is used in prediction. The last row gives the results when w components are used in the construction of the prediction equation. ycvpred is a matrix that contains the cross-validated predicted values for each observation, with the observations in the same order as for x and y. This matrix has w+1 columns. The first column contains the cross-validated means. The next columns contain the cross-validated values when further principal components are added, up to the last column which contains the predicted values when all w components are in the prediction equation. In a similar manner, dd and ddiff contain the cross-validated eigenvalues and the eigenvalue downdates, respectively. dd has w columns, one for each component that is cross-validated. In this way, the first column contains the downdated eigenvalues for the first component, and so on, up to the last column, which contains the results for the w$^{\text{th}}$ component. As for ycvpred, the result are in a single row for each observation, using the same order. ddiff has the same layout for the eigenvalue downdates. Finally, l contains the principal component variances for the complete data, and similarly, q is the principal component coefficient matrix from the initial principal component decomposition on which the efficient cross-validatory computations are based, with the coefficients for the first component in the first column, and so on.

## The Mustard data

Using these procedures, the results from a leave-one-out cross-validation of principal component regression for the Mustard data are as shown in table 2.6. A plot of the cross-validatory index versus the number of components is shown in figure 2.3. The cross-validatory index equals 88.2% for a principal

| $\gamma$ | MSEP | INDEX | $\gamma$ | MSEP | INDEX |
|---|---|---|---|---|---|
| 0 | 3.3466262 | 0 | | | |
| 1 | 3.2601899 | 0.025827909 | 25 | 0.40632839 | 0.87858567 |
| 2 | 2.5588577 | 0.23539185 | 26 | 0.41160549 | 0.87700883 |
| 3 | 1.1609123 | 0.65310965 | 27 | 0.42838069 | 0.87199626 |
| 4 | 0.68795101 | 0.79443447 | 28 | 0.42241545 | 0.87377872 |
| 5 | 0.72123705 | 0.78448832 | 29 | 0.43209750 | 0.87088564 |
| 6 | 0.56175546 | 0.83214275 | 30 | 0.46338557 | 0.86153650 |
| 7 | 0.39447966 | 0.88212617 | 31 | 0.47898626 | 0.85687489 |
| 8 | 0.39490756 | 0.88199831 | 32 | 0.47900025 | 0.85687071 |
| 9 | 0.32686328 | 0.90233051 | 33 | 0.50620231 | 0.84874250 |
| 10 | 0.31741968 | 0.90515234 | 34 | 0.64771512 | 0.80645729 |
| 11 | 0.32714911 | 0.90224510 | 35 | 0.62871990 | 0.81213322 |
| 12 | 0.33641040 | 0.89947775 | 36 | 0.59526733 | 0.82212913 |
| 13 | 0.36743405 | 0.89020762 | 37 | 0.50359747 | 0.84952085 |
| 14 | 0.37193583 | 0.88886245 | 38 | 0.50840066 | 0.84808561 |
| 15 | 0.40468830 | 0.87907574 | 39 | 0.51229011 | 0.84692342 |
| 16 | 0.42650823 | 0.87255576 | 40 | 0.55915348 | 0.83292025 |
| 17 | 0.45367112 | 0.86443926 | 41 | 0.62723446 | 0.81257708 |
| 18 | 0.48900661 | 0.85388072 | 42 | 0.59155810 | 0.82323748 |
| 19 | 0.43636159 | 0.86961150 | 43 | 0.59291498 | 0.82283203 |
| 20 | 0.45729522 | 0.86335635 | 44 | 0.55681674 | 0.83361849 |
| 21 | 0.46513218 | 0.86101460 | 45 | 0.54275339 | 0.83782073 |
| 22 | 0.44417680 | 0.86727624 | 46 | 0.52929253 | 0.84184295 |
| 23 | 0.39996337 | 0.88048759 | 47 | 0.52043810 | 0.84448873 |
| 24 | 0.40344986 | 0.87944579 | 48 | 0.52171454 | 0.84410732 |

Table 2.6: Full leave-one-out cross-validatory analysis for principal component regression of the Mustard data.

Figure 2.3: Cross-validatory index for a full leave-one-out cross-validation of principal component regression of the Mustard data.

component regression prediction equation with seven principal components. The inclusion of the two subsequent principal components increases the cross-validatory index to 90.2%. The maximum is achieved for the cross-validatory index when the first ten components are used for prediction. Although these results are obtained from a cross-validatory analysis, we will take a cautious point of view and restrict ourselves to a seven component prediction equation.

To study the selected principal component prediction equation, a separate GAUSS program has been written which can be used after a suitable prediction equation has been chosen from the previous cross-validatory procedures. The appendix (page 117) contains the code. Within the same GAUSS

session, the statements

$$\{\mathtt{ypred}, \mathtt{yresid}, \mathtt{beta}, \mathtt{q}, \mathtt{u}\} = \mathtt{pcr}(\mathtt{y}, \mathtt{x}, \mathtt{w});$$

generate the first $\mathtt{w}$ principal component regression equations.

The output is as follows. $\mathtt{ypred}$ is a matrix that contains the predicted values, such that each $j^{\mathrm{th}}$ column gives the predicted values for the prediction equation with $j$ components. The layout of $\mathtt{yresid}$ is identical for the residuals. Likewise, $\mathtt{beta}$ gives the regression coefficients, with the same layout for columns, with the coefficients for each $j^{\mathrm{th}}$ variable in the $j^{\mathrm{th}}$ row. $\mathtt{q}$ is the matrix of principal component coefficients which has already been obtained from the cross-validatory computations, and $\mathtt{u}$ is the matrix of principal component scores, such that $\mathtt{u=(x-meanc(x)')*q}$.

Plots of the first seven principal component loadings are shown in figure 2.4. As often, the first component represents the general trend in the observed spectra. An interpretation of the remaining components would be difficult and requires expert knowledge of the near infrared spectroscopy of white mustard seeds. Figure 2.5 shows the regression coefficients for the selected principal component prediction equation.

## 2.4.2 Influence Analysis

The output from the GAUSS procedures $\mathtt{pcrcvq}$, $\mathtt{pcrcvl}$ and $\mathtt{pcr}$ provides the measures necessary to perform an influence analysis.

Figure 2.6 shows plots constructed from the cross-validated predicted values, the cross-validated residuals and the observed values. The top plot considers the cross-validated predicted values versus the observed values. There do not appear to be any consistent deviations from a straight line relation between the observed and the cross-validated predicted values. Neither are there any obvious outliers. The two subsequent plots of the cross-validated residuals versus the observed and cross-validated predicted values suggest some evidence of heterogeneity, although not conclusively. The plots reveal some peculiar features of the sampling design. First of all, the data appear to have been sampled from three distinct groups, according to oil content. Furthermore, within each group, several samples are present which share the same oil content. This suggests that some repeat sampling from the same source may have been employed.

Figure 2.4:  Principal component loadings for the first seven principal components.

Figure 2.5: Regression coefficients for a principal component regression prediction equation with seven components.

Figure 2.6: Scatterplots obtained from the cross-validated predicted values, the cross-validated residuals and the observed values.

Ignoring these problems, the plots appear to preclude the subsequent analysis of influence for the principal component decomposition, as there are no observations with unduly large cross-validated residuals. In practice, this analysis should be carried out for the selected model, as well as for a few submodels, to investigate whether individual observations are unduly inflating the cross-validatory estimate. Once this estimate has been obtained, however, the prediction equation is calibrated from the complete data. Not withstanding the cross-validatory assessment of the prediction equation, it would be undesirable for a single observation to have excessive influence on the decomposition of the complete data. Hence, we would still want to consider the sensitivity of the principal component estimates for leave-one-out perturbations of the data. This is particularly so in applications in chemometrics, where sample sizes are often small.

The most simple principal component influence plot is a plot of the total eigenvalue downdate $\rho$ of each observation versus the observation number, as shown in figure 2.7. This may be complemented with figure 2.8, which show a plot of the cross-validated residuals versus the total downdates. Neither of these plots identify any unduly influential observations.

We should expect that principal components with larger variances will dominate the eigenvalue downdating. Hence, it would make sense to consider the first few principal components separately in a principal component influence analysis. Figure 2.9 shows a plot of the sum of the first seven eigenvalue downdates versus the observation number. This plot is virtually identical to figure 2.7 and as a consequence, figure 2.10 can hardly be distinguished from the corresponding plot in figure 2.8. The reason becomes clear from the scree diagram shown in figure 2.11. The first component accounts for 95.3% of the variation of the data and hence this component dominates the eigenvalue downdating.

The nature of the downdating process for the Mustard data is clarified from a separate plot of the normalized eigenvalue downdates for the first seven components, as shown in figure 2.12. Most of the eigenvalue downdating takes place along the first principal component direction, as the normalized eigenvalue downdates for the first component are between 80% and 100% for most observations. This situation reverses in the remaining principal component directions. The further along we go in the principal component sequence, the less eigenvalue downdating there is. The maximum normalized eigenvalue downdate is smaller than 3% for the seventh principal component,

Figure 2.7: Total eigenvalue downdate $\rho$ versus the observation number.

Figure 2.8: Cross-validated residuals versus the total eigenvalue downdate $\rho$.

Figure 2.9: Sum of the eigenvalue downdates for the first seven principal components versus the observation number.

Figure 2.10: Cross-validated residuals versus the sum of the eigenvalue down-dates for the first seven principal components.

Figure 2.11: Scree diagram.

Figure 2.12: Normalized eigenvalue downdates for the first seven principal components.

while most observations have a normalized eigenvalue downdate smaller than 0.5%.

These plots are appropriate to analyse the influence of observations on the eigenvalues of specific principal component directions. Observation 40 is a clear outlier for the sixth component. Plots of the cross-validated residuals versus each of the seven normalized eigenvalue downdates are shown in figure 2.13. The $40^{\text{th}}$ observation has a cross-validated residual close to zero.

The normalized eigenvalue downdates suggest a summary influence measure for the eigenvalue downdating in the first seven principal components. We may consider the sum of the normalized eigenvalue downdates for the first seven components versus the observation number, as shown in figure 2.14. The $30^{\text{th}}$ observation emerges as a clear outlier from this plot. Plotting the cross-validated residuals versus this sum of normalized eigenvalue downdates, as shown in figure 2.15, we find that this observation has a cross-validated residual close to zero.

The previous analysis is restricted to the influence on the eigenvalues. Figure 2.16 show plots of the downdating angles for the first seven principal components. There are no suspect observations. The size of the downdating angles varies between the components. The downdating angles are smaller than one degree for the first component, whereas the fourth and fifth component have downdating angles of up to 25 degrees. The plots of the downdating norms, shown in figure 2.17, are virtually identical to those of the downdating angles, as expected.

## 2.5  Conclusions

We have demonstrated how an efficient leave-one-out algebra due to Bunch, Nielsen and Sorensen [BNS78] can be applied to the cross-validation of a principal component decomposition. The algebra replicates results from naive implementations, using efficient computations that achieve an order of magnitude reduction in computational cost with respect to the cross-validation of the principal component decomposition. In the application to the calibration of principal component regression equations, this order of magnitude reduction is maintained for the singular case only. Fortunately, in applications in chemometrics, this is precisely the case we are interested in, as spectroscopic analysis will typically generate large numbers of measurements on samples

Figure 2.13: Cross-validated residuals versus normalized eigenvalue downdates for the first seven principal components.

Figure 2.14: Sum of the normalized eigenvalue downdates for the first seven principal components versus the observation number.
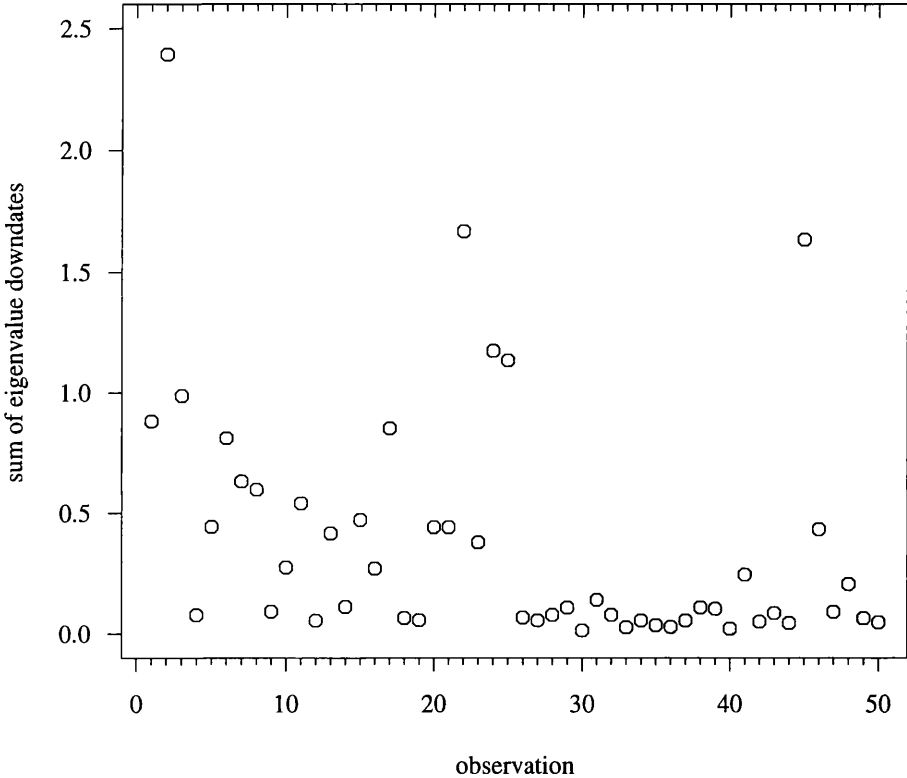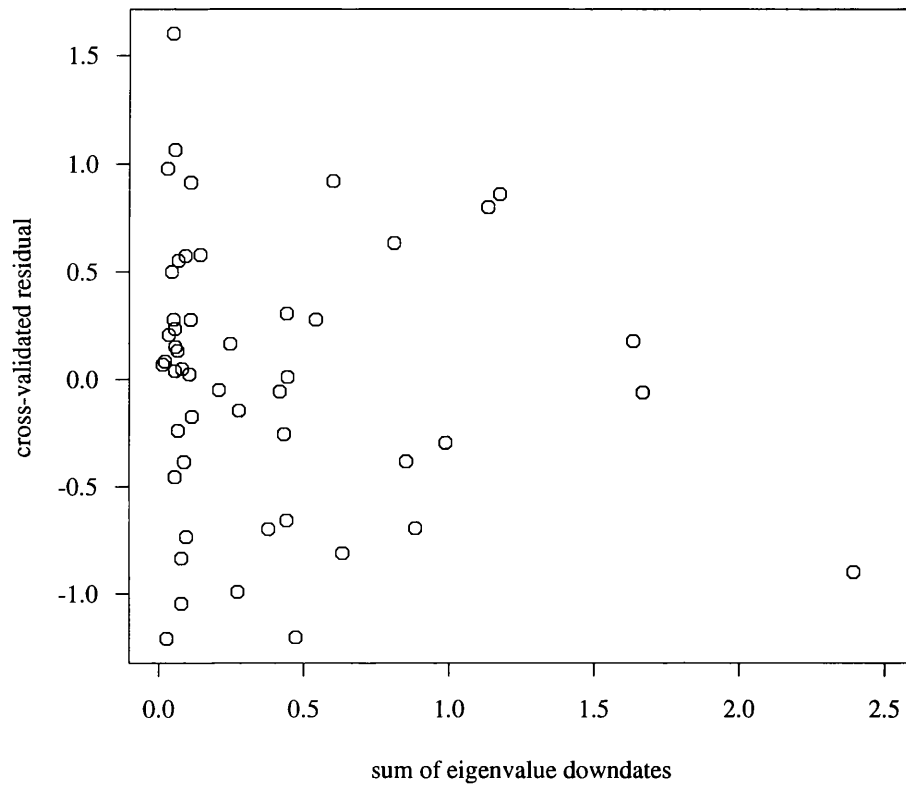
Figure 2.15: Cross-validated residuals versus the sum of the normalized eigenvalue downdates for the first seven principal components.
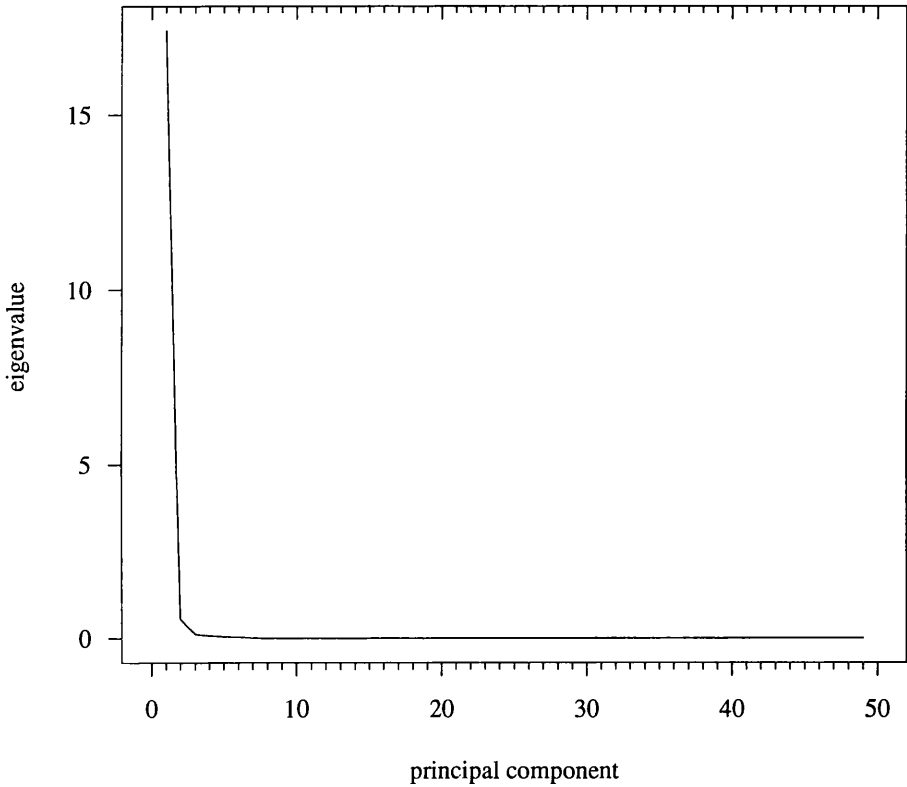
Figure 2.16: Principal component downdating angles.

Figure 2.17: Principal component downdating norms.

of interest.

The efficiency of the leave-one-out computations stems from the reduction of the downdating problem to the cross-validation of the eigenvalues. These computations no longer depend on matrix algebra and can be executed fast and accurately. The algorithm depends crucially on the properties of principal component eigenvalue downdates. These properties may in turn be utilized to formulate a system of principal component influence measures, most of which were already in the literature in some form or another. An application of the influence measures to near infrared spectroscopic data highlights interesting features in the data which would have gone undetected otherwise.

We have only considered the application of the efficient cross-validatory algebra and influence measures to principal component regression. Other applications exist, where the principal component decomposition is of interest in its own right. This may be in terms of its dimension reduction properties or as a model of data, as in the SIMCA [Wol76] approach to classification, for example. The use of principal component influence measures in particular requires more work in these applications. Furthermore, the cross-validatory algebra only applies to the decomposition of the covariance matrix. From a statistical point of view, in many applications of principal component decomposition, the analysis of the correlation matrix is often more useful. Hence, more work is needed to make efficient cross-validatory computations available in this case.

# Chapter 3

# Partial Least Squares Regression

# 3.1 Partial Least Squares Decomposition and Cross-validation

In the Stone-Brooks approach, the full leave-one-out cross-validation of partial least squares components becomes particularly elegant and, to some extent, analogous to the cross-validation of principal components.

Implementations of cross-validation for partial least squares computations are usually based on the naive approach which recomputes the partial least squares decomposition for each construction sample, irrespective of the algorithm used for the computation of the decomposition. The Stone-Brooks [SB90] generalisation of partial least squares regression incorporates an efficient cross-validatory algebra for the leave-one-out case. While the Stone-Brooks algebra generally involves the numerical optimisation of a parameter over a known parameter space, this parameter cancels from the equations in the partial least squares restriction. Thus the Stone-Brooks computations reduce to the sequential application of a computation which is purely matrix-based. For this reason, the Stone-Brooks implementation of partial least squares regression deserves special attention as it leads to efficient cross-validatory computations.

## 3.1.1 Partial Least Squares Decomposition

The Stone-Brooks approach derives partial least squares components sequentially, such that each component maximizes the sample covariance between the corresponding component scores and the observed response, subject to the conditions that the partial least squares component scores are uncorrelated and have positive sample variance, while the component coefficient vectors are of unit length [SB90]. The restriction to uncorrelated partial least squares components is not essential in this definition. As for principal components, we may consider an alternative restriction in the derivation of the partial least squares components, which defines the component coefficient vectors to be mutually orthogonal, as considered by Martens and Naes [MN89, page 119].

## Partial Least Squares Algebra

We will describe the application of the Stone-Brooks approach to the construction of a partial least squares algebra.

The first component must maximize the sample covariance $(s^T q)^2$, or equivalently $(c^T q)^2$, among all $q$ with $q^T q = 1$. Hence, using Lagrange multipliers, we find that $q$ must maximize

$$L_1(q) = (c^T q)^2 - \lambda(q^T q - 1).$$

Therefore, the vector of component loadings for the first component is the solution of the equations

$$\frac{d}{dq}L_1(q) = 2c - 2\lambda q = 0,$$

and we find that the first partial least squares component loadings are proportional to the sample covariance vector $c$. We have

$$q_1 = \frac{c}{\|c\|}.$$

Consider the calculation of the component loadings $q_{k+1}$ for the $k + 1^{\text{th}}$ component, with $1 \leq k < r_x$. The component must maximize the criterion $(c^T q)^2$, subject to the conditions $q^T q = 1$ and $q_j^T C q = 0$, for $j = 1, \ldots, k$.

At this point, we are forced to change the notation for principal components and we will write $C = VEV^T$ for the principal component decomposition of the complete data, such that $V$ is the matrix of principal component coefficients. Using the Stone-Brooks approach to partial least squares decomposition, we derive the partial least squares component from a rotation of the principal component decomposition of the complete data. Stone and Brooks give proof that the principal components span the same space as the partial least squares components. Hence, we may write

$$q = Vr,$$

and we must now find the optimum of the Lagrangian equations

$$L_{k+1}(r) = (d^T r)^2 - \lambda(r^T r - 1) - \alpha^T A^T r,$$

where $\mathbf{d} = \mathbf{V}^T\mathbf{c}$, $\boldsymbol{\alpha}^T = (\alpha_1, \ldots, \alpha_k)$ and $\mathbf{A} = \mathbf{E}\mathbf{V}^T\mathbf{Q}$, with $\mathbf{Q} = (\mathbf{q}_1, \ldots, \mathbf{q}_k)$ the matrix formed by the first $k$ vectors of partial least squares component loadings. The optimum is the solution of

$$\frac{d}{d\mathbf{r}}L_{k+1}(\mathbf{r}) = 2\mathbf{d} - 2\lambda\mathbf{r} - \mathbf{A}\boldsymbol{\alpha} = 0.$$

Hence, absorbing the constant 2 in the Lagrange multipliers, we find that

$$\mathbf{r} \propto \mathbf{d} - \mathbf{A}\boldsymbol{\alpha}$$

with $\mathbf{A}^T\mathbf{r} = 0$ and therefore,

$$\mathbf{r} \propto (\mathbf{I} - \mathbf{A}(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T)\mathbf{d}.$$

We find

$$\mathbf{r} = \frac{\mathbf{M}\mathbf{d}}{\|\mathbf{M}\mathbf{d}\|},$$

with $\mathbf{M} = \mathbf{I} - \mathbf{A}(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T$.

**Partial Least Squares Computations**

The above algebra is inefficient in its present form. This is because the computation of $\mathbf{M}$ involves the inverting of a matrix that must be recomputed for each component, in the sequence of partial least squares components. A naive approach to this computation would involve, for each partial least squares component, the computation of the principal component decomposition of $\mathbf{M}$. This would render the cross-validatory algebra useless. The problem is that the cross-validatory algebra does not yet take full account of the leave-one-out nature of the downdating. Stone and Brooks have shown that an efficient formula can be found for the computation of these inverse matrices, using the Sherman-Morisson theorem, which may be derived as a special case of theorem 2.2. We may employ the updating formula

$$\mathbf{M}_{k+1} = \mathbf{M}_k - \frac{(\mathbf{M}_k\mathbf{a})(\mathbf{M}_k\mathbf{a})^T}{\mathbf{a}^T\mathbf{M}_k\mathbf{a}},$$

for $\mathbf{M}$ at each stage in the cross-validatory computations, where $\mathbf{a}$ is the column appended to $\mathbf{A}$ at stage $k + 1$. The details of this may be found in the paper by Stone and Brooks.

The singular case may be dealt with using the same methods as explained for principal component decomposition. Indeed, we are deriving the partial least squares components from a rotation of the principal component decomposition and hence, we may use the same methods as before to derive the initial principal component decomposition in the singular case.

## 3.1.2  Partial Least Squares and Cross-validation

Partial least squares is more generally known as a class of path modelling techniques in which the interrelation of blocks of observed variables is modelled through sets of latent derived factors. The latent variables are obtained through a sequence of least squares projections. The technique was originated by Herman Wold [Wol66] [Wol85] in the context of econometric prediction modelling and inherits from factor analysis techniques in psychometry. The methodology has been much applied since, with applications in different branches of science.

In chemometrics, partial least squares has been successful in the calibration of linear equations for prediction. In these applications, the number of predictor variables is often huge, as in spectroscopic data for example, while the number of observations is small. The usual formulation of partial least squares in these applications is as an application of the NIPALS algorithm, as described by Wold [WL69].

It is recent studies of the partial least squares approach to regression by Höskuldsson [Hös88], Helland [Hel88] and Stone and Brooks [SB90] that have shown that the partial least squares component coefficient vectors may be interpreted as directions in the measurement space that optimize the covariance with the observed response. This strips partial least squares from its algorithmic formulation. The interpretation was used explicitly by Stone and Brooks [SB90] to derive the partial least squares components from a rotation of the principal components of the predictor data. There has been no work on the efficient cross-validation of partial least squares regression. All implementations of cross-validation for partial least squares are based, essentially, on Wold's approach to the choice of the number of components in factor and principal component models [Wol78]. Thus, they consider the NIPALS algorithm and subsets of data rather than a full leave-one-out cross-validation.

## 3.2 Stone-Brooks Leave-One-Out Cross-validatory Algebra

As for principal component cross-validation we may exploit the efficient leave-one-out downdating formula

$$C_{(i)} = C - \nu x_i^T x_i$$

for the computation of the downdated cross-products matrix. The corresponding leave-one-out downdating formula for the cross-products vector $c$ is

$$c_{(i)} = c - \nu y_i x_i^T.$$

Using the same argument as for the partial least squares decomposition of the complete data, we find that the first cross-validated partial least squares component is proportional to $c_{(i)}$.

Let us again consider the derivation of the $k + 1^{th}$ cross-validated vector of partial least squares component loadings $q_{k+1(i)}$. We must maximize the criterion $(c_{(i)}^T q)^2$, subject to the side-conditions $q_{j(i)}^T C_{(i)} q = 0$, $j = 1, \ldots, k$ and $q^T q = 1$. As for the partial least squares decomposition of the complete data, we may derive the downdated partial least squares decomposition from a rotation of the principal component decomposition $C = V E V^T$ of the complete data, and hence, we write

$$q = Vr,$$

as before. We must now find the maximum of the Lagrangian equation

$$L_{k+1}(r) = (d_{(i)}^T r)^2 - \lambda r^T r - \alpha^T Q_{(i)}^T C_{(i)} Vr,$$

with $d_{(i)} = V^T c_{(i)}$, $\alpha^T = (\alpha_1, \ldots, \alpha_k)$ and $Q_{(i)} = (q_{1(i)}, \ldots, q_{k(i)})$. We have

$$d_{(i)} = V^T c_{(i)} = d - \nu y_i f_i^T,$$

where $f_i = x_i V$ and

$$
\begin{aligned}
Q_{(i)}^T C_{(i)} Vr &= Q_{(i)}^T C Vr - \nu Q_{(i)}^T x_i^T x_i Vr \\
&= Q_{(i)}^T C Vr - \nu Q_{(i)}^T x_i^T x_i Vr \\
&= Q_{(i)}^T V Er - \nu Q_{(i)}^T V f_i^T f_i r \\
&= (Q_{(i)}^T V E - \nu Q_{(i)}^T V f_i^T f_i) r,
\end{aligned}
$$

and hence, writing $\mathbf{A}_{(i)}^T = (\mathbf{Q}_{(i)}^T \mathbf{VE} - \nu \mathbf{Q}_{(i)}^T \mathbf{Vf}_i^T \mathbf{f}_i)$, we have

$$L_{(i)}(\mathbf{r}) = (\mathbf{d}_{(i)}^T \mathbf{r})^2 - \lambda \mathbf{r}^T \mathbf{r} - \boldsymbol{\alpha}^T \mathbf{A}_{(i)}^T \mathbf{r}.$$

With the same argument as for the partial least squares decomposition of the complete data, we find that

$$\mathbf{r} = \frac{\mathbf{M}_{(i)} \mathbf{d}_{(i)}}{\|\mathbf{M}_{(i)} \mathbf{d}_{(i)}\|},$$

with $\mathbf{M}_{(i)} = \mathbf{I} - \mathbf{A}_{(i)}(\mathbf{A}_{(i)}^T \mathbf{A}_{(i)})^{-1} \mathbf{A}_{(i)}^T$.

With this cross-validatory algebra, a matrix $\mathbf{R} = (\mathbf{r}_1, \ldots, \mathbf{r}_{r_x})$ is derived. Each $\mathbf{r}_j$, $j = 1, \ldots, r_x$ is a linear combination that rotates the principal components to the downdated $j^{\text{th}}$ partial least squares component.

As for the computations on the complete data, we must exploit the rank-one nature of the leave-one-out cross-validatory computations to obtain an efficient computation of the projection matrices $\mathbf{M}$. We have, as before

$$\mathbf{M}_{k+1(i)} = \mathbf{M}_{k(i)} - \frac{(\mathbf{M}_{k(i)} \mathbf{a})(\mathbf{M}_{k(i)} \mathbf{a})^T}{\mathbf{a}^T \mathbf{M}_{k(i)} \mathbf{a}},$$

for $\mathbf{M}_{(i)}$ at each stage in the cross-validatory computations, where $\mathbf{a}$ is the column appended to $\mathbf{A}_{(i)}$ at stage $k + 1$.

Detailed proof of the validity of this cross-validatory algebra can be found in the paper from Stone and Brooks [SB90].

## 3.3 Partial Least Squares Regression and Influence Measures

### 3.3.1 Partial Least Squares Regression

The application of the Stone-Brooks leave-one-out cross-validatory algebra to partial least squares regression is completely analogous to the discussion for principal component regression.

As before, we consider the computation of the cross-validated predicted values $\widehat{y}_{i(i,\gamma)}$, for $i = 1, \ldots, n$ and $0 \leq \gamma \leq r_x$. As the partial least squares

components are uncorrelated, we may again exploit the updating formulae

$$\widehat{\widetilde{y}}_{i(i,\gamma)} = \overline{\overline{\mathbf{Y}}}_{(i)} + \sum_{k=1}^{\gamma} (\widehat{\widetilde{y}}_{i(i,k)} - \widehat{\widetilde{y}}_{i(i,k-1)}) \quad \text{for} \quad 0 < \gamma \leq r_x,$$

and

$$\widehat{\widetilde{y}}_{i(i,\gamma)} = \overline{\overline{\mathbf{Y}}}_{(i)} \quad \text{for} \quad \gamma = 0,$$

where

$$\widehat{\widetilde{y}}_{i(i,k)} - \widehat{\widetilde{y}}_{i(i,k-1)} = \frac{\sum_{j \neq i} u_{jk(i)} y_{j(i)}}{\sum_{j \neq i} u_{jk(i)}^2} u_{ik(i)},$$

as before. Therefore, the same procedure may be applied for the cross-validation of partial least squares regression as for principal component regression, and we find

$$u_{ik(i)} = \nu \mathbf{x}_i \mathbf{q}_{k(i)}$$

for $j = i$, and

$$u_{jk(i)} = \mathbf{x}_j \mathbf{q}_{k(i)} + \mathbf{x}_i \mathbf{q}_{k(i)}/(n-1),$$

if $j \neq i$. We have,

$$\mathbf{X} \mathbf{q}_{k(i)} = \mathbf{X} \mathbf{V} \mathbf{r}_k = \mathbf{F} \mathbf{r}_k,$$

with $\mathbf{F}$ the matrix of principal component scores for the complete data. Hence, as for principal component regression, the partial least squares component loadings are not downdated explicitly, but rather, we apply the linear combinations $\mathbf{r}_k$ directly to the principal component scores $\mathbf{F}$.

## 3.3.2 Influence Measures

There is no satisfactory approach in the literature to the analysis of influence in a partial least squares setting. This is probably due to the computational cost of naive implementations of full leave-one-out cross-validatory calculations for partial least squares regression, as well as to the algorithmic approach to partial least squares in most of the literature. Martens and Naes [MN89, page 285] extend the application of influence measures for ordinary least squares to the partial least squares setting. These measures have the advantage that closed form formulae exist which may be employed to reduce the computational cost. However, the validity of this approach in a partial least squares setting is doubtful. This is because it will typically ignore the

contributions of partial least squares components which have not been incorporated in the selected regression equation. Furthermore, this method will not identify observations which are influential for the partial least squares decomposition itself.

The Stone-Brooks cross-validatory algebra provides an ideal opportunity to consider the problem of influence for partial least squares. Full leave-one-out cross-validation for partial least squares regression involves the computation of the cross-validated sums of squared projections

$$h_{k(i)} = \sum_{j \neq i} (\mathbf{x}_j \mathbf{q}_{k(i)} + \mathbf{x}_i \mathbf{q}_{k(i)}/(n-1))^2,$$

for each $k^{th}$ partial least squares component and each $i^{th}$ observation in the cross-validatory calculations. Hence, as for principal component regression, we could consider the influence measures

$$h_k - h_{k(i)},$$

where $h_k$ is the sum of squared projections of the complete data on the $k^{th}$ vector of partial least squares component loadings $\mathbf{q}_k$. However, in contrast to the principal component eigenvalue perturbations, these statistics will reflect the extent to which an observation is outlying in the response space as well as in the predictor space. More importantly, these statistics do not have the elegant properties of the corresponding principal component influence measures. First of all, they may no longer be interpreted as sums of squared projection downdates for the corresponding partial least squares components. The sum of squared projections for a partial least squares component may well increase when an observation is removed from the data. Also, these measures no longer sum to the total principal component eigenvalue downdate $\rho$ and they are not central to the Stone-Brooks approach to partial least squares leave-one-out cross-validation. Hence, they are not appropriate as partial least squares influence measures. Likewise, their sum no longer appears suitable for a possible normalisation of the individual sums of squared partial least squares projection perturbations.

A more natural partial least squares influence measure would be the perturbation of the sample covariance vector $\mathbf{d}$. Indeed, partial least squares decomposition optimizes the sample covariance and hence, we could consider the influence measure

$$\|\mathbf{d}_k - \mathbf{d}_{k(i)}\| = \|\nu y_i \mathbf{f}_i^T\|.$$

This seems to be a summary influence measure that considers the influence of observations in both the predictor and response space simultaneously by weighting the summary principal component influence measure $\rho$ with the square of the observed response for the removed datum.

As always in these cross-validatory computations, we may easily generate the leave-one-out cross-validated residuals

$$y_{i(k)} - \widehat{y}_{i(i,k)}.$$

### 3.3.3   Computational Cost and Performance

We would like to compare the computational cost of the Stone-Brooks cross-validatory computations for partial least squares with naive implementations of full leave-one-out partial least squares cross-validation. For the naive implementations, we will consider the orthogonalized and non-orthogonalized partial least squares algorithms of Martens and Naes [MN89, pages 119 to 125] [Den91, pages 62 and 64], as well as the partial least squares computations described by Helland [Hel88].

### Simulations

The Stone-Brooks approach to full leave-one-out partial least squares regression cross-validation was implemented in the numerical analysis package GAUSS [Apt92] (appendix, page 146). The naive procedures for the cross-validation of partial least squares regression were also implemented in GAUSS (appendix, pages 141, 142 and 143 for the base procedures).

We used the native GAUSS procedure OLSQR [Apt92, volume 2, page 1352] for the computation of the least squares fits in both Martens's non-orthogonalized method as well as in Helland's method for partial least squares decomposition. This procedure is based on the QR decomposition. The initial eigendecomposition necessary for the Stone-Brooks approach was computed with the GAUSS procedure EIGHV [Apt92, volume 2, page 1180].

The simulations were done in the same manner as explained for principal component regression cross-validation. For the regular case ($n - 1 \geq p$), data were generated with the GAUSS pseudo-random number generator RNDU [Apt92, volume 2, page 1452]. We generated predictor matrices with 5, 10, 15, 20, 25, 30, 35, 40, 45 up to 50 predictor variables and 100 observations. For

Figure 3.1: Computation times in seconds for full leave-one-out cross-validation of partial least squares regression in the regular case with 100 observations.

each simulated problem, a response vector of 100 numbers was generated and the four implementations of full leave-one-out cross-validation were compared on each simulated problem. For each simulation, all partial least squares factors were derived and the PRESS statistic was computed for all factors. Figure 3.1 shows the computation times in seconds, with the number of predictor variables plotted on a log scale. There appears to be an order of magnitude difference in the computational cost between the non-orthogonal and the orthogonal methods.

For the singular case $(n - 1 < p)$, the same procedure was repeated, but with matrices of 12 observations and 15, 20, 25,50,75,100,500,1000,2500 up to 3000 predictor variables. For each simulation, all ten factors were derived in the cross-validation. A small number of observations was chosen so that the results will mainly reflect the costs due to the number of predictor variables derived in the cross-validatory computations. Figure 3.2 shows the results plotted on the log scale for the number of predictor variables. As for principal components, the computational cost due to the number of predictor variables appears to cancel out for the Stone-Brooks method. There is a slight increase in the computational cost of this method when the number of predictor variables exceeds 500. This is probably due to the initial eigendecomposition. All other methods are an order of magnitude more expensive.

Let us again consider a theoretical study of the computational costs to analyse these results. To simplify matters, we will only concern ourselves with orders of magnitude descriptions of the computational costs involved in these procedures. Also, we will have special attention for the computational costs due to the number of predictor variables in the problem, as partial least squares computations are often employed in problems with many predictor variables. Therefore, we will describe costs for the case when all partial least squares factors are derived.

A naive implementation of full leave-one-out cross-validation with the orthogonal scores algorithm of Martens and Naes involves an order of $n^2 p^2$ multiplications in the regular case and an order of $pn^3$ multiplications in the singular case. For the non-orthogonalized algorithm described in the same book, as well as for the algorithm due to Helland, we will consider implementations based on the QR decomposition for the ordinary least squares fitting in the partial least squares algorithms. A naive implementation of full

Figure 3.2: Computation times in seconds for full leave-one-out cross-validation of partial least squares regression in the singular case with 12 observations.

| Computation | Case | Regular | Singular |
|---|---|---|---|
| `w=xm'ym` | | $np^2$ | $pn^2$ |
| `w=w/sqrt(w'w)` | | $p^2$ | $pn$ |
| `t=xm*w` | | $np^2$ | $pn^2$ |
| `tt=t't` | | $np$ | $n^2$ |
| `p=xm't/tt` | | $np^2$ | $pn^2$ |
| `q=ym't/tt` | | $np$ | $n^2$ |
| `xm=xm-t*p'` | | $np^2$ | $pn^2$ |
| `ym=ym-t*q` | | $np$ | $n^2$ |
| Subtotal | | $np^2 + np + p^2$ | $pn^2 + pn + n^2$ |
| PRESS | | $p^2$ | $pn$ |
| Total | | $np^2 + np + p^2$ | $pn^2 + pn + n^2$ |

Table 3.1: Contributions of a single observation to the order of the number of multiplications involved in a full leave-one-out cross-validatory analysis of partial least squares regression with a naive implementation based on Martens's orthogonalized algorithm, when all factors are derived in the cross-validation.

| Case<br>Computation | Regular | Singular |
|---|---|---|
| w=xm'ym | $np^2$ | $pn^2$ |
| w=w/sqrt(w'w) | $p^2$ | $pn$ |
| t=xm*w | $np^2$ | $pn^2$ |
| qq=olsqr(y0,tt[.,1:i]) | $np^3$ | $n^4$ |
| xm=xm-t*w' | $np^2$ | $pn^2$ |
| ym=y0-tt[.,1:i]*qq | $np^2$ | $n^3$ |
| Subtotal | $np^3 + np^2 + p^2$ | $pn^2 + pn + n^4$ |
| PRESS | $p^2$ | $pn + n^2$ |
| Total | $np^3 + np^2 + p^2$ | $pn^2 + pn + n^4$ |

Table 3.2: Contributions of a single observation to the order of the number of multiplications involved in a full leave-one-out cross-validatory analysis of partial least squares regression with a naive implementation based on Martens's non-orthogonalized algorithm, when all factors are derived in the cross-validation.

| Case<br>Computation | Regular | Singular |
|---|---|---|
| w=s-xm'(xm*b) | $np^2$ | $pn^2$ |
| w=w/sqrt(w'w) | $p^2$ | $pn$ |
| xw[.,i]=xm*w | $np^2$ | $pn^2$ |
| b=ww[.,1:i]*olsqr(ym,xw[.,1:i]) | $np^3 + p^2$ | $n^4$ |
| Subtotal | $np^3 + np^2 + p^2$ | $pn^2 + pn + n^4$ |
| PRESS | $p^2$ | $pn$ |
| Total | $np^3 + np^2 + p^2$ | $pn^2 + pn + n^4$ |

Table 3.3: Contributions of a single observation to the order of the number of multiplications involved in a full leave-one-out cross-validatory analysis of partial least squares regression with a naive implementation based on Helland's algorithm, when all factors are derived in the cross-validation.

| Computation | Case | Regular | Singular |
|---|---|---|---|
| `ai=e.*zi-nu*sumc(fi.*zi)*fi` | | $p^2$ | $n^2$ |
| `miai=mi*ai` | | $p^3$ | $n^3$ |
| `mi=mi-miai*miai'/ai'miai` | | $p^3$ | $n^3$ |
| `zi=mi*di` | | $p^3$ | $n^3$ |
| `zi=zi/sqrt(zi'zi)` | | $p^2$ | $n^2$ |
| Subtotal | | $p^3$ | $n^3$ |
| PRESS | | $np^2 + np$ | $n^3$ |
| Total | | $np^2 + np + p^3$ | $n^3$ |

Table 3.4: Contributions of a single observation to the order of the number of multiplications involved in a full leave-one-out cross-validatory analysis of partial least squares regression with the Stone-Brooks method, when all factors are derived in the cross-validation.

leave-one-out cross-validation for both these algorithms will involve an order of $n^2 p^3$ multiplications for the regular case, using the Householder method for the QR decomposition [GL89, page 219]. In the singular case, these algorithms require an order of $pn^3 + n^5$ multiplications. The Stone-Brooks cross-validatory computations are based on the eigendecomposition of the cross-products matrix, in the regular case, and of the matrix $\mathbf{XX}^T$, in the singular case. However, as these calculations are only done once, we may largely ignore the computational cost of these decompositions. The remainder of the calculations are of order $n^2 p^2 + n^2 p$ in the regular case and of order $n^4$ in the singular case.

As for principal components, we have summarized these results in tables 3.1 to 3.4, which consider the contributions of a single observation to the computational costs of cross-validating the regression equation, for each observation. Again, the relevant quantities must be multiplied by $n$ for comparison with the text.

A few conclusions may be drawn from this rudimentary analysis of computational cost. First of all, computations are considerably slower in the regular case, not only because the number of observations tends to be larger,

but also because the computational costs due to the number of predictor variables increase. In the singular case, the number of partial least squares components that may be derived is effectively restricted by the number of observations. Furthermore, in the regular case, the computational cost of the non-orthogonal methods relates to that of the orthogonal approaches as $n^2p^3$ relates to $n^2p^2$. Thus, for the same number of observations, we observe in effect an order of magnitude increase in cost due to the number of predictor variables only. This is because both these methods compute correlated scores and hence, a numerically optimized least squares fitting procedure must be employed, like those based on the QR algorithm, to obtain reliable least squares fits. It is the cost of computing these QR decompositions that dominates the total computational cost for the non-orthogonal methods. Our implementation is based on Denham's approach [Den91], who has implemented the algorithms of Martens, Naes and Helland in Splus [Sta92], using the native LS procedure for least squares fitting, which is based on the QR decomposition. Likewise, the tables on computational cost assume a QR decomposition for the non-orthogonal methods. Martens's orthogonalized method and the Stone-Brooks approach do not have these problems with uncorrelated scores.

In the singular case, the computational cost due to the number of predictor variables cancels out for the Stone-Brooks approach. This is in analogy to the efficient cross-validation of principal component regression. The reasons are similar. The cost of computing the partial least squares scores does not depend on the number of predictor variables and neither does the cost of computing the matrices $M_{(i)}$.

## 3.3.4 Accuracy

Applications of the Stone-Brooks algebra, for both partial least squares regression and leave-one-out cross-validation, to partial least squares regression, revealed differences in accuracy when compared to implementations using the existing methods.

Table 3.5 shows maximum absolute differences between the fitted values from a partial least squares regression for the Mustard data with the Stone-Brooks approach, Martens's non-orthogonalized algorithm and Helland's algorithm on the one hand, and those obtained from Martens's orthogonalized algorithm on the other hand. Similarly, table 3.6 shows absolute differences

| $\gamma$ | Stone-Brooks | Martens (*) | Helland |
|---|---|---|---|
| 1 | 2.1094237e-15 | 6.6613381e-16 | 6.6613381e-16 |
| 2 | 4.0772941e-14 | 4.8849813e-15 | 3.9968029e-15 |
| 3 | 1.8185453e-13 | 7.9936058e-15 | 7.0221606e-15 |
| 4 | 4.9138471e-13 | 6.2172489e-15 | 6.6613381e-15 |
| 5 | 1.2800871e-12 | 3.5527137e-15 | 6.5225603e-15 |
| 6 | 3.3240077e-12 | 6.8833828e-15 | 9.7699626e-15 |
| 7 | 4.4608761e-12 | 1.3544721e-14 | 7.2164497e-15 |
| 8 | 1.0283108e-11 | 4.5907722e-14 | 9.7699626e-15 |
| 9 | 2.1280755e-11 | 9.5146113e-14 | 8.8262730e-15 |
| 10 | 5.5720317e-11 | 2.4125146e-13 | 7.2164497e-15 |
| 11 | 1.5374679e-10 | 4.5086157e-13 | 8.1601392e-15 |
| 12 | 3.2566894e-10 | 8.3566487e-13 | 1.2323476e-14 |
| 13 | 4.3988224e-10 | 1.0489387e-12 | 1.3211654e-14 |
| 14 | 6.5851058e-10 | 1.4628299e-12 | 1.9984014e-14 |
| 15 | 8.3024121e-10 | 1.9499957e-12 | 2.2204460e-14 |
| 16 | 1.1768286e-09 | 3.2224223e-12 | 2.1538327e-14 |
| 17 | 1.2683634e-09 | 3.4506842e-12 | 2.2537527e-14 |
| 18 | 1.6198187e-09 | 3.6858294e-12 | 4.1189274e-14 |
| 19 | 1.8266110e-09 | 3.3308911e-12 | 1.1313173e-13 |
| 20 | 1.8995741e-09 | 3.1274983e-12 | 5.9507954e-14 |
| 21 | 1.9695808e-09 | 2.6390001e-12 | 3.7747583e-14 |
| 22 | 2.1251696e-09 | 2.0938806e-12 | 3.2418512e-14 |
| 23 | 2.2867854e-09 | 2.9338754e-12 | 4.4186876e-14 |
| 24 | 2.2528285e-09 | 4.3944848e-12 | 4.3076653e-14 |
| 25 | 2.0772032e-09 | 5.2757798e-12 | 3.7525538e-14 |
| 26 | 1.9582371e-09 | 5.6170624e-12 | 2.6645353e-14 |
| 27 | 1.9700801e-09 | 5.6317173e-12 | 2.9753977e-14 |
| 28 | 1.9536739e-09 | 5.6223914e-12 | 3.9968029e-14 |
| 29 | 2.0686239e-09 | 5.4396487e-12 | 4.9737992e-14 |
| 30 | 2.1287563e-09 | 5.4021232e-12 | 2.9087843e-14 |
| 31 | 2.1090807e-09 | 5.3781424e-12 | 3.5305092e-14 |
| 32 | 2.0819053e-09 | 5.3517191e-12 | 3.4194869e-14 |
| 33 | 2.0552013e-09 | 5.3388405e-12 | 4.1522341e-14 |
| 34 | 1.8951911e-09 | 5.3332894e-12 | 3.8746784e-14 |
| 35 | 1.7918858e-09 | 5.3228533e-12 | 3.1974423e-14 |
| 36 | 1.6417112e-09 | 5.3153038e-12 | 4.7961635e-14 |
| 37 | 1.6355675e-09 | 5.3157478e-12 | 3.2862602e-14 |
| 38 | 1.8173483e-09 | 5.3259619e-12 | 2.3314684e-14 |
| 39 | 1.7308075e-09 | 5.3270721e-12 | 2.9531932e-14 |
| 40 | 1.7773641e-09 | 5.3321791e-12 | 2.9087843e-14 |
| 41 | 1.7830865e-09 | 5.3281823e-12 | 1.7874591e-14 |
| 42 | 1.7169376e-09 | 5.3295146e-12 | 1.2656542e-14 |
| 43 | 1.5133301e-09 | 5.3292926e-12 | 1.8873791e-14 |
| 44 | 1.4821433e-09 | 5.3315130e-12 | 1.3766766e-14 |
| 45 | 1.0802723e-09 | 5.3390625e-12 | 1.9761970e-14 |
| 46 | 9.1030483e-10 | 5.3281823e-12 | 1.8651747e-14 |
| 47 | 8.7347274e-10 | 5.3266280e-12 | 2.2870594e-14 |
| 48 | 9.3170704e-10 | 5.3261839e-12 | 2.2648550e-14 |

Table 3.5: Maximum absolute differences in fitted values for the Mustard data, in a comparison with Martens's orthogonalized algorithm. (* = Martens's non-orthogonalized algorithm)

between mean prediction error sum of squares for the Mustard data and partial least squares regression, in a comparison of leave-one-out cross-validation computed with the same three methods, and a naive implementation based on Martens's orthogonalized algorithm. The choice of the reference algorithm is somewhat arbitrary and we have chosen the orthogonalized algorithm as it has been the earliest implementation, as well as the most simple, of the partial least squares decomposition.

The results suggest that the Stone-Brooks approach performs adequately for the Mustard data in the computation of the partial least squares regression equation. Nevertheless, the maximum absolute differences between fitted values obtained from a comparison with Martens's orthogonalized algorithm tend to be larger than those obtained for Martens's non-orthogonalized algorithm and Helland's algorithm. The best accuracy is obtained for Martens's orthogonalized algorithm and Helland's approach, with a maximum absolute difference in fitted values of $1.13e^{-13}$ between both methods. All methods seem to suffer from a gradual loss of accuracy as further components are added. This decrease is more pronounced for the Stone-Brooks method, although the accuracy stabilizes when 10 components have been added. A maximum absolute difference of $2.29e^{-9}$ is obtained for the Stone-Brooks method.

Things are different in cross-validation. All naive methods perform well for the Mustard data in this comparison. As before, all methods suffer from a gradual loss of accuracy as further partial least squares components are added to the prediction equation. However, in contrast to the Stone-Brooks approach, the accuracy of the naive methods seems to stabilize after a certain number of components has been included in the regression equation. The Stone-Brooks method suffers from an uninterrupted loss of accuracy, and eventually loses all accuracy when most components are in the prediction equation. This loss of accuracy is probably caused by the fact that the Euclidean length of the leave-one-out perturbations in the cross-validatory algebra becomes smaller as more components are added. This may provide an explanation as to the behaviour of the Stone-Brooks algebra in cross-validation, as compared to the computation of the partial least squares decomposition.

The Stone-Brooks algebra performs well for the first 25 partial least squares components. The remaining components are of small variation and we are unlikely to include them in any practical application to calibration.

| $\gamma$ | Stone-Brooks | Martens (*) | Helland |
|---|---|---|---|
| 1 | 1.7763568e-15 | 0 | 0 |
| 2 | 5.9952043e-15 | 6.6613381e-16 | 2.2204460e-16 |
| 3 | 1.4321877e-14 | 3.3306691e-16 | 5.5511151e-16 |
| 4 | 3.7747583e-15 | 0 | 1.1102230e-16 |
| 5 | 6.3116179e-14 | 2.7755576e-16 | 2.7755576e-16 |
| 6 | 1.8823831e-13 | 5.5511151e-16 | 0 |
| 7 | 9.6922470e-14 | 1.6653345e-15 | 1.6653345e-16 |
| 8 | 2.6773028e-13 | 4.3298698e-15 | 3.3306691e-16 |
| 9 | 4.9182880e-13 | 6.1062266e-16 | 3.3306691e-16 |
| 10 | 1.2557733e-12 | 1.3822277e-14 | 6.1062266e-16 |
| 11 | 5.6565863e-13 | 6.1062266e-14 | 1.4432899e-15 |
| 12 | 2.9880320e-11 | 8.2711615e-14 | 1.9428903e-15 |
| 13 | 4.5605630e-11 | 1.1113332e-13 | 1.3877788e-15 |
| 14 | 6.9610817e-11 | 1.9240165e-13 | 1.5543122e-15 |
| 15 | 8.7132079e-11 | 3.0103697e-13 | 8.3266727e-16 |
| 16 | 7.1873729e-11 | 4.5080606e-13 | 3.8857806e-15 |
| 17 | 6.9713457e-11 | 4.4508841e-13 | 5.5511151e-15 |
| 18 | 1.7417290e-11 | 3.2129854e-13 | 3.7747583e-15 |
| 19 | 2.7936542e-12 | 8.0813134e-13 | 7.8825835e-15 |
| 20 | 1.4097845e-10 | 8.7863050e-13 | 9.6589403e-15 |
| 21 | 9.6040953e-11 | 1.0759171e-12 | 1.7763568e-14 |
| 22 | 1.3776080e-10 | 1.0439427e-12 | 1.6209256e-14 |
| 23 | 1.4086821e-10 | 1.0396128e-12 | 1.9095836e-14 |
| 24 | 6.6690653e-10 | 1.0359491e-12 | 2.9642955e-14 |
| 25 | 4.8622550e-10 | 9.1748831e-13 | 3.5749181e-14 |
| 26 | 1.7090758e-09 | 8.9261931e-13 | 4.4853010e-14 |
| 27 | 3.4633689e-09 | 9.0971675e-13 | 5.5067062e-14 |
| 28 | 1.0639859e-08 | 8.9561691e-13 | 6.5836225e-14 |
| 29 | 1.6226559e-08 | 9.0838448e-13 | 7.7049478e-14 |
| 30 | 3.5857015e-08 | 9.1004981e-13 | 8.4598994e-14 |
| 31 | 5.1346823e-08 | 9.1060492e-13 | 9.5923269e-14 |
| 32 | 1.3304011e-07 | 9.1149310e-13 | 1.0480505e-13 |
| 33 | 6.4297527e-08 | 9.1249230e-13 | 1.2234658e-13 |
| 34 | 2.1965565e-07 | 9.1304742e-13 | 1.3000712e-13 |
| 35 | 3.6652035e-06 | 9.1260333e-13 | 1.3855583e-13 |
| 36 | 5.3306525e-06 | 9.1293639e-13 | 1.4599433e-13 |
| 37 | 5.8052278e-06 | 9.1315844e-13 | 1.6264767e-13 |
| 38 | 5.6526496e-05 | 9.1282537e-13 | 1.7053026e-13 |
| 39 | 8.9535199e-05 | 9.1271435e-13 | 1.8829382e-13 |
| 40 | 0.00039045237 | 9.1304742e-13 | 1.9506619e-13 |
| 41 | 0.0010489552 | 9.1260333e-13 | 2.0994317e-13 |
| 42 | 0.0083168697 | 9.1238128e-13 | 2.3081537e-13 |
| 43 | 0.037980110 | 9.1282537e-13 | 2.4757973e-13 |
| 44 | 1.2591844 | 9.1282537e-13 | 2.6489921e-13 |
| 45 | 40.507619 | 9.1271435e-13 | 2.8010927e-13 |
| 46 | 140.25491 | 9.1271435e-13 | 2.9687364e-13 |
| 47 | 201.16291 | 9.1271435e-13 | 3.0431213e-13 |
| 48 | 202.62949 | 9.1249230e-13 | 3.2385206e-13 |

Table 3.6: Absolute differences between mean prediction error sum of squares for the Mustard data, in a comparison with Martens's orthogonalized algorithm. (* = Martens's non-orthogonalized algorithm)

For this reason, the Stone-Brooks algebra may still be of interest. All cross-validatory results on the Mustard data in this chapter have been computed with the Stone-Brooks approach.

## 3.4 Efficient Leave-One-Out Partial Least Squares Regression Cross-validation Applied to Near Infrared Spectroscopy

### 3.4.1 Cross-validation

**GAUSS Implementation**

A complete transcript of a GAUSS implementation of full leave-one-out cross-validation for partial least squares regression with the Stone-Brooks method is in the appendix (page 146). The code is used within a GAUSS session, as follows.

As for the cross-validation of principal component regression, a GAUSS matrix x must have been created for the predictor data, as well as a column vector y for the observed responses, with the same conventions (page 60). The statements

$$\{\mathrm{msep}, \mathrm{dd}, \mathrm{ddi}, \mathrm{ycvpred}\} = \mathrm{plscv}(\mathrm{y}, \mathrm{x}, \mathrm{w});$$

will compute a full leave-one-out cross-validation with the Stone-Brooks method. As for principal component regression, w is the number of partial least squares components that will be derived in the cross-validation.

The output from this procedure is organized with the same conventions as those for the cross-validation of principal component regression. msep is a matrix that contains the PRESS and MSEP statistics. The cross-validated predicted values are in matrix ycvpred. The format of these matrices is exactly the same as that from the corresponding objects that are output for the cross-validation of principal component regression. dd and ddi are matrices of influence measures. The first is a matrix that contains the cross-validated sums of squared projections for each of the partial least squares components derived. The matrix has w columns, such that the first column contains the influence measures $h_{1(i)}$ for the first component, and so on, up

to the last column, which contains the measures $h_{w(i)}$ for the last component derived in the cross-validation. ddi is a column vector that contains the cross-validated influence measures $\|d_k - d_{k(i)}\|$.

## The Mustard data

This procedure was applied to the Mustard data. Table 3.7 shows the results. A plot of the cross-validatory index is shown in figure 3.3. The cross-



number of partial least squares components

Figure 3.3: Cross-validatory index for a full leave-one-out cross-validation of partial least squares regression of the Mustard data, as computed with the Stone-Brooks method.

validatory index drops below zero when more than 44 components are used

for prediction, which is due only to the loss of all significant digits in the Stone-Brooks approach. The plot only considers the first 44 prediction equations. The cross-validatory index equals 89.7% for a six component partial least squares regression prediction equation. This increases to 89.8% when a further partial least squares component is added. Hence, a six component prediction equation seems appropriate.

As for principal component regression, a special GAUSS program has been written to compute the selected partial least squares regression equation within the same GAUSS session, using the Stone-Brooks algebra. The code is in the appendix (page 144). With this code, the statements

$$\{\text{ypred}, \text{yresid}, \text{beta}, \text{q}, \text{u}\} = \text{pls}(\text{y}, \text{x}, \text{w});$$

will compute the prediction equations for each of the first **w** partial least squares regression equations. The output has the same format as described for the corresponding procedure for principal component regression (page 64). Figure 3.4 shows plots of the first six partial least squares component loadings. A plot of the regression coefficients for the selected partial least squares prediction equation is shown in figure 3.5.

## 3.4.2 Influence Analysis

As for principal component regression, the output from `plscv` and `pls` may be used to analyse the influence of observations on the partial least squares regression equations.

A plot of the cross-validated predicted values versus the observed values is shown in figure 3.6. As for principal component regression, this plot shows a straight line relation between the observed and cross-validated predicted values. The second plot in the same figure, which compares the cross-validated residuals with the observed values does not identify any observations with large cross-validated residuals. Hence, the conclusions from this analysis are similar to those found for principal component regression.

We will develop the partial least squares influence analysis, to exhibit its potential and to compare the influence measures with those from principal component regression.

In analogy to principal components, we may study the perturbations of the sums of squared projections for each of the partial least squares components. Figure 3.7 shows a plot of these measures versus the observation

| $\gamma$ | MSEP | INDEX | $\gamma$ | MSEP | INDEX |
|---|---|---|---|---|---|
| 0 | 3.3466262 | 0 | | | |
| 1 | 3.1766852 | 0.050779811 | 25 | 0.51757407 | 0.84534452 |
| 2 | 1.7194746 | 0.48620655 | 26 | 0.52125185 | 0.84424557 |
| 3 | 0.64477544 | 0.80733569 | 27 | 0.52201571 | 0.84401732 |
| 4 | 0.50889399 | 0.84793820 | 28 | 0.52283098 | 0.84377372 |
| 5 | 0.43636637 | 0.86961006 | 29 | 0.52230300 | 0.84393148 |
| 6 | 0.34354859 | 0.89734480 | 30 | 0.52153150 | 0.84416201 |
| 7 | 0.34208401 | 0.89778243 | 31 | 0.52149933 | 0.84417162 |
| 8 | 0.34769427 | 0.89610603 | 32 | 0.52164224 | 0.84412892 |
| 9 | 0.37389645 | 0.88827660 | 33 | 0.52168057 | 0.84411747 |
| 10 | 0.43539933 | 0.86989903 | 34 | 0.52171472 | 0.84410726 |
| 11 | 0.53070963 | 0.84141951 | 35 | 0.52171410 | 0.84410745 |
| 12 | 0.49857490 | 0.85102163 | 36 | 0.52170876 | 0.84410904 |
| 13 | 0.47934740 | 0.85676697 | 37 | 0.52171975 | 0.84410576 |
| 14 | 0.48822256 | 0.85411500 | 38 | 0.52177101 | 0.84409044 |
| 15 | 0.43427971 | 0.87023358 | 39 | 0.52180405 | 0.84408057 |
| 16 | 0.49501340 | 0.85208584 | 40 | 0.52132408 | 0.84422399 |
| 17 | 0.52920586 | 0.84186885 | 41 | 0.52276350 | 0.84379388 |
| 18 | 0.52589643 | 0.84285773 | 42 | 0.53003141 | 0.84162217 |
| 19 | 0.57890167 | 0.82701932 | 43 | 0.55969465 | 0.83275854 |
| 20 | 0.56126072 | 0.83229059 | 44 | 1.7808989 | 0.46785245 |
| 21 | 0.56474361 | 0.83124987 | 45 | 41.029334 | -11.259909 |
| 22 | 0.55657440 | 0.83369090 | 46 | 140.77662 | -41.065236 |
| 23 | 0.52700528 | 0.84252640 | 47 | 201.68462 | -59.265057 |
| 24 | 0.52569122 | 0.84291905 | 48 | 203.15120 | -59.703285 |

Table 3.7: Full leave-one-out cross-validatory analysis for partial least squares regression of the Mustard data, as computed with the Stone-Brooks method.

Figure 3.4:  Partial least squares loadings for the first six partial least squares components.

Figure 3.5: Regression coefficients for a partial least squares regression with six components.

Figure 3.6: Scatterplots obtained from the cross-validated predicted values, the cross-validated residuals and the observed values.

Figure 3.7: Perturbations of the sums of squared projections for the first six partial least squares components.

number for the first six components. Plots of the cross-validated residuals versus these perturbation measures are shown in figure 3.8 for the first six components.

These perturbation measures are not restricted to be positive. This is because they may no longer be interpreted as downdates with respect to the partial least squares regression of the complete data. Nevertheless, for most observations these perturbation measures are positive. There are no clear outliers.

In analogy to the influence measures for principal components, we may normalize these perturbations with the total perturbation of the sum of squared partial least squares projections. This total perturbation is defined as the sum of the perturbations across all partial least squares components.

The effect is shown in figures 3.9 and 3.10. Some of the plots in these figures appear to identify outliers, none of which have large cross-validatory residuals.

The figures 3.11 and 3.12 show plots of the Euclidean length of the covariance perturbations versus the observation number and of the cross-validated residuals versus the Euclidean length of the covariance perturbations respectively. There are some observations with large values for the partial least squares influence measure considered in these plots. As before, none of these have large cross-validated residuals.

## 3.5   Conclusions

The application of the Stone-Brooks algebra to partial least squares regression and leave-one-out cross-validation for partial least squares regression emerges as an efficient computational method in comparison to naive implementations. However, in contrast to the efficient cross-validatory algebra for principal component regression, the present implementation of the Stone-Brooks approach to leave-one-out cross-validation is numerically unstable, at least in the partial least squares case. It is likely that these problems will persist in the full algorithm for continuum regression. This is because the cross-validatory algebra for partial least squares which is presented here is a restriction of the cross-validatory algebra for continuum regression. In contrast, the numerical accuracy of predicted values obtained from the partial least squares regression equation computed with the Stone-Brooks method

Figure 3.8: Cross-validated residuals versus the perturbations of the sums of squared projections for the first six partial least squares components.

Figure 3.9: Normalized perturbations of the sums of squared projections for the first six partial least squares components.

Figure 3.10: Cross-validated residuals versus the normalized perturbations of the sums of squared projections for the first six partial least squares components.

Figure 3.11: Length of the covariance perturbations.

Figure 3.12: Cross-validated residuals versus the length of the covariance perturbation.

was adequate for the analysis of the Mustard data considered here, although consistently lower than that obtained from naive implementations.

Apart from these differences in performance, the fundamental difference between the cross-validatory algebra provided by the Stone-Brooks method and the efficient algebra for the cross-validation of principal components is in the sequential nature of the Stone-Brooks approach. In contrast, the principal component cross-validatory algebra can be modified to cross-validate a single principal component or eigenvalue, ignoring all other components. Furthermore, the restriction of the Stone-Brooks approach to partial least squares leads to an algorithm which is purely algebraic in nature, as a parameter optimization cancels, as compared to the more general implementation for continuum least squares. Efficient computations for principal component cross-validation are based, essentially, on Newton-Raphson optimization, as the computations are reduced to the downdating of the eigenvalues.

The analysis on the accuracy of the Stone-Brooks implementation which is presented here is clearly incomplete. It points to at least two issues which require further investigation. First, the performance of the Sherman-Morisson formula in the downdating of the partial least squares regression equation needs to be investigated. Secondly, given the sequential nature of the Stone-Brooks approach, any comprehensive analysis of the performance of this method should address the accumulation of computational errors in the implementation. Given the results on the Mustard data, it is unlikely that the second point on its own could account for the loss of accuracy observed.

As for principal components, it is relatively easy to generate influence measures from cross-validatory computations for partial least squares regression. In addition to the efficiency gains, it is here that the Stone-Brooks approach has the advantage over naive implementations, as it is based explicitly on the downdating of the sample covariance vector. It is easy to store a measure of the leave-one-out perturbations of this vector during computations.

# Appendix

## A. Principal Component Computations

### A.1 Principal Component Decomposition

```
proc(5)=pcr(y,x,w);  @ univariate principal component regression @
                     @ [Jol86][MN89] @

    local beta,scores,pred,resid,be,re,pr,e,u,v,k,
          meanx,meany,xm,ym,n,p;  @ declarations @

    n=rows(x);p=cols(x);  @ initializations @
    beta=zeros(p,w);pred=zeros(n,w);resid=pred;

    if p==1;print "error:p == 1";end;endif;

    meanx=meanc(x)';meany=meanc(y);  @ mean centering @
    xm=x-meanx;ym=y-meany;

    @ principal component decomposition @

    if n-1 >= p;  @ regular case @
       if p < w;print "error:p < w";end;
       else;
       {e,v}=eighv(xm'xm);e=rev(e);  @ pca from proc EIG @
       v=rev(v')';  @ principal components @
       scores=xm*v;  @ scores @
       endif;
```

```
    endif;
    if n-1 < p;  @ singular case @
        if n-1 < w;print "error:n-2 < w";end;
        else;
        {e,u}=eighv(xm*xm');e=rev(e[2:n]);  @ pca from proc EIG @
        v=xm'rev(u[.,2:n]')'./sqrt(e');  @ principal components @
        scores=rev(u[.,2:n]')'.*sqrt(e');  @ scores @
        endif;
    endif;

    @ principal component regression @

    k=1;
    do while k <= w;
        {be,re,pr}=olsqr2(ym,scores[.,1:k]);
        beta[.,k]=v[.,1:k]*be;
        pred[.,k]=pr+meany;resid[.,k]=re;
        k=k+1;
    endo;

    @ return predicted values @
    @         residuals @
    @         beta coefficients @
    @         principal components @
    @         principal component scores @

    retp(pred,resid,beta,v,scores);
endp;
```

## A.2  Principal Component Cross-validation

**The Quadratic Method**

```
proc(6)=pcrcvq(y,x,w);  @ full leave-one-out cross-validation @
                        @ univariate principal component regression @
                        @ quadratic method @
                        @ [BNS78] @
```

```
    local n,p,msep,dd,ddiff,ypred,d,q;

    n=rows(x);p=cols(x);

    if p==1;print "error:p == 1";end;endif;

    if n-1 >= p; @ regular case @
        if p < w;print "error:p < w";end;
        else;
        {msep,dd,ddiff,ypred,d,q}=pcrcvqr(y,x,w,n,p);
        endif;
    endif;
    if n-1 < p; @ singular case @
        if n-2 < w;print "error:n-2 < w";end;
        else;
        {msep,dd,ddiff,ypred,d,q}=pcrcvqs(y,x,w,n);
        endif;
    endif;

    @ return msep @
    @           cross-validatory index @
    @           downdated eigenvalues @
    @           eigenvalue downdates @
    @           cross-validated predicted values @
    @           initial variances @
    @           initial principal components @

    retp(msep~1-msep/msep[1],dd,ddiff,ypred,d/(n-1),q);
endp;

proc(6)=pcrcvqr(y,x,w,n,p);
@ leave-one-out cross-validation for pcr @
@ quadratic method @
@ regular case @

    local tol1,tol2,toles,xm,u,urt,udf,q,d,nu,i,z,zz,rho,tau,e,pe,
```

```
              press,pressu,modif,dd,ddd,ddff,ddiff,yp,ypred;

toll=10e-16; @ machine epsilon @
tol2=10e-10; @ tolerance relative accuracy @

press=zeros(1,w+1);
ypred=zeros(n,w+1);
ddiff=zeros(n,w);
ddd=ddiff;
modif=ones(p,1);

@ overall principal component decomposition @

xm=x-meanc(x)'; @ pca from proc EIG @
{d,q}=eighv(xm'xm);d=rev(d);q=rev(q')';
if minc(d) < 10*toll;
    print "warning:overall eigenvalues < 10*toll";
endif;
u=xm*q; @ scores @ @ keep u d q @
nu=n/(n-1); @ cross-validatory constant @

@ cross-validation @

i=1;
print "cross-validating";
do while i <= n; @ leave-one-out loop @
    z=u[i,.]';zz=z'z;z=z/sqrt(zz); @ pca scores normalized @
    rho=nu*zz; @ deflation size @
    toles=2*toll*(d[1]+rho); @ tolerance deflation @
    tau=4*toll*maxc(1|(d[1]-rho))/rho; @ tolerance convergence @

    if minc(d[1:(p-1)]-d[2:p]) < toles; @ equal eigenvalues @
        {urt,z}=dfleig(u,d,z,p,toles); @ eigenvalues @
        {udf,e,z,modif,pe}=dflsco(urt,d,z,p,toles); @ scores @
        {pressu,dd,ddff,yp}=
        qmdfpcr(y,i,urt,udf,d,e,z,rho,modif,pe,w,n,nu,tol2,tau);
```

```
        press=press+pressu;
        ddiff[i,.]=ddff;
        ddd[i,.]=dd;
        ypred[i,.]=yp;
        modif=ones(p,1);
    elseif minc(abs(z)) < toles; @ zero scores @
        {udf,e,z,modif,pe}=dflsco(u,d,z,p,toles); @ scores @
        {pressu,dd,ddff,yp}=
        qmdfpcr(y,i,u,udf,d,e,z,rho,modif,pe,w,n,nu,tol2,tau);
        press=press+pressu;
        ddiff[i,.]=ddff;
        ddd[i,.]=dd;
        ypred[i,.]=yp;
        modif=ones(p,1);
    else;
        {pressu,dd,ddff,yp}=
        qmdfpcr(y,i,u,u,d,d,z,rho,modif,p,w,n,nu,tol2,tau);
        press=press+pressu;
        ddiff[i,.]=ddff;
        ddd[i,.]=dd;
        ypred[i,.]=yp;
    endif;

    i=i+1;
  endo;

  retp(press',ddd,ddiff,ypred,d,q);
endp;

proc(6)=pcrcvqs(y,x,w,n);
@ leave-one-out cross-validation for pcr @
@ quadratic method @
@ singular case @

  local tol1,tol2,toles,xm,q,urt,udf,u,d,sqrtd,nu,i,z,zz,rho,tau,e,
       nn,press,pressu,nr,modif,dd,ddd,ddff,ddiff,yp,ypred;
```

```
tol1=10e-16; @ machine epsilon @
tol2=10e-10; @ tolerance relative accuracy @

press=zeros(1,w+1);
ypred=zeros(n,w+1);
nu=n/(n-1);nr=n;n=n-1; @ cross-validatory constant @
ddiff=zeros(nr,w);
ddd=ddiff;
modif=ones(n,1);

@ overall principal component decomposition @

xm=x-meanc(x)'; @ pca from proc EIG @
{d,u}=eighv(xm*xm');d=rev(d[2:nr]);u=rev(u[.,2:nr]')')';
sqrtd=sqrt(d');
q=xm'u./sqrtd;
u=u.*sqrtd; @ scores @ @ keep q d u @
if minc(d) < 10*tol1;
    print "warning :  overall eigenvalues <10*tol1";
endif;

@ cross-validation @

i=1;
print "cross-validating";
do while i <= nr; @ leave-one-out loop @
    z=q'xm[i,.]';zz=z'z;z=z/sqrt(zz); @ pca scores normalized @
    rho=nu*zz; @ deflation size @
    toles=2*tol1*(d[1]+rho); @ tolerance deflation @
    tau=4*tol1*maxc(1|(d[1]-rho))/rho; @ tolerance convergence @

    if minc(d[1:(n-1)]-d[2:n]) < toles; @ equal eigenvalues @
        {urt,z}=dfleig(u,d,z,n,toles); @ eigenvalues @
        {udf,e,z,modif,nn}=dflsco(urt,d,z,n,toles); @ scores @
        {pressu,dd,ddff,yp}=
        qmdfpcr(y,i,urt,udf,d,e,z,rho,modif,nn,w,nr,nu,tol2,tau);
```

```
        press=press+pressu;
        ddiff[i,.]=ddff;
        ddd[i,.]=dd;
        ypred[i,.]=yp;
        modif=ones(n,1);
    elseif minc(abs(z)) < toles; @ zero scores @
        {udf,e,z,modif,nn}=dflsco(u,d,z,n,toles);@ scores @
        {pressu,dd,ddff,yp}=
        qmdfpcr(y,i,u,udf,d,e,z,rho,modif,nn,w,nr,nu,tol2,tau);
        press=press+pressu;
        ddiff[i,.]=ddff;
        ddd[i,.]=dd;
        ypred[i,.]=yp;
        modif=ones(n,1);
    else;
        {pressu,dd,ddff,yp}=
        qmdfpcr(y,i,u,u,d,d,z,rho,modif,n,w,nr,nu,tol2,tau);
        press=press+pressu;
        ddiff[i,.]=ddff;
        ddd[i,.]=dd;
        ypred[i,.]=yp;
    endif;

    i=i+1;
  endo;

  retp(press',ddd,ddiff,ypred,d,q);
endp;

proc(4)=qmdfpcr(y,i,ud,udf,d,ddf,z,rho,modif,p,w,n,nu,tol,tau);
@ rank-one modification with the quadratic method @

  local z2,del,dl,nr,g,gf,gs,j,k,e,t,gf,gs,den,h2,h3,f1,fp1,s1,
        sp1,a,xd,w1,w2,tdel,c,b,f,it,press,count,yh,yhm,ud,dfl,ind,
        ddiff,pupdate,ypred,fe;

  z2=z^2;
```

```
del=(ddf'-ddf)/rho;
nr=seqa(1,1,p);
g=ones(p,1);
count=seqa(1,1,n);
press=zeros(1,w+1);
ypred=press;
dfl=minc(modif)==0;
ddiff=zeros(1,rows(d));
yh=selif(y,count./=i);yhm=meanc(yh);yh=yh-yhm;
press[1]=yhm-y[i];
ypred[1]=yhm;

if maxc(z) /= 1; @ downdate eigenvalues @

    j=1;
    k=1;
    do while j <= w; @ first w components @
        if modif[j] /= 0; @ eigenvalue changed @
            if k < p;

                gf=nr.<=k;
                gs=g-gf;

                @ initialize t @
                if p > 2;
                    e=1+ones(p-2,1)'(selif(z2,nr./=k .and nr./= k+1)./
                    (selif(del[.,k],nr./=k .and nr./= k+1)-
                    del[k+1,k]));
                else;
                    e=1;
                endif;
                f=z2[k]+z2[k+1]+del[k+1,k]*e;
                t=2*z2[k]*del[k+1,k]/
                    (f+sqrt(f^2-4*z2[k]*del[k+1,k]*e));
                if t <= 0 or t >= 1;
                    print "error:initial values k=" k "tinit=" t;
```

```
   end;
endif;

den=del[.,k]-t; @ initialize denominator @
d[j]=d[j]/rho-t; @ initialize normalized eigenvalue @
gf=nr.<= k;
gs=g-gf;
w1=0;
w2=0;
it=1;
do until w1 and w2;
   dl=den[k+1,.];
   h2=z2 ./den;
   h3=h2 ./den;
   f1=gf'h2;
   fp1=gf'h3;
   s1=gs'h2;
   sp1=gs'h3;
   c=1+s1-dl*sp1;
   a=(dl*(1+s1)+f1^2/fp1)/c+f1/fp1;
   w1=1+s1+f1;
   b=(dl*w1*f1)/(fp1*c);
   tdel=2*b/(a+sqrt(a^2-4*b));
   @ convergence eigenvector @
   w1=abs(w1) < tau*sqrt(fp1+sp1);
   @ relative convergence @
   w2=abs(tdel) < tol*abs(t);
   t=t+tdel; @ t update @
   den=den-tdel; @ denominator update @
   d[j]=d[j]-tdel; @ normalized eigenvalue update @
   if it > 30;
      if w1 == 0;
         print "error:no convergence in qmdfpcr";
         end;
      else;
         print
```

```
                    "warning:too many iterations in qmdfpcr";
                endif;
                break;
            endif;
            it=it+1;
        endo;
        ddiff[j]=rho*t;
        d[j]=rho*d[j];  @ downdated eigenvalue @
        xd=z./den;
        ud[.,j]=udf*xd/sqrt(xd'xd);  @ scores downdate @
        k=k+1;

    else;  @ last eigenvalue @

        @ initialize t @
        t=z2[p]/(1+z2/(del[.,p]-1));
        if t <= 0 or t >= 1;
            print "error:initial values k=" k "tinit=" t;
            end;
        endif;

        den=del[.,k]-t;  @ initialize denominator @
        d[j]=d[j]/rho-t;  @ initialize normalized eigenvalue @
        w1=0;
        w2=0;
        it=1;
        do until w1 and w2;
            f1=g'(z2 ./den);
            fp1=g'(z2 ./den^2);
            w1=1+f1;
            tdel=w1*f1/fp1;
            @ convergence eigenvector @
            w1=abs(w1) < tau*sqrt(fp1);
            @ relative convergence @
            w2=abs(tdel) < tol*abs(t);
            t=t+tdel;  @ t update @
```

```
                    den=den-tdel; @ denominator update @
                    d[j]=d[j]-tdel; @ normalized eigenvalue update @
                    if it > 30;
                        if w1 == 0;
                            print "error:no convergence in qmdfpcr";
                            end;
                        else;
                            print
                            "warning:too many iterations in qmdfpcr";
                        endif;
                        break;
                    endif;
                    it=it+1;
                endo;
                ddiff[j]=rho*t;
                d[j]=rho*d[j]; @ downdated eigenvalue @
                xd=z./den;
                ud[.,j]=udf*xd/sqrt(xd'xd); @ scores downdate @

        endif;
    endif;

    if dfl == 0; @ no deflation @
        if d[j]<=tol;
            print "warning PRESS:d[" j "]=" d[j] "<= tol,
                    w may be too large";
        endif;
        pupdate=nu*ud[i,j]*(selif(ud[.,j],count./=i)+
                    ud[i,j]/(n-1))'yh/d[j];
        press[j+1]=press[j]+pupdate;
        ypred[j+1]=ypred[j]+pupdate;
    endif;

    j=j+1;
endo;
```

```
    if dfl == 1; @ sorting for deflation @
        ind=sortind(-d);
        ud=ud[.,ind];
        d=d[ind];
        ddiff=ddiff[ind];
        j=1;
        do while j <= w;
            if d[j]<=tol;
                print "warning PRESS:d[" j "]=" d[j] "<= tol,
                        w may be too large";
            endif;
            pupdate=nu*ud[i,j]*(selif(ud[.,j],count./=i)+
                        ud[i,j]/(n-1))'yh/d[j];
            press[j+1]=press[j]+pupdate;
            ypred[j+1]=ypred[j]+pupdate;
            j=j+1;
        endo;
    endif;

else;

    j=selif(seqa(1,1,rows(d)),modif .== 1);
    d[j]=d[j]-rho;
    ddiff[j]=rho;

    ind=sortind(-d);
    ud=ud[.,ind];
    d=d[ind];
    ddiff=ddiff[ind];
    j=1;
    do while j <= w;
        if d[j]<=tol;
            print "warning PRESS:d[" j "]=" d[j] "<= tol,
                    w may be too large";
        endif;
        pupdate=nu*ud[i,j]*(selif(ud[.,j],count./=i)+
```

```
                        ud[i,j]/(n-1))'yh/d[j];
          press[j+1]=press[j]+pupdate;
          ypred[j+1]=ypred[j]+pupdate;
          j=j+1;
       endo;
    endif;

    retp(press^2/n,d[1:w]',ddiff[1:w],ypred);
endp;
```

## The Linear Method

```
proc(6)=pcrcvl(y,x,w);  @ full leave-one-out cross-validation @
                        @ univariate principal component regression @
                        @ linear method @
                        @ [DeG90] @

    local n,p,msep,dd,ddiff,ypred,d,q;

    n=rows(x);p=cols(x);

    if p==1;print "error:p == 1";end;endif;

    if n-1 >= p;  @ regular case @
       if minc(p|n-2) < w;print "error:minc(p|n-2) < w";end;
       else;
       {msep,dd,ddiff,ypred,d,q}=pcrcvlr(y,x,w,n,p);
       endif;
    endif;
    if n-1 < p;  @ singular case @
       if n-2 < w;print "error:n-2 < w";end;
       else;
       {msep,dd,ddiff,ypred,d,q}=pcrcvls(y,x,w,n);
       endif;
    endif;

    @ return msep @
```

```
@          cross-validatory index @
@          downdated eigenvalues @
@          eigenvalue downdates @
@          cross-validated predicted values @
@          initial variances @
@          initial principal components @


    retp(msep~1-msep/msep[1],dd,ddiff,ypred,d/(n-1),q);
endp;


proc(6)=pcrcvlr(y,x,w,n,p);
@ leave-one-out cross-validation for pcr @
@ linear method @
@ regular case @


    local tol1,tol2,toles,pe,xm,u,d,q,nu,i,z,zz,e,
          rho,tau,urt,udf,press,pressu,modif,dd,ddd,ddff,ddiff,
          ypred,yp;


    tol1=10e-16; @ machine epsilon @
    tol2=10e-10; @ tolerance relative accuracy @


    press=zeros(1,w+1);
    ypred=zeros(n,w+1);
    ddiff=zeros(n,w);
    ddd=ddiff;
    modif=ones(p,1);


    @ overall principal component computations @


    xm=x-meanc(x)'; @ pca from proc EIG @
    {d,q}=eighv(xm'xm);d=rev(d);q=rev(q')';
    if minc(d) < 10*tol1;
        print "warning:overall eigenvalues <10*tol1";
    endif;
    u=xm*q; @ scores @ @ keep u d q @
    nu=n/(n-1); @ cross-validatory constant @
```

```
@ cross-validation @

i=1;
print "cross-validating";
do while i <= n; @ leave-one-out loop @
    z=u[i,.]';zz=z'z;z=z/sqrt(zz); @ pca scores normalized @
    rho=nu*zz; @ deflation size @
    toles=2*tol1*(d[1]+rho); @ tolerance deflation @
    tau=4*tol1*maxc(1|(d[1]-rho))/rho; @ tolerance convergence @

    if minc(d[1:(p-1)]-d[2:p]) < toles; @ equal eigenvalues @
        {urt,z}=dfleig(u,d,z,p,toles); @ eigenvalues @
        {udf,e,z,modif,pe}=dflsco(urt,d,z,p,toles); @ scores @
        {pressu,dd,ddff,yp}=
        lmdfpcr(y,i,urt,udf,d,e,z,rho,modif,pe,w,n,nu,tol2,tau);
        press=press+pressu;
        ddiff[i,.]=ddff;
        ddd[i,.]=dd;
        ypred[i,.]=yp;
        modif=ones(p,1);
    elseif minc(abs(z)) < toles; @ zero scores @
        {udf,e,z,modif,pe}=dflsco(u,d,z,p,toles);
        {pressu,dd,ddff,yp}=
        lmdfpcr(y,i,u,udf,d,e,z,rho,modif,pe,w,n,nu,tol2,tau);
        press=press+pressu;
        ddiff[i,.]=ddff;
        ddd[i,.]=dd;
        ypred[i,.]=yp;
        modif=ones(p,1);
    else;
        {pressu,dd,ddff,yp}=
        lmdfpcr(y,i,u,u,d,d,z,rho,modif,p,w,n,nu,tol2,tau);
        press=press+pressu;
        ddiff[i,.]=ddff;
        ddd[i,.]=dd;
```

```
        ypred[i,.]=yp;
      endif;

      i=i+1;
    endo;

  retp(press',ddd,ddiff,ypred,d,q);
endp;

proc(6)=pcrcvls(y,x,w,n);
@ leave-one-out cross-validation for pcr @
@ linear method @
@ singular case @

  local tol1,tol2,toles,xm,q,urt,udf,u,d,sqrtd,nu,i,z,zz,rho,tau,e,
    nn,press,pressu,nr,modif,dd,ddd,ddff,ddiff,yp,ypred;

  tol1=10e-16; @ machine epsilon @
  tol2=10e-10; @ tolerance relative accuracy @

  press=zeros(1,w+1);
  ypred=zeros(n,w+1);
  nu=n/(n-1);nr=n;n=n-1; @ cross-validatory constant @
  ddiff=zeros(nr,w);
  ddd=ddiff;
  modif=ones(n,1);

  @ overall principal component decomposition @

  xm=x-meanc(x)'; @ pca from proc EIG @
  {d,u}=eighv(xm*xm');d=rev(d[2:nr]);u=rev(u[.,2:nr]')';
  sqrtd=sqrt(d');
  q=xm'u./sqrtd;
  u=u.*sqrtd; @ scores @ @ keep q d u @
  if minc(d) < 10*tol1;
    print "warning :  overall eigenvalues <10*tol1";
  endif;
```

```
@ cross-validation @

i=1;
print "cross-validating";
do while i <= nr; @ leave-one-out loop @
   z=q'xm[i,.]';zz=z'z;z=z/sqrt(zz); @ pca scores normalized @
   rho=nu*zz; @ deflation size @
   toles=2*tol1*(d[1]+rho); @ tolerance deflation @
   tau=4*tol1*maxc(1|(d[1]-rho))/rho; @ tolerance convergence @

   if minc(d[1:(n-1)]-d[2:n]) < toles; @ equal eigenvalues @
      {urt,z}=dfleig(u,d,z,n,toles); @ eigenvalues @
      {udf,e,z,modif,nn}=dflsco(urt,d,z,n,toles); @ scores @
      {pressu,dd,ddff,yp}=
      lmdfpcr(y,i,urt,udf,d,e,z,rho,modif,nn,w,nr,nu,tol2,tau);
      press=press+pressu;
      ddiff[i,.]=ddff;
      ddd[i,.]=dd;
      ypred[i,.]=yp;
      modif=ones(n,1);
   elseif minc(abs(z)) < toles; @ zero scores @
      {udf,e,z,modif,nn}=dflsco(u,d,z,n,toles); @ scores @
      {pressu,dd,ddff,yp}=
      lmdfpcr(y,i,u,udf,d,e,z,rho,modif,nn,w,nr,nu,tol2,tau);
      press=press+pressu;
      ddiff[i,.]=ddff;
      ddd[i,.]=dd;
      ypred[i,.]=yp;
      modif=ones(n,1);
   else;
      {pressu,dd,ddff,yp}=
      lmdfpcr(y,i,u,u,d,d,z,rho,modif,n,w,nr,nu,tol2,tau);
      press=press+pressu;
      ddiff[i,.]=ddff;
      ddd[i,.]=dd;
```

```
        ypred[i,.]=yp;
      endif;

      i=i+1;
    endo;

    retp(press',ddd,ddiff,ypred,d,q);
endp;

proc(4)=lmdfpcr(y,i,ud,udf,d,ddf,z,rho,modif,p,w,n,nu,tol,tau);
@ rank-one modification with the linear method @

    local z2,del,tau,nr,g,j,k,e,t,gf,gs,den,h1,h2,f1,fp1,s1,sp1,a,xd,
          ud,w1,w2,tdel,it,count,press,yh,yhm,dfl,ind,ddiff,
          pupdate,ypred;

    z2=z^2;
    del=(ddf'-ddf)/rho;
    nr=seqa(1,1,p);
    g=ones(p,1);
    count=seqa(1,1,n);
    press=zeros(1,w+1);
    ypred=press;
    dfl=minc(modif)==0;
    ddiff=zeros(1,rows(d));
    yh=selif(y,count./=i);yhm=meanc(yh);yh=yh-yhm;
    press[1]=yhm-y[i];
    ypred[1]=yhm;

    if maxc(z) /= 1; @ downdate eigenvalues @

        j=1;
        k=1;
        do while j <= w; @ first w components @
           if modif[j] /= 0; @ eigenvalue changed @
              if k < p;
```

```
gf=nr.<=k;
gs=g-gf;

@ initialize t@
e=1+sumc(z2[1:k] ./del[1:k,(k+1)])+
    sumc(z2[(k+1):p]./del[(k+1):p,k]);
if e >= 0;
    t=z2[k]*(del[k+1,k]/(z2[k]+del[k+1,k]*e));
    if t <= 0 or t >= 1;
        print
        "error:initial values k=" k "tinit=" t "e>0";
        end;
    endif;
else;
    t=-del[k,k+1]/(z2[k+1]/(del[k,k+1]*e)+1);
    if t <= 0 or t >= 1;
        print
        "error:initial values k=" k "tinit=" t "e<0";
        end;
    endif;
endif;

den=del[.,k]-t; @ initialize denominator @
d[j]=d[j]/rho-t; @ initialize normalized eigenvalue @
w1=0;
w2=0;
it=1;
do until w1 and w2;
    h1=z2 ./den;
    h2=h1 ./den;
    f1=gf'h1;
    fp1=gf'h2;
    s1=gs'h1;
    sp1=gs'h2;
    a=1+s1;
    w1=a+f1;
```

```
tdel=(w1*a*f1)/(a^2 *fp1+f1^2 *sp1);
@ convergence eigenvector @
w1=abs(w1) < tau*sqrt(fp1+sp1);
@ relative convergence @
w2=abs(tdel) < tol*abs(t);
t=t+tdel; @ t update @
den=den-tdel; @ denominator update @
d[j]=d[j]-tdel; @ normalized eigenvalue update @
if it > 30;
    if w1 == 0;
        print "error:no convergence in lmdfpcr";
        end;
    else;
        print
        "warning:too many iterations in lmdfpcr";
    endif;
        break;
    endif;
    it=it+1;
endo;
ddiff[j]=rho*t;
d[j]=rho*d[j]; @ downdated eigenvalue @
xd=z./den;
ud[.,j]=udf*xd/sqrt(xd'xd); @ scores downdate @
k=k+1;

else; @ last eigenvalue @

    @ initialize t @
    t=z2[p]/(1+z2/(del[.,p]-1)+z2[p]);
    if t <= 0 or t >= 1;
        print "error:initial values k=" k "tinit=" t;
        end;
    endif;

    den=del[.,k]-t; @ initialize denominator @
```

```
d[j]=d[j]/rho-t; @ initialize normalized eigenvalue @
w1=0;
w2=0;
it=1;
do until w1 and w2;
    f1=g'(z2 ./den);
    fp1=g'(z2 ./den^2);
    w1=1+f1;
    tdel=w1*f1/fp1;
    @ convergence eigenvector @
    w1=abs(w1) < tau*sqrt(fp1);
    @ relative convergence @
    w2=abs(tdel) < tol*abs(t);
    t=t+tdel; @ t update @
    den=den-tdel; @ denominator update @
    d[j]=d[j]-tdel; @ normalized eigenvalue update @
    if it > 30;
        if w1 == 0;
            print "error:no convergence in lmdfpcr";
            end;
        else;
            print
            "warning:too many iterations in lmdfpcr";
        endif;
        break;
    endif;
    it=it+1;
endo;
ddiff[j]=rho*t;
d[j]=rho*d[j]; @ downdated eigenvalue @
xd=z./den;
ud[.,j]=udf*xd/sqrt(xd'xd); @ scores downdate @

    endif;
endif;
```

```
            if dfl == 0;
                if d[j]<=tol;
                    print "warning PRESS:d[" j "]=" d[j] "<= tol,
                            w may be too large";
                endif;
                pupdate=nu*ud[i,j]*(selif(ud[.,j],count./=i)+
                            ud[i,j]/(n-1))'yh/d[j];
                press[j+1]=press[j]+pupdate;
                ypred[j+1]=ypred[j]+pupdate;
            endif;

            j=j+1;
        endo;

        if dfl == 1; @ sorting for deflation @
            ind=sortind(-d);
            ud=ud[.,ind];
            d=d[ind];
            ddiff=ddiff[ind];
            j=1;
            do while j <= w;
                if d[j]<=tol;
                            print "warning PRESS:d[" j "]=" d[j] "<= tol,
                                    w may be too large";
                endif;
                pupdate=nu*ud[i,j]*(selif(ud[.,j],count./=i)+
                            ud[i,j]/(n-1))'yh/d[j];
                press[j+1]=press[j]+pupdate;
                ypred[j+1]=ypred[j]+pupdate;
                j=j+1;
            endo;
        endif;

    else;

        j=selif(seqa(1,1,rows(d)),modif .== 1);
```

```
        d[j]=d[j]-rho;
        ddiff[j]=rho;

        ind=sortind(-d);
        ud=ud[.,ind];
        d=d[ind];
        ddiff=ddiff[ind];
        j=1;
        do while j <= w;
            if d[j]<=tol;
                print "warning PRESS:d[" j "]=" d[j] "<= tol,
                        w may be too large";
            endif;
            pupdate=nu*ud[i,j]*(selif(ud[.,j],count./=i)+
                        ud[i,j]/(n-1))'yh/d[j];
            press[j+1]=press[j]+pupdate;
            ypred[j+1]=ypred[j]+pupdate;
            j=j+1;
        endo;
    endif;

    retp(press^2/n,d[1:w]',ddiff[1:w],ypred);
endp;
```

## Deflation

```
proc(2)=dfleig(q,d,x,p,tol);
@ deflation @
@ equal eigenvalues @

    local prm,check,dfl,j,tau,c,s,t1,t2;

    print "dfleig:equal eigenvalues";

    prm=seqa(1,1,p);
    check = (d[1:(p-1)]-d[2:p])|1 .< tol;
    dfl=1;
```

```
    do while dfl;
        j=minc(selif(prm,check));
        if abs(x[j]) > tol; @ jacobi rotations @
            if abs(x[j]) < abs(x[j+1]);
                tau=x[j]/x[j+1];c=1/sqrt(1+tau^2);s=c*tau;
            else;
                tau=x[j+1]/x[j];s=1/sqrt(1+tau^2);c=s*tau;
            endif;
            x[j+1]=s*x[j]+c*x[j+1];x[j]=0;
            t1=q[.,j];t2=q[.,j+1];
            q[.,j]=c*t1-s*t2;q[.,j+1]=s*t1+c*t2;
        endif;
        p=p-1;check[j]=0;
        if p == 1;
            dfl=0; @ stop deflation @
        else;
            dfl= maxc(check); @ check deflation @
        endif;
    endo;

    retp(q,x); @ return deflated problem @
endp;

proc(5)=dflsco(q,d,x,p,tol);
@ deflation @
@ zero scores @

    local dfl,j,ind,modif;

    print "dflsco:zero scores ";

    modif=abs(x) .>= tol;
    dfl=1;
    do while dfl;
        if abs(x[p]) < tol;
            p=p-1;d=d[1:p];x=x[1:p];q=q[.,1:p];
        else;
```

```
        j=minindc(abs(x)); @ deflate for jth component @
        q=selif(q',seqa(1,1,p)./=j)';
        d[j]=d[p];
        x[j]=x[p];
        p=p-1;
        d=d[1:p];
        x=x[1:p];
    endif;
    if p == 1;
        dfl=0; @ all eigenvectors unchanged - stop deflating @
    else;
        ind=rev(sortind(d));@ sort deflated results @
        d=d[ind];
        x=x[ind];
        dfl= minc(abs(x)) < tol; @ check deflation @
    endif;
    endo;

    retp(q,d,x,modif,p); @ return deflated problem @
endp;
```

# B.  Partial Least Squares Computations

## B.1  Partial Least Squares Decomposition

### Martens's Orthogonalized Algorithm

```
proc(2)=plsmart1(y,x,k); @ univariate partial least squares @
                         @ orthogonal scores algorithm @
                         @ [MN89, page 121] @


    local i,w,ww,t,tt,p,pp,q,qq,xm,ym; @ declarations @

    xm=x-meanc(x)'; @ initializations @
    ym=y-meanc(y);
    ww=zeros(cols(x),k);
```

```
      pp=ww;
      qq=zeros(k,1);

      @ pls computations @

      i=1;
      do while i <= k;
          w=xm'ym;
          w=w/sqrt(w'w);ww[.,i]=w;
          t=xm*w;
          tt=t't;
          p=xm't/tt;pp[.,i]=p;
          q=ym't/tt;qq[i]=q;
          xm=xm-t*p';
          ym=ym-t*q;
          i=i+1;
      endo;

      @ return partial least squares components and q @
      @ beta coefficients = ww*inv(pp'ww)*qq @

      retp(ww*inv(pp'ww),qq);
endp;
```

## Martens's Non-Orthogonalized Algorithm

```
proc(1)=plsmart2(y,x,k);   @ univariate partial least squares @
                           @ orthogonal loadings algorithm @
                           @ [MN89, page 123] @


      local i,w,ww,t,tt,qq,xm,y0,ym;  @ declarations @

      xm=x-meanc(x)';  @ initializations @
      y0=y-meanc(y);
      ym=y0;
      ww=zeros(cols(x),k);
```

```
    tt=zeros(rows(x),k);

    @ pls computations @

    i=1;
    do while i <= k;
        w=xm'ym;
        w=w/sqrt(w'w);ww[.,i]=w;
        t=xm*w;tt[.,i]=t;
        qq=olsqr(y0,tt[.,1:i]);
        xm=xm-t*w';
        ym=y0-tt[.,1:i]*qq;
        i=i+1;
    endo;

    @ return beta coefficients @

    retp(ww*qq);
endp;
```

## Helland's Algorithm

```
proc(1)=plshell(y,x,k);  @ univariate partial least squares @
                         @ Helland @
                         @ [Hel88] @


    local s,w,ww,xw,b,i,xm,ym;  @ declarations @

    xm=x-meanc(x)';  @ initializations @
    ym=y-meanc(y);
    ww=zeros(cols(x),k);
    xw=zeros(rows(x),k);

    @ pls computations @

    s=xm'ym;  @ first pls component @
```

```
w=s;
w=w/sqrt(w'w);
ww[.,1]=w;
xw[.,1]=xm*w;
b=w*olsqr(ym,xw[.,1]);

i=2; @ further pls components @
do while i <= k;
    w=s-xm'(xm*b);
    w=w/sqrt(w'w);
    ww[.,i]=w;
    xw[.,i]=xm*w;
    b=ww[.,1:i]*olsqr(ym,xw[.,1:i]);
    i=i+1;
endo;

@ return beta coefficients @

retp(b);
endp;
```

## Stone-Brooks's Approach

```
proc(5)=pls(y,x,w); @ univariate partial least squares @
                    @ Stone-Brooks method @
                    @ [SB90] @

    @ y is a column vector that contains the observed responses.@
    @ x is a matrix that contains the predictor variables.@
    @ w is the number of partial least squares factors derived.@

    local c,beta,scores,pred,resid,be,re,pr,s,e,u,v,d,m,k,a,ma,z,
          meanx,meany,xm,ym,n,p; @ declarations @

    n=rows(x);p=cols(x); @ initializations @
    c=zeros(p,w);beta=c;
    scores=zeros(n,w);pred=scores;resid=scores;
```

```
if p==1;print "error :  p == 1";end;endif;

meanx=meanc(x)';meany=meanc(y); @ mean centering @
xm=x-meanx;ym=y-meany;
s=xm'ym; @ covariance @

@ principal component decomposition @

if n-1 >= p; @ regular case @
   if p < w;print "error :  p < w";end;
   else;
       {e,v}=eighv(xm'xm);e=rev(e); @ pca from proc EIG @
       v=rev(v')'; @ principal components @
       m=eye(p); @ initialize m @
   endif;
endif;
if n-1 < p; @ singular case @
   if n-1 < w;print "error :  n-2 < w";end;
   else;
       {e,u}=eighv(xm*xm');e=rev(e[2:n]); @ pca from proc EIG @
       v=xm'rev(u[.,2:n]')'./sqrt(e'); @ principal components @
       m=eye(n-1); @ initialize m @
   endif;
endif;

@ pls computations @

d=v's;
c[.,1]=s/sqrt(s's); @ first pls component @
scores[.,1]=xm*c[.,1];
{be,re,pr}=olsqr2(ym,scores[.,1]);
beta[.,1]=c[.,1]*be;
pred[.,1]=pr+meany;resid[.,1]=re;
k=2;
do while k <= w; @ further pls components @
```

```
      a=e.*v'c[.,(k-1)];
      ma=m*a;
      m=m-ma*ma'/a'ma;
      z=m*d;
      c[.,k]=v*z/sqrt(z'z);
      scores[.,k]=xm*c[.,k];
      {be,re,pr}=olsqr2(ym,scores[.,1:k]);
      beta[.,k]=c[.,1:k]*be;
      pred[.,k]=pr+meany;resid[.,k]=re;
      k=k+1;
   endo;

   @ return predicted values @
   @         residuals @
   @         beta coefficients @
   @         pls components @
   @         pls scores @

   retp(pred,resid,beta,c,scores);
endp;
```

## B.2  Partial Least Squares Cross-validation

```
proc(4)=plscv(y,x,w);  @ full leave-one-out cross-validation @
                       @ univariate partial least squares @
                       @ Stone-Brooks method @
                       @ [SB90] @

   @ y,x and w are as in the partial least squares algorithm.@

   local i,k,s,e,v,u,sqrte,d,fi,di,mi,ai,miai,zi,yh,yhm,fz,
         msep,pressu,hpress,pupdate,dd,ddi,ddii,ypred,xm,ym,
         nr,nu,n,p;

   n=rows(x);p=cols(x);  @ initializations @
   nu=n/(n-1);
   nr=seqa(1,1,n);
```

```
msep=zeros(w+1,1);
pressu=msep;
dd=zeros(n,w);
ddi=zeros(n,1);
ypred=zeros(n,w+1);

if p==1;print "error :  p == 1";end;endif;

xm=x-meanc(x)';ym=y-meanc(y); @ mean centering @
s=xm'ym; @ covariance @

@ principal component decomposition @

if n-1 >= p; @ regular case @
   if p < w;print "error :  p < w";end;
   else;
       {e,v}=eighv(xm'xm); @ pca from proc EIG @
       e=rev(e);v=rev(v')';u=xm*v;
   endif;
endif;
if n-1 < p; @ singular case @
   if n-2 < w;print "error :  n-2 < w";end;
   else;
       {e,u}=eighv(xm*xm'); @ pca from proc EIG @
       e=rev(e[2:n]);u=rev(u[.,2:n]')';sqrte=sqrt(e');
       v=xm'u./sqrte;u=u.*sqrte;
   endif;
endif;

@ cross-validatory computations @

d=v's;

i=1;
do while i <= n;
   yh=selif(y,nr./=i);yhm=meanc(yh);yh=yh-yhm;
```

```
    pressu[1]=yhm-y[i];
    ypred[i,1]=yhm;

    k=1;  @ first pls component @
    mi=eye(minc((n-1)|p));  @ initialize mi @
    fi=u[i,.]';
    ddii=nu*ym[i]*fi;
    di=d-ddii;ddi[i]=sqrt(ddii'ddii);
    zi=di/sqrt(di'di);
    fz=u*zi;
    hpress=selif(fz,nr./=i)+fz[i,.]/(n-1);
    dd[i,1]=hpress'hpress;
    pupdate=nu*fz[i,.]*hpress'yh/dd[i,1];
    pressu[2]=pressu[1]+pupdate;
    ypred[i,2]=ypred[i,1]+pupdate;
    do while k < w;  @ further pls components @
        k=k+1;
        ai=e.*zi-nu*sumc(fi.*zi)*fi;
        miai=mi*ai;
        mi=mi-miai*miai'/ai'miai;
        zi=mi*di;
        zi=zi/sqrt(zi'zi);
        fz=u*zi;
        hpress=selif(fz,nr./=i)+fz[i,.]/(n-1);
        dd[i,k]=hpress'hpress;
        pupdate=nu*fz[i,.]*hpress'yh/dd[i,k];
        pressu[k+1]=pressu[k]+pupdate;
        ypred[i,k+1]=ypred[i,k]+pupdate;
        endo;
    msep=msep+pressu^2/n;  @ update msep @
    i=i+1;
    endo;

@ return msep @
@         dd @
@         ddi @
```

```
    @         ypred @

    retp(msep);
endp;
```

# Bibliography

[And58]   T. W. Anderson. *An Introduction to Multivariate Statistical Analysis*. Wiley, New York, 1958.

[Apt92]   Aptech Systems, Inc., 23804 S.E. Kent-Kangley Road, Maple Valley, Washington 98038. *The GAUSS System, Version 3.0*, Documentation Version 061292 edition, 1992.

[BN78]    J. R. Bunch and C. P. Nielsen. Updating the singular value decomposition. *Numer. Math.*, 31:111–129, 1978.

[BNS78]   J. R. Bunch, C. P. Nielsen, and D. Sorensen. Rank-one modification of the symmetric eigenproblem. *Numer. Math.*, 31:31–48, 1978.

[BR74]    J. R. Bunch and D. J. Rose. Partitioning, tearing and modification of sparse linear systems. *J. Math. Anal. Appl.*, 48:574–593, 1974.

[Bro93]   P. J. Brown. *Measurement, Regression, and Calibration*. Oxford University Press, New York, 1993.

[Cal84]   P. Calder. Influence functions for principal components and their variances. Unpublished, November 1984.

[Cri85]   F. Critchley. Influence in principal components analysis. *Biometrika*, 72(3):627–636, 1985.

[Cup81]   J. J. M. Cuppen. A divide and conquer method for the symmetric tridiagonal eigenproblem. *Numer. Math.*, 36:177–195, 1981.

[DBMS78]  J. J. Dongarra, J. R. Bunch, C. B. Moler, and G. W. Stewart. *LINPACK User's Guide.* Philadelphia, 1978.

[DeG90]  R. D. DeGroat. Efficient, numerically stabilized rank-one eigenstructure updating. *IEEE Transactions on Acoustics, Speech and Signal Processing.*, 38(2):301–316, 1990.

[Den91]  Michael Charles Denham. *Calibration in Infrared Spectroscopy.* PhD thesis, University of Liverpool, January 1991.

[DG87]  J. J. Dongarra and E. Grosse. Distribution of electronic software via electronic mail. *Comm. ACM*, 30:403–407, 1987.

[DS87]  J. J. Dongarra and D. C. Sorensen. A fully parallel algorithm for the symmetric eigenvalue problem. *SIAM J. Sci. Stat. Comput.*, 8(2):139–154, 1987.

[EJH91]  D. G. Evans, K. Jewell, and M. N. Hall. Mathematical interpretation of complex data sets: 1) NIR analysis of mustard seed oil content predicted by regression methods. Technical Report 616, Campden Food and Drink Research Association, Chipping Campden, Gloucestershire, GL55 6LD, UK, May 1991. MAFF project N0. 8495.

[EK82]  H. T. Eastment and W. J. Krzanowski. Crossvalidatory choice of the number of components from a principal component analysis. *Technometrics*, 24(1):73–77, 1982.

[Fea83]  T. Fearn. A misuse of ridge regression in the calibration of a near infrared reflectance instrument. *Applied Statistics*, 32(1):73–79, 1983.

[GBDM72] B. S. Garbow, J. M. Boyle, J. J. Dongarra, and C. B. Moler. *Matrix Eigensystems Routines: EISPACK Guide Extension.* New York, 1972.

[GK65]  G. H. Golub and W. Kahan. Calculating the singular values and pseudo-inverse of a matrix. *SIAM J. Num. Anal.*, 2:205–224, 1965.

[GK72]   R. Gnanadesikan and J. R. Kettenring. Robust estimates, resid-
         uals, and outlier detection with multiple response data. *Biomet-
         rics*, 28:81–124, 1972.

[GL89]   G. H. Golub and C. F. Van Loan. *Matrix Computations (second
         edition)*. John Hopkins University Press, London, 1989.

[Gol73]  G. H. Golub. Some modified matrix eigenvalue problems. *SIAM
         Review*, 15(2):318–334, 1973.

[GR70]   G. H. Golub and C. Reinsch. Singular value decomposition and
         least squares solutions. *Numer. Math.*, 14:403–420, 1970.

[Hel88]  I. S. Helland. On the structure of partial least squares regression.
         *Commun. Statist. Simul.*, 17(2):581–607, 1988.

[Hös88]  A. Höskuldsson. PLS regression methods. *Journal of Chemo-
         metrics*, 2:211–228, 1988.

[Hot33]  H. Hotelling. Analysis of a complex of statistical variables into
         principal components. *J. Educ. Psychol.*, 24:417–441, 1933.

[Hou58]  A. S. Householder. Unitary triangularization of a nonsymmetric
         matrix. *J. Ass. Comp. Mach.*, 5:339–342, 1958.

[Jol82]  I. T. Jolliffe. A note on the use of principal components in re-
         gression. *Applied Statistics*, 31(3):300–302, 1982.

[Jol86]  I. T. Jolliffe. *Principal Component Analysis*. Springer-Verlag,
         New York, 1986.

[Krz83]  W. J. Krzanowski. Crossvalidatory choice in principal compo-
         nent analysis: Some sampling results. *J. Statist Comput. Simul.*,
         18:299–314, 1983.

[Krz87]  W. J. Krzanowski. Crossvalidation in principal component anal-
         ysis. *Biometrics*, 43:575–584, 1987.

[Krz88]  W. J. Krzanowski. *Principles of Multivariate Analysis : a user's
         perspective*. Oxford University Press, New York, 1988.

[Mat92]   Matlab. *High Performance Numeric Computation and Visualization Software.* The Mathworks, Inc., 24 Prime Park Way, Natik, MA 01760, 1992.

[MN89]   H. Martens and T. Naes. *Multivariate Calibration.* Wiley, New York, 1989.

[MT68]   F. Mosteller and J. W. Tukey. Data analysis, including statistics. In G. Lindzey and E. Aronson, editors, *Handbook of social psychology.* Addison-Wesley, Reading, Mass., 1968.

[OFMD84] B. G. Osborne, T. Fearn, A. R. Miller, and S. Douglas. Application of near infrared reflectance spectroscopy to the compositional analysis of biscuits and biscuit doughs. *J.Sci.Food Agric.*, 35:99–105, 1984.

[Pea01]   K. Pearson. On lines and planes of closest fit to systems of points in space. *Phil. Mag. (6)*, 2:559–572, 1901.

[PJM88]   P. Pack, I. T. Jolliffe, and B. J. T. Morgan. Influential observations in principal component analysis: a case study. *Journal of Applied Statistics*, 15(1):39–52, 1988.

[PW79]   G. Peters and J. H. Wilkinson. Inverse iteration, ill-conditioned equations and newton's method. *SIAM Review*, 21(3):339–360, 1979.

[Rao64]   C. R. Rao. The use and interpretation of principal components in applied research. *Sankhya A*, 26:329–358, 1964.

[Rao73]   C. R. Rao. *Linear Statistical Inference and Its Applications (second edition).* Wiley, New York, 1973.

[RK81]   R. Radhakrishnan and A. M. Kshirsagar. Influence functions for certain parameters in multivariate analysis. *Commun. Statist.*, A10(6):515–529, 1981.

[SB90]   M. Stone and R. J. Brooks. Continuum regression: Cross-validated sequentially constructed prediction embracing ordinary

least squares, partial least squares and principal component regression. *Journal of the Royal Statistical Society B*, 52(2):237–269, 1990. (with discussion).

[SB92] M. Stone and R. J. Brooks. Corrigendum: Continuum regression: Cross-validated sequentially constructed prediction embracing ordinary least squares, partial least squares and principal component regression. *Journal of the Royal Statistical Society B*, 54(3):906–907, 1992.

[SBI+70] B. T. Smith, J. M. Boyle, Y. Ikebe, V. C. Klema, and C. B. Moler. *Matrix Eigensystem Routines: EISPACK Guide, 2nd ed.* New York, 1970.

[Seb84] G. A. F. Seber. *Multivariate Observations*. Wiley, New York, 1984.

[Sta92] Statistical Sciences, Inc., Statsci Europe, Statistical Sciences U.K. Ltd., 52 Sandfield Road, Headington, Oxford OX3 7RJ. *S-PLUS for DOS, Version 2.0*, User's Manual, Vol. 1., Version 2.0 edition, January 1992.

[Sto74] M. Stone. Cross-validatory choice and assessment of statistical predictions. *Journal of the Royal Statistical Society B*, 36(2):111–147, 1974. (with discussion).

[Sun93] R. Sundberg. Continuum regression and ridge regression. *Journal of the Royal Statistical Society B*, 54:653–659, 1993.

[Tho76] R. C. Thompson. The behaviour of eigenvalues and singular values under perturbations of restricted rank. *Linear Algebra and Applications*, 13:69–78, 1976.

[Wil65] J. H. Wilkinson. *The Algebraic Eigenvalue Problem*. Oxford University Press, New York, 1965.

[WL69] H. Wold and E. Lyttkens. Nonlinear iterative partial least squares (NIPALS) estimation procedures. *Bull. Internat. Stat. Inst.*, 43:29–51, 1969.

[Wol66]    H. Wold. Nonlinear estimation by iterative least squares procedures. In F. N. David, editor, *Research papers in statistics, Festschrift for J. Neyman.*, pages 411–444. Wiley, New York, 1966.

[Wol76]    S. Wold. Pattern recognition by means of disjoint principal components models. *Pattern Recognition*, 8:127–139, 1976.

[Wol78]    S. Wold. Cross-validatory estimation of the number of components in factor and principal components models. *Technometrics*, 20(4):397–405, 1978.

[Wol85]    H. Wold. Partial least squares. In S. Kotz and N. Johnson, editors, *Encyclopaedia of Statistical Sciences.*, pages 581–591. Wiley, New York, 1985.

[WR71]    J. H. Wilkinson and C. Reinsch, editors. *Handbook for Automatic Computation.*, volume 2, Linear Algebra. Springer-Verlag, New York, 1971.

# Publication

# Principal Component Outlier Detection and SIMCA: A Synthesis

**Bart Mertens and Michael Thompson**

*Department of Chemistry, Birkbeck College, Gordon House, 29 Gordon Square, London, UK WC1H 0PP*

**Tom Fearn**

*Department of Statistical Science, University College London, Gower Street, London, UK WC1E 6BT*

Principal component outlier detection methods are discussed and their application in the soft independent modelling of class analogy (SIMCA) method of pattern recognition is clarified. SIMCA is compared to allocation procedures based on the Mahalanobis distance. Finally, the differences between the SIMCA method and quadratic discriminant analysis are discussed. The discussion is illustrated with an example from spectroscopy.

**Keywords:** *Outliers; principal component; soft independent modelling of class analogy; pattern recognition; spectroscopy*

## Multivariate Outliers

Many applications of chemometrics are concerned with the analysis of multiple measurements on a sample of objects. Typical examples include the problems of pattern recognition and calibration in spectroscopy. In any such application, there is usually a multitude of ways in which erroneous observations can arise. Such aberrant or outlying observations can invalidate the results of an analysis. Hence there is a need for outlier detection methods.[1,2]

Outlier detection poses a particular problem in multivariate statistics. A multivariate measurement may be outlying because the scale of the measurement is inconsistent or because it has an orientation in multidimensional space that is different from that of other observations.[3,4] This contrasts with univariate statistics, where we can readily rank the observations according to their degree of extremity. The problem is exacerbated by the fact that the graphical representation of multivariate data is cumbersome when the number of variates exceeds three.[4,5]

Consider an application of near-infrared spectroscopy in which 13 samples of whole rice were examined and the results stored as transmittance ($T$) measurements at 100 wavelengths, ranging from 850 to 1048 nm at equal intervals (Fig. 1). Let us examine two different ways in which an outlying sample measurement may arise. The transmittance of a sample may be different across the whole measurement range, *i.e.*, by a vertical shift affecting the whole spectrum almost to the same extent, so that the spectrum is clearly separated from other spectra. Alternatively, a spectrum may be affected differently at various parts of the spectrum.

In the first instance, we can usually spot the outliers quite easily, either on a plot of the spectra or from the transmittance at a set of well chosen wavelengths. Such differences in transmittance will often relate to particular physical features in the analysed sample, for instance, moisture content or grain shape and size. The detection of other types of outliers will depend on a careful analysis of the perturbations in the spectrum that are not associated with any major difference in absorption.

For the rice data (Fig. 1), all samples have similar spectra, except for the French Camargue sample, that has a different transmittance pattern. There are three spectra with a lower over-all absorption, while their transmittance pattern seems to be consistent with that from the remaining spectra. These three observations are samples of glutinous rice. The other observations are samples of Basmati rice, except for the odd observation, which is a sample of rice from the French Camargue. Glutinous rice can be distinguished visually from the other samples as its particle size is very different.

The issues illustrated by this example are common to all problems of outlier detection in multivariate statistics. The samples of glutinous rice are outlying because their scale of measurement is inconsistent with that from other samples. Such samples can usually be detected by using methods that are essentially univariate. Samples that are aberrant because of other reasons have a different orientation in multivariate space as compared to other samples. Their outlying position cannot be explained as a mere shift along the general trend or orientation of the data.

With spectral data, we could study scatter plots of the observations for two or three well chosen wavelengths. However, the selection of wavelengths may pose a problem. Furthermore, such a procedure restricts us to work with only three wavelengths at maximum, whereas an anomalous pattern may be owing to a small but consistent aberration across the range of wavelengths. Univariate measures or simple scatter plot techniques will fail to detect such outliers.

As spectra have a natural two dimensional representation, we could study a plot of the raw spectra. However, when the number of features in a spectrum is large, visual comparison will become exceedingly difficult. Hence the need for a statistical evaluation of outliers in spectral data that can accumulate information across all wavelengths.
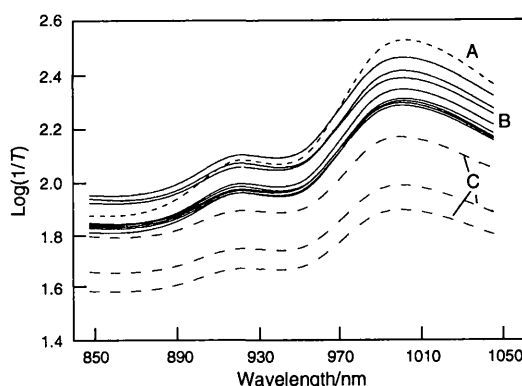


**Fig. 1** Near-infrared transmittance spectra for 13 samples of rice: A, French Camargue; B (all solid lines), Basmati; and C, glutinous.

## Principal Component Analysis and Outliers

Most approaches to outlier detection in multivariate analysis involve some form of dimension reduction. Ideally, we would want to reduce the data to a single univariate statistic that is sensitive to outliers. Alternatively, we may try to reduce the data to a set of uncorrelated univariate statistics which we then analyse independently of one another. Principal component analysis is such a method of dimension reduction. Applications in spectroscopy have been discussed by Cowe and McNichol[6] as well as Cowe, McNichol and Cuthbertson[7,8] and Bertrand, Robert and Tran,[9] for example. We will briefly review the properties of principal component decomposition first and then the applications to outlier detection.

Consider a sample of $n$ observations of $p$ variates $x_i = (x_{i1},..., x_{ip})$, $i = 1,..., n$. We will write

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

To simplify the notation, we will assume that the data are mean centered and that $n - 1 > p$. The generalization to the singular case $(n - 1 < p)$ is straightforward.

Principal components are derived sequentially, as the orthonormal set of vectors $q_j = (q_{1j},..., q_{pj})^T$, $j = 1,..., p$, for which each linear combination $u_j = Xq_j$ has maximum variance among all linear combinations $Xq$ with $q$ orthogonal to the first $j - 1$ principal components.[5,10] This definition was described by Hotelling.[11] The principal component decomposition can be written as

$$X = UQ^T$$

where $Q = (q_1,...,q_p) = [(q_{ij})]$ is the matrix of principal component loadings, with $QQ^T = Q^TQ = I$, and $U = (u_1,..., u_p) = [(u_{ij})]$ is the matrix of principal component scores.

For any $r < p$, the first $r$ principal components give the best representation of the data in $r$ dimensions, in the sense that the total amount of variation explained by these components is maximum in the set of all $r$ dimensional representations.[10] From a purely geometric point of view, the principal components are the set of orthogonal directions in multidimensional space for which the sum of the squared distances of the observations from their projections on the subspace defined by the first $r$ principal components is minimized for each $r < p$.[10,12] This property was originally used by Pearson in 1901 to define the principal component decomposition.[13] This provides the basis for using principal component analysis as a dimension reduction tool.

Principal component decomposition is not scale invariant, and hence we must decide prior to any analysis whether we want to scale the data. In applications to spectroscopy one will often prefer to analyse the original data as the variability is comparable across the different wavelengths. Indeed, in applications to classification or outlier rejection, scaling tends to amplify noise at wavelengths with small between-group variability. Similarly, scatter correction may remove some of the between-group variability.

The adequacy of the principal component representation can be assessed from the variances $\lambda_j$, $j = 1,..., p$ of the principal component scores $u_j$, $j = 1,..., p$. We have

$$\lambda_j = \frac{\sum_{i=1}^{n} u_{ij}^2}{n - 1}$$

$$= u_j^T u_j/(n - 1)$$

$$= q_j^T X^T X q_j/(n - 1)$$

$$= q_j^T S q_j$$

where $S = X^T X/(n - 1)$ is the sample covariance matrix. It

follows that the principal component decomposition may also be written as

$$S = Q\Lambda Q^T$$

with $\Lambda$ a diagonal matrix with the variances $\lambda_1 > ... > \lambda_p$ on the diagonal. Rao[12] proposed the criterion $\lambda_{r+1} + ... + \lambda_p$ to assess the $r$ dimensional representation. We have

$$\lambda_{r+1} + ... + \lambda_p = \frac{\sum_{i=1}^{n} \sum_{j=r+1}^{p} u_{ij}^2}{n - 1}$$

which is proportional to the sum of the squared distances of the observations from their projections on the fitted principal component subspace. We may refer to this criterion as the absolute lack-of-fit of the principal component representation, as it equals the total amount of variation that is left unexplained by the principal component representation. The total amount of variation in the measured data is $\lambda_1 + ... + \lambda_p$, and hence, we may use the equivalent lack-of-fit ratio

$$\frac{\lambda_{r+1} + ... + \lambda_p}{\lambda_1 + ... + \lambda_p}$$

or alternatively, the goodness-of-fit ratio

$$\frac{\lambda_1 + ... + \lambda_r}{\lambda_1 + ... + \lambda_p}$$

If the first $k < p$ principal components provide an adequate representation of the data, then the remaining principal components are often interpreted as representing the near-constant linear relations in the data. The variation of the corresponding principal component scores will be small.

Consider the rice data. The first principal component accounts for 98.9% of the variation in the data. This variation is caused by differences in packing density in the sample cell. As can be seen from Fig. 2, the weights on this component are virtually constant, as compared to the second component. Hence, this component accounts for the vertical separation of the spectra as it represents the average transmittance spectrum of the rice samples. As a consequence, the first component separates the samples of glutinous rice, as shown in the principal component scores plot in Fig. 3. This is because their spectra are characterized by a substantial and consistent difference in absorption, across all wavelengths. As can be seen on this plot, there is a further separation within the group of Basmati rice. This separation is owing to three samples that have a smaller transmittance compared with the remaining samples. This difference is masked somewhat on the plot of the raw spectra, as the spectrum of rice from the French Camargue is overlaying. In comparison, the second component contrasts the responses at the two ends of the wavelength spectrum. This may be compared with the appearance of the spectrum of the sample from the French Camargue on Fig. 1 and the position of this sample in the principal component scores plot. The mean transmittance of
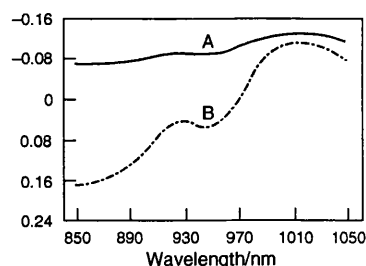


Fig. 2 Principal component loadings for the first two principal components for the rice data: A, first and B, second component.

this sample compares to that from Basmati rice, and hence this sample is not outlying on the first principal component. However, this sample is extreme on the second component, because of an increase in transmission at the higher wavelengths, compared with a decrease in response at the lower end of the spectrum.

This example illustrates some of the general principles of principal component outlier detection.

Observations that are extreme on the larger variance components inflate the variances and covariances of the original variables.[3,4,10] Such outliers can often be detected with univariate methods, as they are usually caused by a large perturbation in one or a few of the variates. Our example is extreme in this respect, as the first principal component causes a vertical shift of the spectra, and hence, almost any wavelength could be used to test for such outliers. Observations that are outlying on the smaller variance components are typically incompatible with the orientation of the observed data and may obscure the linear relations in the data.[3,10] We may not always be able to detect these observations as outliers on individual variates or even from scatter plots of the measured data. They arise from small but consistent perturbations in the measured variates. As a consequence, the smaller variance principal components are especially useful in multivariate outlier analysis. In spectroscopy, the smaller variance components are essential, as they will detect small pattern aberrations that could be obscured by large variation trends in the data. These properties of smaller variance components have been exploited in applications to spectroscopy by several authors (*e.g.*, Cowe and McNichol[6]).

### Principal Component Outlier Detection and SIMCA

The principal component scores may be used individually in univariate outlier detection tests,[14] as well as in scatter plots,[3,15-17] for the detection of outliers. This implicitly uses the property that the principal component scores are uncorrelated. A natural progression of these ideas is to combine the information from several principal components in a single statistic.

Suppose that the first $k$ components give an adequate representation of the data. To assess whether an observation $x_i$ is an outlier, we may consider the contribution of the observation to the absolute lack-of-fit statistic. Hence, we define the statistic

$$D^2_{1(k)}(x_i) = \sum_{j=k+1}^{p} u^2_{ij}$$

which is the squared length of the perpendicular from the observation to the principal component subspace defined by the first $k$ components. This statistic was originally proposed in 1964 for outlier detection by Rao in his extensive paper on applied principal component analysis.[12] The statistic reflects
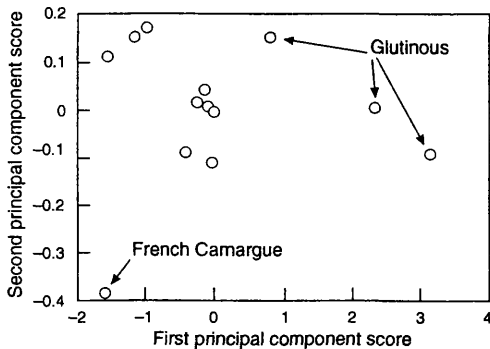
only the small variance perturbations in the data and ignores the major sources of variation. Hence, the purpose of this statistic is to accumulate information on the small variation aberrations. Gnanadesikan and Kettenring[3] used this statistic in a gamma probability plot for outlier detection. The statistic has also been used in quality control problems by Hawkins,[15] Hawkins and Fatti[16] and Jackson.[18,19]

Consider the one-dimensional principal component representation of the rice data. The absolute lack-of-fit of this representation,

$$\lambda_2 + \ldots + \lambda_p = \frac{\sum_{i=1}^{n} \sum_{j=2}^{p} u^2_{ij}}{n-1} = \frac{\sum_{i=1}^{n} D^2_{1(1)}(x_i)}{n-1}$$

is 0.023 and, hence, the outlier statistic $D^2_{1(1)}(x_i)$ for an individual observation should be of the same magnitude, if the principal component representation fits the data well and if the observation is not an outlier.

Fig. 4 contains a histogram of statistic $D^2_{1(1)}$. The average value is 0.021. All values are smaller than 0.03, except for the sample from the French Camargue, which has a value of 0.149 on this statistic. The outlier statistic detects this sample, as it is outlying on the second component. The glutinous samples, however, can not be detected with this statistic, as it ignores all information from the first component. The values for the glutinous samples were 0.024, 0.001 and 0.0001.

This procedure has been applied by Wold to problems of pattern recognition.[20] Wold studied the problem of comparing a new observation $x$ with a set of observations $X$ from a single homogeneous group. It is implicit in this notation that $x$ has been mean centered for the mean of the group with which it is compared. Consider a $k$-dimensional principal component representation. The absolute lack-of-fit of the representation is

$$\frac{\sum_{i=1}^{n} \sum_{j=k+1}^{p} u^2_{ij}}{n-1} = \frac{\sum_{i=1}^{n} D^2_{1(k)}(x_i)}{n-1}$$

To compare the observation $x$ with the group, we calculate the distance of the observation from the principal component representation. Hence, we compute Rao's statistic

$$D^2_{1(k)}(x) = \sum_{j=k+1}^{p} v \cdot r_j^2$$

where $v_j = xq_j$, $j = 1,\ldots, p$. We compare the value of the observation on this statistic with the spread of the values of the observations within the group. We reject the observation if this value is well separated.

This procedure is referred to by Wold as the SIMCA method of pattern recognition, where SIMCA stands for 'Soft Independent Modelling of Class Analogy'. Wold's derivation of the SIMCA procedure differs from the approach outlined
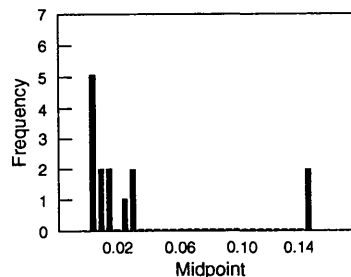


**Fig. 3** Plot of the rice data *versus* the first two principal components.



**Fig. 4** Histogram of the values of Rao's statistic for the rice data.

above. Wold writes the $k$-dimensional principal component representation as

$$X = U_{(k)} Q_{(k)}^T + \epsilon,$$

with $\epsilon$ a matrix of residuals from the $k$-dimensional representation, $U_{(k)} = (u_1,..., u_k)$ and $Q_{(k)} = (q_1,..., q_k)$. The principal component representation defines a $k$-dimensional subspace that fits best to the mean-centered data in a least squares sense.[21] The absolute lack-of-fit of this representation may be defined as

$$\frac{\sum_{i=1}^{n} \sum_{j=1}^{p} \epsilon_{ij}^2}{n - 1}$$

We can prove that

$$\frac{\sum_{i=1}^{n} \sum_{j=1}^{p} \epsilon_{ij}^2}{n - 1} = \frac{\sum_{i=1}^{n} \sum_{j=k+1}^{p} u_{ij}^2}{n - 1}$$

which shows that we obtain Rao's lack-of-fit test in this derivation.

To compare the observation $x$ with the group, we fit the observation to the principal component representation and calculate the residual:

$$x = \beta Q_{(k)}^T + \eta$$

where $\beta = (\beta_1,..., \beta_k)$ may be interpreted as a vector of regression coefficients and $\eta$ is the vector of residuals: Rao's statistic may now be calculated as

$$D_{1(k)}^2 (x) = \eta\eta^T$$

This approach is reminiscent of analysis of variance in regression. Indeed, in analogy to least squares residuals, Gnanadesikan and Kettenring refer to the orthogonal distance $D_{1(k)}(x)$ of an observation from the principal component representation as the orthogonal residual. The corresponding principal component scores $v_j, j = k + 1,..., p$ are referred to as the principal component residuals[3]. This analogy is not too far fetched, as we have seen that the principal component representation may be interpreted in terms of a least squares fit. Hence, we will introduce some new terminology by referring to the principal components $q_j, j = 1,..., k$ as the model or primary components and similarly, we may refer to the complementary set of components as the residual components.

This analogy between least squares residuals and principal component residuals is quite fundamental in the paper by Gnanadesikan and Kettenring. It is also exploited by Wold, who interprets the principal component representation as a model of the data. The residuals from this model are used to detect aberrant samples, in a similar way as residuals are analysed for the presence of outliers in analysis of variance.

Wold's implementation of SIMCA is essentially an application of principal component outlier analysis to classification. A principal component outlier statistic is defined and the extremity of an observation with respect to a particular group is evaluated with this statistic. Wold proposed the statistic

$$D_{W(k)}^2 (x) = \frac{D_{1(k)}^2 (x)/(p - k)}{\sum_{i=1}^{n} D_{1(k)}^2 (x_i)/((p - k)(n - k - 1))}$$

to assess a new observation $x$. The numerator in this expression is interpreted as the variance of the deviations $\eta$ of the observation $x$ from the principal component model and the denominator as the variance of the deviations $\epsilon$. This statistic

is compared to some critical value of an $F$-distribution. It is important to note that this formula is for the regular case. In the singular case, $p$ should be replaced by $n - 1$, which is the number of principal component dimensions that can be estimated from the data.

If several groups are present in the data, then a separate principal component representation is derived for each group and the observation is compared with each group, using the above procedure. As a consequence of Wold's approach, an observation may be left unclassified if it is extreme for all the groups, as the allocation procedure consists of a set of outlier tests that are evaluated independently of one another. Alternatively, an observation may be assigned to several groups.

There is further information in the position of an observation in the primary space. Hence, the above classification procedure may be complemented by an additional outlier analysis on the primary principal component scores.[20] As a consequence, some authors have described SIMCA as a 'box' model.[22] Indeed, if we decide on the appropriate rejection values for the principal component outlier tests, then a box is defined for each group. The previously described procedure is then equivalent to assigning an observation to a group if it is inside the corresponding box.

This extension of SIMCA is not general in the literature. In his original paper on the SIMCA method,[20] Wold takes the view that the classification procedure becomes 'rather complicated' if we incorporate the primary scores in the analysis. His original approach was to generate a misfit indicator' from the primary component scores that is reported together with the actual classification as obtained from the residual scores, 'if one wishes'. This may go some way towards explaining why there appears to be no consensus in the literature on SIMCA as to how the primary information should be used.

Another approach would be to ignore the information from the primary components altogether and assign an observation to the group for which the smallest value of $D_{W(k)}^2$ (x) is observed. This is the approach of Frank,[23] Frank and Friedman[24] and McLachlan[25]. It defines SIMCA as a classification procedure, in a manner similar to some more conventional allocation procedures in statistics. Furthermore, it restricts the SIMCA method to the bare essentials and hence, we will restrict attention to this formulation from now on.

A problem with the SIMCA method that we have not yet discussed is the choice of the dimensionality of the principal component representations. The same problem poses itself in outlier detection. For the rice data, for example, it is by no means clear whether we should regard the second component as redundant information, even though its variation is small. However, the performance of Rao's statistic will decrease if we add this component to the model components.

To choose the number of components in the SIMCA procedure, Wold used a cross-validated estimate of the absolute goodness-of-fit

$$\frac{\sum_{i=1}^{n} \sum_{j=1}^{p} \epsilon_{ij}^2}{n - 1}$$

Unfortunately, while this aproach may be safeguarded against influential observations, it does not optimize the SIMCA classification rule directly. Dröge and Van't Klooster[26,27] studied an empirical criterion that is referred to as Malinowski's indicator function:[27]

$$\frac{\rho}{(p - k)^2}$$

with

$$\rho^2 = \frac{\sum\limits_{j=k+1}^{p} \lambda_j}{n(p-k)}$$

Essentially, this criterion is a version of Rao's absolute lack-of-fit test that takes the number of deleted components into account. According to Dröge, this criterion is more reliable than cross-validation when there are less than 15 observations available to calibrate the procedure. Similar caution should be applied here with respect to the degrees of freedom to be used, as we explained for Wold's $F$-test.

Hawkins prefers a more *ad hoc* approach. He standardizes the measured variates prior to analysis and discusses three different criteria to choose $k$. The first of these is to treat all components with variance less than one as residual components. The second criterion is to postulate a normal distribution for the measured data and calculate the power of Rao's statistic for different values of $k$. We may then select that value for which the maximum power is achieved. As a last procedure, he considers the principal component weights and chooses $k$ so that all measured variates are adequately represented in the residual components. His final decision is a compromise value based on the results of these separate analyses. The aim of this approach is not so much to obtain an adequate representation of the data, but rather to calibrate the outlier statistic optimally with respect to outlier detection. Standardization of the measured data is not essential to this procedure.

### SIMCA and the Mahalanobis Distance

We have demonstrated how SIMCA is based, essentially, on a lack-of-fit test. This statistic has a serious drawback, however. Indeed, Rao's statistic merely adds all residual principal component scores together, irrespective of their variability. As a consequence, a smaller variance component that separates groups well may be obscured by an irrelevant larger variance component. More generally, we may wonder whether Rao's statistic gives enough weight to the smaller variance components, as we have already pointed out that it is especially these components that may provide the directions in which groups separate.[10]

The problem stems essentially from the fact that we have implicitly assumed that the measured data are represented in a Euclidean space[12] and, as a consequence, the orthogonal distance from the principal component representation was measured as a Euclidean distance. Indeed, we have

$$D^2_{1(k)}(\mathbf{x}_i) = u^2_{ik+1} + \ldots + u^2_{ip}$$
$$= u^2_{i1} + \ldots + u^2_{ip} - u^2_{i1} - \ldots - u^2_{ik}$$
$$= \mathbf{x}_i Q(\mathbf{x}_i Q)^T - (\mathbf{x}_i \mathbf{q}_1)^2 - \ldots - (\mathbf{x}_i \mathbf{q}_k)^2$$

and hence

$$D^2_{1(k)}(\mathbf{x}_i) = \mathbf{x}_i \mathbf{x}_i^T - (\mathbf{x}_i \mathbf{q}_1)^2 - \ldots - (\mathbf{x}_i \mathbf{q}_k)^2$$

The first term in this equation is the Euclidean distance of the observation from the sample mean. The remaining terms remove the contributions from the primary components from the distance measure.

Fig. 5 contains an illustration of the problem associated with the Euclidean metric. Observations A and B are plotted at equal Euclidean distances from the mean center of a group of observations. However, the position of observation A is along the first principal component direction while the orientation of observation B is orthogonal to the first component. Hence, the orientation of most observations in the group will be closer to the orientation of A than to the orientation of B.

Furthermore, the variation of the data is substantially larger along the first principal component direction than along the second component. For these reasons, A is closer to the group in statistical terms. These arguments show why a proper statistical distance measure must take the scale and orientation of the data into account.

This is achieved if we compute the standardized principal component scores and use the Euclidean distance for these transformed measures. Indeed, consider Fig. 6, which contains the same data as shown in Fig. 5, plotted *versus* the standardized components. Observation A is much closer to the group in the transformed space. This is because the standardized principal component scores are an uncorrelated set of variates with unit variance. Essentially, the principal component decomposition has removed the correlations and differences of scale from the data so that we may again use the Euclidean distance. This procedure defines a distance measure that is fundamental in multivariate statistics. The distance is referred to as the Mahalanobis distance and may be written as

$$\mathbf{x}_i S^{-1} \mathbf{x}_i^T = u^2_{i1}/\lambda_1 + \ldots + u^2_{ip}/\lambda_p$$

We may apply these ideas to define the Mahalanobis analogue to Rao's statistic:

$$D^2_{2(k)}(\mathbf{x}_i) = \mathbf{x}_i S^{-1} \mathbf{x}_i^T - (\mathbf{x}_i \mathbf{q}_1)^2/\lambda_1 - \ldots - (\mathbf{x}_i \mathbf{q}_k)^2/\lambda_k$$

This statistic was proposed by Hawkins.[15] As with Rao's statistic, we have

$$D^2_{2(k)}(\mathbf{x}_i) = u^2_{i1}/\lambda_1 + \ldots + u^2_{ip}/\lambda_p - u^2_{i1}/\lambda_1 - \ldots - u^2_{ik}/\lambda_k$$
$$= u^2_{ik+1}/\lambda_{k+1} + \ldots + u^2_{ip}/\lambda_p$$
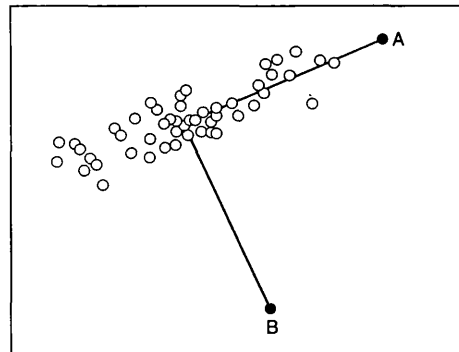


Fig. 5 Plot of two observations, A and B, at equal Euclidean distances from a group of observations. Observation A is positioned along the first principal component, while B is along the second principal component.
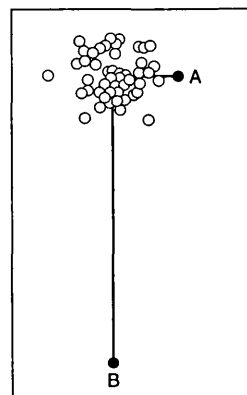


Fig. 6 Plot of the observations, A and B, *versus* the standardized principal components of the group of observations.

and hence Hawkins's statistic may be seen as a lack-of-fit test that is defined as the sum of the squared residual standardized principal component scores, in a manner analogous to the definition of Rao's statistic. Hawkins proposed an interesting alternative derivation of this statistic.[15-17]

As Hawkins states, there is no obvious advantage involved in working with the Euclidean distance. We should expect that implementations of SIMCA that are based on the Mahalanobis distance will separate groups better in most applications. This is confirmed in applications to quality control, as discussed by Hawkins,[15] Hawkins and Fatti[16] and Fellegi.[28] Some chemometricians might argue that Rao's statistic can be computed conveniently as a residual sum of squares with the nonlinear iterative partial least squares (NIPALS) algorithm.[29] Given the efficient algorithms for principal component decomposition that are available today, however, this would be a very weak argument in favour of this statistic. Furthermore, the above formula implies that we may calculate Hawkins's statistic from the primary components and hence we do not need to calculate the residual scores explicitly. According to Jackson,[30] there is a qualitative difference between Hawkins's and Rao's statistic, in that Rao's statistic may be interpreted 'in terms of residuals' while Hawkins' statistic must be interpreted 'in terms of the deleted components'. However, the principal component decomposition of the data is the same, irrespective of whether we are using the Euclidean or the Mahalanobis distance[12] and hence, despite Jackson's statements, we may interpret both statistics as lack-of-fit tests for the same principal component representation.

The Mahalanobis distance upweights the smaller variance components so that all components have equal weight. The opposite effect is achieved if we multiply all residual scores with their respective standard deviations. This leads to the statistic

$$D_{3(k)}^2 (\mathbf{x}_i) = \lambda_{k+1} u_{ik+1}^2 + \ldots + \lambda_p u_{ip}^2$$
$$= \lambda_1 u_{i1}^2 + \ldots + \lambda_p u_{ip}^2 - \lambda_1 u_{i1}^2 - \ldots - \lambda_k u_{ik}^2$$
$$= \mathbf{x}_i \mathbf{S} \mathbf{x}_i^T - \lambda_1 (\mathbf{x}_i \mathbf{q}_1)^2 - \ldots - \lambda_k (\mathbf{x}_i \mathbf{q}_k)^2$$

which was studied by Gnanadesikan and Kettenring in 1972.[3] This statistic gives more weight to the larger variance components and reduces the contribution of the smaller variance components. As a consequence, we will often find that this statistic does not separate groups well in discriminant problems.

Hawkins[15] has proposed a further statistic for outlier detection that is based on the normalized components. It is defined as

$$D_{4(k)} (\mathbf{x}_i) = \max_{k < j} |u_{ij}/\sqrt{\lambda_j}|$$

In contrast to the statistics that were defined above, this statistic is not a lack-of-fit test for a principal component representation. It is proposed more on empirical grounds, because of the optimal properties of the normalized principal components in outlier detection. If we postulate a normal distribution for the data, then the principal components are independent and the above statistic can be interpreted as the largest test result from a set of independent test statistics. Similar methods to those discussed above may be used to determine the number of components to delete.

Let us consider the application of these statistics in a SIMCA discriminant analysis of the rice data. We add three more samples of rice from the French Camargue (F), as well as a new group of 5 samples of white rice (W) and a group of 5 samples of rice from Laos (L) (samples of Basmati rice are indicated with B).

Using a full leave-one-out cross-validated estimate of the absolute goodness-of-fit statistics, we find that the first

component accounts for more than 97% of the cross-validated variation in each of the groups. Hence, we use a one-dimensional principal component representation for each group. This essentially removes the grain size effects from the SIMCA outlier statistics, while retaining all other small pattern aberrations. Tables 1 to 5 give the full leave-one-out cross-validated SIMCA classification tables for each of the lack-of-fit statistics that we have discussed.

The results confirm the previous discussion. The Mahalanobis distance (Table 2) reduces the total cross-validated error rate, as compared to the Euclidean metric (Table 1). The total cross-validated error rate obtained for Wold's statistic (42.5%) (Table 5) is smaller than that obtained for Rao's statistic (52.5%) but larger than the cross-validated error rate obtained for Hawkins's statistic (31.8%). The total cross-validated error rate for distance $D_4$ is 35.3%.

**Table 1** Classification table using Rao's statistic $D_{1(1)}$

| From | To | | | | Total |
|------|---|---|---|---|-------|
|      | B | W | F | L |       |
| B | 9 | 0 | 0 | 0 | 9 |
| W | 0 | 2 | 0 | 3 | 5 |
| F | 0 | 1 | 2 | 1 | 4 |
| L | 0 | 4 | 1 | 0 | 5 |
| Total | 9 | 7 | 3 | 4 | 23 |
| Total cross-validated error rate | | | | | 52.5% |

**Table 2** Classification table using Hawkins's statistic $D_{2(1)}$

| From | To | | | | Total |
|------|---|---|---|---|-------|
|      | B | W | F | L |       |
| B | 7 | 0 | 0 | 2 | 9 |
| W | 0 | 1 | 0 | 4 | 5 |
| F | 0 | 0 | 3 | 1 | 4 |
| L | 0 | 0 | 0 | 5 | 5 |
| Total | 7 | 1 | 3 | 12 | 23 |
| Total cross-validated error rate | | | | | 31.8% |

**Table 3** Classification table using Gnanadesikan's statistic $D_{3(1)}$

| From | To | | | | Total |
|------|---|---|---|---|-------|
|      | B | W | F | L |       |
| B | 7 | 2 | 0 | 0 | 9 |
| W | 0 | 4 | 0 | 1 | 5 |
| F | 0 | 2 | 2 | 0 | 4 |
| L | 0 | 4 | 1 | 0 | 5 |
| Total | 7 | 12 | 3 | 1 | 23 |
| Total cross-validated error rate | | | | | 48.1% |

**Table 4** Classification table using statistic $D_{4(1)}$

| From | To | | | | Total |
|------|---|---|---|---|-------|
|      | B | W | F | L |       |
| B | 8 | 0 | 0 | 1 | 9 |
| W | 0 | 1 | 0 | 4 | 5 |
| F | 0 | 0 | 2 | 2 | 4 |
| L | 0 | 0 | 0 | 5 | 5 |
| Total | 8 | 1 | 2 | 12 | 23 |
| Total cross-validated error rate | | | | | 35.3% |

## SIMCA and Biased Discriminant Analysis

We have discussed the application of the Mahalanobis distance to principal component outlier detection and SIMCA methods of pattern recognition. The Mahalanobis distance is a fundamental concept in statistical discriminant analysis. We would like to compare the SIMCA approach to classification with some classical methods of discriminant analysis in statistics. We will assume that the prior probabilities of group membership are equal.

Consider the minimum-distance classification procedure in which an observation is assigned to the group for which the smallest Mahalanobis distance is observed. In such applications where it is reasonable to assume that the dispersion of the data is similar in the distinct groups, a common covariance matrix may be postulated for the groups. Some pooled estimate of the covariance matrices may then be used in the calculation of the Mahalanobis distances. As a consequence of pooling the covariance estimates, the variability of the discriminant rule decreases, due to a drastic reduction in the number of parameters that are estimated from the data. The minimum-distance discriminant rule itself reduces to classical linear discriminant analysis,[4] first derived in 1936 by Fisher[31] for the two-group case, from a purely data-based point of view.[5]

In many classification problems, however, an assumption of equal dispersion is not warranted and we must allow for the differences in the spatial distribution of the data in the distinct groups. Quadratic discriminant analysis generalizes linear discrimination by postulating different covariance matrices for the distinct groups. As for linear discriminant analysis, the quadratic allocation rule may be derived from an assumption of multivariate normality for each group, assigning a new observation $x$ to the group for which the posterior probability of group membership is maximum. Hence we assign $x$ to the group for which we observe the smallest value of the discriminant score

$$D_Q(x) = xS^{-1}x^T + \ln[\det(S)]$$

where $\det(S)$ refers to the determinant of $S$.[4,5] The last term in this equation is a measure of the size of the covariance matrix, as $\det(S) = \prod_{j=1}^{p} \lambda_j$. It reflects the scale of the spread of the data in the group. The first term is the Mahalanobis distance of $x$ from the group. Needless to say, this allocation rule will reduce to linear discriminant analysis if the covariance matrices are pooled in a single estimate of covariance. Hence, it is essential that the estimation of the covariance matrix $S$ is specific to each group.

Quadratic discriminant analysis is probably the standard statistical discriminant technique when the data have unequal within-group dispersion. However, a serious problem associated with the quadratic rule is that it may not perform well when the number of observations is limited. Many statisticians have reported applications where the linear rule performed better, even though the assumptions for pooling the data were

not met (Seber gives an interesting discussion on minimal sample sizes,[4] as well as Friedman[32]). This is because the quadratic rule requires the estimation of a much larger number of parameters and hence the variability of the quadratic discriminant function increases rapidly when the number of observations decreases. We may employ a linear discriminant analysis in such a situation as a form of regularization, in which estimates are biased to achieve smaller variability.

In contrast to these statements, it is claimed in the chemometrics literature that the SIMCA method performs well, even for exceptionally small calibration sets.[33] This should be of particular interest, as both quadratic discriminant analysis and SIMCA postulate distinct covariance matrices $S$ for each group and therefore, we may think of both these procedures as methods that model the within-group variation separately for each group. As the SIMCA approach retains a separate modelling of the within-group variation, we are particularly interested in a comparison of the SIMCA method with the quadratic rule and we may wonder whether the SIMCA rule can improve on the performance of the quadratic rule when there are few observations.

There is an interpretation of SIMCA as another method of regularization for discriminant problems with unequal class covariance matrices and a small number of observations. To discuss this, we will first discuss a reformulation of the outlier statistics that we have discussed.

We may rewrite Hawkins's statistic as

$$D^2_{2(k)}(x) = sS^c_{(k)}{}^{-1}x^T$$

where $S^c_{(k)} = Q^c_{(k)}\Lambda^c_{(k)}Q^c_{(k)}{}^T$, with $Q^c_{(k)} = (q_{k+1},\ldots,q_p)$ and $\Lambda^c_{(k)}$ is a diagonal matrix with the values $\lambda_{k+1} > \ldots > \lambda_p$ on the diagonal. In a similar fashion, Gnanadesikan's statistic can be written as

$$D^2_{3(k)}(x) = xS^c_{(k)}x^T$$

These statistics are of the type

$$xDx^T$$

where $D = Q\Lambda Q^T$ and $\Lambda$ is a diagonal matrix with non-negative weights $\delta_1,\ldots,\delta_p$ on the diagonal. $D$ is a symmetric matrix and hence the above statistics define a set of distance measures. These measures are variants of the Mahalanobis distance, which itself assigns equal importance to each of the principal component directions. We may change the relative importance of the principal components by changing the weights in the distance measure.

In the SIMCA method, we always have $\delta_1 = \ldots = \delta_k = 0$. Hawkins's statistic is then obtained for $\delta_j = 1/\lambda_j$, $j = k + 1,\ldots,p$. Likewise, Gnanadesikan's statistic is obtained for $\delta_j = \lambda_j$, $j = k + 1,\ldots,p$, while Rao's statistic is obtained from $\delta_j = 1$, $j = k + 1,\ldots,p$.

Wold's version of the SIMCA method fits into the same framework. We have

$$\sum_{i=1}^{n} D^2_{1(k)}(x_i) = (n-1)\sum_{j=k+1}^{p} \lambda_j.$$

and hence, Wold's statistic is obtained for

$$\delta_{k+1},\ldots,\delta_p = \frac{(n-k-1)}{(n-1)\sum_{j=k+1}^{p} \lambda_j}$$

As we can see, Wold's version of the SIMCA method downweights the smaller variance components, as compared to the Mahalanobis distance, by using a uniform weight that is inversely proportional to the average residual principal component variation, instead of $\delta_j = 1/\lambda_j$, $j = k + 1,\ldots,p$.

**Table 5** Classification table using Wold's statistic $D_{W(1)}$

| From | To | | | | |
| --- | --- | --- | --- | --- | --- |
| | B | W | F | L | Total |
| B | 9 | 0 | 0 | 0 | 9 |
| W | 0 | 1 | 0 | 4 | 5 |
| F | 0 | 0 | 2 | 2 | 4 |
| L | 0 | 1 | 1 | 3 | 5 |
| Total | 9 | 2 | 3 | 9 | 23 |
| Total cross-validated error rate | | | | | 42.5% |

Wold's implementation of SIMCA is not identical with that which is obtained for Rao's statistic, as the principal component variances are calculated from a within-group principal components decomposition and are therefore group dependent.

From these results, we find that an implementation of SIMCA with any of the above statistics is a quadratic allocation procedure. A further interpretation of these results has been considered by Frank and Friedman, in the following manner.[24] The small eigenvalues contribute substantially to the variation of the Mahalanobis distance in the quadratic allocation rule. This effect is worsened by the biased estimation of the principal component variances in the principal component decompostion of the covariance matrix S. The small variances are biased to values that are too low, while the larger variances are biased to values that are too high. This bias is more severe when the poplulation size decreases.[4,10]

Wold's variant of the SIMCA method tries to rectify this problem by down-weighting the small variance components, as compared to the weighting in the Mahalanobis distance. Hence, SIMCA is a biased form of quadratic discriminant analysis. A similar discussion can also be found in Friedman's paper on regularized discriminant analysis, although he does not discuss SIMCA explicitly.[32]

Some critical comments on this interpretation are in order. Although Wold's version of the SIMCA method may reduce the variability of the distance measure, we may lose some relevant information from the smaller variance principal component directions. More importantly, it is not at all clear how the particular reweighting employed by Wold's version of the SIMCA method should be justified. Originally, Wold's statistic was introduced as a statistic that could easily be obtained from the NIPALS algorithm and there was no intention of optimizing a bias-variance trade-off in any way. Indeed, Wold's argument makes no reference to any interpretation of SIMCA as a biased discriminant analysis.

Let us now turn to a final comparison between quadratic discriminant analysis and the SIMCA allocation rule.[24] First of all, SIMCA differs from quadratic discriminant analysis in that it ignores the last term $\{\ln[\det(S)]\}$ in the quadratic discriminant rule. As a consequence, the allocation rule may perform badly when the scale of the variation is very different in the distinct groups. Unfortunately, this is precisely one of the situations in which we would want to use a SIMCA type variant of the quadratic rule, when the number of observations is small. Indeed, if there are no substantial differences in within-group variation, we may wonder why we should bother at all with a method that postulates separate models for the within-group variation. Secondly, as we have pointed out before, a minimal implementation of SIMCA will ignore all information from the primary principal components and there is no generally accepted method to incorporate this information in the analysis. We may therefore expect poor performance from such SIMCA classifiers when the groups are separated mainly along the primary components. The last difference is concerned with the biased estimation of the residual principal component variances. SIMCA biases the estimation of the group covariance matrices to reduce the variability of the allocation rule. Hence, SIMCA may be regarded as a form of regularized quadratic discriminant analysis.[24,25]

Frank and Friedman[24] used simulations and real data to compare the SIMCA method with some new methods of discriminant analysis in which the bias-variance trade-off is optimized explicitly. Given the above analysis, it is not surprising that the results did not favour the SIMCA method.

## References

1   *Proceedings of the IUFoST Symposium on Food Research and Data Analysis, Oslo, Norway, September 1982.* IUFoST, eds. Martens, H., and Russwurm, Jr., H., Applied Science Publishers, London 1982.
2   Martens, H., and Naes, T., *Multivariate Calibration.* Wiley, New York, 1989, pp. 267–296.
3   Gnanadesikan, R., and Kettenring, J. R., *Biometrics,* 1972, **28**, 81.
4   Seber, G. A. F., *Multivariate Analysis: a user's perspective,* Oxford University Press, New York, 1988, pp. 1–176.
5   Krzanowski, W. J., *Principles of Multivariate Analysis: a user's perspective,* Oxford University Press, New York, 1988, pp. 127–174.
6   Cowe, I. A., and McNichol, J. W., *App. Spectrosc.,* 1985, **39**, 257.
7   Cowe, I. A., McNichol, J. W., and Cuthbertson, D. C., *Analyst,* 1985, **110**, 1227.
8   Cowe, I. A., McNichol, J. W., and Cuthbertson, D. C., *Analyst,* 1985, **110**, 1233.
9   Bertrand, D., Robert, P., and Tran, V., *Rep. Int. Assoc. Cereal Sci. Technol.,* 1984, **11**, 93.
10  Jolliffe, I. T., *Principal Component Analysis,* Springer-Verlag, New York, 1986.
11  Hotelling, H., *J. Educ. Psychol.,* 1933, **24**, 417.
12  Rao, C. R., *Sankhya, Ser. A,* 1964, **26**, 329.
13  Pearson, K., *Phil. Mag.,* 1901, **2**, 559.
14  Barnett, V., and Lewis, T., *Outliers in Statistical Data.* Wiley, New York, 1978, pp. 223–227.
15  Hawkins, D. M., *J. Am. Statist. Assoc.,* 1974, **69**, (346, 340.
16  Hawkins, D. M., and Fatti, L. P., *Statistician,* 1984, **33**, 325.
17  Hawkins, D. M., *Identification of Outliers,* Chapman and Hall, London, 1980, pp. 110–114.
18  Jackson, J. E., and Mudholkar, G. S., *Technomet.,* 1979, **21**, 341.
19  Jackson, J. E., and Hearne, F. T., *Technomet.,* 1979, **21**, 253.
20  Wold, S., *Pattern Recognit.,* 1976, **8**, 127.
21  Eckart, C., and Young, G., *Psychometrika,* 1936, **1**, 211.
22  Massart, D. L., Vandeginste, B. G. M., Deming, S. N., Michotte, Y., and Kaufman, L., *Chemometrics: a textbook,* Elsevier, Amsterdam, 1988, pp. 403–407.
23  Frank, I. E., *Chemom. Intell. Lab. Syst.,* 1989, **5**, 247.
24  Frank, I. E., and Friedman, J. H., *J. Chemom.,* 1989, **3**, 463.
25  McLachlan, G. F., *Discriminant Analysis and Pattern Recognition,* Wiley, New York, 1992, pp. 141–144.
26  Droge, J. B. M., and Van't Klooster, H. A., *J. Chemom.,* 1987, **1**, 221.
27  Droge, J. B. M., and Van't Klooster, H. A., *J. Chemom.,* 1987, **1**, 231.
28  Fellegi, I. P., *Bull. Int. Stat. Inst.,* 1975, **46**, 249.
29  Wold, H., in *Research papers in statistics, Festschrif for J. Neyman* ed. David, F. N., Wiley, New York, 1966, pp. 4 1–444.
30  Jackson, J. E., *A User's Guide to Principal Components.* Wiley, New York, 1991.
31  Fisher, R. A., *Ann. Eugenics,* 1936, **7**, 179.
32  Friedman, J. H., *J. Am. Statist. Assoc.,* 1989, **84** (405), 165.
33  Albano, C., Blomquist, G., Coomans, D., Dunn, W. J., Edlund, U., Eliasson, B., Hellberg, S., Johansson, E., Nordén, E., Johnels, D., Sjörström, M., Söderström, B., Wold, H., and Wold, S., in *Proc. Symposium i anvent statistik., Köperhamm, January 22, 1981,* eds. Höskuldsson, A., Sloth Jensen, B., and Esbensen, K., NEUCC, RECAU and RECKU.

# Errata

- In the publication, page 2779: change the last displayed equation to

$$D^2_{1(k)}(\mathbf{x}) = \sum_{j=k+1}^{p} v_j^2.$$