



Mathematical properties of Consensus PCA

J.J. Jansen, H.F.M. Boelens, A.K. Smilde

Process Analysis and Chemometrics, Department of Chemical Engineering,
University of Amsterdam, The Netherlands

website: <http://www-its.chem.uva.nl/> ✉ jjansen@science.uva.nl

Introduction

A multiblock extension for PCA has been developed by Wold et al.¹ called Consensus PCA (CPCA). Westerhuis et al.² modified this method. Westerhuis' CPCA version is on the super level equivalent to ordinary PCA^{2,3}.

CPCA generates several models of the data in each block. The properties of these models in terms of orthogonality and splitting of variance are determined. Knowledge of these properties gives insight into the interpretation of the obtained block scores and loadings.

Algorithm

The algorithm developed by Westerhuis et al. is given by:

- (1) $\mathbf{X}^{(t)} = \mathbf{T}_{\text{sup}} \mathbf{P}_{\text{sup}}^T + \mathbf{X}^{(R+t)}$
 - (2) $\mathbf{T}_{\text{sup}} = [\mathbf{t}_{\text{sup},1} \ \dots \ \mathbf{t}_{\text{sup},R}]; \mathbf{P}_{\text{sup}}^T = [\mathbf{P}_{\text{sup},1}^T \ \dots \ \mathbf{P}_{\text{sup},B}^T];$
 $\mathbf{P}_{\text{sup},b} = [\mathbf{p}_{\text{sup},b1} \ \dots \ \mathbf{p}_{\text{sup},bR}]$
 - (3) for every dimension r , for $r = 1 \dots R$
 - a) $\mathbf{p}_{br} = \frac{\mathbf{P}_{\text{sup},br}}{\|\mathbf{P}_{\text{sup},br}\|}$
 - b) $\mathbf{t}_{br} = \mathbf{X}_b^{(r)} \mathbf{p}_{br}$
 - c) $\mathbf{T}_r = [\mathbf{t}_{1r} \ \dots \ \mathbf{t}_{Br}]$
 - d) $\mathbf{w}_r = \frac{\mathbf{T}_r^T \mathbf{t}_{\text{sup},r}}{\|\mathbf{T}_r^T \mathbf{t}_{\text{sup},r}\|}$
 - e) Deflate: $\mathbf{X}_b^{(r+1)} = \mathbf{X}_b^{(r)} - \mathbf{t}_{\text{sup},r} \mathbf{P}_{\text{sup},br}^T$
- end

Models

Model A1: Regular PCA

$\mathbf{X}^{(t)} = \mathbf{T}_{\text{sup}} \mathbf{P}_{\text{sup}}^T + \mathbf{X}^{(R+t)}$				
Symbol	Dimensions		Symbol	
$\mathbf{X}^{(t)}$	$I \times \sum_{b=1}^B J_b$	Raw Data Matrix	$1 \dots b \dots B$	Block Index
\mathbf{T}_{sup}	$I \times R$	Super Scores	$1 \dots j_b \dots J_b$	Variable index
\mathbf{P}_{sup}	$\sum_{b=1}^B J_b \times R$	Super Loadings	I	Number of samples
$\mathbf{X}^{(R+t)}$	$I \times \sum_{b=1}^B J_b$	Residuals of the R -component Model A1	$1 \dots r \dots R$	PC index

Properties of Model A1

Variance can be split in PCA. Orthogonality properties in PCA are already well known.

$$\|\mathbf{X}^{(t)}\|^2 = \|\mathbf{t}_{\text{sup},1} \mathbf{P}_{\text{sup},1}^T\|^2 + \dots + \|\mathbf{t}_{\text{sup},R} \mathbf{P}_{\text{sup},R}^T\|^2 + \|\mathbf{X}^{(R+t)}\|^2$$

Model A2: Model on super level for each block

$\mathbf{X}_b^{(t)} = \mathbf{T}_{\text{sup}} \mathbf{P}_{\text{sup},b}^T + \mathbf{X}_b^{(R+t)}$		
Symbol	Dimensions	
$\mathbf{X}_b^{(t)}$	$I \times J_b$	Raw data matrix for block b
$\mathbf{P}_{\text{sup},b}$	$J_b \times R$	Super Loadings for block b
$\mathbf{X}_b^{(R+t)}$	$I \times J_b$	Residuals of the R -component Model A2

Properties of Model A2

$$\mathbf{P}_{\text{sup},b}^T = (\mathbf{T}_{\text{sup}}^T \mathbf{T}_{\text{sup}})^{-1} \mathbf{T}_{\text{sup}}^T \mathbf{X}_b^{(t)}; \mathbf{P}_{\text{sup},b}^T \mathbf{P}_{\text{sup},b} \neq \mathbf{0};$$

$$\mathbf{X}_b^{(R+t)} \mathbf{P}_{\text{sup},b} \neq \mathbf{0}; \mathbf{T}_{\text{sup}}^T \mathbf{X}_b^{(R+t)} = \mathbf{0}$$

Variance can be split into the different PCs in Model A2, due to the fact that

$\mathbf{T}_{\text{sup}}^T \mathbf{T}_{\text{sup}} = \text{diag}$ and the orthogonality of \mathbf{T}_{sup} with the residuals.

$$\|\mathbf{X}_b^{(t)}\|^2 = \|\mathbf{t}_{\text{sup},1} \mathbf{P}_{\text{sup},b1}^T\|^2 + \dots + \|\mathbf{t}_{\text{sup},R} \mathbf{P}_{\text{sup},bR}^T\|^2 + \|\mathbf{X}_b^{(R+t)}\|^2$$

Model B: Model on block level for each block

$\mathbf{X}_b^{(t)} = \mathbf{T}_b \mathbf{P}_b^T + \mathbf{Z}_b^{(R+t)}$		
Symbol	Dimensions	
\mathbf{T}_b	$I \times R$	Block Scores
\mathbf{P}_b	$J_b \times R$	Block Loadings
$\mathbf{Z}_b^{(R+t)}$	$I \times J_b$	Residuals of the R -component Model B

Properties of Model B

$$\mathbf{t}_{br} = \mathbf{X}_b^{(r)} \mathbf{p}_{br}; \mathbf{T}_b^T \mathbf{T}_b \neq \text{diag}; \mathbf{P}_b^T \mathbf{P}_b \neq \text{diag}; \mathbf{T}_b^T \mathbf{Z}_b^{(R+t)} \neq \mathbf{0}; \mathbf{Z}_b^{(R+t)} \mathbf{P}_b \neq \mathbf{0}$$

Variance can not be split in this model, due to the non-orthogonality of \mathbf{T}_b and \mathbf{P}_b and their non-orthogonality with the residuals. To be able to split variance on the block level, Model C is required.

Model C: Updated model on block level for each block

When the block data $\mathbf{X}_b^{(t)}$ is projected on the block scores \mathbf{T}_b , updated block loadings are obtained. Calculation of $\tilde{\mathbf{P}}_b$ can only take place for all PCs simultaneously, due to the fact that \mathbf{T}_b is not orthogonal.

$\mathbf{X}_b^{(t)} = \mathbf{T}_b \tilde{\mathbf{P}}_b^T + \tilde{\mathbf{Z}}_b^{(R+t)}$		
Symbol	Dimensions	
$\tilde{\mathbf{P}}_b$	$J_b \times R$	Updated Block Loadings
$\tilde{\mathbf{Z}}_b^{(R+t)}$	$I \times J_b$	Residuals of the R -component Model C

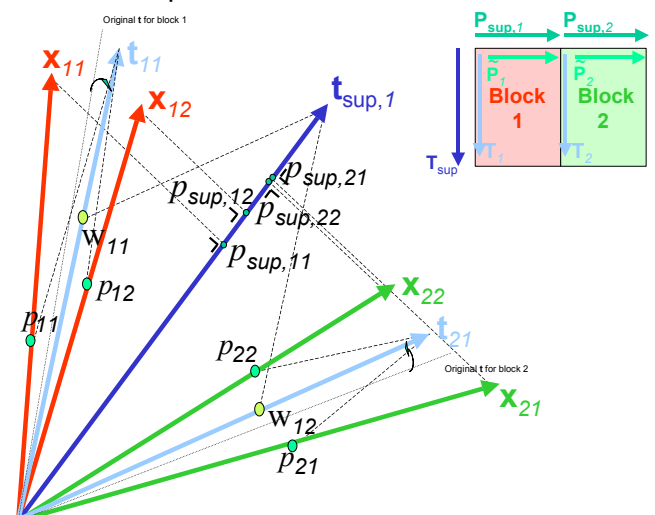
Properties of Model C

$$\tilde{\mathbf{P}}_b^T = (\mathbf{T}_b^T \mathbf{T}_b)^{-1} \mathbf{T}_b^T \mathbf{X}_b^{(t)}; \mathbf{T}_b^T \tilde{\mathbf{Z}}_b^{(R+t)} = \mathbf{0}$$

Variance can be split into explained and unexplained variance in Model C, due to the orthogonality of \mathbf{T}_b with the residuals. It is not possible to split variance per component.

$$\|\mathbf{X}_b^{(t)}\|^2 = \|\mathbf{T}_b \tilde{\mathbf{P}}_b^T\|^2 + \|\tilde{\mathbf{Z}}_b^{(R+t)}\|^2$$

Geometric Description



¹ S. Wold, S. Hellberg, T. Lundstedt, M. Sjostrom, H. Wold, *Proc. Symp. on PLS Model Building: Theory and Application*, Frankfurt am Main 1987

² Westerhuis JA, Kourti T and MacGregor JF, Analysis of multiblock and hierarchical PCA and PLS models, *Journal of Chemometrics*, 12, (1998), 301-321.

³ S.J. Qin, S. Valle, M.J. Piovoso, Unifying Multiblock Analysis, *Journal of Chemometrics*, 15, 715-742