

AN ONLINE NIPALS ALGORITHM FOR PARTIAL LEAST SQUARES

Alexander E. Stott, Sithan Kanna, Danilo P. Mandic, William T. Pike

Electrical and Electronic Engineering Department, Imperial College London, SW7 2AZ, UK

E-mails: {alexander.stott10, shri.kanagasabapathy08, d.mandic, w.t.pike}@imperial.ac.uk

ABSTRACT

Partial Least Squares (PLS) has been gaining popularity as a multivariate data analysis tool due to its ability to cater for noisy, collinear and incomplete data-sets. However, most PLS solutions are designed as block-based algorithms, rendering them unsuitable for environments with streaming data and non-stationary statistics. To this end, we propose an online version of the nonlinear iterative PLS (NIPALS) algorithm, based on a recursive computation of covariance matrices and gradient-based techniques to compute eigenvectors of the relevant matrices. Simulations over synthetic data show that the regression coefficients from the proposed online PLS algorithm converge to those of its block-based counterparts.

Index Terms— Latent variables, NIPALS, regression, Rayleigh quotient, online learning.

1. INTRODUCTION

Partial Least Squares (PLS) regression is a data analysis tool used to find relationships between blocks of multivariate data. The key attribute of PLS is its ability to deal with collinear and noisy data, an ill-posed problem for traditional ordinary least squares (OLS) solutions [1]. Importantly, unlike standard regression algorithms for collinear data, such as principal component regression (PCR), the PLS does not require the user to choose principal components in the data before finding the relationships between the independent and dependent variables. These benefits have led to successful applications of PLS in many areas such as chemometrics, social sciences and bioinformatics [2], however, the PLS is still rather unexplored in the signal processing community.

The main drawback of conventional PLS algorithms is that their block-based nature requires simultaneous processing of all the available observations. This grossly limits the use of PLS in applications involving streaming data with non-stationary statistics. In addition, block-based methods tend to be computationally prohibitive for large scale problems, where gradient-based and online estimators are becoming standard [3].

Although finding a unified solution to the online PLS algorithm is still an open problem, several efforts in this direction have been made. These include, a recursive solution based on a sliding-window type technique [4], incremental solutions specifically designed for dimensionality reduction [5] or solutions limited to univariate data [6]. None of these methods have exploited ideas from stochastic gradient-based learning which has been shown to routinely outperform other methods in large scale settings [3, 7].

To this end, we propose an online PLS algorithm which is capable of finding a common subspace for multivariate data which are collected in an online fashion. The proposed online PLS algorithm employs well-known ideas, such as recursive computation of empirical cross-covariance matrices and stochastic gradient maximisation

of the Rayleigh quotients, to estimate the scores and loading factors of the data matrices. The proposed result has similarities with online manifold learning methods [8].

The rest of this paper is organised as follows. In Section 2, a background of the PLS solution is provided and its extension to an online format is derived in Section 3. Simulation over synthetic data for an online prediction setting is included in Section 4 to verify the performance for the proposed algorithm.

2. BACKGROUND

Consider a multivariate regression problem of predicting dependent variables (output), $\mathbf{Y} \in \mathbb{R}^{N \times p}$, from independent variables (input), $\mathbf{X} \in \mathbb{R}^{N \times m}$, where N denotes the number of observations of the p -dimensional dependent variables and m -dimensional independent variables. The linear model for this regression problem is given by

$$\mathbf{Y} = \mathbf{X}\mathbf{B} \quad (1)$$

where $\mathbf{B} \in \mathbb{R}^{m \times p}$ is the matrix of regression coefficients. The ordinary least squares solution to (1) is given by

$$\hat{\mathbf{B}}_{\text{LS}} = \mathbf{X}_{\text{LS}}^+ \mathbf{Y} \quad (2)$$

and is performed through calculating the matrix pseudoinverse defined as $\mathbf{X}_{\text{LS}}^+ = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$. Crucially, this requires the matrix $\mathbf{X}^T \mathbf{X}$ to be invertible, so that for collinear independent variables in \mathbf{X} , the least squares solution in (2) is ill-posed due to the singularity of the matrix $\mathbf{X}^T \mathbf{X}$.

The partial least squares (PLS) method was introduced by Wold [9] to solve this class of problems. Several methodologies to compute a block-based PLS solution exist, and without loss in generality, we consider the classical nonlinear iterative PLS (NIPALS) algorithm owing to its widespread use and simplicity [1].

The principle of PLS is to jointly decompose the input matrix, \mathbf{X} , and the output matrix, \mathbf{Y} , into the matrices which contain only the mutually relevant components [9]. This is achieved through coupled decompositions in the form

$$\mathbf{X} = \mathbf{T}\mathbf{P}^T \quad \mathbf{Y} = \mathbf{T}\mathbf{C}^T \quad (3)$$

where $\mathbf{T} \in \mathbb{R}^{N \times r}$ is an orthogonal matrix common to both \mathbf{X} and \mathbf{Y} , and is referred to as the matrix of “scores”. The matrices $\mathbf{P} \in \mathbb{R}^{m \times r}$ and $\mathbf{C} \in \mathbb{R}^{p \times r}$ are the so-called “loadings” of the individual input and output matrices,¹ while r refers to the rank of the input matrix and is upper bounded as $r \leq m$. Notice that the PLS

¹In the full NIPALS solution, \mathbf{Y} is also decomposed in terms of another matrix \mathbf{U} which are the equivalent scores in the column space of \mathbf{Y} , but are omitted here for brevity as they are not required for calculations.

decomposition in (3), constrains \mathbf{X} and \mathbf{Y} to lie in the same column space spanned by the score matrix, \mathbf{T} , and in that sense it establishes an input-output relationship between \mathbf{X} and \mathbf{Y} which can be exploited for prediction.

The decomposition in (3) is accomplished through an iterative computation of the columns of the matrices in \mathbf{T} , \mathbf{P} and \mathbf{C} . This is achieved by first calculating a component, $\mathbf{w} \in \mathbb{R}^{m \times 1}$, that maximises the empirical cross-covariance between the \mathbf{X} and \mathbf{Y} blocks defined as

$$\mathbf{S}_1 \stackrel{\text{def}}{=} \mathbf{X}_1^\top \mathbf{Y}_1 \quad (4)$$

where the subscript “1” in \mathbf{S}_1 is used to indicate that the PLS is an iterative procedure in which subsequent covariance matrices are calculated using “deflated” versions of the original data matrices, with the initial values $\mathbf{X}_1 \leftarrow \mathbf{X}$ and $\mathbf{Y}_1 \leftarrow \mathbf{Y}$.

The first principal vector, $\mathbf{w}_1 \leftarrow \mathbf{w}$, is found as the solution to the following optimisation problem

$$\mathbf{w}_1 = \arg \max_{\|\mathbf{w}\|=1} \mathbf{w}^\top \mathbf{S}_1 \mathbf{S}_1^\top \mathbf{w} \quad (5)$$

and is given by the maximum eigenvector of the sample cross-covariance matrix (in the same iteration) post-multiplied by its transpose, that is

$$\mathbf{w}_1 = \text{Eig}_{\max}\{\mathbf{S}_1 \mathbf{S}_1^\top\} \quad (6)$$

where the $\text{Eig}_{\max}\{\cdot\}$ operator extracts the maximum eigenvector of the matrix in its argument [2]. The so obtained \mathbf{w}_1 is then used to produce the first column of the score matrix, \mathbf{T} , in (3), in the form

$$\mathbf{t}_1 = \mathbf{X}_1 \mathbf{w}_1 \quad (7)$$

To compute the loading matrices, \mathbf{C} and \mathbf{P} , it is convenient to express the PLS decomposition in (3) as a sum of rank-1 matrices, produced through outer products of the columns of the matrix of scores, \mathbf{T} , and the corresponding columns of the loading matrices, \mathbf{C} and \mathbf{P} , that is

$$\mathbf{X} = \sum_{d=1}^r \mathbf{t}_d \mathbf{p}_d^\top, \quad \mathbf{Y} = \sum_{d=1}^r \mathbf{t}_d \mathbf{c}_d^\top \quad (8)$$

Recall that equation (7) produces only the first score, \mathbf{t}_1 , so that the initial rank-1 approximation of (8) is given by

$$\mathbf{X} \approx \mathbf{t}_1 \mathbf{p}_1^\top, \quad \mathbf{Y} \approx \mathbf{t}_1 \mathbf{c}_1^\top \quad (9)$$

where the vectors \mathbf{p}_1 and \mathbf{c}_1 from (9) are respectively the first columns of the loading matrices \mathbf{P} and \mathbf{C} in (3) which are yet to be found. Notice that (9) represents nothing else than a variant of the standard least squares problem in (1) – (2), from which the solutions for \mathbf{p}_1 and \mathbf{c}_1 can be found by projecting \mathbf{X}_1 and \mathbf{Y}_1 on to the score \mathbf{t}_1 in (7), to give

$$\mathbf{p}_1 = \frac{\mathbf{X}_1^\top \mathbf{t}_1}{\mathbf{t}_1^\top \mathbf{t}_1}, \quad \mathbf{c}_1 = \frac{\mathbf{Y}_1^\top \mathbf{t}_1}{\mathbf{t}_1^\top \mathbf{t}_1} \quad (10)$$

Finally, the remaining columns of \mathbf{T} , \mathbf{C} and \mathbf{P} are found by removing the impact of the score \mathbf{t}_1 from the original data \mathbf{X} and \mathbf{Y} to produce the “deflated” matrices

$$\mathbf{X}_2 \leftarrow \mathbf{X} - \mathbf{t}_1 \mathbf{p}_1^\top, \quad \mathbf{Y}_2 \leftarrow \mathbf{Y} - \mathbf{t}_1 \mathbf{c}_1^\top \quad (11)$$

The process in (4) – (11) is repeated r times to find all \mathbf{w}_d , \mathbf{t}_d , \mathbf{p}_d , \mathbf{c}_d ,

for $d = 1, \dots, r$, starting from the deflated covariance matrix, $\mathbf{S}_d = \mathbf{X}_d^\top \mathbf{Y}_d$ in (4). The eigenvector matrix, \mathbf{W} , score matrix, \mathbf{T} , and loading matrices, \mathbf{P} and \mathbf{C} , then become

$$\begin{aligned} \mathbf{W} &= [\mathbf{w}_1, \dots, \mathbf{w}_r], & \mathbf{T} &= [\mathbf{t}_1, \dots, \mathbf{t}_r], \\ \mathbf{P} &= [\mathbf{p}_1, \dots, \mathbf{p}_r], & \mathbf{C} &= [\mathbf{c}_1, \dots, \mathbf{c}_r]. \end{aligned} \quad (12)$$

The eigenvectors, \mathbf{w}_d , of the iteration wise matrix $\mathbf{S}_d \mathbf{S}_d^\top$, $d = 1, 2, \dots, r$, are obtained from the NIPALS given in Algorithm 1, and span the subspace of \mathbf{X} which is relevant to the regression with \mathbf{Y} . Hence, they provide a low-rank PLS decomposition of the input matrix \mathbf{X} in the form

$$\tilde{\mathbf{X}} = \mathbf{T}(\mathbf{P}^\top \mathbf{W}) \mathbf{W}^\top \quad (13)$$

The regression coefficients, \mathbf{B} , are now calculated as

$$\hat{\mathbf{B}}_{\text{PLS}} = \tilde{\mathbf{X}}^+ \mathbf{Y} \quad (14)$$

where, owing to its structure, $\tilde{\mathbf{X}}^+$ is straightforwardly calculated as

$$\tilde{\mathbf{X}}^+ = \mathbf{W}(\mathbf{P}^\top \mathbf{W})^{-1} \mathbf{T}^\top \quad (15)$$

Algorithm 1 The NIPALS Algorithm

- 1: **Initialise:** $\mathbf{X}_1 \leftarrow \mathbf{X}$, $\mathbf{Y}_1 \leftarrow \mathbf{Y}$,
 - 2: **for** $d = 1, \dots, r$ **do**
 - 3: $\mathbf{S}_d = \mathbf{X}_d^\top \mathbf{Y}_d$
 - 4: $\mathbf{w}_d = \text{Eig}_{\max}\{\mathbf{S}_d \mathbf{S}_d^\top\}$
 - 5: $\mathbf{t}_d = \mathbf{X}_d \mathbf{w}_d$
 - 6: $\mathbf{c}_d = \mathbf{Y}_d^\top \mathbf{t}_d / \mathbf{t}_d^\top \mathbf{t}_d$
 - 7: $\mathbf{p}_d = \mathbf{X}_d^\top \mathbf{t}_d / \mathbf{t}_d^\top \mathbf{t}_d$
 - 8: $\mathbf{X}_{d+1} = \mathbf{X}_d - \mathbf{t}_d \mathbf{p}_d^\top$, $\mathbf{Y}_{d+1} = \mathbf{Y}_d - \mathbf{t}_d \mathbf{c}_d^\top$
 - 9: **end for**
 - 10: Create matrices \mathbf{W} , \mathbf{T} , \mathbf{P} and \mathbf{C} from (12)
 - 11: Find \mathbf{B}_{PLS} from (14).
-

3. AN ONLINE PARTIAL LEAST SQUARES

The NIPALS algorithm outlined in Algorithm 1 requires block data inputs, $\mathbf{X} \in \mathbb{R}^{N \times m}$ and $\mathbf{Y} \in \mathbb{R}^{N \times p}$, where N refers to the number of observations. However, in many practical scenarios all the observations, N , cannot be obtained simultaneously but are observed as streaming values, that is, at every time instant, n , only one row of \mathbf{X} and \mathbf{Y} is observed.

Therefore, the key difference between the NIPALS formulation in (4) – (11) and the proposed online-NIPALS (OL-PLS) is that in the online setting, at each time instant n , we only have access to the n -th row of the matrices \mathbf{X} and \mathbf{Y} , denoted by $\mathbf{x}_n \in \mathbb{R}^{1 \times m}$ and $\mathbf{y}_n \in \mathbb{R}^{1 \times p}$, given as

$$\mathbf{x}_n \stackrel{\text{def}}{=} \mathbf{X}_{(n,1:m)}, \quad \mathbf{y}_n \stackrel{\text{def}}{=} \mathbf{Y}_{(n,1:p)} \quad (16)$$

The OL-PLS algorithm is next developed through recursive solutions to the steps in (4) – (11). To this end, the empirical cross-covariance matrix in (4) is first expressed as a sum of rank-1 matrices $\mathbf{x}_n^\top \mathbf{y}_n$, for

$n = 1, \dots, N$, to give

$$\mathbf{S}_{1,N} = \sum_{n=1}^N \mathbf{x}_n^\top \mathbf{y}_n \quad (17)$$

where the subscript N in $\mathbf{S}_{1,N}$ indicates the number of observations used to compute the sample cross-covariance matrix in (17). Then, (17) can be computed recursively as

$$\mathbf{S}_{1,n} = \mathbf{S}_{1,n-1} + \mathbf{x}_n^\top \mathbf{y}_n \quad (18)$$

and the same procedure can also be used for the deflated cross-covariance matrices, to give

$$\mathbf{S}_{d,n} = \mathbf{S}_{d,n-1} + \mathbf{x}_{d,n}^\top \mathbf{y}_{d,n}, \quad d = 1, \dots, r \quad (19)$$

where $\mathbf{x}_{d,n}$ and $\mathbf{y}_{d,n}$ denote the deflated data vectors which constitute the deflated data matrices in (11). A forgetting factor, $0 \ll \lambda < 1$, can be introduced in (18) and (19) to account for non-stationary data statistics, which yields

$$\mathbf{S}_{d,n} = \lambda \mathbf{S}_{d,n-1} + (1 - \lambda) \mathbf{x}_{d,n}^\top \mathbf{y}_{d,n} \quad (20)$$

The recursion in (20) provides a reliable estimate owing to the inherent averaging of the cross-covariance, which reduces the impact of outliers.

To explain the principle of OL-PLS, recall that the block-based PLS solution in steps (4) – (11) uses the empirical cross-covariance matrix, \mathbf{S}_1 , to identify a basis to describe a subspace in \mathbf{X} that exhibits co-variation with \mathbf{Y} . This basis is used to produce the corresponding latent variables, which are the columns of the score matrix \mathbf{T} , these in turn provide a robust prediction of the output \mathbf{Y} from the input \mathbf{X} . The most important step for an online PLS algorithm is to produce an estimate of the eigenvectors $\mathbf{w}_{d,n}$, $d = 1, 2, \dots, r$, at each time instant n . With an increase in n , these estimates should converge to a set of vectors equivalent to the NIPALS solution that describe the same joint subspace. Upon calculating the eigenvector, $\mathbf{w}_{d,n}$, the remaining variables can be straightforwardly computed in an online fashion.

The online PLS is now fully equipped to produce estimates of the covariance matrix $\mathbf{S}_{d,n}$, which is then used to produce the eigenvectors $\mathbf{w}_{d,n}$, at each time instant, n , and for every deflation iteration, $d = 1, 2, \dots, r$. The estimation of $\mathbf{w}_{d,n}$ from (6), when performed in an online setting, can therefore be considered as a problem of adaptively estimating the largest eigenvectors of the matrix $\mathbf{S}_{d,n} \mathbf{S}_{d,n}^\top$.

3.1. Adaptive Estimation of the OL-PLS Variables

Consider the Rayleigh quotient, which for a matrix and vector pair (\mathbf{A}, \mathbf{w}) is defined as [10]

$$\mathcal{R}(\mathbf{A}, \mathbf{w}) \stackrel{\text{def}}{=} \frac{\mathbf{w}^\top \mathbf{A} \mathbf{w}}{\mathbf{w}^\top \mathbf{w}} \quad (21)$$

The vector \mathbf{w} that maximises (21) represents the maximum eigenvector of the matrix \mathbf{A} . Therefore, the estimation of the maximum eigenvector, $\mathbf{w}_{d,n}$, of the matrix $\mathbf{S}_{d,n} \mathbf{S}_{d,n}^\top$ in (6) can be accomplished using the maximisation of the Rayleigh quotient in (21) for the matrix-vector pair $(\mathbf{S}_{d,n} \mathbf{S}_{d,n}^\top, \mathbf{w})$ along with the constraint

$\|\mathbf{w}\| = 1$. This yields the cost function

$$\mathcal{J}_w \stackrel{\text{def}}{=} \frac{1}{2} \mathcal{R}(\mathbf{S}_{d,n} \mathbf{S}_{d,n}^\top, \mathbf{w}) = \frac{\mathbf{w}^\top \mathbf{S}_{d,n} \mathbf{S}_{d,n}^\top \mathbf{w}}{2 \mathbf{w}^\top \mathbf{w}} + \beta (1 - \mathbf{w}^\top \mathbf{w}) \quad (22)$$

where β is a Lagrange multiplier. The Rayleigh cost function in (22) can be maximised recursively, using gradient ascent methodology [11], in the form

$$\mathbf{w}_{d,n+1} = \mathbf{w}_{d,n} + \mu \frac{\partial \mathcal{J}_w}{\partial \mathbf{w}} \bigg|_{\mathbf{w}=\mathbf{w}_{d,n}} \quad (23)$$

where μ is a step-size which controls the trade-off between convergence speed, stability and steady state error. The gradient $\partial \mathcal{J}_w / \partial \mathbf{w}$ is calculated as

$$\frac{\partial \mathcal{J}_w}{\partial \mathbf{w}} = - \left(\frac{\|\mathbf{S}_{d,n}^\top \mathbf{w}\|^2}{\|\mathbf{w}\|^2} \mathbf{I} - \mathbf{S}_{d,n} \mathbf{S}_{d,n}^\top \right) \frac{\mathbf{w}}{\|\mathbf{w}\|^2} - \beta \mathbf{w}$$

The role of the Lagrange optimisation is to keep $\|\mathbf{w}\| = 1$ and to ensure the uniqueness of the vector that maximises the Rayleigh quotient. It can be shown that the Lagrange multiplier, β , is in the form

$$\beta = \frac{\mathbf{w}^\top \mathbf{S}_{d,n} \mathbf{S}_{d,n}^\top \mathbf{w}}{\mathbf{w}^\top \mathbf{w}} \left(1 - \frac{1}{\|\mathbf{w}\|} \right)$$

The OL-PLS algorithm follows the direction of this gradient, estimated from the available variables at the time instant n , to give the OL-PLS update

$$\mathbf{w}_{d,n+1} = \mathbf{w}_{d,n} - \mu \beta \mathbf{w}_{d,n} + \mu \left(\frac{\|\mathbf{S}_{d,n}^\top \mathbf{w}_{d,n}\|^2}{\|\mathbf{w}_{d,n}\|^2} \mathbf{I} - \mathbf{S}_{d,n} \mathbf{S}_{d,n}^\top \right) \frac{\mathbf{w}_{d,n}}{\|\mathbf{w}_{d,n}\|^2} \quad (24)$$

Remark 1. The analysis in [11] shows that the gradient ascent step in (24) converges to the largest eigenvector of $\mathbf{S}_d \mathbf{S}_d^\top$. This implies that if the estimate of $\mathbf{S}_{d,n}$ converges to the cross-covariance matrices \mathbf{S}_d obtained by NIPALS in (4) as $n \rightarrow \infty$, then the online estimate, $\mathbf{w}_{d,n}$ will also converge to the eigenvectors, \mathbf{w}_d in (6).

Upon estimating the eigenvectors, $\mathbf{w}_{d,n}$, the scores $t_{d,n}$ from (7) are computed as

$$t_{d,n} = \mathbf{x}_{d,n} \mathbf{w}_{d,n} \quad (25)$$

The next step is to find the corresponding loading vectors, \mathbf{p}_d and \mathbf{c}_d , necessary for the deflation. Note that the computation in the block-based NIPALS algorithm in (10) gives the loading vectors

$$\mathbf{p}_{d,n} = \frac{\mathbf{x}_{d,n}}{t_{d,n}}, \quad \mathbf{c}_{d,n} = \frac{\mathbf{y}_{d,n}}{t_{d,n}} \quad (26)$$

However, the deflation of the input and output data using the loadings in (26) will yield a null-vector at each time instant. This is obvious from the deflation step in (11), modified for the online case, where

$$\mathbf{x}_{d+1,n} = \mathbf{x}_{d,n} - t_{d,n} \mathbf{p}_{d,n}^\top = \mathbf{x}_{d,n} - t_{d,n} \frac{\mathbf{x}_{d,n}}{t_{d,n}} = \mathbf{0}$$

and is a consequence of the fact the available data at each time instant are rank-1. Our solution to this issue is through an online estimation of optimal loading vectors, these are found via gradient descent min-

imisation of the following cost functions

$$\mathcal{J}_p = \|\mathbf{x}_{d,n} - t_{d,n}\mathbf{p}^\top\|^2, \quad \mathcal{J}_c = \|\mathbf{y}_{d,n} - t_{d,n}\mathbf{c}^\top\|^2 \quad (27)$$

for which the gradients are given by

$$\begin{aligned} \frac{\partial \mathcal{J}_p}{\partial \mathbf{p}} &= -2t_{d,n}(\mathbf{x}_{d,n}^\top - t_{d,n}\mathbf{p}) \\ \frac{\partial \mathcal{J}_c}{\partial \mathbf{c}} &= -2t_{d,n}(\mathbf{y}_{d,n}^\top - t_{d,n}\mathbf{c}) \end{aligned}$$

The online estimates for the loading vectors now take the form

$$\mathbf{p}_{d,n+1} = \mathbf{p}_{d,n} + \mu t_{d,n}(\mathbf{x}_{d,n}^\top - t_{d,n}\mathbf{p}_{d,n}) \quad (28)$$

$$\mathbf{c}_{d,n+1} = \mathbf{c}_{d,n} + \mu t_{d,n}(\mathbf{y}_{d,n}^\top - t_{d,n}\mathbf{c}_{d,n}) \quad (29)$$

which are standard least mean square adaptive filters. This solution is equivalent to the block NIPALS algorithm which produces rank-1 least squares estimates of \mathbf{X} and \mathbf{Y} , denoted by $\hat{\mathbf{X}} = t\mathbf{p}^\top$ and $\hat{\mathbf{Y}} = t\mathbf{c}^\top$, for each score on each iteration. Recall that $t\mathbf{p}^\top$ and $t\mathbf{c}^\top$ are then removed from the data blocks \mathbf{X} and \mathbf{Y} . To this end, the instantaneous score, $t_{d,n}$, within OL-PLS is used to deflate each $\mathbf{x}_{d,n}$ and $\mathbf{y}_{d,n}$ through the estimates $\hat{\mathbf{x}}_{d,n} = t_{d,n}\mathbf{p}_{d,n}^\top$ and $\hat{\mathbf{y}}_{d,n} = t_{d,n}\mathbf{c}_{d,n}^\top$, with the full deflation scheme given by

$$\begin{aligned} \mathbf{x}_{d+1,n} &= \mathbf{x}_{d,n} - t_{d,n}\mathbf{p}_{d,n}^\top \\ \mathbf{y}_{d+1,n} &= \mathbf{y}_{d,n} - t_{d,n}\mathbf{c}_{d,n}^\top \end{aligned} \quad (30)$$

In other words, the OL-PLS performs an adaptive estimation of $\mathbf{p}_{d,n}$ and $\mathbf{c}_{d,n}$ in order to converge to the least squares solution.

As a result, the proposed online extension of PLS adaptively estimates $\mathbf{w}_{d,n}$, $\mathbf{p}_{d,n}$ and $\mathbf{c}_{d,n}$ from the inputs \mathbf{x}_n and \mathbf{y}_n , which provides the instantaneous estimates of the PLS components $t_{d,n}$. To this end, the matrices $\mathbf{S}_{d,n}$ are updated recursively based on the deflated inputs $\mathbf{x}_{d,n}$ and $\mathbf{y}_{d,n}$, which are also estimated adaptively as shown in (30).

The OL-PLS algorithm summarised in Algorithm 2, where the tuning parameters are μ and λ .

Algorithm 2 Proposed online PLS (OL-PLS)

```

1: Inputs at time instant  $n$ :  $\mathbf{x}_n \in \mathbb{R}^{1 \times m}$  and  $\mathbf{y}_n \in \mathbb{R}^{1 \times p}$ 
2: Initialise:  $\mathbf{S}_{d,1} \leftarrow \mathbf{x}_1^\top \mathbf{y}_1$ 
3: Initialise:  $\mathbf{p}_{d,1}$ ,  $\mathbf{q}_{d,1}$  and  $\mathbf{c}_{d,1} \leftarrow \mathbf{0}$ 
4: Initialise:  $\mathbf{w}_{d,1} \leftarrow$  first column of  $\mathbf{S}_{d,1}$ 
5: for  $d = 1, \dots, r$  do
6:    $\mathbf{S}_{d,n} = \lambda \mathbf{S}_{d,n} + (1 - \lambda) \mathbf{x}_{d,n}^\top \mathbf{y}_{d,n}$ 
7:    $\mathbf{w}_{d,n+1} = \mathbf{w}_{d,n} + \mu \frac{\partial \mathcal{J}_{\mathbf{S}_{d,n}}}{\partial \mathbf{w}_{d,n}}$ , using (24)
8:    $t_{d,n} = \mathbf{x}_{d,n} \mathbf{w}_{d,n}$ 
9:    $\mathbf{p}_{d,n+1} = \mathbf{p}_{d,n} + \mu t_{d,n}(\mathbf{x}_{d,n}^\top - t_{d,n}\mathbf{p}_{d,n})$ 
10:   $\mathbf{c}_{d,n+1} = \mathbf{c}_{d,n} + \mu t_{d,n}(\mathbf{y}_{d,n}^\top - t_{d,n}\mathbf{c}_{d,n})$ 
11:   $\mathbf{x}_{d+1,n} = \mathbf{x}_{d,n} - t_{d,n}\mathbf{p}_{d,n}^\top$ ,  $\mathbf{y}_{d+1,n} = \mathbf{y}_{d,n} - t_{d,n}\mathbf{c}_{d,n}^\top$ 
12: end for
```

4. SIMULATIONS AND ANALYSIS

Simulations on synthetic data were conducted to verify the performance of the proposed OL-PLS algorithm. A particular emphasis was on the accuracy of OL-PLS in an online prediction setting. To

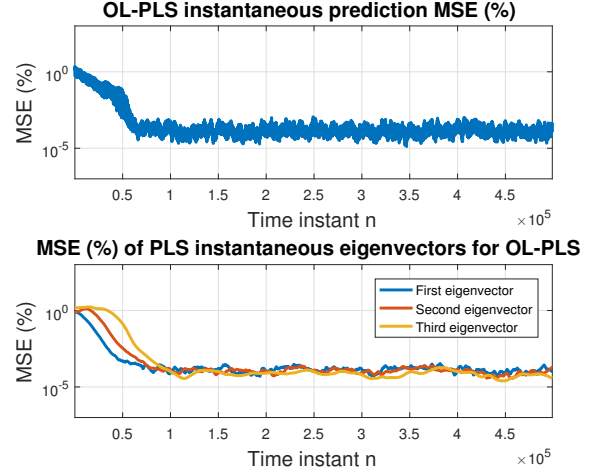


Fig. 1: OL-PLS error in predicting \mathbf{Y}

this end, we generated $M = 10$ realisations of a low-rank input input matrix, $\mathbf{X} \in \mathbb{R}^{N \times 5}$, with rank $r = 3$, drawn from a unit variance Gaussian distribution. This input matrix was multiplied by a vector $\mathbf{b} \in \mathbb{R}^{5 \times 1}$ to produce a univariate output $\mathbf{y} = \mathbf{X}\mathbf{b}$. The OL-PLS algorithm in Algorithm 2 was then used to estimate the output \mathbf{y} at each time instant. The performance measure was the empirical mean square error (MSE) computed as

$$\text{MSE}_n = \frac{1}{M} \sum_{\ell=1}^M |\mathbf{y}_n^{(\ell)} - \hat{\mathbf{y}}_n^{(\ell)}|^2$$

where $\hat{\mathbf{y}}_n^{(\ell)}$ denotes the ℓ -th realisation of the OL-PLS output estimate at time n . Figure 1 (top panel) shows the ensemble average prediction MSE, calculated as a percentage of the total variance in \mathbf{y} , for $\lambda = 0.9999$ and $\mu = 0.0001$. For successful performance, the OL-PLS solution should be equivalent to that of the NIPALS algorithm, that is, the OL-PLS vectors $\mathbf{w}_{d,n}$ should converge to those produced by NIPALS, \mathbf{w}_d , as $n \rightarrow \infty$. To verify this, Figure 1 (bottom panel) shows the ensemble average of the MSE of each vector $\mathbf{w}_{d,n}$ (as a percentage of its total variance) and the corresponding vector calculated by NIPALS.

5. CONCLUSION

We have introduced an Online PLS (OL-PLS) algorithm as an extension of the popular NIPALS method for PLS-regression. It has been shown that a recursive estimation of the empirical cross-covariance matrices, \mathbf{S}_d , provides a suitable instantaneous estimate from which the eigenvectors of $\mathbf{S}_d \mathbf{S}_d^\top$ can be adaptively calculated through gradient ascent of the Rayleigh quotient. This makes it possible for the online version of PLS to adaptively calculate the scores, while the corresponding loadings are estimated through a stochastic gradient descent scheme. We have validated the proposed OL-PLS through illustrative simulations which confirm an accurate, online, estimate of the output \mathbf{Y} and the NIPALS vectors \mathbf{w}_d . The convergence analysis and simulations over a larger number of case studies are the subject of ongoing work.

6. REFERENCES

- [1] S. Wold, M. Sjöström, and L. Eriksson, “PLS-regression: A basic tool of chemometrics,” *Chemometrics and Intelligent Laboratory Systems*, vol. 58, no. 2, pp. 109–130, 2001.
- [2] R. Rosipal and N. Krämer, “Overview and recent advances in partial least squares,” in *Subspace, Latent Structure and Feature Selection*, C. Saunders, M. Grobelnik, S. Gunn, and J. Shawe-Taylor, Eds. Springer Berlin Heidelberg, 2006, pp. 34–51.
- [3] O. Bousquet and L. Bottou, “The tradeoffs of large scale learning,” in *Advances in Neural Information Processing Systems 20*, J. C. Platt, D. Koller, Y. Singer, and S. T. Roweis, Eds. Curran Associates, Inc., 2008, pp. 161–168.
- [4] K. Helland, H. E. Berntsen, O. S. Borgen, and H. Martens, “Recursive algorithm for partial least squares regression,” *Chemometrics and Intelligent Laboratory Systems*, vol. 14, no. 1, pp. 129 – 137, 1992.
- [5] X. Q. Zeng and G. Z. Li, “Incremental partial least squares analysis of big streaming data,” *Pattern Recognition*, vol. 47, no. 11, pp. 3726 – 3735, 2014.
- [6] L. Qin, H. Snoussi, and F. Abdallah, “Online learning partial least squares regression model for univariate response data,” in *Proc. of the 22nd European Signal Processing Conference (EUSIPCO)*, Sept 2014, pp. 1073–1077.
- [7] L. Bottou, *Online learning and stochastic approximations*, D. Saad, Ed. Cambridge Univ Pr, 1998.
- [8] A. B. Goldberg, M. Li, and X. Zhu, *Online Manifold Regularization: A New Learning Setting and Empirical Study*, W. Daelemans, B. Goethals, and K. Morik, Eds. Springer Berlin Heidelberg, 2008.
- [9] H. Wold, “Estimation of principal components and related models by iterative least squares,” in *Multivariate Analysis*. Academic Press, New York, 1966, vol. 6, pp. 391–420.
- [10] R. A. Horn and C. R. Johnson, *Matrix analysis*. Cambridge University Press, 2012.
- [11] R. Arablouei, S. Werner, and K. Dogancay, “Analysis of the gradient-descent total least-squares adaptive filtering algorithm,” *IEEE Transactions on Signal Processing*, vol. 62, no. 5, pp. 1256–1264, 2014.